Математические методы для компьютера Зрение, робототехника и графика

Примечания к курсу CS 205A, осень 2013 г.

Джастин Соломон Департамент компьютерных наук Стэндфордский Университет Machine Translated by Google

Содержание

я предварительные	Ġ
0 Обзор математики	11
0.1 Предварительные занятия: числа и наборы	 . 1
0.2 Векторные пространства.	 . 12
0.2.1 Определение векторных пространств · · · · · · · · · · · · · · · ·	 . 12
0.2.2 Размах, линейная независимость и основания	 . 13
0.2.3 Наш фокус: Rn	 . 14
0.3 Линейность.	 . 16
0.3.1 Матрицы.	 . 17
0.3.2 Скаляры, векторы и матрицы.	 . 19
0.3.3 Модельная задача: Ax =b	 . 21
0.4 Нелинейность: дифференциальное исчисление	 . 22
0.4.1 Дифференциация	 . 22
0.4.2 Оптимизация	 . 25
0,5 Проблемы.	 . 27
1 Числа и анализ ошибок 1.1 Хранение	29
чисел с дробными частями .	 . 29
1.1.1 Представления с фиксированной точкой.	 . 30
1.1.2 Представления с плавающей запятой . · · · · · · · · · · · · · · · · · ·	 . 31
1.1.3 Дополнительные экзотические опции.	 . 32
1.2 Понимание ошибки	 . 33
1.2.1 Ошибка классификации.	 . 34
1.2.2. Кондиционирование, стабильность и точность.	 . 35
1.3 Практические аспекты.	 . 36
1.3.1 Пример большего масштаба: суммирование	 . 37
1.4 Проблемы.	 . 39
	4.
II Линейная алгебра	41
2. Линейные системы и LU-разложение 2.1. Разрешимость	43
линейных систем	 . 43
2.2 Специальные стратегии решения	 . 45
2.3 Колирование операций со строками	 . 46

2.3.1 Перес	становка.						 . 46
2.3.2 Масшт	абирование строк.						 . 47
2.3.3 Устра	нение.						 . 47
2.4 Исключение Г	aycca						 . 49
2.4.1 Замен	на вперед.						 . 49
2.4.2 Обра ⁻	тная замена.						 . 50
2.4.3 Анали	из исключения Гау	cca					 . 50
2.5 LU-факториза	ция						 . 52
2.5.1 Постр	ооение факторизац	ции.	• • •				 . 53
2.5.2 Реали	ıзация ЛЕ.						 . 54
2.6 Проблемы.							 . 55
3 Проектирование и а	анализ линейных с	истем 3.1 Реш	іение				57
квадратных систе							 . 57
3.1.1 Регре	ессия.						 . 57
3.1.2 Метод	наименьших квадрат	ов • • • • • • • • • • • • • • • • • •					 . 59
3.1.3 Доп	олнительные пр	имеры					 . 60
3.2 Специальные	свойства линейнь	іх систем.					 . 61
3.2.1. Пол	ожительно опр	еделенные і	иатрицы и с	факториза	ция Холецк	ого. · · · ·	 . 61
3.2.2 Разр	еженность.						 . 64
•	твительности . оичные и векторн						 . 66 . 66
	ра состояний .						 . 68
3.4 Проблемы.	•						 . 70
4 K	OD						73
4 Колоночные простр		.4.4					 . 73
4.1 Структура нор 4.2 Ортогонально	омальных уравнен						 . 74
•	егия для неортого						 . 75
·	егия для неортого ация Грама-Шмидт	·	иц.				 . 76
4.3.1 Прогн		a. 					 . 76
•	тозы. гонализация Грама	LUMARTO					 . 77
	ния домохозяев						 . 78
	 ная QR-факториз						 . 80
4.6 Проблемы.							 . 81
5 Собственные							83
	тивация						 . 83
' 5.1.1 Стати	•						. 83
5 1 2 Лиф.	ференциальные						. 84
5.1.2 диф 5.2 Спектральное							 . 85
	: вложение обственных вект						 . 86
	лоственных вект иетричные и поло	-					. 87
	летричные и поло циализированны						. 89
5.4 REIUMCROUMO C	циализированны собственных значе	с своиства.					 . 90
	ооственных значе над итепания	лиш,					90

5.4.2 Обратная итерация.					 	. 92
5.4.3 Переключение . • •					 	. 92
5.4.4 Поиск нескольких с	обственных значе	ний. • • • • •			 	. 93
5.5 Чувствительность и конди	іционирование				 	. 97
5.6 Проблемы					 	. 98
						00
6 Разложение по сингулярным числ	ам 6.1					99
Получение SVD .					 	. 99
6.1.1 Расчет СВД.	•••				 	. 101
6.2 Применение СВД.					 	. 101
6.2.1. Решение линей		•			 	. 101
6.2.2. Разложение на вне	шние произведения	•	-	кого ранга.	• • •	. 102
6.2.3 Нормы матрицы					 	. 104
6.2.4 Проблема Прокруст	а и выравнивание.		• • • •		 	. 105
6.2.5 Анализ основных ко	мпонентов (АПК) .		• • •		 	. 106
6.3 Проблемы					 	. 107
III Нелинейные методы						109
пт пелипеиные методы						
7 Нелинейные системы 7.1						111
Задачи с одной переменной .					 	. 111
7.1.1 Характеристика про	блем				 	. 111
7.1.2 Непрерывность и бы					 	. 112
 7.1.3 Анализ поиска корн					 	. 112
7.1.4 Итерация с фиксирова					 	. 113
7.1.5 Метод Ньютона.					 	. 114
7.1.6 Метод секущих					 	. 115
7.1.7 Гибридные методы.					 	. 116
7.1.8 Случай с одной		одка			 	. 116
7.2. Многомерные задачи					 	. 117
7.2.1 Метод Ньютона.					 	. 117
7.2.2 Делаем Ньютона бы	істрее: Квазиньютої	н и Бройен.			 	. 117
7.3 Кондиционирование	•				 	. 119
					 	. 119
8 Неограниченная оптимизация 8.1						121
Неограниченная оптимизация:	мотивация				 	. 121
8.2 Оптимальность					 	. 122
8.2.1 Дифференциальная	оптимальность				 	. 123
8.2.2 Оптимальность чер	ез функциональнь	іе свойства.			 	. 125
8.3 Одномерные стратегии.					 	. 126
8.3.1 Метод Ньютона.					 	. 126
8.3.2 Поиск золотого сече	ения				 	. 127
8.4 Многовариантные стратеги	1				 	. 129
8.4.1 Градиентный спуск.					 	. 129

8.4.2 Метод Ньютона.	 	 . 129
8.4.3 Оптимизация без производных: BFGS	 	 . 130
8.5 Проблемы	 	 . 132
		135
9 Ограниченная оптимизация 9.1 Мотивация		 . 135
·		 . 137
3.2 Теория оттимизации с ограничениями.		 . 140
э.э Алторитмы оттимизации.		 . 140
9.3.1 Последовательное квадратичное программирование (SQP) 9.3.2 Барьерные методы.		 . 141
		 . 141
9.4 Выпуклое программирование		 . 143
3.3 проолемы.	 	 . 143
10 итерационных линейных решателей		145
10.1 Градиентный спуск.	 	 . 146
10.1.1 Получение итерационной схемы.	 	 . 146
10.1.2 Конвергенция	 	 . 147
10.2 Сопряженные градиенты	 	 . 149
	 	 . 149
10.2.2. Субоптимальность градиентного спуска • • • • •	 	 . 151
10.2.3 Генерация А-сопряженных направлений.	 	 . 152
10.2.4 Формулировка алгоритма сопряженных градиентов.	 	 . 154
10.2.5 Условия сходимости и остановки .	 	 . 155
10.3 Предварительная подготовка.	 	 . 156
10.3.1 CG с предварительным кондиционированием.	 	 . 157
10.3.2 Общие предварительные условия.	 	 . 158
10.4 Другие итерационные схемы.	 	 . 159
10.5 Проблемы	 	 . 160
		1.61
IV функции, производные и интегралы		161
11 Интерполяция		163
11.1 Интерполяция в одной переменной.	 	 . 163
11.1.1 Полиномиальная интерполяция.		 . 164
11.1.2 Альтернативные базы.		 . 165
·	 	 . 167
11.1.3 Кусочная интерполяция		 . 168
11.2 Многопараметрическая интерполяция		 . 168
11.3 Теория интерполяции		 . 171
11.3.1 Линейная алгебра функций.		 . 171
11.3.2 Аппроксимация кусочными полиномами		 . 173
11.4 Проблемы.	 	 . 174

12 Численное интегрирова	ание и дифферен	цирова	ние				1/5
12.1 Мотивация				• • •		 	 . 176
12.2 Квадратура						 	 . 177
12.2.1 Интерпол	тяционная квадр	атура.				 	 . 177
12.2.2 Правила кв	вадратуры. •					 	 . 178
12.2.3 Квадратура	а Ньютона-Котеса.					 	 . 179
12.2.4 Квадратура	а Гаусса					 	 . 182
12.2.5 Адаптивная	я квадратура					 	 . 183
12.2.6 Нескольк	о переменных .					 	 . 183
12.2.7 Кондицио 12.3 Дифференциация.	нирование					 	 . 184 . 185
	нциация основны		ZI 11414			 	 . 185
12.3.2. Конечные		ыл фупг 				 	 . 185
12.3.3 Выбор разм						 	 . 187
• •	мера шага занные количества					 	 . 187
12.4 Проблемы.						 	 . 188
13 Обыкновенные диффере	нциальные						189
уравнения 13.1 Мотива	•					 	 . 190
13.2 Теория ОДУ.						 	 . 190
13.2.1 Основные г	понятия					 	 . 191
13.2.2 Существо	вание и уникаль	ность.				 	 . 192
13.2.3 Уравнения	_					 	 . 193
13.3 Схемы временного						 	 . 194
13.3.1 Форвард Эі						 	 . 194
13.3.2 Эйлер в обратн	ом направлении					 	 . 195
13.3.3 Трапециеві	идный метод.					 	 . 196
13.3.4 Методы Рун	нге-Кутты					 	 . 197
13.3.5 Экспоненц	иальные интегра	торы				 	 . 198
13.4 Многозначные мет						 	 . 199
13.4.1 Схемы Нью	омарка.					 	 . 199
13.4.2 Ступенчата	я сетка					 	 . 202
13.5 Сделать.						 	 . 203
13.6 Проблемы.						 	 . 203
14 Уравнения с частными производными 14.1 М						 	 205 . 205
14.2 Основные определе	ения					 	 . 209
14.3 Уравнения модели						 	 . 209
	неские УЧП					 	 . 210
14.3.2 Параболич						 	 . 211
14.3.3 Гипербол	ические УЧП. ·					 	 . 211
14.4. Производные ка						 	 . 212
14.5 Численное решени						 	 . 214
·	липтических урав	нений .				 	 . 214
	араболических и				павнений	 	 . 215

14.6 Метод конечных элементов.	 	. 217
14.7 Практические примеры	 	. 217
14.7.1 Обработка изображений в области градиента . • • • • • • • • • • • • • • • • • •	 	. 217
14.7.2 Фильтрация с сохранением границ	 	. 217
14.7.3 Жидкости на основе сетки		
14.8 Сделать	 	. 217
14.9 Проблемы.	 	. 218

Machine Translated by Google

Часть I

Предварительные



Глава 0

Обзор математики

В этой главе мы рассмотрим соответствующие понятия из линейной алгебры и исчисления с несколькими переменными, которые будут фигурировать в нашем обсуждении вычислительных методов. Он задуман как обзор справочного материала с уклоном в сторону идей и интерпретаций, обычно встречающихся на практике; эта глава может быть пропущена или использована в качестве справочного материала студентами с более сильным математическим образованием.

0.1 Предварительные игры: числа и наборы

Вместо того чтобы рассматривать алгебраические (а иногда и философские) рассуждения вроде «Что такое число?», мы будем полагаться на интуицию и математический здравый смысл, чтобы определить несколько

множеств: • Натуральные числа N = $\{1, 2, 3, \ldots\}$

- Целые числа Z = $\{\dots, 2, 1, 0, 1, 2, \dots\}$
- Рациональные числа Q = {a/b : a, b Z} •

Действительные числа R , охватывающие Q , а также иррациональные числа, такие как π и $\qquad 2$ •

Комплексные числа $C = \{a + bi : a, b \in R\}$, где мы думаем, что і удовлетворяет i = 1.

Стоит признать, что наше определение R далеко не строгое. Построение действительных чисел может быть важной темой для практиков криптографических методов, использующих альтернативные системы счисления, но эти тонкости не имеют отношения к данному обсуждению.

Как и с любыми другими наборами, N, Z, Q, R и C можно манипулировать с помощью общих операций для создания новых наборов чисел. В частности, напомним, что мы можем определить «евклидово произведение» двух множеств A и B как

$$A \times B = \{(a, b) : a \quad A \cup b \quad B\}.$$

Мы можем взять степени множеств, написав

$$A^{H} = A \times \underbrace{A \times \cdots \times A}_{\text{n pas}} .$$

¹Это первый из многих случаев, когда мы будем использовать обозначение {A : B}; фигурные скобки должны обозначать набор, а двоеточие можно читать как «такой, что». Например, определение Q можно прочитать как «множество дробей a/b , таких что a и b — целые числа». В качестве второго примера мы могли бы написать N = {n Z : n > 0}.

Эта конструкция дает то, что станет нашим любимым набором чисел в следующих главах:

0.2 Векторные пространства

Вводные курсы по линейной алгебре вполне можно было бы назвать «Введение в конечномерные векторные пространства». Хотя определение векторного пространства может показаться абстрактным, мы найдем множество конкретных приложений, удовлетворяющих формальным аспектам и, таким образом, извлекающих пользу из разработанного нами механизма.

0.2.1 Определение векторных пространств

Начнем с определения векторного пространства и приведем ряд примеров:

Определение 0.1 (Векторное пространство). Векторное пространство — это множество V , замкнутое относительно скалярного умножения и сложения.

Для наших целей скаляр — это число в R, а операции сложения и умножения удовлетворяют обычным аксиомам (коммутативность, ассоциативность и т. д.). Как правило, определить векторные пространства в дикой природе несложно, включая следующие примеры:

Пример 0.1 (Rn как векторное пространство). Наиболее распространенным примером векторного пространства является Rn . Здесь сложение и скалярное умножение происходят покомпонентно:

$$(1, 2) + (3, 4) = (1 3, 2 + 4) = (2, 6)$$

 $10 \cdot (1, 1) = (10 \cdot 1, 10 \cdot 1) = (10, 10)$

Пример 0.2 (многочлены). Вторым важным примером векторного пространства является «кольцо» многочленов с вещественными входными параметрами, обозначаемое R[x]. Многочленом р R[x] называется функция р : R R, принимающая вид2

$$p(x) = \kappa$$

Сложение и скалярное умножение выполняются обычным образом, например, если p(x) = x $^2 + 2x - 1$ и $q(x) = x + 3x^2$, Втбуду фикурь Буб братите $^2 + 6x - 3$, который является еще одним полиномом. Кроме того, внимание, что функции типа p(x) = (x - 1)(x + 1) + x, хотя они явно не 2 (икс 3 for 3 bce еще полиномы, даже записаны в приведенной выше форме.

Элементы v V векторного пространства V называются векторами, а взвешенная сумма вида і aivi , где ai R и vi V, известна как линейная комбинация элементов vi . В нашем втором примере «векторы» — это функции, хотя мы обычно не используем этот язык для обсуждения R[x]. Один k с последовательным способом зрения состоял бы в том, чтобы идентифицировать многочлен k akx (a0, a1, a2, ···); связать эти две точки помните, что многочлены имеют конечное число членов, так что последовательность в конце концов будет заканчиваться цепочкой нулей.

²Обозначение f: A В означает, что f — это функция, которая принимает на вход элемент множества A и выводит элемент множества B. Например, f: R Z принимает на вход вещественное число из R и выводит целое число Z., как может быть в случае f(x) = x, функции «округления в меньшую сторону».

0.2.2 Размах, линейная независимость и основания

Предположим, мы начинаем с векторов v1, . . . , vk V для векторного пространства V. По определению 0.1 у нас есть два способа начать с этих векторов и построить новые элементы V: сложение и скалярное умножение. Идея span заключается в том, что он описывает все векторы, которых вы можете достичь с помощью этих двух операций:

Определение 0.2 (Размах). Оболочкой множества S V векторов называется множество

span S
$$\{a1v1 + \cdots + akvk : k 0, vi V$$
 для всех і и аі R для всех і $\}$.

Обратите внимание, что span S является подпространством V, то есть подмножеством V , которое само по себе является векторным пространством. Мы можем приведите несколько примеров:

Пример 0.3 (Миксология). Типичный «колодец» в коктейль-баре содержит как минимум четыре ингредиента в распоряжении бармена: водка, текила, апельсиновый сок и гренадин. Предполагая, что у нас есть этот простой колодец, мы можем представить напитки в виде точек в R4 с одним слотом для каждого ингредиента. Например, типичный «текила санрайз» может быть представлен точкой (0, 1,5, 6, 0,75), представляющей количество водки, текилы, апельсинового сока и гренадина (в унциях), соответственно.

Набор напитков, которые можно приготовить из типового колодца, содержится в

то есть все комбинации четырех основных ингредиентов. Однако бармен, стремящийся сэкономить время, может не заметить, что во многих напитках одинаковое соотношение апельсинового сока и гренадина, и смешать бутылки. Новый упрощенный колодец может быть проще для розлива, но может приготовить меньше напитков:

Пример 0.4 (многочлены). Определите pk(x) x K. Тогда легко увидеть, что

$$R[x] = диапазон {pk : k 0}.$$

Убедитесь, что вы достаточно хорошо понимаете нотацию, чтобы понять, почему это так.

Добавление другого элемента в набор векторов не всегда увеличивает размер его диапазона. Для Например, в R2 очевидно, что

диапазон
$$\{(1, 0), (0, 1)\}$$
 = диапазон $\{(1, 0), (0, 1), (1, 1)\}$.

В этом случае говорят, что множество {(1, 0), (0, 1), (1, 1)} линейно зависимо:

Определение 0.3 (Линейная зависимость). Мы даем три эквивалентных определения. Множество S V векторов линейно зависимо. если:

- 1. Один из элементов S может быть записан как линейная комбинация других элементов, или S содержит нуль.
- 2. Существует непустая линейная комбинация элементов vk S, дающая ck = 0 для bcex k. dec M=1 ckvk = 0, где
- 3. Существует v S такой, что span S = span S\ $\{v\}$. То есть мы можем удалить вектор из S без влияет на его размах.

Если S не является линейно зависимым, то говорят, что оно линейно независимо.

Предоставление доказательств или неофициальных свидетельств того, что каждое определение эквивалентно своим аналогам (по принципу «если и только если»), является полезным упражнением для студентов, менее знакомых с обозначениями и абстрактной математикой.

Концепция линейной зависимости приводит к идее «избыточности» набора векторов. В этом смысле естественно спросить, насколько большой набор мы можем выбрать, прежде чем добавление другого вектора не может увеличить диапазон. В частности, предположим, что у нас есть линейно независимое множество S V, и теперь мы выбираем дополнительный вектор V V. Добавление v к S приводит к одному из двух возможных результатов:

- 1. Размах S {v} больше размаха S.
- 2. Добавление v к S не влияет на диапазон.

Размерность V — это не что иное, как максимальное количество раз, когда мы можем получить результат 1, добавить v к S и повторить.

Определение 0.4 (размерность и базис). Размерность V — это максимальный размер |S| линейно независимого множества S V такого, что span S = V. Любое множество S, удовлетворяющее этому свойству, называется базисом для V.

Пример 0.5 (Rn). Стандартным базисом для Rn является множество векторов вида

То есть ek имеет все нули, кроме одного в k-м слоте. Ясно, что эти векторы линейно независимы и образуют базис; например, в R3 любой вектор (a, b, c) может быть записан как ae1 + be2 + ce3.

Таким образом, размерность Rn равна n, как и следовало ожидать.

Пример 0.6 (многочлены). Ясно, что множество $\{1, x, x, \ldots\}$ — линейно 2 не 3 а́висимый набор полиномов, охватывающий R[x]. Обратите внимание, что это множество бесконечно велико, и поэтому размерность R[x] равна .

0.2.3 Наш фокус: Rn

Особое значение для наших целей имеет векторное пространство Rn клидово , так называемый n-мерный Eu пространство. Это не что иное, как набор координатных осей, встречающихся на уроках математики в средней школе:

- R1 R числовая прямая
- R2 двумерная плоскость с координатами (x, y) R3 трехмерное

пространство с координатами (x, y, z)

Почти все методы этого курса будут иметь дело с преобразованиями и функциями на Rn.

Для удобства мы обычно записываем векторы в Rn в «столбцовой форме» следующим образом

Эта нотация будет включать векторы как частные случаи матриц, обсуждаемые ниже.

В отличие от некоторых векторных пространств, Rn имеет не только структуру векторного пространства, но и одну дополнительную конструкцию, которая имеет все значение: скалярное произведение.

Определение 0,5 (Скалярный продукт). Скалярное произведение двух векторов a = (a1, ..., an) и b = (b1, ..., bn) в Rn определяется выражением

Пример 0.7 (R2). Скалярное произведение (1, 2) и (2 , 6) равно 1 · 2 + 2 · 6 = 2 + 12 = 10.

Скалярное произведение является примером метрики, и его существование дает Rn понятие геометрии. Например, мы можем использовать теорему Пифагора, чтобы определить норму или длину вектора как квадратный корень

$$\frac{22+\cdots+aa2}{a1} = a \cdot a.$$

Тогда расстояние между двумя точками a,b Rn равно просто b a2.

Скалярные произведения дают не только представления о длинах и расстояниях, но и об углах. Напомним следующее тождество из тригонометрии для a,b R3:

$$a \cdot b = ab \cos \theta$$

где θ — угол между а и b. Однако для n 4 понятие «угол» гораздо труднее визуализировать для Rn . Мы могли θ определить угол θ между а и b как значение θ , заданное выражением

$$\theta$$
 arccos $\frac{a \cdot 6}{a6}$.

Мы должны сделать нашу домашнюю работу, прежде чем давать такое определение! В частности, вспомните, что косинус выводит значения в интервале [1, 1], поэтому мы должны проверить, что входные данные арккосинуса (также обозначаемого как соз 1) находятся в этом интервале; к счастью, известное неравенство Коши-Шварца а ·b аb гарантирует именно это свойство.

Когда a = cb для некоторого с R, мы имеем $\theta = \arccos 1 = 0$, как и следовало ожидать: угол между параллельными векторами равен нулю. Что означает перпендикулярность векторов? Подставим $\theta = 90$ Тогда у нас есть

$$0 = \cos 90$$

$$= \frac{a \cdot 6}{a6}$$

Умножение обеих сторон на аb мотивирует определение:

Определение 0.6 (ортогональность). Два вектора перпендикулярны или ортогональны, когда а \cdot b = 0.

Это определение несколько неожиданно с геометрической точки зрения. В частности, нам удалось определить, что значит быть перпендикулярным, без явного использования углов. Эта конструкция облегчит решение некоторых задач, для которых нелинейность синуса и косинуса могла бы создать головную боль в более простых условиях.

В сторону 0,1. Здесь нужно обдумать множество теоретических вопросов, некоторые из которых мы рассмотрим в следующих главах, когда они станут более мотивированными:

- Все ли векторные пространства допускают скалярные произведения или аналогичные структуры?
- Все ли конечномерные векторные пространства допускают скалярное произведение?
- Каким может быть разумное скалярное произведение между элементами R[x]?

Заинтересованные студенты могут ознакомиться с текстами по реальному и функциональному анализу.

0.3 Линейность

Функция между векторными пространствами, которая сохраняет структуру, известна как линейная функция:

Определение 0.7 (Линейность). Предположим, что V и V — векторные пространства. Тогда L : V — V линейна, если она удовлетворяет следующим двум критериям для всех v,v1,v2 — V и с — R:

- L сохраняет суммы: L[v1 +v2] = L[v1] + L[v2]
- L сохраняет скалярные произведения: L[cv] = cL[v]

Легко генерировать линейные карты между векторными пространствами, как мы можем видеть в следующих примерах:

$$f(x, y) = (3x, 2x + y, y)$$

Мы можем проверить линейность следующим образом:

• Сохранение суммы:

$$f(x1 + x2, y1 + y2) = (3(x1 + x2), 2(x1 + x2) + (y1 + y2), (y1 + y2)) =$$

 $(3x1, 2x1 + y1, y1) + (3x2, 2x2 + y2, y2)$
 $= f(x1, y1) + f(x2, y2)$

• Сохранение скалярного произведения:

$$f(cx, cy) = (3cx, 2cx + cy, cy) = c(3x, 2x + y, y)$$

= c f(x, y)

Напротив, g(x, y) ху2 не является линейным. Например, g(1, 1) = 1, но $g(2, 2) = 8 = 2 \cdot g(1, 1)$, так что эта форма не сохраняет скалярные произведения.

Пример 0.9 (Интеграция). Следующий «функционал» L от R[x] до R является линейным:

Этот несколько более абстрактный пример отображает многочлены p(x) в действительные числа L[p(x)]. Например, мы можем написать

$$_{\Pi[3x]}^2 + x - 1] = \int_0^1 (3x)^2 (1 + x - 1) \, \mu x = 2$$
 -.

Линейность исходит из следующих хорошо известных фактов из исчисления:

Мы можем написать особенно красивую форму для линейных отображенийна Rn. Напомним, что вектор а = (a1, ..., an) равен сумме k akek, где ek — k-й стандартный базисный вектор. Тогда, если L линейно, мы знаем:

Этот вывод показывает следующий важный факт:

L полностью определяется его действием на стандартной основе vectorsek .

То есть для любого вектора а , мы можем использовать приведенную выше сумму, чтобы определить L [а] путем линейного комбинирования Rn L[e1], . . . , L[en].

Пример 0.10 (Расширение линейной карты). Вспомним карту из примера 0.8, заданную формулой f(x, y) = (3x, 2x + y, y). Имеем f(e1) = f(1, 0) = (3, 2, 0) и f(e2) = f(0, 1) = (0, 1, 1). Таким образом, приведенная выше формула показывает:

$$f(x, y) = x f(e1) + y f(e2) = x$$

$$3 2 + y 1$$

$$0 1$$

0.3.1 Матрицы

Расширение линейных карт выше предлагает один из многих контекстов, в которых полезно хранить несколько векторов в одной и той же структуре. В более общем случае, скажем, у нас есть п векторов v1, . . . ,vn Rm. Мы можем записать каждый как вектор-столбец:

$$v1 = \begin{cases} v11 & v12 & v1n \\ v21 & v22 & \cdots, BH = \\ \vdots & \vdots & \ddots & BH = \\ vm1 & v2 & BM2 & BMH \end{cases}$$

Носить их по отдельности может быть громоздко, поэтому для упрощения мы просто объединяем их в единую матрицу размера $m \times n$:

Будем называть пространство таких матриц Rm×n

Пример 0.11 (Идентификационная матрица). Мы можем сохранить стандартный базис для Rn в «единичной матрице» размера n × n. In×n определяется как:

Поскольку мы построили матрицы как удобные способы хранения наборов векторов, мы можем использовать умножение, чтобы выразить, как их можно комбинировать линейно. В частности, матрицу в Rm×n можно умножить на вектор-столбец в Rn следующим образом:

Расширение этой суммы дает следующую явную формулу для произведений матрицы на вектор:

Пример 0.12 (умножение единичной матрицы). Ясно, что для любого х $Rn x = In \times nx$, где $In \times n -$, мы можем написать единичная матрица из примера 0.11.

Пример 0.13 (Линейная карта). Вернемся еще раз к выражению из примера 0.8, чтобы показать еще одну альтернативную форму:

Аналогичным образом мы определяем произведение между матрицей в M Rm×n и другой матрицей в Rn×p путем объединения отдельных произведений матриц-векторов:

Пример 0.14 (Миксология). Продолжая пример 0.3, предположим, что мы делаем текилу восход солнца и вторую смесь с равными частями двух ликеров в нашем упрощенном колодце. Чтобы узнать, сколько базовых ингредиентов содержится в каждом заказе, мы могли бы объединить рецепты для каждого столбца и использовать матрицу. умножение:

	Скважина	1 Скважина 2	Скважина 3	Hagusay 1	Напиток 1 Напиток 2			Напиток 1 Напиток 2		
Водка	1	0	0	панинок			0 0,75 1	,5 6	Водка	
Текила	0	1	0	. 1,5	0,75 0.75	" <u>-</u> "		0,75	Текила	
ОЖ 0		0	6	1,3	2			12	ОЖ	
Гренадин 0		0	0,75	'	۷		0,75	1,5	Гренадин	

В общем, мы будем использовать заглавные буквы для обозначения матриц, например, А Rm×n^{. Мы будем использовать} обозначение Aij R для обозначения элемента A в строке i и столбце j.

0.3.2 Скаляры, векторы и матрицы

Неудивительно, что мы можем записать скаляр как вектор 1 × 1 с $R1 \times 1$. Аналогично, как мы уже предложенный в \$0.2.3, если мы запишем векторы в Rn в виде столбцов, их можно рассматривать как матрицы размера $n \times 1.$ V $Rn \times 1$. Обратите внимание, что в этом контексте можно легко интерпретировать произведения матрицы-вектора; например, если A $Rm \times n$, X Rn, U D Rm, то мы можем записать выражения вида

A
$$_{\text{MKC}} = 6$$

M×П $\Pi \times 1$ $M\times 1$

Мы введем один дополнительный оператор на матрицах, полезный в этом контексте:

Определение 0.8 (Транспонирование). Транспонированием матрицы A $Rm \times n$ называется матрица A $Rn \times m$ с элементами (A)ij = Aji.

Пример 0.15 (Транспозиция). Транспонирование матрицы

дан кем-то

Геометрически транспонирование можно представить как переворачивание матрицы по диагонали.

Эта унифицированная обработка скаляров, векторов и матриц в сочетании с такими операциями, как транспозиция и умножение, может привести к гладким выводам хорошо известных тождеств. Например,

мы можем вычислить скалярное произведение векторов sa,b Rn, выполнив следующую серию шагов:

Многие важные тождества из линейной алгебры могут быть получены путем объединения этих операций с помощью нескольких правил:

$$(A) = A$$

 $(A + B) = A + B$
 $(AB) = BA$

Пример 0.16 (Остаточная норма). Предположим, у нас есть матрица A и два вектора x и b. Если мы хотим знать, насколько хорошо Ax аппроксимирует b, мы можем определить невязку r b Ax; эта невязка равна нулю ровно тогда, когда Ax = b. В противном случае мы могли бы использовать норму r2 в качестве показателя отношения между Ax и b. Мы можем использовать приведенные выше тождества для упрощения:

$$p = \frac{2}{2} = 6 - Ax$$
 $\frac{2}{2}$ $= (b - Ax) \cdot (b - Ax)$, как объяснено в $\S 0.2.3 = (b - Ax)$ $(b - Ax)$ с помощью нашего выражения для скалярного произведения выше $= (b - x - A)$ $(b - Ax)$ по свойствам транспонирования $= bb - bAx - x - AAx$ после умножения

Все четыре члена в правой части являются скалярами или, что эквивалентно, матрицами 1 × 1. Скаляры, рассматриваемые как матрицы, тривиально обладают одним дополнительным приятным свойством с = с, поскольку транспонировать нечего! Таким образом, мы можем написать

$$x A b = (x A b) = b Ax$$

Это позволяет нам еще больше упростить наше выражение:

$$p = \frac{2}{2} = bb = 2b Ax + x A Ax$$

= $TORIOP = \frac{2}{2} = 26 Ax + 6 = \frac{2}{2}$

Мы могли бы получить это выражение, используя тождества скалярного произведения, но промежуточные шаги, описанные выше, окажутся полезными в нашем дальнейшем обсуждении.

0.3.3 Модельная задача: Ax = b

На вводном уроке алгебры студенты тратят много времени на решение линейных систем, таких как следующие для троек (x, y, z):

$$3x + 2y + 5z = 0$$

 $4x + 9y$ $3z = 72x$
 $3y$ $3z = 1$

Наши конструкции в §0.3.1 позволяют нам кодировать такие системы более чистым способом:

В более общем смысле мы можем записать линейные системы уравнений в форме Ах = b, следуя той же схеме, что и выше; здесь вектор х неизвестен, а A и b известны. Такая система уравнений не всегда имеет решение. Например, если A содержит только нули, то очевидно, что ни один х не будет удовлетворять Ax = b всякий раз, когда b = 0. Мы отложим общее рассмотрение того, когда решение существует, до нашего обсуждения линейных решателей в следующих главах.

Ключевая интерпретация системы Ax =b заключается в том, что она решает задачу:
Запишите b как линейную комбинацию столбцов матрицы A.

Почему? Напомним из §0.3.1, что произведение Ах кодирует линейную комбинацию столбцов А с весами, содержащимися в элементах х. Итак, уравнение Ах = b требует, чтобы линейная комбинация Ах равнялась заданному вектору b. Учитывая эту интерпретацию, мы определяем пространство столбца А как пространство правых частей b, для которых система имеет решение:

Определение 0.9 (столбцовое пространство). Пространство столбцов матрицы A Rm×n является оболочкой столбцов матрицы A. Мы можем записать как

$$col A \{Ax:x Rn\}.$$

Несколько легче рассмотреть один важный случай. Предположим, что A квадрат, поэтому мы можем написать A Rn×n · Кроме того, предположим, что система Ax = b имеет решение для всех вариантов b. Таким образом, согласно единственное условие на b состоит в том, что он является членом · нашей вышеприведенной интерпретации Ax =b, мы можем Rn , чтобы сделать вывод, что столбцы A охватывают Rn ·

В этом случае, поскольку линейная система всегда разрешима, предположим, что мы подставляем стандартный базис e1, . . . ,en для получения векторов x1, . . . ,xn , удовлетворяющих условию Axk = ek для каждого k. Затем мы можем «сложить» эти выражения, чтобы показать:

где In×n — единичная матрица из примера 0.11. Будем называть матрицу со столбцами xk обратной A 1

$$AA-1 = A -1A = In \times n$$
.

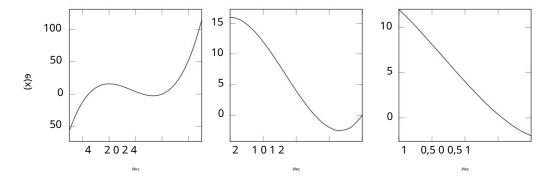


Рисунок 1: Чем ближе мы приближаемся к f(x) = x 8x + 4, тем больше это похоже на линию.

Также легко проверить, что (A = A. Когда такая обратная существует, легко решить систему Ax = b. В частности, мы находим:

0.4 Нелинейность: дифференциальное исчисление

В то время как красота и применимость линейной алгебры делает ее ключевым объектом изучения, нелинейности изобилуют в природе, и нам часто приходится проектировать вычислительные системы, которые могут справиться с этим фактом жизни. В конце концов, на самом базовом уровне квадрат в знаменитом соотношении E = mc2 делает его мало поддающимся линейному анализу.

0.4.1 Дифференциация

Хотя многие функции глобально нелинейны, локально они демонстрируют линейное поведение. Эта идея «локальной линейности» является одним из основных мотивов дифференциального исчисления. Например, на рисунке 1 показано, что если вы достаточно приблизите гладкую функцию, в конечном итоге она станет похожа на линию. Производная f (x) функции f (x): R R есть не что иное, как наклон аппроксимирующей линии, вычисляемый путем нахождения наклона линий, проходящих через все более и более близкие к х точки:

= lim
$$y \times x \frac{f(y) - f(x) f(x)}{y - x}$$

Мы можем выразить локальную линейность, записав $f(x + x) = f(x) + x \cdot f(x) + O(x)$

Если функция f принимает несколько входных данных, то ее можно записать как f(x): Rn R для x Rn; иными словами, каждой точке x = (x1, ..., xn) в n-мерном пространстве f ставит в соответствие единственное число f(x1, ..., xn). Наше представление о локальной линейности здесь несколько нарушается, потому что линии — это одномерные объекты. Однако фиксация всех переменных, кроме одной, сводится к случаю исчисления одной переменной. Например, мы могли бы написать g(t) = f(t, x2, ..., xn), где мы просто фиксируем константы x2, ..., xn. Тогда g(t) — дифференцируемая функция одной переменной. Конечно, мы могли бы поместить t в любой из входных слотов для t, поэтому в общем случае мы даем следующее определение частной производной от t:

f Определение 0.10 (частная производная). k-я частная производная f , обозначенная как xk , — дается по-разному связывающая f в своей k-й входной переменной:

$$\frac{\Gamma}{dt}$$
 (x1, ..., xn) xk $\frac{\Gamma}{-}$ f(x1, ..., xk 1, t, xk+1, ..., xn)|t=xk

Обозначение « t=xk » следует читать как «оценивается при t=xk ».

Пример 0.17 (Относительность). Соотношение E = mc2 можно рассматривать как функцию от m и c к E. Таким образом, мы могли бы написать E(m, c) = mc2, что дало бы производные

$$\frac{E}{M} = 2 = C$$

$$\frac{E}{M} = 2MKC$$

Используя исчисление с одной переменной, мы можем написать:

$$f(x + x) = f(x1 + x1, x2 + x2, ..., xn + xn)$$

$$= f(x1, x2 + x2, ..., xn + xn) + f(x + x1 + x1) = f(x1, x2 + x2, ..., xn + xn) + f(x1 + x1 + x1) = f(x1, ..., xn) + f(x1 + x1 + x1) = f(x1, ..., xn) + f(x1 + x1 + x1) = f(x1, ..., xn) + f(x1 + x1 + x1) = f(x1, ..., xn) + f(x1 + x1 + x1) = f(x1 + x1 +$$

где мы определяем градиент f как

$$\varphi \qquad \frac{}{f \quad x1}, \frac{}{f \quad x2}, \cdots \quad , \frac{}{xn_{-}} \quad e \; R \quad ^H$$

Из этого соотношения легко видеть, что f можно дифференцировать по любому направлению v; мы можем оценить эту производную Dv f следующим образом:

Dv f(x)
$$f(x-tv) | t=0 dt =$$

$$f(x) \cdot v$$

Пример 0.18 (R2). Возьмем f(x, y) = x 2y 3. Затем,

$$-f = 2xy3$$

$$\frac{x}{3x} \quad f = 2y^2$$

Таким образом, мы можем написать $f(x, y) = (2xy3 \ 3x \ 2y \ 2)$. Производная f в точке (1, 2) в направлении (1, 4) определяется выражением (1, 4) · $f(1, 2) = (1, 4) \cdot (16, 12) = 32$.

Пример 0.19 (Линейные функции). Очевидно, но стоит отметить, что градиент функции f(x) $a \cdot x + c = (a1x1 + c1, ..., anxn + cn)$ равен .

Пример 0.20 (Квадратичные формы). Возьмите любую матрицу A , и определить f(x) х Ax. Расширение $Rn \times n$, эта функция поэлементно показывает

$$e(x) = Aijxixj;$$

Стоит расширить f и явно проверить это отношение. Возьмем некоторое k $\{1, \ldots, n\}$. Затем мы можем выделить все термины, содержащие xk:

С помощью этой факторизации легко увидеть

$$\frac{-\phi}{xk}$$
 = 2Akkxk + $\underset{g=\kappa}{\underset{g=\kappa}{\text{Aйкси}}}$ Akjxj

Эта сумма есть не что иное, как определение умножения матрицы на вектор! Таким образом, мы можем написать

$$f(x) = Ax + Ax$$
.

Мы обобщили f: R R на f: Rn R. Чтобы достичь полной общности, мы хотели бы рассмотреть f: Rn Rm. Другими словами, f принимает n чисел и выводит m чисел. K счастью, это расширение является простым, потому что мы можем думать о f как о наборе однозначных функций f1, . . . , fm: Rn R слились в один вектор. То есть пишем:

$$e(x) = \begin{cases} f1(x) \\ f2(x) \\ \vdots \\ \phi M(x) \end{cases}$$

Каждую fk можно дифференцировать, как и раньше, поэтому в итоге мы получаем матрицу частных производных, называемую якобианом f:

Определение 0.11 (якобиан). Якобианом f : Rn — Rm называется матрица D f — Rm×n с элементами

(Df)ij
$$\frac{\Phi^{N}}{Xj}$$
.

Пример 0.21 (Простая функция). Предположим, что f(x, y) = (3x, xy2, x + y). Затем,

D f (x, y) =
$$\begin{pmatrix} 3 \\ y^2 \\ 1 \end{pmatrix}$$
 0 2xy 1

Убедитесь, что вы можете произвести это вычисление вручную.

Пример 0.22 (умножение матриц). Неудивительно, что якобиан функции f(x) = Ax для матрицы A равен D f(x) = A.

Здесь мы сталкиваемся с общей точкой путаницы. Предположим, что функция имеет векторный вход и скалярный выход, то есть f: Rn R. Мы определили градиент f как вектор-столбец, поэтому, чтобы согласовать это определение с определением якобиана, мы должны написать

D ж знак равно ж.

0.4.2 Оптимизация

Напомним, что из исчисления одной переменной минимумы и максимумы f: R R должны встречаться в точках x, удовлетворяющих f (x) = 0. Конечно, это условие скорее необходимое, чем достаточное: могут существовать точки x c f (x) = 0 которые не являются максимумами или минимумами. Тем не менее, нахождение таких критических точек f может быть шагом алгоритма минимизации функции, если следующий шаг гарантирует, что результирующий x действительно является минимумом/максимумом.

Если f : Rn R минимизируется или максимизируется в точке x, мы должны гарантировать, что не существует единственного направления x от x, в котором f уменьшается или увеличивается, соответственно. Согласно обсуждению в \$0.4.1, это означает, что мы должны найти точки, для которых f = 0.

Пример 0.23 (Простая функция). Предположим, что f(x, y) = x 2x + 8y. 2 + 2xy + 4y $2_{\text{и x y}} - x \times x \cdot \text{Тогда} = 2x + 2y$. Таким образом, критические точки f удовлетворяют:

$$2x + 2y = 0 2x +$$

$$8y = 0$$

Ясно, что эта система решается при (x, y) = (0, 0). В самом деле, это минимум f, что можно увидеть более четко, написав f(x, y) = (x + y) 2 + 3 roda 2.

Пример 0.24 (Квадратичные функции). Предположим, что f(x) = x Ax + bx + c. Тогда из примеров в предыдущем разделе мы можем написать f(x) = (A + A)x + b. Таким образом, критические точки x функции f удовлетворяют условию (A + A)x + b = 0.

В отличие от исчисления с одной переменной, когда мы выполняем исчисление на Rn , мы можем добавлять ограничения к нашей оптимизации. Наиболее общий вид такой задачи выглядит так:

минимизировать

f(x) так, чтобы g(x) = 0

Пример 0.25 (прямоугольные области). Предположим, прямоугольник имеет ширину w и высоту h. Классическая задача геометрии состоит в том, чтобы максимизировать площадь с фиксированным периметром 1:

максимизировать wh

такое, что 2w + 2h - 1 = 0

Когда мы добавляем это ограничение, мы больше не можем ожидать, что критические точки удовлетворяют f(x) = 0, поскольку эти точки могут не удовлетворять g(x) = 0.

Пока предположим, что g: Rn R. Рассмотрим множество точек S0 $\{x: g(x) = 0\}$. Очевидно, для любых двух x,у S0 выполняется соотношение g(y) g(x) = 0 0 = 0. Предположим, что y = x + x для малых x. Тогда g(y) $g(x) = g(x) \cdot x + O(x)$. Другими словами, если мы начнем c x, удовлетворяющего g(x) = 0, затем, если мы сместим в направлении x, $g(x) \cdot x = 0$, чтобы продолжать удовлетворять этому соотношению.

Теперь вспомним, что производная от f по направлению v в точке x задается как $f \cdot v$. Если x является минимумом задачи оптимизации c ограничениями, описанной выше, то любое небольшое смещение x на x + v должно вызывать увеличение от f(x) до f(x + v). Поскольку нас интересуют только перемещения v, сохраняющие ограничение g(x + v) = c, из нашего рассуждения выше мы хотим, чтобы $f \cdot v = 0$ для всех v, удовлетворяющих $g(x) \cdot v = 0$. Другими словами, f и g должны быть параллельны, условие, которое мы можем записать как $f = \lambda$ g для некоторого λ R.

Определять

$$\Lambda(x, \lambda) = f(x)$$
 $\lambda g(x)$.

Тогда критические точки Л без ограничений удовлетворяют условию:

$$0 = \frac{\Lambda}{\lambda} = -\Gamma(x)$$

$$0 = x\Lambda = f(x) \quad \lambda \quad g(x)$$

Другими словами, критические точки Λ удовлетворяют условиям g(x) = 0 и $f(x) = \lambda$ g(x) — именно тем условиям оптимальности, которые мы получили!

Расширение до многомерных ограничений дает следующее:

Теорема 0.1 (метод множителей Лагранжа). Критические точки приведенной выше задачи оптимизации с ограничениями являются неограниченными критическими точками функции множителя Лагранжа.

$$\Lambda(x,\lambda)$$
 $f(x)$ $\lambda \cdot g(x)$,

как по х, так и по λ.

Пример 0.26 (Максимизация площади). Продолжая пример 0.25, определим функцию множителя Лагранжа Λ (w, h, λ) = wh λ (2w + 2h 1). Дифференцируя, находим:

$$w \frac{\Lambda}{\Lambda} = h \quad 2\lambda = 0$$

$$= w \frac{0}{2\lambda}$$

$$0 = \frac{h}{\lambda} \quad \Lambda = 1 \quad 2w \quad 2h$$

Итак, критические точки системы удовлетворяют

Решение системы показывает w = h = 1/4 и $\lambda = 1/8$. Другими словами, при фиксированном периметре прямоугольник с максимальной площадью является квадратом.

Пример 0.27 (собственные задачи). Предположим, что A — симметричная положительно определенная матрица, что означает A = A (симметрия) и x Ax > 0 для всех x $Rn\{0\}$ (положительно определенная). Часто мы хотим минимизирова $\frac{2}{2}$ ь x Ax при x = 1 для данной матрицы A $Rn \times n$; обратите внимание, что без ограничения минимум тривиально имеет место при x = 0. Определим функцию множителя Лагранжа

$$\Lambda(x, \lambda) = x Ax \qquad \lambda(x \qquad \qquad \frac{2}{2}$$

$$1) = x Ax \qquad \lambda(x x \qquad 1).$$

Дифференцируя по х, находим

$$0 = x\Lambda = 2Ax \quad 2\lambda x$$

Другими словами, х является собственным вектором матрицы А:

$$Ax = \lambda x$$
.

0.5 Проблемы

Проблема 0.1. В качестве C1 (R) возьмем множество функций f: R R, допускающих первую производную f(x). Почему $C_{+}(R)$ векторное пространство? Докажите, что C1 (R) имеет размерность .

Задача 0.2. Предположим, что строки A Rm×n заданы транспонированием r1, . . . ,rm Rn , а столбцы A Rm×n задаются как c1, . . . ,cn Rm. To есть,

Выразите элементы АА и АА через эти векторы.

Задача 0.3. Приведите линейную систему уравнений, которой удовлетворяют минимумы энергии f(x) = Ax b относительно x для x Rn и b Rn ВЗТА истема называется «нормальными уравнениями» и будет встречаться в других местах в этих примечаниях; даже в этом случае стоит проработать и полностью понять вывод.

Задача 0.4. Предположим, что A, B Rn×n . Сформулируйте условие того, что векторы х Rn являются критическими $\frac{1}{2}$ точками при Bx = $\frac{2}{3}$. Также приведите альтернативную форму для оптимальных значений Ax $\frac{2}{2}$.

Задача 0.5. Зафиксируем некоторый вектор а $Rn\{0}$ и определим $f(x) = a \cdot x$. Напишите выражение для максимума f(x) при x = 1.

Machine Translated by Google

Глава 1

Числа и анализ ошибок

Изучая численный анализ, мы переходим от работы с целыми и длинными числами к числам с плавающей запятой и удвоениям. Этот, казалось бы, невинный переход включает в себя огромный сдвиг в том, как мы должны думать о разработке и реализации алгоритма микрофона. В отличие от основ дискретных алгоритмов, мы больше не можем ожидать наши алгоритмы дают точные решения во всех случаях. «Большое О» и подсчет операций не всегда царствовать безраздельно; вместо этого, даже в понимании самых основных методов, мы вынуждены изучать компромисс между временем, ошибкой аппроксимации и так далее.

1.1 Хранение чисел с дробными частями

Напомним, что компьютеры обычно хранят данные в двоичном формате. В частности, каждая цифра положительного целое число соответствует другой степени двойки. Например, мы можем преобразовать 463 в двоичный код. используя следующую таблицу:

Другими словами, эта запись кодирует тот факт, что 463 можно разложить на степени двойки. однозначно как:

$$463 = 2 ^{87} \frac{10}{9} \frac{2}{2} + 2 + 2 + 2 + 2 + 2 + 2$$
$$= 256 + 128 + 64 + 8 + 4 + 2 + 1$$

Помимо проблем с переполнением, все положительные целые числа могут быть записаны в этой форме с использованием конечного числа цифры. Отрицательные числа также могут быть представлены таким образом, либо путем введения ведущего знака немного или с помощью трюка «дополнение до двух».

Такое разложение вдохновляет на простое расширение чисел, содержащих дроби: просто включают отрицательные степени двойки. Например, разложить 463,25 так же просто, как сложить два слоты:

1	1	10	013	2	1	1	1. 0	0 2	1	
8 2	72	62	5 2	42	22	12		2	2	2

Однако, как и в десятичной системе, представление дробных частей чисел таким образом не так хорошо себя ведет, как представление целых чисел. Например, запись дроби 1/3 в двоичном формате дает выражение:

$$\frac{1}{3}$$
 = 0,0101010101 . . .

Такие примеры показывают, что во всех масштабах существуют числа, которые нельзя представить с помощью конечная двоичная строка. На самом деле такие числа, как π = 11.00100100001...2 имеют бесконечно длинные расширения независимо от того, какую (целочисленную) базу вы используете!

По этой причине при разработке вычислительных систем, выполняющих математические операции с R вместо Z, мы вынуждены делать приближения почти для любого достаточно эффективного числового представления.

Это может привести ко многим путаницам при кодировании. Например, рассмотрим следующий C++
фрагмент:

```
двойной x = 1,0;
двойной y = x/3,0;
if (x == y *3.0) cout << Они равны! ";
иначе cout << Они НЕ равны. ;
```

Вопреки интуиции, эта программа выводит «Они НЕ равны». Почему? Определение у делает приближение к 1/3, поскольку его нельзя записать в виде завершающей двоичной строки, округляя к ближайшему числу, которое он может представлять. Таким образом, у*3.0 больше не умножает 3 на 1/3. Один из способов исправить эта проблема ниже:

```
двойной x = 1,0;
двойной y = x/3,0;
if ( fabs (x - y *3.0) < numeric_limits < double >:: epsilon ) cout << else cout << Они НЕ равны. "Они равны! ";
```

Здесь мы проверяем, что х и у*3.0 находятся в пределах некоторого допуска друг друга, а не проверяем точное равенство. Это пример очень важного момента:

Редко, если вообще когда-либо, оператор == и его эквиваленты должны использоваться для дробных значений.

Вместо этого следует использовать некоторый допуск, чтобы проверить, равны ли числа.

Конечно, здесь есть компромисс: размер допуска определяет грань между равенством и «близкий, но не такой же», который необходимо тщательно выбирать для данного приложения.

Ниже рассмотрим несколько вариантов представления чисел на компьютере.

1.1.1 Представления с фиксированной точкой

Самый простой вариант хранения дробей — добавить фиксированную десятичную точку. То есть как в В приведенном выше примере мы представляем значения, сохраняя коэффициенты 0/1 перед степенями двойки, которые варьироваться от 2 k до 2 для некоторого k, Z. Например, представление всех неотрицательных значений между 0 и 127,75 с шагом 1/4 так же просто, как взять k = 2 и = 7; в этой ситуации мы представляем эти значения используют 9 двоичных разрядов, два из которых идут после десятичной точки.

Основное преимущество этого представления состоит в том, что почти все арифметические операции могут быть осуществляется по тем же алгоритмам, что и с целыми числами. Например, легко видеть, что

$$a + b = (a \cdot 2 \qquad \qquad ^{K} + 6 \cdot 2 \qquad ^{K}) \cdot 2 \qquad ^{K}.$$

Умножение нашего фиксированного представления на 2k гарантирует, что результат будет целым, поэтому это наблюдение по существу показывает, что сложение может быть выполнено с использованием целочисленного сложения, по существу, «игнорируя» десятичную точку. Таким образом, вместо того, чтобы использовать специализированное оборудование, ранее существовавшее целочисленное арифметико-логическое устройство (АЛУ) быстро выполняет математические операции с фиксированной точкой.

Арифметика с фиксированной точкой может быть быстрой, но она может страдать от серьезных проблем с точностью. В частности, часто бывает так, что вывод двоичной операции, такой как умножение или деление, может потребовать битов больше, чем операндов. Например, предположим, что мы включили одну десятичную точку точности и

хотите выполнить произведение $1/2 \cdot 1/2 = 1/4$. Мы пишем $0,12 \times 0,12 = 0,012$, что усекается до 0. В этой системе довольно просто разумно комбинировать числа с фиксированной точкой и получить неразумный результат.

Из-за этих недостатков большинство основных языков программирования по умолчанию не включают десятичный тип данных с фиксированной запятой. Однако скорость и регулярность арифметических операций с фиксированными точками может быть значительным преимуществом для систем, которые отдают предпочтение времени, а не точности. Фактически, некоторые младшие графические процессоры (GPU) реализуют только эти операции, поскольку для многих графических приложений достаточно точности в несколько десятичных знаков.

1.1.2 Представления с плавающей запятой

Одной из многих числовых проблем при написании научных приложений является разнообразие масштабов, которые могут появиться. Только химики имеют дело со значениями где-то между 9,11 × 10–31 и 6,022 × 1023. Такая невинная операция, как смена единиц измерения, может вызвать внезапный переход между этими режимами: одно и то же наблюдение, записанное в килограммах на световой год, будет выглядеть значительно иначе в мегатоннах. в секунду. Наша работа как числовых аналитиков состоит в том, чтобы написать программное обеспечение, которое может изящно переходить между этими шкалами, не накладывая на клиента неестественные ограничения на их методы.

обсуждение. Во-первых, очевидно, что одно из следующих представлений более компактно, чем другое:

Несколько понятий и наблюдений из искусства научных измерений имеют отношение к такому

$$6,022 \times 1023 = 602, 200, 000, 000, 000, 000, 000, 000$$

Кроме того, в отсутствие исключительного научного оборудования разница между $6,022 \times 1023$ и $6,022 \times 1023 + 9,11 \times 10$ –31 незначительна. Один из способов прийти к этому заключению — сказать, что $6,022 \times 1023$ имеет только три цифры точности и, вероятно, представляет некоторый диапазон возможных измерений $[6,022 \times 1023 - \epsilon, 6,011 \times 1023 + \epsilon]$ для некоторого $\epsilon = 0,001 \times 1023$.

Наше первое наблюдение позволило компактифицировать наше представление 6,022 × 1023 , записав его в экспоненциальной записи. Эта система счисления отделяет «интересные» цифры числа от его порядка величины, записывая его в виде а × 10b для некоторых а 1 и b Z. Мы называем этот формат формой числа с плавающей запятой, потому что в отличие от установки с фиксированной запятой в §1.1.1, здесь десятичная точка «плавает» вверх. Мы можем описать системы с плавающей запятой, используя несколько параметров (СІТЕ):

- База β N; для научного обозначения, описанного выше, основание равно 10.
- Точность р N , представляющая количество цифр в десятичной записи.
- Диапазон показателей [L, U] , представляющих возможные значения b

Такое расширение выглядит так:

$$\pm$$
 (d0 + d1 · β 1 + d2 · β 2 + · · · + dp 1 · β 1 1 р) × β 6 экспонент.

где каждая цифра dk лежит в диапазоне [0, β 1] и b [L, U].

Представления с плавающей запятой обладают любопытным свойством, которое может неожиданным образом повлиять на программное обеспечение: интервалы между ними неравномерны. Например, количество значений, которыем кеутковыть в которыем кеутковы по пределим машинную точность стак как

наименьшее $\epsilon m > 0$ такое, что 1 + ϵm представимо. Тогда такие числа, как β + ϵm , невозможно выразить в системе счисления, потому что ϵm слишком мало!

На сегодняшний день наиболее распространенным стандартом для хранения чисел с плавающей запятой является стандарт IEEE 754. Этот стандарт определяет несколько классов чисел с плавающей запятой. Например, число с плавающей запятой двойной точности записывается с основанием β = 2 (как и большинство чисел на компьютере), с одним битом знака ± , 52 цифрами для d и диапазоном показателей степени от -1022 до 1023. Стандарт также определяет, как хранить ± и такие значения, как NaN или «не-число», зарезервированные для результатов вычислений, таких как 10/0. Дополнительную точность можно получить, записав нормализованные значения с плавающей запятой и предполагая, что старшая значащая цифра d0 равна 1, а не записывая ее.

Стандарт IEEE также включает в себя согласованные варианты работы с конечным числом значений, которые могут быть представлены заданным конечным числом битов. Например, распространенная несмещенная стратегия округления вычислений — округление до ближайшего, привязки до четности, что нарушает равноудаленные привязки путем округления до ближайшего значения с плавающей запятой с четным наименее значащим (самым правым) битом. Обратите внимание, что существует много одинаково законных стратегий округления; выбор одного из них гарантирует, что научное программное обеспечение будет одинаково работать на всех клиентских машинах, реализующих один и тот же стандарт.

1.1.3 Дополнительные экзотические опции

Двигаясь вперед, мы будем предполагать, что десятичные значения хранятся в формате с плавающей запятой, если не указано иное. Это, однако, не означает, что других систем счисления не существует, и для конкретных приложений может потребоваться альтернативный выбор. Мы признаем некоторые из этих ситуаций здесь.

Головная боль, связанная с добавлением допусков для учета ошибок округления, может оказаться неприемлемой для некоторых приложений. Такая ситуация возникает в приложениях вычислительной геометрии, например, когда различие между почти и полностью параллельными линиями может быть трудно определить. Одним из решений может быть использование арифметики произвольной точности, то есть реализация арифметики без округления или ошибок любого рода.

Арифметика произвольной точности требует специальной реализации и тщательного рассмотрения того, какие типы значений вам нужно представлять. Например, может быть так, что для данного приложения достаточно рациональных чисел Q , которые можно записать в виде отношений a/b для a, b

Z.

Основные арифметические операции можно выполнять в Q без потери точности. Например, легко увидеть

Арифметика в рациональных числах исключает существование оператора квадратного корня, поскольку такие значения, как 2, иррациональны. Кроме того, это представление неоднозначно, поскольку, например, a/b = 5a/5b.

В других случаях может быть полезно заключить ошибку в скобки, представив значения вместе с оценками ошибки в виде пары а,є R; мы думаем о паре (a,є) как о диапазоне a ± є. Затем арифметические операции также обновляют не только значение, но и оценку ошибки, как в

$$(x \pm \epsilon 1) + (y \pm \epsilon 2) = (x + y) \pm (\epsilon 1 + \epsilon 2 + ошибка(x + y)),$$

где последний член представляет собой оценку ошибки, вызванной добавлением х и у.

1.2 Понимание ошибки

За исключением систем произвольной точности, описанных в § 1.1.3, почти каждое компьютеризированное представление действительных чисел с дробными частями вынуждено использовать округление и другие схемы приближения. Эта схема представляет собой один из многих источников приближений, обычно встречающихся в числовых системах:

- Ошибка усечения возникает из-за того, что мы можем представить только конечное подмножество всего возможного набора значений в R; например, мы должны обрезать длинные или бесконечные последовательности после десятичной точки до количества битов, которое мы хотим сохранить.
- Ошибка дискретизации возникает из-за нашей компьютеризированной адаптации исчисления, физики и других аспекты непрерывной математики. Например, мы делаем аппроксимацию

$$\frac{y(x+e)-y(x)}{\epsilon}$$
.

Мы узнаем, что это приближение является законным и полезным, но в зависимости от выбора ε оно может быть не совсем правильным.

- Ошибка моделирования возникает из-за неполного или неточного описания проблем, которые мы хотим решить. Например, моделирование для прогнозирования погоды в Германии может пренебречь коллективным взмахом крыльев бабочек в Малайзии, хотя смещения воздуха этими бабочками может быть достаточно, чтобы несколько нарушить погодные условия в других местах.
- Ошибка эмпирической константы возникает из-за плохого представления физических или математических констант. Например, мы можем вычислить π , используя последовательность Тейлора, которую мы обрываем раньше, и даже ученые могут даже не знать скорость света с точностью до некоторого числа цифр.
- Ошибка ввода может возникать из-за аппроксимаций параметров данной системы, созданных пользователем (и из-за опечаток!). Моделирование и численные методы могут использоваться для ответа на вопросы типа «что, если», в которых исследовательский выбор входных настроек выбирается только для того, чтобы получить некоторое представление о том, как ведет себя система.

Пример 1.1 (Вычислительная физика). Предположим, мы разрабатываем систему для имитации вращения планет вокруг Земли. Система по существу решает уравнение Ньютона F = ma путем интегрирования сил вперед во времени. Примеры источников ошибок в этой системе могут включать:

- Ошибка усечения: использование IEEE с плавающей запятой для представления параметров и выходных данных системы и усечение при вычислении результата.
- Ошибка дискретизации: замена ускорения а разделенной разностью
- Ошибка моделирования: Пренебрежение моделированием влияния Луны на движение Земли внутри планетарной орбиты. система
- Эмпирическая ошибка: ввод массы Юпитера только до четырех цифр.
- Ошибка ввода: пользователь может захотеть оценить стоимость отправки мусора в космос, а не рисковать накоплением в стиле ВАЛЛ-И на Земле, но может только оценить количество мусора, которое правительство готово выбросить за борт таким образом.

1.2.1 Ошибка классификации

Учитывая наше предыдущее обсуждение, следующие два числа можно рассматривать как имеющие одинаковую величину потенциальной ошибки:

 1 ± 0.01

 105 ± 0.01

Хотя он имеет размер в виде диапазона [1 - 0,01, 1 + 0,01], диапазон [105 - 0,01, 105 + 0,01], по-видимому, кодирует более надежное измерение, поскольку ошибка 0,01 намного меньше по отношению к 105, чем к 1.

Различие между этими двумя классами ошибок описывается путем разграничения между ab растворимая ошибка и относительная ошибка:

Определение 1.1 (Абсолютная ошибка). Абсолютная погрешность измерения определяется разницей между приблизительным значением и лежащим в его основе истинным значением.

Определение 1.2 (Относительная ошибка). Относительная ошибка измерения определяется делением абсолютной ошибки на истинное значение.

Один из способов различить эти два вида ошибок — использовать единицы измерения, а не проценты.

Пример 1.2 (Абсолютная и относительная погрешность). Вот два эквивалентных утверждения в противоположных формах:

Абсолют: 2 дюйма ± 0,02 дюйма

Относительный: 2 дюйма ± 1%

В большинстве приложений истинное значение неизвестно; в конце концов, если бы это было не так, использование приближения вместо истинного значения могло бы быть сомнительным предложением. Есть два популярных способа решить эту проблему. Во-первых, просто быть консервативным при проведении вычислений: на каждом шаге брать максимально возможную оценку ошибки и распространять эти оценки вперед по мере необходимости. Такие консервативные оценки хороши тем, что, когда они малы, мы можем быть уверены, что наше решение полезно.

Альтернативное разрешение связано с тем, что вы можете измерить. Например, предположим, что мы хотим решить уравнение f(x) = 0 для x с заданной функцией f : R R. Мы знаем, что где-то существует корень x0, точно удовлетворяющий f(x0) = 0, но если бы мы знали это гооt наш алгоритм не был бы необходим в первую очередь. На практике наша вычислительная система может дать некоторый xest , удовлетворяющий $f(xest) = \varepsilon$ для некоторого ε c |ε| 1. Возможно, мы не сможем оценить разность x0 xest, поскольку x0 неизвестно. С другой стороны, просто оценивая f, мы можем вычислить f(xest) f(x0) f(xest), поскольку мы знаем, что f(x0) = 0 по определению. Это значение дает некоторое представление об ошибке для нашего расчета.

Этот пример иллюстрирует различие между прямой и обратной ошибкой. Прямая ошибка, создаваемая аппроксимацией, скорее всего, определяет нашу интуицию для анализа ошибок как разницу между аппроксимированным и фактическим решением, но, как мы обсуждали, ее не всегда возможно вычислить. Однако обратная ошибка отличается тем, что ее можно вычислить, но она не является нашей точной целью при решении данной проблемы. Мы можем корректировать наше определение и интерпретацию обратной ошибки по мере того, как мы подходим к различным проблемам, но одно подходящее, хотя и расплывчатое, определение выглядит следующим образом:

Определение 1.3 (обратная ошибка). Обратная ошибка определяется величиной, которую пришлось бы изменить в постановке задачи, чтобы реализовать заданное приближение к ее решению.

Это определение несколько бестолковое, поэтому мы проиллюстрируем его использование на нескольких примерах.

Пример 1.3 (Линейные системы). Предположим, мы хотим решить линейную систему размера n × n Ax = b. Назовем истинным решением x0 A 1b. На самом деле из-за ошибки усечения и других проблем наша система дает почти решение xest. Прямая ошибка этого приближения, очевидно, измеряется с помощью разности xest x0; на практике это значение невозможно вычислить, так как мы не знаем x0. На самом деле xest — это точное решение модифицированной системы Ax = лучший из лучших Axest; таким образом, мы могли бы измерить обратную ошибку с точки зрения разности b-best. В отличие от прямой ошибки, эту ошибку легко вычислить, не обращая A, и легко видеть, что xest является решением проблемы именно тогда, когда обратная (или прямая) ошибка равна нулю.

Пример 1.4 (Решение уравнений, СІТЕ). Предположим, мы пишем функцию для нахождения квадратных корней из положительных чисел, которая выдает 2 1,4. Прямая ошибка |1,4 1,41421 ··· | 0,0142. Обратите внимание, что 1,42 = 1,96, поэтому обратная ошибка равна |1,96 2 | = 0,04.

Два приведенных выше примера демонстрируют более широкую закономерность, согласно которой обратная ошибка может быть намного проще вычислена, чем прямая. Например, для оценки прямой ошибки в примере 1.3 требовалось инвертировать матрицу А, тогда как для оценки обратной ошибки требовалось только умножение на А. Аналогично, в примере 1.4 переход от прямой ошибки к обратной ошибке заменил вычисление квадратного корня умножением.

1.2.2 Кондиционирование, стабильность и точность

Почти в любой численной задаче нулевая обратная ошибка подразумевает нулевую прямую ошибку, и наоборот.

Таким образом, часть программного обеспечения, предназначенного для решения такой проблемы, безусловно, может прекратить работу, если обнаружит, что решение-кандидат имеет нулевую обратную ошибку. Но что, если обратная ошибка не равна нулю, но мала?

Обязательно ли это означает небольшую ошибку вперед? Такие вопросы мотивируют анализ большинства численных методов,

целью которых является минимизация прямой ошибки, но на практике они могут измерять только обратную ошибку.

Мы хотим проанализировать изменения обратной ошибки по сравнению с прямой ошибкой, чтобы наши алгоритмы могли с уверенностью сказать, используя только обратную ошибку, что они дали приемлемые решения.
Это соотношение может быть разным для каждой проблемы, которую мы хотим решить, поэтому в итоге мы делаем следующую грубую классификацию:

- Проблема нечувствительна или хорошо обусловлена, когда небольшое количество обратных ошибок подразумевает небольшое количество прямых ошибок. Другими словами, небольшое изменение постановки хорошо обусловленной задачи дает лишь небольшое изменение истинного решения.
- Проблема чувствительна или плохо обусловлена, если это не так.

Пример 1.5 (ах = 6). Предположим в качестве игрушечного примера, что мы хотим найти решение х0 b/а линейного уравнения ах = b для а, x, b R. Прямая ошибка потенциального решения х определяется как x x0, а обратная ошибка равна определяется как b ax = a(x x0). Итак, когда |a| 1, задача хорошо обусловлена, так как малые значения обратной ошибки a(x x0) влекут за собой еще меньшие значения x x0; и наоборот, когда |a| 1 проблема плохо обусловлена, так как даже если a(x x0) мало, прямая ошибка x x0 1/a · a(x x0) может быть большой, учитывая коэффициент 1/a.

Мы определяем число условий как меру чувствительности проблемы:

Определение 1.4 (Номер условия). Число обусловленности задачи — это отношение того, насколько изменится ее решение, к тому, насколько изменится ее формулировка при малых возмущениях. В качестве альтернативы, это отношение прямой ошибки к обратной для небольших изменений в постановке задачи.

Пример 1.6 (ах = b, часть вторая). Продолжая пример 1.5, мы можем точно вычислить номер условия:

$$c = \frac{\text{ошибка вперед}}{\text{ошибка назад}} = \frac{x - x0}{a(x - x0)} = \frac{1}{a}$$

В общем, вычисление чисел условий почти так же сложно, как вычисление прямой ошибки, и поэтому их точное вычисление, вероятно, невозможно. Тем не менее, во многих случаях можно найти границы или приближения для чисел условий, чтобы помочь оценить, насколько можно доверять решению.

Пример 1.7 (Поиск корня). Предположим, что нам дана гладкая функция f: R R и мы хотим найти значения x с f(x) = 0. Обратите внимание, что f(x + y) = f(x) + f(x). Таким образом, аппроксимация числа условий для нахождения x может быть

изменение прямой ошибки изменение обратной ошибки
$$\frac{f(x+)-x}{f(x+)-f(x)}$$

Обратите внимание, что это приближение совпадает с приближением в примере 1.6. Конечно, если мы не знаем x, мы не можем вычислить f(x), но если мы можем посмотреть на форму f и оценить | * x | вблизи x y нас есть представление о наихудшей ситуации.

Прямая и обратная ошибка являются мерой точности решения. Ради научной воспроизводимости мы также хотим вывести стабильные алгоритмы, которые дают самосогласованные решения для класса задач. Например, алгоритм, который генерирует очень точные решения только в одной пятой части времени, может не стоить реализовывать, даже если мы можем вернуться, используя описанные выше методы, чтобы проверить, является ли решение-кандидат хорошим.

1.3 Практические аспекты

Бесконечность и плотность вещественных чисел R могут вызывать опасные ошибки при реализации численных алгоритмов. Хотя теория анализа ошибок, представленная в § 1.2, в конечном счете поможет нам гарантировать качество численных методов, представленных в следующих главах, прежде чем мы продолжим, стоит отметить ряд распространенных ошибок и «подводных камней», которыми сопровождаются реализации численных методов.

В начале §1.1 мы преднамеренно представили самый большой нарушитель, который мы повторяем более крупным шрифтом для

заслуженного внимания: Редко, если вообще когда-либо, оператор == и его эквиваленты должны использоваться для дробных значений.

Поиск подходящей замены == и соответствующих условий завершения численного метода зависит от рассматриваемого метода. Пример 1.3 показывает, что метод решения Ax = b может закончиться, когда невязка b — Ax равна нулю; поскольку мы не хотим явно проверять, является ли A*x==b, на практике реализации будут проверять norm(A*xb)<epsilon. Обратите внимание, что этот пример демонстрирует два метода:

- Использование обратной ошибкиь Ах, а не прямой ошибки, чтобы определить, когда прекратить, и
- Проверка того, меньше ли обратная ошибка эпсилон, чтобы избежать запретного предиката ==0.

Параметр эпсилон зависит от того, насколько точным должно быть желаемое решение, а также от разрешения используемой системы счисления.

Программист, использующий эти типы данных и операции, должен проявлять бдительность, когда речь идет об обнаружении и предотвращении некорректных числовых операций. Например, рассмотрим следующий фрагмент кода для вычисления нормы x2 для вектора x Rn, представленного в виде одномерного массива x[]:

```
двойная нормаКвадрат = 0; for (int i = 0; i < n; i ++) normSquared += x [i]* x [i]; вернуть sqrt (normSquared);
```

Легко видеть, что теоретически mini $|x| = x2/\bar{n} = maxi |xi|$, т. е. норма х порядка значений элементов, содержащихся в х. Однако при вычислении x2 скрыто выражение x[i]*x[i]. Если существует такое i, что x[i] порядка DOUBLE MAX, произведение x[i]*x[i] переполнится, даже если x2 все еще находится в диапазоне удвоений. Такое переполнение легко предотвратить, разделив x на его максимальное значение, вычислив норму и умножив обратно:

```
двойной maxElement = эпсилон; // не нужно делить на ноль! для ( int я = 0; я < п; я ++)

maxElement = max ( maxElement for ( int i , потрясающие ( x [ я ]));

= 0; i < n; i ++) {

двойной масштаб = x[i]/maxElement; normSquared

+= в масштабе * в масштабе;
}
вернуть sqrt (normSquared) * maxElement;
```

Коэффициент масштабирования устраняет проблему переполнения, гарантируя, что суммируемые элементы не превышают 1.

Этот небольшой пример показывает одно из многих обстоятельств, при котором один символ кода может привести к неочевидной числовой проблеме. Хотя нашей интуиции из непрерывной математики достаточно для создания многих численных методов, мы всегда должны перепроверять, что используемые нами операции верны с дискретной точки зрения.

1.3.1 Пример большего масштаба: суммирование

Теперь мы приведем пример числовой проблемы, вызванной арифметикой с конечной точностью, которую можно решить с помощью менее чем очевидного алгоритмического трюка.

Предположим, что мы хотим суммировать список значений с плавающей запятой, что легко требуется для систем бухгалтерского учета, машинного обучения, графики и почти любой другой области. Фрагмент кода для выполнения этой задачи, которая, без сомнения, используется в бесчисленных приложениях, выглядит следующим образом:

```
двойная сумма = 0;
для (int i = 0; i < n; i ++) sum += x [i];
```

Прежде чем мы продолжим, стоит отметить, что для подавляющего большинства приложений это совершенно стабильный и, безусловно, математически обоснованный метод.

Но что может пойти не так? Рассмотрим случай, когда п велико, а большинство значений х[i] малы и положительны. В этом случае, когда і достаточно велико, сумма переменных будет велика относительно х[i]. В конце концов сумма может быть настолько велика, что х[i] влияет только на младшие биты суммы, а в крайнем случае сумма может быть достаточно большой, чтобы добавление х[ii] не имело никакого эффекта. Хотя одна такая ошибка может не иметь большого значения, накопленный эффект от повторного совершения этой ошибки может превысить сумму, которой мы вообще можем доверять.

Чтобы понять этот эффект математически, предположим, что вычисление суммы a + b может быть ошибочным на величину, равную $\varepsilon > 0$. Тогда явно описанный выше метод может вызвать ошибку порядка $n\varepsilon$, которая линейно растет с ростом n. На самом деле, если большинство элементов x[i] имеют порядок ε , то этой сумме вообще нельзя доверять! Это неутешительный результат: ошибка может быть такой же большой, как и сама сумма.

К счастью, есть много способов сделать лучше. Например, добавление сначала наименьших значений может помочь объяснить их кумулятивный эффект. Парные методы, рекурсивно добавляющие пары значений из х[] и создающие сумму, также более стабильны, но их может быть сложно реализовать так же эффективно, как описанный выше цикл for. К счастью, алгоритм Кахана (СІТЕ) предоставляет легко реализуемый метод «компенсированного суммирования», который работает почти так же быстро.

Полезное наблюдение здесь заключается в том, что мы действительно можем отслеживать приближение ошибки в сумме во время данной итерации. В частности, рассмотрим выражение

$$((a + 6) - a) - 6.$$

Очевидно, это выражение алгебраически равно нулю. Однако численно это может быть не так.

В частности, сумма (a + b) может округлять результат, чтобы он оставался в пределах области значений с плавающей запятой. Вычитание а и b по одному за раз дает приблизительное значение ошибки, вызванной этой операцией; обратите внимание, что операции вычитания, вероятно, лучше обусловлены, поскольку переход от больших чисел к маленьким добавляет цифры точности из-за отмены.

Таким образом, техника Кахана действует следующим образом:

```
двойная сумма = 0;
двойная компенсация = 0; // приближение ошибки

for ( int i = 0; i < n ; i ++) { // пытаемся
    прибавить и к x [ i ] и к недостающей части double nextTerm = x [ i ] +
    компенсация ;

    // вычисляем результат суммирования этой итерации double nextSum =
    sum + nextTerm ;

    // вычислить компенсацию как разницу между термином, который вы хотели // добавить, и фактическим
    результатом компенсация = nextTerm -
    ( nextSum - sum );

    сумма = следующая сумма;
}
```

Вместо того, чтобы просто поддерживать сумму, теперь мы отслеживаем сумму, а также аппроксимацию компенсации разницы между суммой и желаемым значением. Во время каждой итерации мы пытаемся добавить обратно

эту компенсацию в дополнение к текущему элементу x[], а затем мы пересчитываем компенсацию для учета последней ошибки.

Анализ алгоритма Каана требует более тщательного учета, чем анализ более простого метода нарастания. В конце этой главы вы познакомитесь с одним выводом выражения ошибки; окончательный математический результат будет заключаться в том, что ошибка улучшится с пє до $O(\varepsilon + n\varepsilon$ значительное улучшение, когда²), а $O<\varepsilon$

Реализовать суммирование Кэхана несложно, но оно более чем удваивает количество операций результирующей программы. Таким образом, существует неявный компромисс между скоростью и точностью, на который инженеры-программисты должны пойти при выборе наиболее подходящего метода.

В более широком смысле алгоритм Кахана является одним из нескольких методов, позволяющих избежать накопления численной ошибки в ходе вычислений, состоящих более чем из одной операции. Другие примеры включают алгоритм Брезенхема для растеризации линий (СІТЕ), который использует только целочисленную арифметику для рисования линий, даже если они пересекают строки и столбцы пикселей в нецелочисленных местах, и быстрое преобразование Фурье (СІТЕ), которое эффективно использует двоичное разбиение. трюк с суммированием, описанный выше.

1.4 Проблемы

Проблема 1.1. Вот проблема.

Machine Translated by Google

Часть II

Линейная алгебра



Глава 2

Линейные системы и LU Разложение

В главе 0 мы обсуждали различные ситуации, в которых линейные системы уравнений Ax = b появляются в математической теории и на практике. В этой главе мы решаем основную проблему и исследуем численные методы решения таких систем.

2.1 Разрешимость линейных систем

Как было введено в §0.3.3, системы линейных уравнений типа

$$3x + 2y = 6$$
$$4x + y = 7$$

можно записать в матричной форме, как в

В более общем случае мы можем записать системы вида Ax =b для A Rm×n , x Rn , и b Rm. Разрешимость системы должна попадать в один из трех случаев:

1. Система может не допускать никаких решений, например:

Эта система требует, чтобы x = 1 и x = 1 одновременно, очевидно, два несовместимых условия.

2. Система может допускать единственное решение; например, система в начале этого раздела решается как (x, y) = (8/11, 45/11).

3. Система может иметь бесконечно много решений, например, 0x = 0. Заметим, что если система Ax =b допускает два решения x0 и x1, то она автоматически имеет бесконечно много решений вида cx0 + (1 c)x1 при c R, так как

$$A(cx0 + (1 c)x1) = cAx0 + (1 c)Ax1 = cb + (1 c)b = b.$$

Эту линейную систему можно было бы назвать недоопределенной.

В общем случае разрешимость системы зависит как от A, так и от b. Например, если мы изменим приведенную выше неразрешимую систему, чтобы она была

тогда система переходит от отсутствия решений к бесконечному множеству решений в форме (1, у). Фактически, каждая матрица А допускает правую часть b такую, что Ax = b разрешима, поскольку Ax = 0 всегда может быть решена x 0 независимо от A. Напомним из \$0.3.1, что умножение матрицы на вектор можно рассматривать как линейное объединение столбцов A с весами из x. Таким образом, как упоминалось в \$0.3.3, мы можем ожидать, что Ax = b будет разрешимым именно тогда, когда b находится в пространстве столбцов A.

В широком смысле «форма» матрицы А Rm×n существенно влияет на разрешимость Ax = b. Напомним, что столбцы A являются m-мерными векторами. Во-первых, рассмотрим случай, когда A «широкий», то есть столбцов в нем больше, чем строк (n > m). Каждый столбец является вектором в Rm, поэтому пространство столбцов может иметь размерность не более m. Поскольку n > m, n столбцов A должны быть линейно зависимы; это означает, что существует x = 0 такой, что на самом деле существует бесконечно много решений x = 0 другими словами, мы показали, что ни одна широкая матричная система не допускает единственного решения.

Когда A «высокий», то есть когда в нем больше строк, чем столбцов (m > n), тогда n столбцов не могут охватывать Rm. Таким образом, существует некоторый вектор b0 Rm\col A. По определению этот вектор b0 не может удовлетворять условию Ax = b0 ни для какого x. Другими словами, любая длинная матрица A допускает неразрешимые системы Ax = b0.

Обе описанные выше ситуации далеко не благоприятны для построения численных алгоритмов. Например, если линейная система допускает множество решений, мы должны сначала определить, какое решение требуется от пользователя: в конце концов, решение x + 1031x0 может быть не таким осмысленным, как x 0,1x0. С другой стороны, в длинном случае, даже если Ax = b разрешимо для конкретного b, любое малое возмущение Ax =b + ɛb0 уже не разрешимо; такая ситуация может возникнуть просто потому, что процедуры округления, обсуждавшиеся в предыдущей главе, могут, во-первых, только аппроксимировать A и b.

Учитывая эти сложности, в этой главе мы сделаем несколько упрощающих предположений: • Мы будем рассматривать только квадрат A Rn×n

• Будем считать, что A неособо, т. е. что Ax =b разрешима для любого b.

Напомним из \$0.3.3, что условие невырожденности эквивалентно требованию, чтобы столбцы A охватывает Rn и влечет существование матрицы A 1 такое, что A 1 1 = In×n.

Вводящее в заблуждение наблюдение состоит в том, чтобы думать, что решение Ах =b эквивалентно явному А ¹ вычислению матрицы и последующему умножению, чтобы найти х А 1b. Хотя эта стратегия решения, безусловно, действительна, она может представлять собой значительное количество излишеств: в конце концов, нас интересуют только значения n 2 d ҳначениянов ме того, даже когда А ведет себя хорошо, может случиться так, что написание А ¹ дает числовые трудности, которые можно обойти.

2.2 Стратегии специальных решений

Во вводной алгебре мы часто подходим к проблеме решения линейной системы уравнений как к искусству. Стратегия состоит в том, чтобы «изолировать» переменные, итеративно записывая альтернативные формы линейной системы до тех пор, пока каждая строка не будет иметь форму x = const.

При формулировании систематических алгоритмов решения линейных систем полезно привести пример этого процесса решения. Рассмотрим следующую систему:

$$y - z = -1$$

 $3x - y + z = 4x$
 $+ y - 2z = -3$

Параллельно мы можем поддерживать матричную версию этой системы. Вместо того, чтобы явно записывать Ax = b, мы можем сэкономить немного места, написав «расширенную» матрицу ниже:

Мы всегда можем писать линейные системы таким образом, если согласны с тем, что переменные остаются в левой части уравнений, а константы — в правой.

Возможно, мы хотим сначала разобраться с переменной х. Для удобства мы можем переставить строки системы так, чтобы третье уравнение появилось первым:

$$x + y - 2z = -3 y - z$$

= -1 3x - y +
z = 4
1 1 2 3
0 1 1 1
3 1 1 4

Затем мы можем подставить первое уравнение в третье, чтобы исключить член 3x. Это то же самое, что масштабировать отношение x + y - 2z = -3 на -3 и добавлять результат к третьему уравнению:

Точно так же, чтобы исключить у из третьего уравнения, мы можем умножить второе уравнение на 4 и добавить результат к третьему:

$$x + y - 2z = -3 y - z$$

= -1 3z = 9
1 1 2 3 0
1 1 1 1
0 0 3 9

Теперь мы изолировали z! Таким образом, мы можем масштабировать третью строку на 1/3 , чтобы получить выражение для z:

$$x + y - 2z = -3y - z$$

= -1 z = 3 1 2 3 0
1 1 1 1
0 0 1 3

Теперь мы можем подставить z = 3 в два других уравнения, чтобы удалить z из всех строк, кроме последней:

$$x + y = 3y =$$
 1103
2r=3 0102
0013

Наконец, мы делаем аналогичную замену для у, чтобы завершить решение:

Этот пример может показаться несколько педантичным, но оглядываясь назад на нашу стратегию, можно сделать несколько важных наблюдений о том, как мы можем решать линейные системы:

- Мы написали последовательные системы Aix = bi , которые можно рассматривать как упрощения исходной системы. Ax = 6.
- Мы решили систему, даже не записывая А
- Мы неоднократно использовали несколько простых операций: масштабирование, добавление и перестановку строк таблицы. система.
- Те же операции были применены к A и b. Если мы масштабируем k-ю строку A, мы также масштабируем k-й ряд b. Если мы добавили строки k и из A, мы добавили строки k и из b.
- Менее очевидно, что шаги решения не зависели от b. То есть все наши решения о том, как решать, были мотивированы устранением ненулевых значений в A, а не изучением значений в b; b, которые просто попались на пути.
- Мы остановились, когда свели систему к In×nx =b.

Мы будем использовать все эти общие наблюдения о решении линейных систем в наших интересах.

2.3 Кодирование операций со строками

Возвращаясь к примеру в § 2.2, мы видим, что решение линейной системы на самом деле включало в себя только три операции: перестановку, масштабирование строк и добавление масштаба одной строки к другой. На самом деле, таким образом мы можем решить любую линейную систему, поэтому стоит изучить эти операции более подробно.

2.3.1 Перестановка

Нашим первым шагом в §2.2 было поменять местами две строки в системе уравнений. В более общем случае мы могли бы индексировать строки матрицы, используя числа 1, . . . , м. Тогда перестановку этих строк можно записать в виде функции σ такой, что список σ(1), . . . , σ(m) покрывает тот же набор индексов.

То есть произведение $P\sigma A$ — это в точности матрица A со строками, переставленными в соответствии с σ .

Пример 2.1 (матрицы перестановок). Предположим, мы хотим переставить строки матрицы в R3×3 с $\sigma(1) = 2$, $\sigma(2) = 3$ и $\sigma(3) = 1$. Согласно нашей формуле мы бы имели

$$010$$
 $P\sigma = 001$
 100

Из примера 2.1 видно, что Рσ имеет единицы в позициях вида (k, σ(k)) и нули в остальных местах. Пара (k, σ(k)) представляет утверждение: «Мы хотели бы, чтобы строка k выходной матрицы была строкой σ(k) из входной матрицы». Основываясь на этом описании матрицы перестановок, легко увидеть, что обратным Рσ является транспонированное P, поєкольку это просто меняет местами роли строк и столбцов — теперь мы берем строку σ(k) входных данных и помещаем ее в строка k вывода. Другими сдорвам мпл.

2.3.2 Масштабирование строк

Предположим, мы записываем список констант a1, . . . , am и попытаться масштабировать k-ю строку некоторой матрицы A на ak . Очевидно, это достигается путем применения масштабной матрицы Sa, определяемой следующим образом:

Предполагая, что все ak удовлетворяют ak = 0, легко инвертировать Sa, «уменьшив масштаб»:

$$C_a^{1} = S1/a$$

$$\begin{array}{c} 1/a1 \ 0 \ 0 \cdots 1/a2 \ 0 \cdots \\ \vdots \ \vdots \ \ddots \ \vdots \\ 0 \ 0 \cdots 1/HOYLD \end{array}$$

2.3.3 Устранение

Наконец, предположим, что мы хотим масштабировать строку k с помощью константы с и добавить результат к строке . Эта операция может показаться менее естественной, чем две предыдущие, но на самом деле она весьма практична.

нужно объединить уравнения из разных строк линейной системы! Мы реализуем эту операцию, используя «матрицу исключения» М, так что произведение МА применяет эту операцию к матрице А. Напомним, что произведение е_к А выбирает k-ю строку А. Затем предварительное умножение на е дает матрица к _{А, который равен нулю, за исключением того, что -th строка равна k-й строке А.}

Пример 2.2 (Построение матрицы исключения). Брать

Предположим, мы хотим изолировать третью строку матрицы A R3×3 и переместить ее во вторую строку. Как обсуждалось выше, эта операция выполняется записью:

Конечно, выше мы умножали справа налево, но так же легко могли бы сгруппировать произведение, как (e2e₃)A. Структуру этого продукта легко увидеть:

Нам удалось изолировать строку k и переместить ее в строку . Наша первоначальная операция исключения хотела добавить с раз строку, k к 6троку ее еперы Троку к беевыполнить как сумму A + сее

Пример 2.3 (Решение системы). Теперь мы можем закодировать каждую из наших операций из раздела 2.2, используя матрицы, которые мы построили выше:

1. Переставьте строки, чтобы переместить третье уравнение в первую строку:

- 2. Масштабируйте первую строку на -3 и прибавьте результат к третьей строке: E1 = I3×3 3e3e.1
- 3. Масштабируйте вторую строку на 4 и добавьте результат к третьей строке: $E2 = I3 \times 3 + 4e3e$. 2
- 4. Масштабируйте третью строку на 1/3: S = diag(1, 1, 1/3)

- 5. Масштабируйте третью строку на 2 и добавьте ее к первой строке: E3 = I3×3 + 2e1e.3
- 6. Добавьте третью строку ко второй: E4 = I3×3 +e2e.

2

7. Масштабируйте третью строку на -1 и прибавьте результат к первой строке: E5 = I3×3 e1e 3

Таким образом, обратный А в разделе 2.2 удовлетворяет

$$A^{-1} = E5E4E3SE2E1P.$$

Убедитесь, что вы понимаете, почему эти матрицы отображаются в обратном порядке!

2.4 Исключение Гаусса

Последовательность шагов, выбранная в разделе 2.2, отнюдь не уникальна: существует множество различных путей, которые могут привести к решению Ax = b. Наши шаги, однако, следовали стратегии исключения Гаусса, известного алгоритма решения линейных систем уравнений.

В более общем смысле, скажем, наша система имеет следующую «форму»:

Алгоритм работает поэтапно, как описано ниже.

2.4.1 Замена вперед

Рассмотрим верхний левый элемент нашей матрицы:

Мы назовем этот элемент нашей первой точкой опоры и будем считать, что он не равен нулю; если он равен нулю, мы можем переставить строки так, чтобы это было не так. Сначала мы применяем матрицу масштабирования, чтобы точка опоры равнялась единице:

Теперь мы используем строку, содержащую сводную точку, чтобы исключить все остальные значения ниже в том же столбце:

Теперь мы перемещаем нашу точку опоры на следующую строку и повторяем аналогичную серию операций:

Обратите внимание, что здесь происходит приятная вещь. После того, как первая точка опоры удалена из всех остальных строк, первый столбец равен нулю под строкой 1. Это означает, что мы можем безопасно добавлять числа, кратные второй строке, к строкам под ней, не затрагивая нули в первом столбце.

Повторяем этот процесс до тех пор, пока матрица не станет верхнетреугольной:

2.4.2 Обратная замена

Устранение оставшихся × из системы теперь является прямым процессом. Теперь мы действуем в обратном порядке строк и удаляем назад. Например, после первой серии шагов обратной замены у нас остается следующая форма:

Точно так же вторая итерация дает:

После нашего последнего шага исключения у нас осталась желаемая форма:

Правая часть теперь является решением линейной системы Ах = b.

2.4.3 Анализ исключения Гаусса

Каждая операция со строками в методе исключения Гаусса — масштабирование, исключение и замена двух строк — очевидно, занимает O(n) времени для завершения, так как вам нужно выполнить итерацию по всем n элементам строки (или

два) из А. Как только мы выбираем точку опоры, мы должны сделать n прямых или обратных замен в строках). В общей сложности ниже или выше этой точки опоры, соответственно; это означает, что работа для одного поворота в целом составляет О (n мы выбираем по одной опорной точке на строку, добавляя последний коэффициент n. Таким образом, достаточно легко видеть, что гауссовский исключение выполняется за O(n) время.

Одно решение, которое принимается во время исключения Гаусса и которое мы не обсуждали, — это выбор опорных точек. Напомним, что мы можем переставлять строки линейной системы по своему усмотрению перед выполнением обратной или прямой замены. Эта операция необходима, чтобы иметь возможность работать со всеми возможными матрицами А. Например, рассмотрим, что произошло бы, если бы мы не использовали поворот на следующих матрица:

Обратите внимание, что элемент, обведенный кружком, равен нулю, поэтому мы не можем ожидать деления первой строки на любое число, чтобы заменить этот 0 на 1. Это не означает, что система неразрешима, это просто означает, что мы должны выполнить поворот, достигнутый путем замены элементов местами. первую и вторую строки, чтобы поместить в этот слот ненулевое значение.

В более общем случае предположим, что наша матрица выглядит так:

$$A = \begin{pmatrix} \varepsilon & 1 \\ 10 & 7 \end{pmatrix}$$

где 0 < ε

1. Если мы не поворачиваемся, то первая итерация исключения Гаусса дает:

$$A^{\sim} = \frac{1}{1/\epsilon} 0$$

Мы преобразовали матрицу A, которая выглядит почти как матрица перестановок (на самом деле это ^{1 A} , а очень простой способ решения системы!) в систему с потенциально огромными значениями 1/ε.

Этот пример показывает, что бывают случаи, когда мы можем захотеть развернуться, даже если в этом, строго говоря, нет необходимости. Поскольку мы масштабируем по обратному значению опорного значения, очевидно, что наиболее численно устойчивые варианты — это иметь большое опорное значение: маленькие опорные значения имеют большие обратные значения, масштабируя числа до больших значений в режимах, которые, вероятно, потеряют точность. Есть две хорошо известные стратегии разворота:

- 1. Частичная сводка просматривает текущий столбец и переставляет строки матрицы так, что наибольшее абсолютное значение появляется на диагонали.
- 2. Полный поворот выполняет итерацию по всей матрице и переставляет как строки, так и столбцы, чтобы получить максимально возможное значение по диагонали. Обратите внимание, что перестановка столбцов матрицы является допустимой операцией: она соответствует изменению маркировки переменных в системе или последующему умножению А на перестановку.

Пример 2.4 (Поворот). Предположим, что после первой итерации исключения Гаусса у нас осталась следующая матрица:

Если мы реализуем частичный поворот, то мы будем смотреть только во второй столбец и поменяем местами вторую и третью строки; обратите внимание, что мы оставляем 10 в первой строке, так как алгоритм уже посещал ее:

Если мы реализуем полный поворот, то мы переместим 9:

Очевидно, что полная сводка дает наилучшие числовые значения, но цена — более дорогой поиск больших элементов в матрице.

2.5 Факторизация LU

Много раз мы хотим решить последовательность задач Ax1 =b1, Ax2 =b2, Как мы уже обсуждали, шаги исключения Гаусса для решения Ax =bk зависят в основном от структуры A, а не от значений в конкретном bk . Поскольку A здесь сохраняется постоянным, мы можем захотеть «запомнить» шаги, которые мы предприняли для решения системы, чтобы каждый раз, когда мы сталкиваемся с новым b, нам не приходилось начинать с нуля.

Укрепляя это подозрение, что мы можем переместить некоторые из O(n³) для исключения Гаусса во время предварительного вычисления вспомните верхнюю треугольную систему, полученную после этапа прямой подстановки:

На самом деле, решение этой системы обратной подстановкой ²) время! Почему? Обратная замена требует только O(n, в этом случае это намного проще благодаря структуре нулей в системе. Например, в первой серии обратных подстановок мы получаем следующую матрицу:

Поскольку мы знаем, что значения (обведены кружком) слева от опорной точки равны нулю по построению, нам не нужно копировать их явно. Таким образом, этот шаг занял всего O(n) времени, а не O(n), ² затрачиваемого на прямую замену.

Теперь наша следующая точка поворота делает аналогичную замену:

Опять же, нули по обе стороны от 1 не нужно копировать явно.

Таким образом, мы нашли:

З Наблюдение. В то время как исключение Гаусса принимает O (n) времени решение треугольных систем занимает O(n) время.

2.5.1 Построение факторизации

Напомним из §2.3, что все операции исключения Гаусса можно рассматривать как предварительное умножение Ax = b на разные матрицы Ax = b на разные матрици Ax = b на разные матрици Ax = b на разнение Ax = b на разнения Ax = b на разнен

После фазы прямой замены исключения Гаусса у нас остается верхняя треугольная матрица, которую мы можем назвать U Rn×n . С точки зрения умножения матриц мы можем писать:

$$M\kappa \cdot \cdot \cdot M1A = U$$
,

или, что то же самое,

$$A = (Mk \cdots M1)$$
 1U
$$= (M_{11} M_{2} 1 \cdots M_{K} 1)U$$
 LU, если сделать определение L $M_{11} M_{2} 1 \cdots M_{K} 1$

Мы еще ничего не знаем о структуре L, но мы знаем, что системы вида Uy = d легче решать, поскольку U является верхнетреугольным. Если L одинаково хорошо, мы могли бы решить Ax = b в два шага, записав (LU)x = b или x = U 1L 1b: 1. Решить Ly = b относительно

```
y, получив y = L 1b.
```

2. Теперь, когда у нас есть у, решите Ux = y, что даст x = U 1y = U 1 (L 1b) = (LU) 1b = A 1b. Мы уже знаем, что этот шаг занимает всего $O(n^2)$ время.

Наша оставшаяся задача — убедиться, что L имеет хорошую структуру, которая облегчит решение Ly = b, чем решение Ax = b. К счастью — и это неудивительно — мы обнаружим, что L является нижнетреугольным и, следовательно, может быть решена с использованием O (r) прямая замена.

Чтобы убедиться в этом, предположим, что мы не реализуем поворот. Тогда каждая наша матрица Мк либо является масштабирующей матрицей, либо имеет структуру

$$M\kappa = B \times n + c. 9.9. \kappa$$

где > k, так как мы выполнили только прямую замену. Помните, что эта матрица служит определенной цели: масштабируйте строку k по с и добавляйте результат к строке . Эту операцию, очевидно, легко отменить: масштабируйте строку k на с и вычтите результат из строки . Мы можем проверить это формально:

$$(In\times n + B.п._{K})(In\times n$$
 сее $_{K}) = In\times n + ($ сее $_{K} + CUK$ $) - C$ $2 \ni 3 \ni KK$ $= B \times \Pi - C$ 2 e (e_{K} Д) $_{AK}$ $= In\times n$, так как $e = ek \cdot ek$, $UK =$

Таким образом, матрица L является произведением масштабирующих матриц и матриц вида M 1 = In×n _K ченесторование на в.п. нижняя треугольная при > k. Масштабирующие матрицы диагональные, а матрица M нижняя треугольная. Вы покажете в упражнении 2.1, что произведение нижних треугольных матриц является нижним треугольным, показав, в свою очередь, что L является нижнетреугольным по мере необходимости.

Мы показали, что если возможно выполнить гауссово исключение А без использования поворота, вы можете разложить A = LU на произведение нижне- и верхнетреугольных матриц. Прямая и обратная подстановка каждая занимает O(n) времени, поэтому, если эту факторизацию можно вычислить заранее, линейное решение может ³) Исключение Гаусса. Ты покажешь в быть выполнено быстрее, чем полное O(n) необходимы.

2.5.2 Реализация ЛЕ

Простая реализация исключения Гаусса для решения Ax = b достаточно проста, чтобы ее можно было сформулировать. В частности, как мы обсуждали ранее, мы можем сформировать расширенную матрицу (A | b) и применять операции со строками по одной к этому n × (n + 1) блоку, пока он не будет выглядеть как (In×nA |b). Этот процесс, однако, является деструктивным, то есть, в конце концов, нас интересует только последний столбец расширенной матрицы, и мы не сохранили никаких свидетельств нашего пути решения. Такое поведение явно неприемлемо для факторизации LU.

Давайте посмотрим, что происходит, когда мы умножаем две матрицы исключения:

$$(In \times n - cee \ \kappa)(In \times n - cee \ \kappa) = In \times n$$
 ceek cpepe k

Как и в нашей конструкции обратной матрицы исключения, произведение двух членов еі равно нулю, поскольку стандартный базис ортогонален. Эта формула показывает, что после масштабирования опорной точки до 1 произведение матриц исключения, используемых для прямой замены этой опорной точки, имеет вид:

$$M = \begin{pmatrix} 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & \times 1 & 0 \\ 0 & \times 0 & 1 \end{pmatrix}$$

где значения × — это значения, используемые для исключения остальной части столбца. Перемножение матриц этой формы вместе показывает, что элементы под диагональю L просто происходят из коэффициентов, используемых для выполнения подстановки.

Мы можем принять окончательное решение оставить элементы по диагонали L в факторизации LU равными 1. Это решение является законным, поскольку мы всегда можем пост-умножить L на масштабирующую матрицу S, переводя эти элементы в 1. и запишем LU = (LS)(S $\,$ 1U), не затрагивая треугольный шаблон L или U. Приняв это решение, мы можем сжать наше хранилище как L, так и U в единую матрицу $\,$ n $\,$ x $\,$ n, верхний треугольник которой равен U и который равен L ниже диагонали; все отсутствующие диагональные элементы L равны 1.

Теперь мы готовы написать псевдокод для простейшей стратегии факторизации LU, в которой мы не переставляем строки или столбцы для получения опорных точек:

```
// Принимает на вход матрицу размера n на n A [i, j]
// Редактируем A на месте, чтобы получить описанную выше компактную факторизацию LU
для поворота от 1 до n {
    pivotValue = A [поворот, поворот]; // Плохое предположение, что это не ноль!
```

2.6 Проблемы

Проблема 2.1. Произведение нижних треугольных вещей есть нижнее треугольное; произведение сводных матриц выглядит правильно

Задача 2.2. Реализовать LU с поворотом

Задача 2.3. Неквадратный LU

Machine Translated by Google

Глава 3

Проектирование и анализ линейных Системы

Теперь, когда у нас есть некоторые методы решения линейных систем уравнений, мы можем использовать их для решения множества задач. В этой главе мы рассмотрим несколько таких приложений и сопутствующих аналитических методов, чтобы охарактеризовать типы решений, которые мы можем ожидать.

3.1 Решение квадратных систем

В начале предыдущей главы мы сделали несколько предположений о типах линейных систем, которые собирались решать. Хотя это ограничение было нетривиальным, на самом деле многие, если не большинство приложений линейных решателей можно представить в терминах квадратных обратимых линейных систем. Ниже мы рассмотрим несколько контрастирующих приложений.

3.1.1 Регрессия

Мы начнем с простого приложения, используемого в анализе данных, известного как регрессия. Предположим, что мы проводим научный эксперимент и хотим понять структуру наших экспериментальных результатов. Одним из способов моделирования такой связи может быть запись независимых переменных эксперимента в некотором векторе х Rn и рассмотрение зависимой переменной как функции f(x): Rn R. Наша цель — предсказать результат f без проведения полного эксперимента.

Пример 3.1 (Биологический эксперимент). Предположим, мы хотим измерить влияние удобрений, солнечного света и воды на рост растений. Мы могли бы провести ряд экспериментов, применяя различное количество удобрений (в см3), солнечного света (в ваттах) и воды (в мл) и измеряя высоту растения через несколько дней. Мы могли бы смоделировать наши наблюдения как функцию f: R3 R, принимая три параметра, которые мы хотим проверить, и выводя высоту растения.

В параметрической регрессии мы делаем упрощающее предположение о структуре f. Например, предположим, что f линейна:

$$f(x) = a1x1 + a2x2 + \cdots + anxn.$$

Тогда наша цель становится более конкретной: оценить коэффициенты ак.

Предположим, мы проводим серию экспериментов, которые показывают f(x(k)). Подключившись к нашему форму х у для f, мы получаем ряд утверждений:

$$_{\Gamma}^{(1)} = f(x(1)) = a1x$$
 $_{1}^{(1)(1)(1)+a2x+\cdots+anx}$
 $_{2}^{(2)} = f(x(2)) = a1x$
 $_{3}^{(2)} = f(x(2)) = a1x$
 $_{4}^{(2)} = f(x(2)) = a1x$
 $_{5}^{(2)} = f(x(2)) = a1x$
 $_{7}^{(2)} = f(x(2)) = a1x$
 $_{1}^{(2)} = f(x(2)) = a1x$

Обратите внимание, что вопреки нашему обозначению Ax = b, неизвестными здесь являются переменные ak , а не переменные x. Если мы сделаем n наблюдений, мы можем написать:

Другими словами, если мы проведем n испытаний нашего эксперимента и запишем их в столбцы матрицы X Rn×n , а зависимые переменные запишем в вектор у Rn , то коэффициенты а можно восстановить, решив X а = y .

На самом деле, мы можем обобщить наш подход на другие, более интересные нелинейные формы функции f. Здесь важно то, что f является линейной комбинацией функций. В частности, предположим, что f(x) принимает следующий вид:

$$f(x) = a1 f1(x) + a2 f2(x) + \cdots + am fm(x),$$

где fk : Rn — R и мы хотим оценить параметры ak . Тогда путем параллельного вывода (k) (k) по m наблюдениям мы можем найти параметры, решив: — у формы х

То есть, даже если f нелинейны, мы можем узнать веса ak, используя чисто линейные методы.

Пример 3.2 (Линейная регрессия). Система Xa = y может быть восстановлена из общей формулировки, если взять fk(x) = xk.

Пример 3.3 (полиномиальная регрессия). Предположим, что мы наблюдаем функцию одной переменной f(x) и хотим записать ее в виде полинома n-й степени

$$f(x) = a0 + a1x + a2x$$
 $2 + \cdots + \tau$

Даны n пар x(k) у $^{(\kappa)}$, мы можем решить для параметрова через систему

Другими словами, мы берем fk(x) = x в нашей общей форме выше. Между прочим, матрица в левой части этого соотношения известна как матрица Вандермонда, которая обладает многими особыми свойствами, характерными для ее структуры.

Пример 3.4 (Колебание). Предположим, мы хотим найти а и ϕ для функции $f(x) = a \cos(x + \phi)$. Напомним из тригонометрии, что мы можем написать $\cos(x + \phi) = \cos x \cos \phi$ $\sin x \sin \phi$. Таким образом, учитывая две точки выборки, мы можем использовать описанную выше технику, чтобы найти $f(x) = a1 \cos x + a2 \sin x$, и применяя это тождество, мы можем написать

Эта конструкция может быть расширена до нахождения f(x) = k ak $cos(x + \phi k)$, что дает один из способов мотивировать дискретное преобразование Фурье f.

3.1.2 Метод наименьших квадратов

Методы, описанные в §3.1.1, предоставляют ценные методы для нахождения непрерывного f, точно соответствующего набору пар данных xk yk . У этого подхода есть два связанных недостатка:

- Возможна ошибка измерения значений xk и yk . В этом случае приближенное соотношение f(xk) уk может быть приемлемым или даже предпочтительнее точного f(xk) = yk .
- Заметьте, что если бы всего было m функций fk , то нам понадобилось ровно m наблюдений xk yk . Дополнительные наблюдения пришлось бы отбросить, или нам пришлось бы изменить форму f.

Обе эти проблемы связаны с более крупной проблемой переобучения: подбор функции с n степенями свободы к n точкам данных не оставляет места для ошибки измерения.

В более общем случае предположим, что мы хотим решить линейную систему Ax = b относительно x. Если мы обозначим строку k матрицы A как r toгда наша система выглядит так



по определению умножения матриц.

Таким образом, каждой строке системы соответствует наблюдение вида $rk \cdot x = bk$. То есть еще один способ интерпретировать линейную систему Ax = b - 3 о утверждений вида «Скалярное произведение x на rk равно bk ».

С этой точки зрения, длинная система Ax = b c A $Rm \times n u m > n$ просто кодирует более n этих наблюдений скалярного произведения. Однако когда мы делаем более n наблюдений, они могут оказаться несовместимыми; как объяснялось в §2.1, высокие системы, скорее всего, не допускают решения. В нашей «экспериментальной» установке, описанной выше, такая ситуация может соответствовать ошибкам в измерении пар xk yk.

Когда мы не можем точно решить Ax = b, мы можем немного ослабить задачу, чтобы аппроксимировать Ax b. В частности, мы можем потребовать, чтобы невязка b Ax была как можно меньше, минимизируя

норма b — Ax. Заметим, что если существует точное решение линейной системы, то эта норма минимизируется в нуле, так как в этом случае мы имеем b — Ax = b — b = 0. Минимизация b — Ax аналогична минимизации b — Ax $\frac{2}{2}$, который мы расширили в примере 0.16 до:

$$6$$
 - топор $2 = x A Ax - 26 Ax + 6 $2.1$$

Градиент этого выражения по х должен быть равен нулю в его минимуме, что дает следующую систему:

$$0 = 2Ax - 2A6$$

Или, что то же самое: Ax = Ab.

Это знаменитое соотношение достойно теоремы:

Теорема 3.1 (нормальные уравнения). Минимумы невязки b Ах для A Rm×n (без ограничений на m или n) удовлетворяют условию AAx = A b.

Если хотя бы n строк матрицы A линейно независимы, то матрица AA Rn×n обратима. В этом случае минимальная невязка возникает (единственным образом) при (AA) 1A b или, что то же самое, решение задачи наименьших квадратов так же просто, как решение квадратичной линейной системы A Ax = A b из теоремы 3.1. Таким образом, мы расширили наш набор стратегий решения до A Rm×n c m n, применяя только методы для квадратных матриц.

С недоопределенным случаем m < n справиться значительно труднее. В частности, мы теряем возможность единственного решения задачи Ax = b. В этом случае мы должны сделать дополнительное предположение относительно x, чтобы получить единственное решение, например, что оно имеет малую норму или что оно содержит много нулей. Каждое такое регуляризирующее допущение приводит к разным стратегиям решения; мы рассмотрим некоторые из них в упражнениях, сопровождающих эту главу.

3.1.3 Дополнительные примеры

Важным навыком является умение идентифицировать линейные системы «в дикой природе». Здесь мы быстро перечислим еще несколько примеров.

Выравнивание

Предположим, мы делаем две фотографии одной и той же сцены с разных позиций. Одна общая задача в компьютерном зрении и графике — сшить их вместе. Для этого пользователь (или автоматическая система) может отметить некоторое количество точек xk ,yk R2 так, что xk на первом изображении соответствует yk на втором изображении. Конечно, при сопоставлении этих точек были допущены ошибки, поэтому мы хотим найти стабильное преобразование между двумя изображениями путем передискретизации количества необходимых пар (x, y).

Предполагая, что у нашей камеры стандартный объектив, проекции камеры линейны, поэтому разумно, как предположение состоит в том, что существуют некоторый A R2×2 и вектор сдвига b R2 такие, что

¹На этом этапе может быть полезно вернуться к предварительным этапам главы 0 для повторения.

Наши неизвестные переменные здесь — это A и b, а не xk и yk .

В этом случае мы можем найти преобразование, решив:

Это выражение снова является суммой квадратов линейных выражений для наших неизвестных A и b, и путем вывода, аналогичного нашему обсуждению задачи наименьших квадратов, оно может быть решено линейно.

Деконволюция

Часто мы случайно делаем фотографии, которые немного не в фокусе. Хотя фотография, которая полностью размыта, может быть потерянной причиной, если есть локальное или мелкомасштабное размытие, мы можем восстановить более четкое изображение с помощью вычислительных методов. Одной из простых стратегий является деконволюция, описанная ниже.

Мы можем думать о фотографии как о точке в Rp. , где p — количество пикселей; конечно, если Фотография находится в цвете, нам могут понадобиться три значения (RGB) на пиксель , что приводит к аналогичной технике в R3p . соседи на картинке. При обработке изображений эти линейные операции часто обладают другими особыми свойствами, такими как инвариантность к сдвигу, но для наших целей мы можем думать о размытии как о некотором линейном операторе х G х.

мин
$$x0 - G * x$$
 2.

То есть мы просим, чтобы, когда вы размываете х с помощью G, вы получали наблюдаемое фото x0. Конечно, многие четкие изображения могут давать такой же размытый результат при G, поэтому мы часто добавляем дополнительные условия к приведенной выше минимизации, прося, чтобы x0 не менялось слишком сильно.

3.2 Специальные свойства линейных систем

Наше обсуждение исключения Гаусса и факторизации LU привело к совершенно общему методу решения линейных систем уравнений. Хотя эта стратегия всегда работает, иногда мы можем получить преимущество в скорости или количественных показателях, исследуя конкретную систему, которую решаем. Здесь мы обсудим несколько распространенных примеров, когда дополнительные сведения о линейной системе могут упростить стратегии решения.

3.2.1 Положительно определенные матрицы и факторизация Холецкого

Как показано в теореме 3.1, решение задачи наименьших квадратов Ах b дает решение x, удовлетворяющее квадратичной линейной системе (A A)x = A b. Независимо от A, матрица AA обладает несколькими особыми свойствами, которые делают эту систему особенной.

Во-первых, легко видеть, что АА симметричен, так как

$$(A A) = A (A) = A A.$$

Здесь мы просто использовали тождества (AB) = BA и (A) = A. Мы можем выразить эту симметрию по индексам, записав (AA)іі = (AA)іі для всех индексов і, j. Это свойство означает, что достаточно хранить только значения AA на диагонали или над ней, так как остальные элементы могут быть получены по симметрии.

Кроме того, АА является положительно полуопределенной матрицей, как определено ниже:

Определение 3.1 (Положительное (полу)определенное). Матрица В Rn×n является положительно полуопределенной, если для всех х Rn x Bx 0. Матрица В является положительно определенной, если x Bx > 0 при x = 0.

Легко показать, что АА положительно полуопределенно, так как:

$$x A Ax = (Ax) (Ax) = (Ax) \cdot (Ax) = Ax$$
 2. 0.

На самом деле, если столбцы А линейно независимы, то АА положительно определена.

В более общем случае предположим, что мы хотим решить симметричную положительно определенную систему Сх = d. Как мы уже выяснили, мы можем LU-факторизовать матрицу C, но на самом деле мы можем сделать несколько лучше. Запишем C Rn×n как блочную матрицу:

$$C = \frac{c11 \text{ vv}}{C^{\sim}}$$

где v Rn 1 и C R(n 1)×(n 1) . Благодаря особой структуре С мы можем сделать следующее наблюдение:

> 0, так как С положительно определена и е1 =0.

Это показывает, что — игнорируя числовые проблемы — нам не нужно использовать поворот, чтобы гарантировать, что c11 = 0 для исключения Гаусса на первом этапе.

Продолжая исключение Гаусса, мы можем применить матрицу прямой замены E, которая в общем случае имеет вид

$$\exists = \begin{array}{ccc} 1/&\overline{c11} & 0 \\ & p & I(n-1)\times(n-1) \end{array} .$$

Здесь вектор r Rn 1 содержит числа, кратные строке 1, чтобы сократить остальную часть первого столбца С. Мы также масштабируем строку 1 на 1/ c11 по причинам, которые вскоре станут очевидными!

По замыслу, после прямой замены мы знаем продукт

EC =
$$\frac{\text{c11 в}}{\text{c11}}$$
 c11

для некоторого D $R(n 1) \times (n 1)$.

Вот где мы отклоняемся от исключения Гаусса: вместо того, чтобы переходить ко второй строке, мы можем пост-умножить на E, чтобы получить произведение ECE:

То есть мы исключили первую строку и первый столбец C! Кроме того, легко проверить, что матрица D^{\sim} также является положительно определенной.

Мы можем повторить этот процесс, чтобы удалить все строки и столбцы С симметрично. Уведомление что мы использовали как симметрию, так и положительную определенность для вывода факторизации, поскольку

- симметрия позволяла нам применять одну и ту же Е к обеим сторонам, и
- положительная определенность гарантирует, что с11 > 0, что означает существование с11 .

В конце концов, аналогично LU-факторизации, мы теперь получаем факторизацию C = LL для нижней треугольной матрицы L. Это известно как факторизация Холецкого C., где D — диагональная матрица, позволяет избежать этой проблемы, и его легко вывести из приведенного выше обсуждения.

Факторизация Холецкого важна по ряду причин. Наиболее заметно то, что для хранения L требуется половина памяти, чем LU-разложение C или даже самого C, поскольку элементы над диагональю равны нулю, и, как и в LU, решение Cx = d так же просто, как прямая и обратная замена .
В упражнениях вы изучите другие свойства факторизации.

В конце концов, код для факторизации Холецкого может быть очень кратким. Чтобы получить особенно ком пакта, предположим, что мы выбираем произвольную строку k и пишем L в блочной форме, изолируя эту строку:

Здесь L11 и L33 — нижние треугольные квадратные матрицы. Тогда выполнение продукта дает:

Мы опускаем значения произведения, которые не нужны для нашего вывода.

В конце концов, мы знаем, что можем написать C = LL. Средний элемент продукта показывает:

$$KK = CKK - K$$

где k Rk 1 содержит элементы k-й строки L слева от диагонали. Более того, средний левый элемент продукта показывает

 $L11k = c\kappa$

где ck содержит элементы C в той же позиции ask . Поскольку L11 является нижнетреугольным, это система может быть решена прямой заменой!

Обратите внимание, что наше обсуждение выше дает алгоритм вычисления факторизации Холецкого. сверху вниз, так как L11 уже будет вычислена к моменту, когда мы достигнем строки k. Мы предоставляем приведенный ниже псевдокод, адаптированный из CITE:

```
// Принимает на вход матрицу размера n на n A [i, j]
// Редактирует А на месте, чтобы получить факторизацию Холецкого в его нижнем треугольнике
для k от 1 до n {
      // Назад - подставить, чтобы найти l_k
      for i from 1 to k -1 { // элемент i из l_k
            сумма = 0;
            для јот 1 до і -1
                   сумма += A [i, j]* A [k, j];
            A[K, g] = (A[K, i] - cymma)/A[i, g];
      }
      // Применяем формулу для l_kk
      нормаКвадрат = 0
      для ј от 1 до і -1
            normSquared += A[k, j]^2;
      A[K, k] = sqrt(A[k]
                               , k] - нормаКвадрат);
```

Как и в случае факторизации LU, этот алгоритм работает за O(n

³) время.

3.2.2 Разреженность

Многие линейные системы уравнений естественным образом обладают свойствами разреженности, а это означает, что большинство элементы A в системе Ax = b ровно равны нулю. Разреженность может отражать определенную структуру в данной проблемы, включая следующие варианты использования:

- При обработке изображений многие системы для редактирования и понимания фотографий выражают взаимосвязь между значениями пикселей и значениями их соседей в сетке изображения. Изображение может быть точкой в Rp для p пикселей, но при решении Ax =b для нового изображения размера p A Rp×p может иметь только O(p), а не O(p 2) отличен от нуля, так как каждая строка включает только один пиксель. и его верхние/нижние/левые/правые соседи.
- В машинном обучении графическая модель использует структуру графа G (V, E) для выражения распределения вероятностей по нескольким переменным. Каждая переменная представлена с помощью узла v V график, а ребро е E представляет собой вероятностную зависимость. Линейные системы, возникающие в этот контекст часто имеет одну строку на вершину v V с ненулевыми значениями только в столбцах, содержащих v и его соседи.
- В вычислительной геометрии формы часто выражаются с помощью наборов связанных треугольников. вместе в сетку. Уравнения для сглаживания поверхности и другие задачи снова связывают положения и другие значения в данной вершине с теми, которые находятся в их соседях по сетке.

Пример 3.5 (Гармоническая параметризация). Предположим, мы хотим использовать изображение для текстурирования треугольной сетки. Сетку можно представить как набор вершин V R3, соединенных ребрами E V × V в треугольники. Поскольку вершины геометрии находятся в R3, мы должны найти способ сопоставить их с плоскостью изображения, чтобы сохранить текстуру как изображение. Таким образом, мы должны присвоить текстурные координаты t(v) R2 на плоскости изображения каждому элементу v V. См. иллюстрацию на рисунке HOMEP.

Одна из стратегий создания этой карты включает одно линейное решение. Предположим, что сетка имеет дисковую топологию, то есть ее можно сопоставить с внутренней частью круга на плоскости. Для каждой вершины vb на границе сетки мы можем указать положение vb , поместив ее на окружность. Внутри мы можем попросить, чтобы позиция карты текстуры была средним значением соседних позиций:

1
$$t(v) = | T(u)$$

Здесь n(v) V — множество соседей v V на сетке. Таким образом, каждому элементу v V ставится в соответствие линейное уравнение, которое либо фиксирует его на границе, либо требует, чтобы его положение равнялось среднему значению соседних положений. Это |V| × |B| система уравнений приводит к устойчивой стратегии параметризации, известной как гармоническая параметризация; матрица системы имеет только O(|V|) отличных от нуля слотов, соответствующих вершинам и их соседям.

Конечно, если A Rn×n разрежена до такой степени, что содержит значения O(n), а не O(n), нет причин тричин хранить A как матрицу размера n × n. Вместо этого методы хранения разреженных матриц сохраняют только O (n) отличных от нуля в более разумной структуре данных, например, в списке троек строк/столбцов/значений Выбор матричной структуры данных включает в себя рассмотрение вероятных операций, которые будут выполняться с матрицей, возможно, включая умножение, итерацию по ненулевым значениям. или перебор отдельных строк или столбцов.

К сожалению, легко видеть, что LU-факторизация разреженного А может не привести к разреженным матрицам L и U; эта потеря структуры сильно ограничивает применимость использования этих методов для решения Ax = b, когда A большое, но разреженное. К счастью, есть много прямых разреженных решателей, адаптирующих LU к разреженным матрицам, которые могут производить LU-подобную факторизацию, не вызывая большого заполнения или дополнительных ненулевых значений; обсуждение этих методов выходит за рамки этого текста. В качестве альтернативы для получения приближенных решений линейных систем использовались итерационные методы; мы отложим обсуждение этих методов до будущих глав.

Некоторые матрицы не только разрежены, но и структурированы. Например, трехдиагональная система линейные уравнения имеют следующий набор ненулевых значений:

В упражнениях, следующих за этой главой, вы получите специальную версию исключения Гаусса для работы с этой простой ленточной структурой.

В других случаях матрицы не могут быть разреженными, но могут допускать разреженное представление. Для экзамена ple, рассмотрим циклическую матрицу:

abcddabccdabbcda

Очевидно, что эту матрицу можно сохранить, используя только значения a, b, c, d. Специализированные методы для этого и других классов матриц хорошо изучены и часто более эффективны, чем общее исключение Гаусса.

3.3 Анализ чувствительности

Как мы видели, важно исследовать матрицу линейной системы, чтобы выяснить, обладает ли она особыми свойствами, которые могут упростить процесс решения. Разреженность, положительная определенность, симметрия и т. д. могут дать ключи к правильному решателю для использования в конкретной задаче.

Однако даже если данная стратегия решения может работать теоретически, не менее важно понимать, насколько хорошо мы можем доверять ответу на линейную систему, данному конкретным решателем. Например, из-за округления и других дискретных эффектов может случиться так, что реализация исключения Гаусса для решения Ах = b дает решение x0 такое, что 0 < Ax0 b 1; другими словами, x0 лишь приблизительно решает систему.

Одним из способов понять вероятность этих эффектов аппроксимации является анализ чувствительности. В этом подходе мы спрашиваем, что может произойти с x, если вместо решения Ax = b в действительности мы решим возмущенную систему уравнений $(A + \delta A)x = b + \delta b$. Существует два способа рассмотрения выводов, сделанных в результате этого типа анализа:

- 1. Вероятно, мы допускаем ошибки при представлении A и b из-за округления и других эффектов. Затем этот анализ показывает наилучшую возможную точность, которую мы можем ожидать для x с учетом допущенных ошибок, представляющих проблему.
- 2. Если наш решатель генерирует приближение x0 к решению Ax = b, это точное решение системы Ax0 = b0, если мы определяем b0 Ax0 (убедитесь, что вы понимаете, почему это предложение не является тавтологией!). Понимание того, как изменения в x0 влияют на изменения в b0, показывает, насколько чувствительна система к слегка неправильным ответам.

Обратите внимание, что наше обсуждение здесь похоже на наши определения прямой и обратной ошибки в предыдущих главах и действительно мотивировано ими.

3.3.1 Матричные и векторные нормы

Прежде чем мы сможем обсудить чувствительность линейной системы, мы должны быть несколько осторожны, чтобы определить, что означает, что изменение δх является «малым». Как правило, мы хотим измерить длину или норму вектора х. Мы уже встречались с двумя нормами вектора:

$$2 \times 2 \times 1$$
 $+ \times \begin{array}{c} 2 \\ 2 \end{array} + \cdots + \times \begin{array}{c} H \end{array}$

для х Rn · Эта норма популярна благодаря своей связи с евклидовой геометрией, но она ни в коем случае не является единственной нормой на Rn · Обычно мы определяем норму следующим образом:

Определение 3.2 (Векторная норма). Векторной нормой называется функция : Rn = [0,], удовлетворяющий следующим \cdot условия:

• x = 0 тогда и только тогда, когда x = 0.

```
• cx = |c|x для всех скаляров c R и векторов x Rn
• x + y x + y для всех x,y Rn
```

Хотя мы используем два нижних индекса · 2 для обозначения двойной нормы вектора, если не указано иное, мы будем использовать обозначение х для обозначения двойной нормы х. Помимо этой нормы, есть много других примеров: • р- норма хр для р 1,

определяемая как:

xp (
$$|x| + |x| +$$

Особое значение имеет 1-норма или норма «такси», заданная

Эта норма получила свое прозвище, потому что она представляет собой расстояние, которое такси проезжает между двумя точками в городе, где дороги проходят только с севера на юг и с востока на запад.

• -норма х определяется как:

$$x = max(|x1|, |x2|, \dots, |xn|).$$

В некотором смысле многие нормы на Rn совпадают. В частности, предположим, что мы говорим, что две нормы эквивалентны, если они удовлетворяют следующему свойству:

Определение 3.3 (Эквивалентные нормы). Две нормы · и · эквивалентны, если существуют константы clow и chigh такие, что

для всех х Rn ·

Это условие гарантирует, что с точностью до некоторых постоянных множителей все нормы согласуются в отношении того, какие векторы бывают «маленькие» и «большие». Фактически мы сформулируем без доказательства известную теорему из анализа:

Теорема 3.2 (Эквивалентность норм на Rn). Все нормы в Rn эквивалентны.

Этот несколько неожиданный результат подразумевает, что все векторные нормы ведут себя примерно одинаково, но выбор нормы для анализа или постановки конкретной проблемы может иметь огромное значение.

Например, на R3 --норма считает, что вектор (1000, 1000, 1000) имеет ту же норму, что и (1000, 0, 0), тогда как на 2-норму, безусловно, влияют дополнительные ненулевые значения.

Так как мы возмущаем не только векторы, но и матрицы, мы также должны уметь брать норму матрицы. Конечно, основное определение нормы не меняется на Rn×m. По этой причине мы можем «развернуть» любую матрицу из Rm×n в вектор из Rnm , чтобы принять любую норму вектора к матрицам. Одной из таких норм является норма Фробениуса, заданная формулой

Однако такая адаптация векторных норм не всегда имеет большое значение. В частности, приоритетом для понимания структуры матрицы А часто является ее действие на векторы, то есть вероятные результаты при умножении А на произвольный х. С этой мотивацией мы можем определить норму, индуцированную векторной нормой, следующим образом:

Определение 3.4 (Индуцированная норма). Норма на Rm×n , индуцированная нормой · на Rn , определяется выражением

A
$$Makc {Ax : x = 1}.$$

То есть индуцированная норма — это максимальная длина изображения единичного вектора, умноженная на А.

Поскольку векторные нормы удовлетворяют сх = |с|х, легко видеть, что это определение эквивалентно требованию

$$\begin{array}{ccccc} A & \max & Rn & & \frac{\text{Tonop}}{} \\ & & & & \\ & & & & \\ & & & & \\ & & & & \\ & & & & \\ \end{array}.$$

С этой точки зрения норма А, индуцированная · , является наибольшим достижимым отношением нормы Ах к норме входа х.

Это общее определение несколько затрудняет вычисление нормы А при заданной матрице А и выборе · . К счастью, матричные нормы, индуцированные многими популярными векторными нормами, можно упростить. Приведем некоторые такие выражения без доказательства:

• Индуцированная одна норма А – это максимальная сумма любого одного столбца А:

• Индуцированная -норма А – это максимальная сумма любой одной строки А:

• Индуцированная двойная норма, или спектральная норма, A Rn×n – это квадратный корень из наибольшего собственного значения A A. То есть,

A
$$\frac{2}{2}$$
 = max{ λ : существует x R $\frac{1}{2}$ где Ax = λ x}

По крайней мере, первые две нормы относительно легко вычислить; к третьему мы еще вернемся при обсуждении проблем с собственными значениями.

3.3.2 Номера условий

Теперь, когда у нас есть инструменты для измерения действия матрицы, мы можем определить число обусловленности линейной системы, адаптировав наше общее определение чисел обусловленности из главы 1. Мы следуем развитию, представленному в СІТЕ.

Предположим, что мы возмущаем δА матрицы А и соответствующим возмущением δb. Для каждого

 ϵ 0, игнорируя технические аспекты обратимости, мы можем записать вектор $x(\epsilon)$ как решение

$$(A + \varepsilon \cdot \delta A)x(\varepsilon) = b + \varepsilon \cdot \delta b.$$

Если мы продифференцируем обе части по ε и применим правило произведения, мы получим следующий результат:

$$\delta A \cdot x + (A + \epsilon \cdot \delta A) = \delta b d\epsilon$$
 $\frac{dx}{}$

В частности, при ϵ = 0 находим

$$\delta A \cdot x(0) + A = \delta b \cdot d\epsilon \epsilon = 0$$

или, что то же самое,

$$\frac{Ax}{d\epsilon}$$
 = A $\frac{1}{(\delta b + \delta A \cdot x(0))}$.

Используя разложение Тейлора, мы можем написать

$$x(\varepsilon) = x + \varepsilon x (0) + O(\varepsilon^{2}),$$

где мы определяем x (0) = $\frac{dx}{d\epsilon}_{\epsilon=0}$. Таким образом, мы можем расширить относительную ошибку, допущенную путем решения возмущенная система:

$$\frac{x (\epsilon) - x (0) x}{(0)} = \frac{ex (0) + O(ex^{-2})}{(0)} - \text{по разложению Тейлора}$$

$$\frac{1 \epsilon A (\delta b - \delta A \cdot x(0)) + O(\epsilon x(0)^{-2})}{\epsilon |} - \text{по производной, которую мы вычислили}$$

$$\frac{1}{\epsilon |} - (A - \frac{1}{1} - 1 \delta b + A - \delta A \cdot x(0))) + O(\epsilon^{-2})$$

$$x(0) \text{ по неравенству треугольника } A + B - A + B$$

$$|\epsilon|A - \frac{\delta b}{x(0)} + \delta A + O(\epsilon^{-2}) \text{ тождеством } AB - AB$$

$$= |\epsilon|A - \frac{\delta b}{1} - \frac{\delta b}{1} - \frac{\delta A}{1} - \frac{\delta A}{1} + O(\epsilon^{-2})$$

$$|\epsilon|A - \frac{\delta b}{1} - \frac{\delta A}{1} - \frac{\delta A}{1} - \frac{\delta A}{1} - \frac{\delta A}{1} + O(\epsilon^{-2}) , \text{ так как } Ax(0) - Ax(0)$$

$$= |\epsilon|A - \frac{\delta b}{1} - \frac{\delta A}{1} - \frac{\delta A}{1} - \frac{\delta A}{1} - \frac{\delta A}{1} - O(\epsilon^{-2}) , \text{ так как } Ax(0) - Ax(0)$$

поскольку по определению Ax(0) =b

Здесь мы применили некоторые свойства матричной нормы, вытекающие из соответствующих свойств векторов. Обратите внимание, что сумма D $\delta b/b + \delta A/A$ кодирует относительные возмущения A и b. С этой точки зрения мы в первом порядке ограничили относительную погрешность возмущения системы по ϵ с помощью множителя ϵ AA 1:

$$\frac{x(\varepsilon) - x(0) x}{(0)} \qquad \varepsilon \cdot D \cdot \kappa + O(\varepsilon^{2})$$

Таким образом, величина к ограничивает обусловленность линейных систем с участием А. По этой причине мы даем следующее определение:

Определение 3.5 (Матричное число обусловленности). Число обусловленности A Rn×n для данной матричной нормы · есть

конд
$$A$$
 AA 1 .

Если A необратима, мы берем cond A

Легко видеть, что cond A 1 для всех A, что масштабирование A не влияет на его число обусловленности и что число обусловленности единичной матрицы равно 1. Эти свойства контрастируют с определителем, который может увеличиваться и уменьшаться. по мере масштабирования A.

Если \cdot индуцируется векторной нормой и A обратим, то имеем

$$A = MAKC.$$
 $X = MAKC.$ $X =$

В этом случае номер условия А определяется как:

конд. A = макс.
$$x=0$$
 $x=0$ мин $x=0$ $y=0$ $x=0$

Другими словами, cond A измеряет максимальное или минимально возможное растяжение вектора х при A.

В более общем смысле желательное свойство устойчивости системы Ax = b состоит в том, что если A orb возмущена, решение x существенно не изменится. Наша мотивация для cond A показывает, что, когда число обусловленности мало, изменение x мало по сравнению с изменением A orb, как показано на рисунке HOMEP. В противном случае небольшое изменение параметров линейной системы может вызвать большие отклонения x; эта нестабильность может привести к тому, что линейные решатели будут делать большие ошибки в x из-за округления и других приближений в процессе решения. может быть столь же сложным, как

Норма А ¹ вычисление полной обратной границы А, число условия ¹. Один из способов снизить состоит в том, чтобы применить тождество А 1х А 1х. Таким образом, для любого х = 0 А 1х/х. Таким образом, 1 мы можем написать А

конд
$$A = AA$$
 1 $AA = AA$ 1.

Итак, мы можем ограничить число обусловленности, решив A 1х для некоторых векторов х; конечно, необходимость линейного решателя для нахождения A 1х создает циклическую зависимость от числа обусловленности для оценки качества оценки! Когда · индуцируется двукратной нормой, в следующих главах мы дадим более надежные оценки.

3.4 Проблемы

Что-то вроде:

- Регрессия ядра на примере §3.1.1
- Решение по методу наименьших норм для Ах =b, в противном случае матрица наименьших квадратов обратима
- Вариационные версии тихоновской регуляризации/«гребневой» регрессии (не обычный подход к этому, но что угодно); завершение недоопределенной истории таким образом
- 1 L подходы к регуляризации для контраста нарисуйте картину, почему следует ожидать разреженности, нарисуйте единичных кругов, покажите, что норма p не является нормой при p < 1, примите предел при p 0
- Mini-Riesz выведите матрицу для внутреннего продукта, используйте, чтобы показать, как вращать пространство
- трехдиагональное решение
- свойства номера условия

Machine Translated by Google

Глава 4

Пространства столбцов и QR

Один из способов интерпретации линейной задачи Ах = b для х состоит в том, что мы хотим записать b как линейную комбинацию столбцов A с весами, заданными в х. Эта перспектива не меняется, когда мы допускаем, что A Rm×n не является квадратом, но решение может не существовать или быть уникальным в зависимости от структуры пространства столбцов. По этим причинам некоторые методы факторизации матриц и анализа линейных систем ищут более простые представления пространства столбцов для устранения неоднозначности и более точного охвата, чем факторизация на основе строк, такая как LU.

4.1 Структура нормальных уравнений

Как мы показали, необходимым и достаточным условием для того, чтобы х было решением задачи наименьших квадратов Ах b, является то, что х удовлетворяет нормальным уравнениям (AA)х = A b. Эта теорема предполагает, что решение метода наименьших квадратов является достаточно простым расширением линейных методов. Такие методы, как факторизация Холецкого, также показывают, что особая структура задач наименьших квадратов может быть использована в интересах решателя.

Однако есть одна большая проблема, ограничивающая использование этого подхода. А пока предположим, что А квадрат; то мы можем написать:

КОНД
$$AA = A A(A A)$$

1
1
AAA
(A)
= A 2 A 12
= (КОНД A)²

То есть число обусловленности АА примерно равно квадрату числа обусловленности А! Таким образом, в то время как общие линейные стратегии могут работать с АА, когда задача наименьших квадратов является «легкой», когда столбцы А почти линейно зависимы, эти стратегии, вероятно, будут генерировать значительную ошибку, поскольку они не имеют дело с А напрямую.

Интуитивно основная причина того, что cond A может быть большим, заключается в том, что столбцы A могут выглядеть «похожими». Думайте о каждом столбце A как о векторе в Rm. Если два столбца аі и ај удовлетворяют аі , то остаточная длина ај методом наименьших квадратов b Ах, вероятно, не сильно пострадает, если мы заменим числа, кратные аі , на числа, кратные аі , или наоборот. Этот широкий диапазон почти, но не полностью эквивалентных решений приводит к плохой обусловленности. Однако, хотя результирующий вектор х нестабилен, произведение

Топор остается практически без изменений, по замыслу нашей замены. Следовательно, если мы хотим решить Ax b просто для записи b в пространстве столбцов A, будет достаточно любого решения.

Чтобы решить такие плохо обусловленные задачи, мы будем использовать альтернативную стратегию, уделяя больше внимания пространству столбцов матрицы A, а не прибегая к операциям со строками, как при исключении Гаусса. Таким образом, мы можем явно идентифицировать такие близкие зависимости и работать с ними численно стабильным образом.

4.2 Ортогональность

Мы определили, когда задача наименьших квадратов сложна, но мы могли бы также спросить, когда она наиболее проста. Если мы сможем свести систему к простому случаю, не вызывая при этом проблем с обусловленностью, мы найдем более устойчивый способ обойти проблемы, описанные в § 4.1.

Очевидно, что проще всего решить линейную систему In×nx = b: решение просто x b! Мы вряд ли введем эту конкретную линейную систему в наш решатель явно, но мы можем сделать это случайно при решении метода наименьших квадратов. В частности, даже когда A = In×n — на самом деле, A не обязательно должна быть квадратной матрицей — мы можем в особенно удачных обстоятельствах обнаружить, что нормальная матрица AA удовлетворяет AA = In×n. Чтобы не путать с общим случаем, мы будем использовать букву Q для обозначения такой матрицы.

Простая молитва о том, что QQ = In×n , вряд ли приведет к желаемой стратегии решения, но мы можем изучить этот случай, чтобы увидеть, как он становится таким благоприятным. Запишем столбцы матрицы Q в виде векторов $q1, \cdots, qn$ Rm. Тогда нетрудно проверить, что произведение QQ имеет следующую структуру:

Установка выражения справа равным In×n дает следующее соотношение:

Другими словами, столбцы Q имеют единичную длину и ортогональны друг другу. Мы говорим, что они образуют ортонормированный базис для пространства-столбца Q:

Определение 4.1 (Ортонормированный; ортогональная матрица). Набор векторов $\{v1, \cdots, vk\}$ ортонормирован, если vi=1 для всех i и $vi \cdot vj=0$ для всех i=j. Квадратная матрица, столбцы которой ортонормированы, называется ортогональной матрицей.

Мы мотивировали нашу дискуссию вопросом, когда мы можем ожидать, что $QQ = In \times n$. Теперь легко видеть, что это происходит, когда столбцы матрицы Q ортонормированы. Кроме того, если Q является квадратным и обратимым с $QQ = In \times n$, то, просто умножая обе части этого выражения на Q = 1, мы находим Q = 1 = Q. Таким образом, решение Qx = b в этом случае так же просто, как умножение обе стороны транспонированным Q.

Ортонормированность также имеет сильную геометрическую интерпретацию. Напомним из главы 0, что мы можем рассматривать два ортогональных вектора a и b как перпендикулярные. Итак, ортонормированный набор векторов

просто представляет собой набор перпендикулярных векторов единичной • Если Q ортогонален, то его действие длины в Rn , не влияющих на длину векторов:

$$Qx^{2} = x Q Qx = x In \times nx = x \cdot x = x$$

Точно так же Q не может влиять на угол между двумя векторами, поскольку:

$$(Qx) \cdot (Qy) = x Q Qy = x In \times ny = x \cdot y C$$
 этой точки

зрения, если Q ортогонален, то Q представляет собой изометрию Rn, т. е. сохраняет длины и углы. Другими словами, он может вращать или отражать векторы, но не может масштабировать или сдвигать их. На высоком уровне линейная алгебра ортогональных матриц проще, потому что их действие не влияет на геометрию основного пространства каким-либо нетривиальным образом.

4.2.1 Стратегия для неортогональных матриц

За исключением особых обстоятельств, большинство наших матриц А при решении Ах = b или соответствующей задачи наименьших квадратов не будут ортогональны, поэтому механизм § 4.2 не применяется напрямую. По этой причине мы должны сделать некоторые дополнительные вычисления, чтобы связать общий случай с ортогональным.

Возьмем матрицу A Rm×n, и обозначим его пространство столбцов как соl A; напомним, что соl A представляет диапазон столбцов A. Теперь предположим, что матрица B Rn×n обратима. Мы можем сделать простое наблюдение о пространстве столбцов AB относительно

пространства столбцов А: Лемма 4.1 (инвариантность пространства столбцов). Для любого A Rm×n и обратимого B Rn×n столбец A = столбец AB.

Доказательство. Предположим, что b соl A. Тогда по определению умножения на A существует x такой, что Ax = b. Тогда (AB) \cdot (B 1x) = Ax =b, sob col AB. Обратно, пусть с col AB, значит, существует y, для которого (AB)y = c. Тогда A \cdot (By) = c, что показывает, что c находится в столбце A.

для R = E. задачи квадратов Ax b могут упроститься. B частности, когда A = QR, мы можем записать решение A Ax = A b следующим образом:

Или, что то же самое, Rx = Q b

Таким образом, если мы создадим R как треугольную матрицу, то решение линейной системы Rx = Qb будет таким же простым, как обратная подстановка.

Наша задача в оставшейся части главы — разработать стратегии такой факторизации.

4.3 Ортогонализация Грама-Шмидта

Наш первый подход к нахождению QR-факторизаций является самым простым для описания и реализации, но может страдать числовыми проблемами. Мы используем его здесь в качестве первоначальной стратегии, а затем усовершенствуем его, улучшив операции.

4.3.1 Прогнозы

Предположим, у нас есть два вектора а и b. Тогда мы могли бы легко спросить: «Какое число, кратное a, ближе всего к tob?» Математически эта задача эквивалентна минимизации са b по всем возможным с R. Если представить a и b как матрицы размера $n \times 1$, a c — как матрицу размера 1×1 , то это не более чем нетрадиционная задача наименьших квадратов. \cdot c 6. В этом случае нормальные уравнения показывают $a \cdot c = ab$, или

$$c = \frac{a \cdot 6}{a \cdot a} = \frac{a \cdot 6}{a^2}.$$

Обозначим эту проекцию b на а как:

$$_{\text{проект}_a} = \text{ca} = \qquad \frac{\text{a} \cdot \text{bb}}{\text{a} \cdot \text{a}} = \frac{\text{a} \cdot \text{6}}{\text{a} - 2}$$

Очевидно, а б параллелен а. А как насчет остатка b — proj вычисление проекта, чтобы узнать:

аб? Мы можем сделать простой

$$a \cdot (6 \quad проект_a 6) = a \cdot 6 \quad a \cdot \frac{a \cdot 6}{a - 2} 2$$

$$= a \cdot 6 - \frac{a \cdot 6}{a \cdot 6} \cdot a \cdot 6$$

$$= 0$$

Таким образом, мы разложили b на компоненту, параллельную а, и другую компоненту, ортогональную тоа.

Предположим теперь, что a^1, a^2, · · · , a^k ортонормированы; мы наденем шляпы на векторы единичной длины. Тогда для любого одиночного і мы можем увидеть:

Член нормы не появляется, потому что aˆi = 1 по определению. Мы могли бы спроецировать b на промежуток {aˆ1, · · · , aˆk} , минимизируя следующую энергию по c1, . . . , ск € R:

Обратите внимание, что второй шаг здесь допустим только из-за ортонормированности. Как минимум, производная по сі равна нулю для каждого сі , что дает:

$$ci = a^i \cdot b$$

Таким образом, мы показали, что когда a^1, \cdots , a^k ортонормированы, выполняется следующее соотношение:

projspan
$$\{a^1, \dots, a^k\} = (a^1 \cdot b)a^1 + \dots + (a^k \cdot b)a^k$$

Это просто расширение нашей проекционной формулы, и с помощью аналогичного доказательства легко увидеть, что

$$a^i \cdot (b \quad projspan \{a^1, \dots, a^k\}) = 0.$$

То есть мы разделили b на составляющую, параллельную размаху a^i, и перпендикулярную невязку.

4.3.2 Ортогонализация Грама-Шмидта

Наши наблюдения выше приводят к простому алгоритму ортогонализации или нахождению ортогонального базиса $\{a^1, \dots, a^k\}$, размах которого такой же, как у набора линейно независимых входных векторов $\{v1, \dots, vk\}$:

1 комплект

То есть мы принимаем a¹ за единичный вектор, параллельный v1.

- 2. Для і от 2 до k,
 - (а) Вычислите проекцию

По определению a^1, \cdots , a^i 1 ортонормированы, поэтому применима наша вышеприведенная формула.

(б) Определить

Этот метод, известный как «ортогонализация Грама-Шмидта», представляет собой прямое применение нашего обсуждения выше. Ключ к доказательству этой техники состоит в том, чтобы заметить, что span $\{v1, \dots, vi\}$ = span $\{a^1, \dots, a^i\}$ для каждого і $\{1, \dots, k\}$. Шаг 1 явно делает это случайным для i = 1, а для i > 1 определение a^i на шаге 2b просто удаляет проекцию на векторы, которые мы уже видели.

Если мы начнем с матрицы A, столбцами которой являются v1, · · · ,vk , то мы можем реализовать Грэма Шмидта как серию операций со столбцами над A. Деление столбца i матрицы A на его норму эквивалентно последующему умножению A на аk × k диагональная матрица. Точно так же вычитание проекции столбца на ортонормированные столбцы слева от него, как на шаге 2, эквивалентно последующему умножению на верхнетреугольную матрицу: обязательно поймите, почему это так! Таким образом, наше обсуждение в § 4.2.1 применимо, и мы можем использовать Грамма-Шмидта, чтобы получить факторизацию A = QR.

К сожалению, алгоритм Грама-Шмидта может вносить серьезные числовые нестабильности из-за шага вычитания. Например, предположим, что мы предоставили векторы v1 = (1, 1) и v2 = (1 + ϵ , 1) в качестве входных данных для Грамма-Шмидта для некоторого 0 < ϵ 1. Обратите внимание, что очевидным основанием для span $\{v1,v2\}$ является $\{(1,0),(0,1)\}$. Но если применить Грамма-Шмидта, то получим:

$$^{2}1 = \frac{v1}{v1} = \frac{1}{2} = 11$$

$$p2 = 2\frac{2+e}{2} = 1$$

$$v2 \quad p2 = \frac{1+e1}{2} = \frac{2+e}{2} = 1$$

$$v3 \quad p4 = \frac{1}{2} = \frac{1}{2}$$

Обратите внимание, что v2 $p2 = (2/2) \cdot \epsilon$, поэтому для вычисления a^2 потребуется деление на скаляр порядка ϵ . Деление на маленькие числа — нестабильная числовая операция, которой следует избегать.

4.4 Преобразования домохозяев

В §4.2.1 мы мотивировали построение QR-факторизации пост-умножением и операциями со столбцами. Эта конструкция разумна в контексте анализа пространств столбцов, но, как мы видели в нашем выводе алгоритма Грама-Шмидта, полученные в результате численные методы могут быть нестабильными.

Вместо того, чтобы начинать с A и пост-умножать на операции со столбцами, чтобы получить Q = AE1 · · · Ek , однако мы можем сохранить нашу высокоуровневую стратегию от исключения Гаусса. То есть мы можем начать с A и предварительно умножить на ортогональные матрицы Qi , чтобы получить Qk · · · Q1A = R; эти Q будут действовать как операции со строками, исключая элементы A до тех пор, пока результирующий продукт R не станет верхнетреугольным. Тогда, благодаря ортогональности Q, мы можем записать A = QR, получив коэф ищиент QR , так как произведение ортогональных матриц ортогонально.

Матрицы операций со строками, которые мы использовали в исключении Гаусса и LU, не будут достаточными для факторизации QR, поскольку они не ортогональны. Было предложено несколько альтернатив; мы представим одну общую стратегию, представленную в 1958 году Алстоном Скоттом Хаусхолдером.

Пространство ортогональных матриц размера n × n очень велико, поэтому мы должны найти меньшее пространство Qi , c которым легче работать. Из наших геометрических обсуждений в § 4.2 мы знаем, что ортогональные матрицы должны сохранять углы и длины, поэтому интуитивно они могут вращать и отражать только векторы.

К счастью, отражения можно легко записать в виде проекций, как показано на рисунке НОМЕР. Предположим, у нас есть вектор b, который мы хотим отразить над вектором v. Мы показали, что невязка r b proj 2proj Мы в b перпендикулярна v. Как и на рисунке НОМЕР, разность в b отражает b над v.

можем расширить нашу формулу отражения следующим образом:

$$_{2 \text{проект}_{B}} \cdot \text{bb}$$
 $b = 2 \frac{V}{B \cdot B} V$ b по определению проекции $= 2 B \cdot \frac{B 6}{B B}$ b с использованием матричных обозначений $= \frac{2 B B}{B B}$ $B \times n 6$

Hv b, где отрицательный знак вводится для согласования с другими методами лечения

Таким образом, мы можем думать об отражении b над v как о применении линейного оператора Hv к b! Конечно, Hv без отрицания по-прежнему ортогонален, поэтому с этого момента мы будем использовать его.

Предположим, мы делаем первый шаг прямой замены во время исключения Гаусса. Затем мы хотим предварительно умножить A на матрицу, которая переводит первый столбец A, который мы будем обозначать, в некоторое число, кратное первому единичному вектору e1. Другими словами, мы хотим для некоторого с R:

ce1 = XBa
$$= B \times \Pi - \frac{2BB}{BB} \quad a$$

$$= a - 2B \quad \frac{Ba}{BB}$$

Перемещение терминов по шоу

$$v = (a ce1) \cdot \frac{BB}{2va}$$

Другими словами, v должно быть параллельно разности а се1. На самом деле масштабирование v не влияет на формулу для Hv , поэтому мы можем выбрать v = a се1. Затем, чтобы наши отношения сохранялись, мы должны иметь

$$1 = \frac{BB}{2Ba}$$

$$= \frac{a^{2} 2ce1 \cdot a + c^{2}}{2(a \cdot a ce1 \cdot a)}$$
Или 0 = a
$$= c = \pm a$$

При таком выборе с мы показали:

$$\mathsf{HvA} = \begin{array}{c} \mathsf{c} \times \times \times \\ \mathsf{0} \times \times \times \\ \vdots & \vdots & \vdots \\ \mathsf{0} \times \times \times \end{array}$$

Мы только что выполнили шаг, похожий на прямое исключение, используя только ортогональные матрицы! Исходя из этого, в обозначениях СІТЕ на k-м шаге триангуляризации имеем вектор а, который мы можем разделить на две составляющие:

Здесь a1 Rk и a2 Rm k. Мы хотим найти такое v, что

Следуя выводу, параллельному приведенному выше, легко показать, что

совершает именно это преобразование при с = ±a2; мы обычно выбираем знак с, чтобы избежать отмены, заставляя его иметь знак, противоположный знаку k-го значения ina.

Таким образом, алгоритм Хаусхолдера QR довольно прост. Для каждого столбца матрицы A мы вычисляем v, аннулирующее нижние элементы столбца, и применяем Hv к A. Конечным результатом является верхняя треугольная матрица R = Hvn··· Hv1 A. Ортогональная матрица Q задается произведением H, которое может быть не выше, сохранен как список векторов v, который соответствует нижнему треугольнику, как v1, показанному выше.

4.5 Сокращенная QR-факторизация

Мы завершаем наше обсуждение, возвращаясь к наиболее общему случаю Ax b, когда A Rm×n не является квадратом. Обратите внимание, что оба алгоритма, которые мы обсуждали в этой главе, могут разлагать неквадратные матрицы A на произведения QR, но результат несколько отличается:

- При применении Грама-Шмидта мы выполняем операции со столбцами над A, чтобы получить Q путем ортогонализации. По этой причине размерность A равна размерности Q, что дает Q Rm×n и R Rn×n:
- При использовании отражений Хаусхолдера мы получаем Q как произведение числа m \times m матрицы отражения, оставляя R Rm \times n ·

Предположим, что мы находимся в типичном случае для метода наименьших квадратов, для которого m n. Мы по-прежнему предпочитаем использовать метод Хаусхолдера из-за его численной стабильности, но теперь матрица Q размером m × m может оказаться слишком большой для хранения! К счастью, мы знаем, что R является верхнетреугольным. Например, рассмотрим структуру матрицы 5 × 3 R:

Легко видеть, что все, что находится ниже верхнего квадрата размера n × n матрицы R, должно быть равно нулю, что дает упрощенную формулу. **КАТИОН:**

$$A = QR = Q1 Q2$$
 $P1 = Q1R1$

Здесь Q1 Rm×n и R1 Rn×n по-прежнему содержит верхний треугольник R. Это называется «приведенным» QR-факторизация A, поскольку столбцы Q1 содержат базис для пространства столбцов A, а не для всего Rm; он занимает гораздо меньше места. Обратите внимание, что обсуждение в § 4.2.1 все еще применимо, поэтому сокращенная факторизация QR может использоваться для метода наименьших квадратов аналогичным образом.

4.6 Проблемы

- тридиагонализация с Хаусхолдером
- Даны
- Недоопределенный QR

Machine Translated by Google

Глава 5

Собственные векторы

Обратимся теперь к нелинейной задаче о матрицах: нахождение их собственных значений и собственных векторов. Собственные векторы х и соответствующие им собственные значения λ квадратной матрицы A определяются уравнением Ax = λx. Есть много способов увидеть, что эта задача нелинейна.

Например, существует произведение неизвестных λ и х, и, чтобы избежать тривиального решения х = 0, мы ограничиваем х = 1; это ограничение является круговым, а не линейным. Благодаря такой структуре наши методы нахождения собственных пространств будут существенно отличаться от методов решения и анализа линейных систем уравнений.

5.1 Мотивация

Несмотря на произвольный вид уравнения Ax = λx, задача нахождения собственных векторов и собственных значений естественно возникает во многих случаях. Мы мотивируем нашу дискуссию несколькими примерами ниже.

5.1.1 Статистика

Предположим, у нас есть механизм для сбора нескольких статистических наблюдений о наборе элементов. Например, в медицинском исследовании мы можем узнать возраст, вес, артериальное давление и частоту сердечных сокращений 100 пациентов. Тогда каждый пациент і может быть представлен точкой хі в R4, хранящей эти четыре значения.

Конечно, такая статистика может демонстрировать сильную корреляцию. Например, пациенты с более высоким кровяным давлением могут иметь более высокий вес или частоту сердечных сокращений. По этой причине, хотя мы собирали наши данные в R4, в действительности они могут — в некоторой приблизительной степени — жить в пространстве более низкого измерения, лучше фиксируя отношения между различными переменными.

А пока предположим, что на самом деле существует одномерное пространство, аппроксимирующее наш набор данных. Затем мы ожидаем, что все точки данных будут почти параллельны некоторому вектору v, так что каждую можно записать как xi civ для различных ci R. Ранее мы знали, что наилучшее приближение xi параллельно v есть proj

Здесь мы определяем v^* v/v. Конечно, величина v не имеет значения для рассматриваемой задачи, поэтому разумно искать в пространстве единичных векторов v^* .

Следуя схеме наименьших квадратов, у нас есть новая задача оптимизации:

Мы можем немного упростить нашу цель оптимизации:

2
 = xi (xi · v^)v^2 по определению проекции

= 8 xi 2 (xi · v^) 2 поскольку v^ = 1 и w 2 = константа 2 (xi · v^) 2

Этот вывод показывает, что мы можем решить эквивалентную задачу оптимизации:

максимизировать
$$X v^2$$

такой, что $V^2 = 1$.

где столбцы X являются векторами xi . Обратите внимание, что X $v^2 = v^2 \chi \chi v^2$, поэтому в примере 0.27 вектор v соответствует собственному вектору XX с наибольшим собственным значением. Вектор v^2 известен как первый главный компонент набора данных.

5.1.2 Дифференциальные уравнения

Многие физические силы могут быть записаны как функции положения. Например, сила пружины между двумя частицами в положениях x и y в R3 может быть записана как k(x y) по закону Гука; такие силы пружины используются для аппроксимации сил, скрепляющих ткань во многих системах моделирования. Хотя эти силы не обязательно линейны по положению, мы часто аппроксимируем их линейным образом. В частности, в физической системе с n частицами положения всех частиц одновременно кодируются вектором X R3n . Тогда, если мы примем такое приближение, мы можем написать, что силы в системе приблизительно равны F AX для некоторой матрицы A.

Вспомните второй закон Ньютона F = ma, или сила равна массе, умноженной на ускорение. В нашем контексте мы можем написать диагональную матрицу масс M R3n×3n, содержащую массу каждой частицы в системе. Тогда мы знаем, что F = MX, где штрих означает дифференцирование по времени. Конечно, X = (X), поэтому в итоге мы имеем систему уравнений первого порядка:

Здесь мы одновременно вычисляем как положение в X R3n , так и скорости V R3n всех n частиц как функции времени.

В более общем смысле дифференциальные уравнения вида x = Ax появляются во многих контекстах, включая моделирование ткани, пружин, тепла, волн и других явлений. Предположим, мы знаем собственные векторы

x1, . . . ,xk группы A, такие что Axi = λixi . Если мы запишем начальное условие дифференциального уравнения через собственные векторы, как

$$x(0) = c1x1 + \cdots + ckxk$$
,

то решение уравнения можно записать в замкнутом виде:

$$x(T) = c1e$$
 $\lambda 1t \lambda k t x 1 + \cdots + cke x k$,

Это решение легко проверить вручную. То есть, если мы запишем начальные условия этого дифференциального уравнения в терминах собственных векторов A, то мы знаем его решение для всех моментов времени t 0 бесплатно. Конечно, эта формула не является концом истории моделирования: нахождение полного набора собственных векторов A дорого обходится, а A может меняться со временем.

5.2 Спектральное вложение

Предположим, у нас есть набор из n элементов в наборе данных и мера wij 0 того, насколько похожи каждая пара элементов i и j; будем считать wij = wji. Например, может быть, нам дана коллекция фотографий, и мы используем wij , чтобы сравнить сходство их цветовых распределений. Возможно, мы пожелаем отсортировать фотографии на основе их сходства, чтобы упростить просмотр и изучение коллекции.

Одна из моделей упорядочения коллекции может состоять в том, чтобы присвоить номер хі каждому элементу і, попросив, чтобы подобным объектам были присвоены одинаковые номера. Мы можем измерить, насколько хорошо задание группирует похожие объекты, используя энергию

$$E(x) = \sup_{ij} wij(xi xj)^{2}.$$

То есть E(x) требует, чтобы элементы і и j с высокими показателями подобия wij были сопоставлены с соседними значениями. Конечно, минимизация E(x) без ограничений дает очевидный минимум: xi = const. для всех я. Добавление ограничения x = 1 не удаляет это постоянное решение! В частности, взятие xi = 1/ п для всех і дает x = 1 и E(x) = 0 неинтересным образом. Таким образом, мы должны удалить и этот случай:

минимизировать
$$E(x)$$
 такой, что x $= 1$ $1 \cdot x = 0$

Обратите внимание, что наше второе ограничение требует, чтобы сумма х была равна нулю.

Еще раз мы можем упростить энергию:

Легко проверить, что 1 является собственным вектором 2A 2W с собственным значением 0. Что еще интереснее, собственный вектор, соответствующий второму наименьшему собственному значению, соответствует решению нашей задачи минимизации выше! (TODO: добавьте доказательство ККТ из лекции)

5.3 Свойства собственных векторов

Мы установили множество приложений, нуждающихся в вычислении собственного пространства. Однако прежде чем мы сможем изучить алгоритмы для этой цели, мы более подробно рассмотрим структуру проблемы собственных значений.

Мы можем начать с нескольких определений, которые, вероятно, очевидны на данный момент:

Определение 5.1 (собственное значение и собственный вектор). Собственным вектором x=0 матрицы A Rn×n называется любой вектор, удовлетворяющий условию $Ax=\lambda x$ для некоторого λ R; соответствующее λ известно как собственное значение. Комплексные собственные значения и собственные векторы удовлетворяют тем же соотношениям с λ C и x Cn.

Определение 5.2 (Спектр и спектральный радиус). Спектр A — это множество собственных значений A. Спектральный радиус ρ(A) — это собственное значение λ, максимизирующее |λ|.

Масштаб собственного вектора не важен. В частности, масштабирование собственного вектора x с помощью с дает A (cx) = cAx = cA

Лемма 5.1 (теорема ЦИТЭ 2.1). Каждая матрица А Rn×n имеет хотя бы один (комплексный) собственный вектор.

Доказательство. Возьмем любой вектор х Rn\{0}. Множество $\{x, Ax, A 2x, \cdots Anx \}$ д \emptyset лжно быть линейно зависимым, поскольку оно содержит n + 1 вектор в n измерениях. Итак, существуют константы c0, . . . , cn R, cn = 0 такие, что

$$0 = c0x + c1Ax + \cdots + cnA$$

Мы можем записать многочлен

$$f(z) = c0 + c1z + \cdots + cnz$$

По основной теореме алгебры существует n корней zi C таких, что f(z) = cn(z)

$$z1)(z z2) \cdots (z zn).$$

Тогда у нас есть:

$$0 = c0x + c1Ax + \cdots + cnA$$
 H_{Mac} $= (c0 In \times n + c1A + \cdots + cnA n)x =$ $cn(A z1 In \times n) \cdots (A zn In \times n)x$ по нашей факторизации

Таким образом, по крайней мере один A zi In×n имеет нулевое пространство, что показывает, что существует v с Av = ziv, что и требовалось.

Есть еще один факт, который стоит проверить, чтобы мотивировать наше обсуждение вычисления собственных векторов. нация:

Лемма 5.2 (предложение 2.2 ЦИТЭ). Собственные векторы, соответствующие разным собственным значениям, должны быть линейно независимыми.

Доказательство. Предположим, что это не так. Тогда существуют собственные векторы х1, \cdots , хk с различными собственными значениями λ 1, \cdots , λ k, линейно зависимые. Отсюда следует, что существуют коэффициенты с1, \ldots , ck не все нули с 0 = c1x1 + \cdots + ckxk. Если предварительно умножить на матрицу (A λ 2 In×n) \cdots (A λ k In×n), найдем:

$$0 = (A \quad \lambda 2 \text{ In} \times \text{n}) \cdot \cdot \cdot \cdot (A \quad \lambda k \text{ In} \times \text{n}) (c1x1 + \cdot \cdot \cdot + ckxk)$$

= c1(λ1 \ \lambda 2)\cdot \cdot \((\lambda 1 \) \lambda k)x1 , τακ κακ Axi = \lambda ixi

Поскольку все λі различны, это показывает, что c1 = 0. Аналогичное доказательство показывает, что остальные сі должны быть равны нулю, что противоречит линейной зависимости.

Эта лемма показывает, что матрица размера n × n может иметь не более n различных собственных значений, поскольку набор из n собственных значений дает n линейно независимых векторов. Максимальное количество линейно независимых собственных векторов, соответствующих одному собственному значению λ , известно как геометрическая кратность λ .

Однако неверно, что матрица должна иметь ровно n линейно независимых собственных векторов. Это имеет место для многих матриц, которые мы будем называть недефектными:

Определение 5.3 (Бездефектный). Матрица A Rn×n называется недефектной или диагонализируемой, если ее собственные векторы порождают Rn ·

Мы называем такую матрицу диагонализируемой по следующей причине: если матрица диагонализируема, то она имеет п собственных векторов x1, . . . , xn Rn с соответствующими (возможно, не единственными) собственными значениями λ1, . . . , λн. Возьмем столбцы X в качестве векторов xi и определим D как диагональную матрицу с собственными значениями λ1, . . . , λn по диагонали. Тогда по определению собственных значений имеем AX = XD; это просто «сложенная» версия Axi = λixi . Другими словами,

$$Д = X - 1AX,$$

это означает, что А диагонализуется преобразованием подобия А Х 1АХ: Определение 5.4

(Подобные матрицы). Две матрицы A и B подобны, если существует такая T, что B = T 1AT.

Подобные матрицы имеют одинаковые собственные значения, так как если $Bx = \lambda x$, то T 1AT $x = \lambda x$. Эквивалентно, $A(Tx) = \lambda(Tx)$, показывая, что Tx является собственным вектором с собственным значением λ .

5.3.1 Симметричные и положительно определенные матрицы

Учитывая наше специальное рассмотрение нормальных матриц АА, неудивительно, что симметричные и/или положительно определенные матрицы имеют особую структуру собственных векторов. Если мы сможем проверить любое из этих свойств, можно будет использовать специализированные алгоритмы для более быстрого извлечения их собственных векторов.

Во-первых, мы можем доказать свойство симметричных матриц, которое устраняет необходимость в сложных арифметика. Начнем с обобщения симметричных матриц на матрицы из Cn×n:

Определение 5.5 (комплексно сопряженное). Комплексно-сопряженным числом z=a+bi=C является z=a-bi.

Определение 5.6 (сопряженное транспонирование). Сопряженное транспонирование А - Cm $^{\times}$ n равно AH - A $^{-}$.

Определение 5.7 (эрмитова матрица). Матрица A Сn×n эрмитова, если A = AH.

Обратите внимание, что симметричная матрица A Rn×n автоматически является эрмитовой, поскольку не имеет комплексной части. Имея это небольшое обобщение, мы можем доказать свойство симметрии для собственных значений.

Наше доказательство будет использовать скалярное произведение векторов в Сп , заданное выражением

$$x, y = xiy^i$$
,

где x,y Сn . Заметим, что это определение снова совпадает с $x \cdot y$, когда x,y Rn . По большей части свойства этого скалярного произведения совпадают со свойствами скалярного произведения на Rn , за исключением того, что v, w = w ,v. , заметный

Лемма 5.3. Все собственные значения эрмитовых матриц действительны.

Доказательство. Предположим, что A $Cn \times n$ эрмитово с $Ax = \lambda x$. Масштабируя, мы можем принять х 1. Тогда $^2 = X$, X = M мы имеем:

^{Таким образом, $\lambda = \lambda^-$,} что может произойти, только если λ R, что и требовалось.

Симметричные и эрмитовы матрицы также обладают особым свойством ортогональности собственных векторов. Торы:

Лемма 5.4. Собственные векторы, соответствующие различным собственным значениям эрмитовых матриц, должны быть ортогональны.

П

Доказательство. Предположим, что A $Cn \times n$ эрмитово, и пусть $\lambda = \mu$ с $Ax = \lambda x$ и $Ay = \mu y$. По предыдущей лемме мы знаем λ , μ R. Тогда $Ax,y = \lambda x,y$. Но поскольку A эрмитово, мы можем также написать Ax,y = x, $Ay = \mu x,y$. Таким образом, $\lambda x,y = \mu x,y$. Поскольку $\lambda = \mu$, мы должны иметь x,y = 0.

Наконец, мы можем сформулировать без доказательства венчающий результат линейной алгебры — спектральную теорему. Эта теорема утверждает, что ни одна симметричная или эрмитова матрица не может быть дефектной, а это означает, что матрица размера n × n, удовлетворяющая этому свойству, имеет ровно n ортогональных собственных векторов.

Теорема 5.1 (спектральная теорема). Предположим, что A Сп×п эрмитово (если A Rn×n , то пусть оно симметрично). Тогда А имеет ровно п ортонормированных собственных векторов x1, · · · ,xn с (возможно, повторяющимися) собственными значениями λ1, . . . , λн . Другими словами, существуют ортонормированная матрица собственных векторов X и диагональная матрица собственных значений D такие, что D = X AX.

Из этой теоремы следует, что любой вектор у Rn можно разложить в линейную комбинацию собственных векторов эрмитовой матрицы A. В этом базисе многие вычисления проще, как показано ниже:

Пример 5.1 (Вычисление с использованием собственных векторов). Возьми x1, . . . , xn Rn — собственные векторы единичной длины симметричной матрицы A Rn×n . Предположим, мы хотим решить Ay = b. Мы можем написать

$$b = c1x1 + \cdots + cnxn$$

где $ci = b \cdot xi$ по ортонормированности. Нетрудно догадаться о следующем решении:

$$\frac{c1}{m}$$
 cn x1 + · · · + y = $\frac{c1}{m}$ xn. λ 1 λ n

В частности, мы находим:

$$A$$
й = A $\frac{c1}{\lambda} \frac{cn}{\lambda} \frac{n}{\lambda} \frac{n}{\lambda} \frac{n}{\lambda} \frac{n}{c1} \frac{n}{cn}$

"=" — $Ax1 + \dots +$ — AxH
 $= c1x1 + \dots + cnxn$

= b , по желанию.

Приведенный выше расчет является как положительным, так и отрицательным результатом. Он показывает, что при заданных собственных векторах симметричного A такие операции, как инверсия, выполняются просто. С другой стороны, это означает, что найти полный набор собственных векторов симметричной матрицы A «как минимум» так же сложно, как решить Ax = b.

Возвращаясь от нашего экскурса в комплексные числа, мы вернемся к действительным числам, чтобы доказать один последний полезный, хотя и простой факт о положительно определенных матрицах:

Лемма 5.5. Все собственные значения положительно определенных матриц неотрицательны.

Доказательство. Возьмем положительно определенное A Rn×n и предположим, что $Ax = \lambda x$, где x = 1. По положительной определенности мы знаем, что x Ax = 0. Но x Ax = x (λx) = λx = λx , что и требовалось.

5.3.2 Специализированные свойства1

Характеристический полином

Напомним, что определитель матрицы det A удовлетворяет соотношению det A = 0 тогда и только тогда, когда A обратима. Таким образом, одним из способов нахождения собственных значений матрицы является нахождение корней характеристического многочлена

$$pA(\lambda) = det(A \quad \lambda In \times n).$$

Мы не будем определять определители в нашем обсуждении здесь, но упрощение рА показывает, что это полином n-й степени от λ. Это дает альтернативную причину, по которой существует не более n различных собственных значений, поскольку у этой функции не более n корней.

Исходя из этой конструкции, мы можем определить алгебраическую кратность собственного значения как его кратность как корня рА. Легко видеть, что алгебраическая кратность не меньше геометрической

¹Этот раздел можно пропустить, если у читателей недостаточно знаний, но он включен для полноты картины.

множественность. Если алгебраическая кратность равна 1, корень называется простым, поскольку он соответствует одному собственному вектору, который линейно зависит от любых других. Собственные значения, у которых алгебраическая и геометрическая кратности не равны, называются дефектными.

В численном анализе мы избегаем обсуждения определителя матрицы. Хотя это удобная теоретическая конструкция, ее практическое использование ограничено. Детерминанты трудно вычислить. На самом деле алгоритмы собственных значений не пытаются найти корни рА , поскольку для этого потребовалось бы вычисление определителя. Кроме того, определитель det A не имеет ничего общего с обусловленностью A, поэтому близкий к нулю определитель det(A λIn×n) может не показывать, что λ почти является собственным значением A.

Нормальная форма Джордана

Мы можем диагонализовать матрицу только тогда, когда она имеет полное собственное пространство. Однако все матрицы подобны матрице в жордановой нормальной форме, которая имеет следующий вид: • Ненулевые

значения находятся на диагональных элементах аіі и на «наддиагональных» элементах аі(і+1).

- Диагональные значения это собственные значения, повторяющиеся столько раз, сколько их кратно; матрица диагональ блока вокруг этих кластеров.
- Внедиагональные значения равны 1 или 0.

Таким образом, форма выглядит примерно так:

Нормальная форма Жордана привлекательна теоретически, потому что она всегда существует, но структура 1/0 дискретна и неустойчива при численном возмущении.

5.4 Вычисление собственных значений

Вычисление и оценка собственных значений матрицы — хорошо изученная задача, имеющая множество возможных решений. Каждое решение настраивается для разных ситуаций, и для достижения максимальной подготовки или скорости требуется экспериментировать с несколькими методами. Здесь мы рассмотрим несколько наиболее популярных и простых решений проблемы собственных значений, часто встречающихся на практике.

5.4.1 Итерация мощности

Пока предположим, что A Rn×n симметрична. Тогда по спектральной теореме мы можем записать собственные векторы x1, . . . , xn Rn; мы сортируем их так, чтобы их соответствующие собственные значения удовлетворяли |λ1| |λ2| · · · · |λn|.

Предположим, мы взяли произвольный вектор v. Поскольку собственные векторы матрицы A порождают Rn, мы можем написать:

$$v = c1x1 + \cdots + cnxn$$
.

Затем,

Av = c1Ax1 + ··· + cnAxn
= c1\lambda1x1 + ··· + cn\lambdanxn , tak kak Axi = \lambdaixi \lambda2 \lambdan = \lambda1 c1x1

$$\frac{+ c2x2 + \cdots + cnxn \lambda1}{\lambda1}$$
A 2 v = λ 1 2 c1x1 + $\frac{\lambda^2}{\lambda_1}$ 2 c2x2 + ··· + $\frac{\lambda n}{\lambda_1}$ 2 cnxn

$$\vdots$$

$$A^{KK}_{v} = \lambda$$
1 2 c1x1 + $\frac{\lambda^2}{\lambda_1}$ 2 c2x2 + ··· + $\frac{\lambda n}{\lambda_1}$ 2 cnxn

Обратите внимание, что при k отношение ($\lambda i/\lambda 1$) 0, если только $\lambda i = \lambda 1$, поскольку $\lambda 1$ по определению имеет наибольшую величину из всех собственных значений. Таким образом, если х является проекцией вектора v на пространство собственных векторов v собственными значениями v все более точно выполняется следующее приближение:

Это наблюдение приводит к чрезвычайно простому алгоритму вычисления собственного вектора х оператора A, соответствующее наибольшему собственному значению λ 1:

- 1. Возьмем v1 Rn произвольный ненулевой вектор.
- 2. Итерировать до сходимости для увеличения k:

Этот алгоритм, известный как степенная итерация, будет создавать векторы vk, все более и более параллельные желаемому x1. Он гарантированно сходится, даже если A асимметричен, хотя доказательство этого факта сложнее, чем приведенный выше вывод. Единственный случай, когда этот метод может дать сбой, — это если мы случайно выберем v1 так, что c1 = 0, но шансы на то, что это произойдет, ничтожны.

Конечно , если |\(\)1 | > 1, то vk при k , нежелательное свойство для арифметики с плавающей запятой. Напомним, что нас интересует только направление собственного вектора, а не его величина, поэтому масштабирование не влияет на качество нашего решения. Таким образом, чтобы избежать этой ситуации расхождения, мы можем просто нормализовать на каждом шаге, создавая алгоритм итерации нормализованной мощности:

- 1. Возьмем v1 Rn произвольный ненулевой вектор.
- 2. Итерировать до сходимости для увеличения k:

$$k = Avk \quad 1 \text{ W}$$

$$BK = \frac{\Box W K}{\Box W K}$$

5.4.2 Обратная итерация

Теперь у нас есть стратегия нахождения собственного значения λ1 наибольшей величины . Предположим, что A обратим, так что мы можем вычислить y = A 1v, решив Ay = v с помощью методов, описанных в предыдущих главах.

Если $Ax = \lambda x$, то $x = \lambda A$ 1x или, что то же самое,

A
$$1 x = \frac{1}{\lambda} x$$
.

- 1. Возьмем v1 Rn произвольный ненулевой вектор.
- 2. Итерировать до сходимости для увеличения k:
 - (a) Решите для w k : Aw k = vk 1
 - (б) Нормализация: vk = ______

Мы постоянно решаем системы уравнений, используя одну и ту же матрицу A, что является прекрасным применением методов факторизации из предыдущих глав. Например, если мы напишем A = LU, то мы могли бы сформулировать эквивалентную, но значительно более эффективную версию обратной степенной итерации:

- Фактор A = ЛУ
- 2. Возьмем v1 Rn произвольный ненулевой вектор.
- 3. Итерировать до сходимости для увеличения k:
 - (a) Решить для ук с помощью прямой подстановки: Lyk = vk 1 (b)

Решить для w k c помощью обратной подстановки: Uw k = yk (c)

Hормировать: vk =
$$\frac{\underline{\mathsf{w}\,\mathsf{K}}}{\underline{\mathsf{w}\,\mathsf{K}}}$$

5.4.3 Переключение

Предположим, что λ2 является собственным значением со второй по величине величиной А. Учитывая наш первоначальный вывод степенной итерации, легко увидеть, что степенная итерация сходится быстрее всего, когда |λ2/λ1| мала, так как в этом случае мощность (λ2/λ1) быстро убывает. Напротив, если это отношение близко к 1, может потребоваться много итераций мощности, прежде чем будет выделен один собственный вектор.

Если собственные значения A равны $\lambda 1, \ldots, \lambda n$, то легко видеть, что собственные значения A σ In×n равны $\lambda 1$ $\sigma, \ldots, \lambda n$ σ . Тогда одна из стратегий быстрой сходимости степенной итерации состоит в том, чтобы выбрать σ таким образом, чтобы:

$$\frac{\lambda 2}{\lambda 1} \frac{\sigma}{\sigma} < \frac{\lambda 2}{\lambda 1}$$

Конечно, угадывание такого σ может быть искусством, так как собственные значения A, очевидно, изначально неизвестны. Точно так же, если мы считаем, что σ близко к собственному значению A, то A σIn×n имеет собственное значение, близкое к 0, которое мы можем выявить с помощью обратной итерации.

Одна стратегия, в которой используется это наблюдение, известна как итерация по фактору Рэлея. Если у нас есть фиксированное предположение о собственном векторе х матрицы A, то с помощью ЧИСЛА аппроксимация методом наименьших квадратов соответствующего собственного значения о определяется выражением

$$\sigma = \frac{\underset{\text{MKC}}{\sum} \frac{2}{2}}{\frac{2}{2}}.$$

Эта дробь известна как частное Рэлея. Таким образом, мы можем попытаться увеличить сходимость, итерируя следующим образом:

- 1. В качестве v1 Rn возьмем произвольный ненулевой вектор или начальное приближение собственного вектора.
- 2. Итерировать до сходимости для увеличения k:
 - (а) Запишите текущую оценку собственного значения

$$\sigma k = \frac{v k \quad 1Avk \quad 1}{BK-1 \quad 2}$$

Эта стратегия сходится намного быстрее при хорошем начальном предположении, но матрица А ок In×n меняется на каждой итерации и не может быть предварительно факторизована с использованием LU или большинства других стратегий. Таким образом, требуется меньше итераций, но каждая итерация занимает больше времени!

5.4.4 Поиск нескольких собственных значений

До сих пор мы описывали методы нахождения одной пары собственное значение/собственный вектор: итерация по мощности для нахождения наибольшего собственного значения, обратная итерация для нахождения наименьшего и переход к промежуточным целевым значениям. Конечно, для многих приложений одного собственного значения недостаточно. К счастью, мы можем расширить наши стратегии, чтобы справиться и с этим случаем.

дефляция

Вспомните нашу стратегию степенной итерации: выберите произвольное v1 и итеративно умножайте его на A, пока не останется только самое большое собственное значение λ1 . Возьмем x1 за соответствующий собственный вектор.

Однако мы быстро отбросили маловероятную неисправность этого алгоритма, когда v1 ·x1 = 0.
В этом случае, независимо от того, сколько раз вы предварительно умножаете на A, вы никогда не восстановите вектор.

параллельно x1, так как вы не можете усилить нулевую составляющую. Вероятность выбора такой v1 ровно равна нулю, поэтому во всех случаях, кроме самых пагубных, итерация мощности остается безопасной.

Мы можем перевернуть этот недостаток с ног на голову, чтобы сформулировать стратегию нахождения более чем одного собственного значения, когда A симметрично. Предположим, что мы находим x1 и λ1 с помощью степенной итерации, как и раньше. Теперь мы перезапускаем итерацию мощности, но перед тем, как начать проект x1 из v1. Затем, поскольку собственные векторы A ортогональны, степенная итерация восстановит второе по величине собственное значение!

Из-за числовых проблем может случиться так, что применение A к вектору вводит небольшую компоненту, параллельную х1. На практике мы можем избежать этого эффекта, проецируя каждую итерацию. В итоге эта стратегия дает следующий алгоритм вычисления собственных значений в порядке убывания величины:

- Для каждого желаемого собственного значения = 1, 2, . . .
 - 1. Возьмем v1 Rn произвольный ненулевой вектор.
 - 2. Итерировать до сходимости для увеличения k: (a)

Проецировать собственные векторы, которые мы уже вычислили:

$$uk = vk \quad 1 \qquad projspan\{x1,...,x \quad 1\} \ vk \quad 1$$
 (b) Умножить Auk = w k (c)
 Нормировать: vk =
$$\frac{ux}{}$$

3. Добавьте результат итерации к множеству хі

Внутренний цикл эквивалентен итерации мощности на матрице AP, где P проецирует x1, . . . , x-1 . Легко видеть, что AP имеет те же собственные векторы, что и A; его собственные значения равны λ , . . . , λn с обнулением остальных собственных значений.

В более общем плане стратегия дефляции включает изменение матрицы А таким образом, чтобы степенная итерация выявляла собственный вектор, который вы еще не вычислили. Например, АР — это модификация А, так что большие собственные значения, которые мы уже вычислили, обнуляются.

Наша проекционная стратегия терпит неудачу, если A асимметрична, так как в этом случае ее собственные векторы могут не быть ортогональными. В этом случае могут работать и другие менее очевидные стратегии дефляции. Например, предположим, что Ax1 = λ1x1 с x1 = 1. Возьмем H в качестве матрицы Хаусхолдера, такой что Hx1 = е1, первый стандартный базисный вектор. Преобразования подобия снова не влияют на набор собственных векторов, поэтому мы можем попробовать выполнить сопряжение с помощью H. Рассмотрим, что происходит, когда мы умножаем HAH на e1:

НАНе1 = HAHe1 , так как H симметричен
$$= HAx1 \ , \text{ так как Hx1 =e1 } \text{ и H} = \lambda 1 \text{Hx1 }, ^2 = \text{B} \times \text{П}$$

$$\text{так как Ax1 = } \lambda 1 \text{x1} \qquad \qquad = \lambda 1 \text{e1 по определению H}$$

Таким образом, первый столбец НАН равен λ1e1, показывая, что НАН имеет следующую структуру (СITE 3ДОРОВьE):

Матрица В $R(n-1)\times(n-1)$ имеет собственные значения $\lambda 2, \ldots, \lambda h$. Таким образом, другая стратегия дефляции состоит в том, чтобы строить все меньшие и меньшие В-матрицы, где каждое собственное значение вычисляется с использованием степенной итерации.

QR-итерация

Дефляция имеет тот недостаток, что мы должны вычислять каждый собственный вектор отдельно, что может быть медленным и может накапливать ошибку, если отдельные собственные значения неточны. Наши оставшиеся стратегии пытаются найти более одного собственного вектора за раз.

Напомним, что подобные матрицы A и B = T 1AT должны иметь одинаковые собственные значения. Таким образом, алгоритм, пытающийся найти собственные значения A, может свободно применять преобразования подобия к A. Конечно, применение T в целом может быть трудным предложением, поскольку фактически потребовало бы инвертирования T, поэтому мы ищем T-матрицы, обратные матрицы которых легко найти. применять.

Одним из наших мотивов для получения QR-факторизации было то, что матрица Q ортогональна, удовлетворяя Q-1 = Q. Таким образом, Q и Q-1 одинаково просты в применении, что делает ортогональные матрицы сильным выбором для преобразований подобия.

Но какую ортогональную матрицу Q выбрать? В идеале Q должен включать в себя структуру A, но при этом быть простым для вычислений. Неясно, как стратегически применять простые преобразования, такие как матрицы Хаусхолдера, для выявления множества собственных значений2, но мы знаем, как сгенерировать одно такое Q, просто разложив на множители A = QR. Тогда мы могли бы связать A с Q, чтобы найти: Q 1AQ = Q (QR)Q = (Q Q)RQ = RQ

Удивительно, но сопряжение A = QR ортогональной матрицей Q идентично записи произведения RQ!

Основываясь на этом рассуждении, в 1950-х годах несколько групп европейских математиков выдвинули гипотезу тот же элегантный итерационный алгоритм для нахождения собственных значений матрицы A:

- 1. Возьмем А1 = А.
- 2. При к = 1, 2, . . .
 - (a) Фактор Ak = QkRk . (б)
 - Запишите Ak+1 = RkQk.

Согласно приведенному выше выводу, все матрицы Ak имеют те же собственные значения, что и A. Кроме того, предположим, что матрицы Ak сходятся K некоторому A. Тогда мы можем разложить A = Q R, и по сходимости мы знаем, что A = Q R = R Q. По ЧИСЛУ собственные значения R — это просто значения вдоль диагонали R, а по ЧИСЛУ произведение R Q = A, в свою очередь, должно иметь те же собственные значения. Наконец, по построению A имеет те же собственные значения, что и A. Итак, мы показали, что если QR-итерация сходится, она выявляет собственные значения A прямым способом.

Конечно, приведенный выше вывод предполагает, что существует A с Ak A при k . На самом деле итерация QR — это стабильный метод, который гарантирует сходимость во многих важных ситуациях, и сходимость можно даже улучшить, изменив стратегии. Мы не будем здесь выводить точные условия, но вместо этого можем дать некоторое интуитивное представление о том, почему эта, казалось бы, произвольная стратегия должна сходиться. Ниже мы приводим некоторую интуицию для симметричного случая A = A, который легче анализировать благодаря ортогональности собственных векторов в этом случае.

Предположим, что столбцы A заданы a1, . . . ,an, и рассмотрим матрицу A можно написать:

к для больших к. Мы

2Более продвинутые методы, однако, делают именно это!

K

Согласно нашему выводу степенной итерации, первый столбец A в общем случае параллелен собственному вектору x1 с наибольшей величиной |\(\begin{align*} \begin{align*} \begin{align*} \begin{align*} \left(\begin{align*} \left(\begin{align*} \left(\begin{align*} \left(\begin{align*} \begin{align*} \left(\begin{align*} \left(

к Конечно, вычисление А для большого k приводит число обусловленности A к k-й степени, поэтому QR на результирующей матрице, скорее всего, не удастся; это ясно видно, поскольку все столбцы A должны выглядеть как x1 для больших k. Однако мы можем сделать следующее наблюдение:

```
A = Q1R1
A^2 = (Q1R1)(Q1R1)
= Q1(R1Q1)R1
= Q1Q2R2R1, используя приведенное выше обозначение итерации QR, поскольку A2 = R1Q1
\vdots
A^K = Q1Q2 \cdots QkRkRk 1 \cdots R1
```

Группировка переменных Qi и переменных Ri по отдельности обеспечивает QR-факторизацию A, и мы ожидаем, что столбцы Q1 · · · Qk будут сходиться к собственным векторам A.

Аналогичным образом можно найти

где Ak - k-я матрица из итерации QR. Таким образом, Ak+1 — это просто матрица A, сопряженная произведением Q ^-k Q1 \cdots Qk . Ранее мы утверждали, что столбцы матрицы Q ^-k сходятся к собственным векторам матрицы A. Таким образом, поскольку сопряжение матрицей собственных векторов дает диагональную матрицу собственных значений, мы знаем, что $Ak+1 = Q^{-k}AQ^-$ будет иметь приближенные собственные значения A вдоль его диагонали при k .

Методы подпространств Крылова

Наше обоснование итерации QR включало анализ столбцов A итерации при k как расширение мощности. В более общем случае для вектора b Rn , мы можем рассмотреть так называемую матрицу Крылова

Методы анализа Кk для нахождения собственных векторов и собственных значений обычно известны как методы подпространств Крылова. Например, итерационный алгоритм Арнольди использует ортогонализацию Грама-Шмидта для поддержания ортогонального базиса {q1,...,qk} для пространства столбца Kk:

1. Начнем с того, что возьмем q1 в качестве произвольного вектора единичной нормы.

2. При к = 2, 3, . . .

(a) Takeak =
$$Aqk 1$$
 (b)

Спроецируйте значения q, которые вы уже вычислили:

(c) Перенормируйте, чтобы найти следующее qk = bk/bk.

Матрица Qk, столбцами которой являются векторы, найденные выше, является ортогональной матрицей с тем же пространством столбцов, что и Кk, и оценки собственных значений могут быть восстановлены из структуры Q AQk. Использование метода Грама-Шмидта делает этот метод нестабильным, и синхронизация становится все хуже по мере увеличения k, однако для того, чтобы сделать его возможным, требуется так много расширений. Например, одна стратегия включает в себя выполнение нескольких итераций Arnoldi, использование выходных данных для получения лучшего предположения для начального q1 и перезапуск. Методы этого класса подходят для задач, требующих нескольких собственных векторов на одном из концов спектра без вычисления полного набора.

5.5 Чувствительность и кондиционирование

Как и предупреждали, мы обрисовали в общих чертах только несколько методов собственных значений из богатой и давней литературы. Для нахождения спектров были опробованы почти любые алгоритмические методы, от итерационных методов до поиска корней характеристического многочлена и методов, которые делят матрицы на блоки для параллельной обработки.

Как и в линейных решателях, мы можем оценить обусловленность проблемы собственных значений независимо от метода решения. Этот анализ может помочь понять, будет ли успешной упрощенная итерационная схема для нахождения собственных векторов данной матрицы или же необходимы более сложные методы; важно отметить, что обусловленность задачи на собственные значения — это не то же самое, что число обусловленности матрицы для решения систем, поскольку это отдельные задачения — это не то же самое, что число обусловленности матрицы для решения систем, поскольку это отдельные задачения — это не то же самое, что число обусловленности матрицы для решения систем, поскольку это отдельные задачения — это не то же самое, что число обусловленности матрицы для решения систем, поскольку это отдельные задачения систем.

Предположим, что матрица A имеет собственный вектор x с собственным значением λ. Анализ обусловленности проблемы собственных значений включает анализ устойчивости x и λ к возмущениям в A. С этой целью мы могли бы возмущать A малой матрицей δA, тем самым изменяя набор собственных векторов. В частности, мы можем записать собственные векторы A + δA как возмущения собственных векторов A, решив задачу

$$(A + \delta A)(x + \delta x) = (\lambda + \delta \lambda)(x + \delta x).$$

Расширение обеих сторон дает:

$$Ax + A\delta x + \delta A \cdot x + \delta A \cdot \delta x = \lambda x + \lambda \delta x + \delta \lambda \cdot x + \delta \lambda \cdot \delta x$$

Предполагая δА малым, будем считать3, что δх и δλ также малы. Произведения между этими переменными тогда незначительны, что дает следующее приближение:

$$Ax + A\delta x + \delta A \cdot x$$
 $\lambda x + \lambda \delta x + \delta \lambda \cdot x$

Поскольку $Ax = \lambda x$, мы можем вычесть это значение с обеих сторон, чтобы найти:

$$A\delta x + \delta A \cdot x$$
 $\lambda \delta x + \delta \lambda \cdot x$

$$y (A\delta x + \delta A \cdot x)$$
 $y (\lambda \delta x + \delta \lambda \cdot x)$

Поскольку A $y = \lambda y$, мы можем упростить:

Перегруппировка выходов:

$$\frac{y (\delta A)x \delta \lambda}{vx}$$

Предположим, что x = 1 и y = 1. Тогда, если мы возьмем нормы с обеих сторон, мы найдем:

$$|\delta\lambda| = \frac{\delta A2 |}{|y \cdot x|}$$

Таким образом, в целом условие задачи на собственные значения зависит от размера возмущения δA , как и ожидалось, и угла между левым и правым собственными векторами x и y. Мы можем использовать $1/x \cdot y$ как приблизительное число обусловленности. Обратите внимание, что x = y, когда A симметричен, что дает число условия 1; это отражает тот факт, что собственные векторы симметричных матриц ортогональны и, следовательно, максимально разделены.

5.6 Проблемы

³Это предположение следует проверить более строго!

Глава 6

Разложение по сингулярным значениям

В главе 5 мы вывели ряд алгоритмов для вычисления собственных значений и собственных векторов матриц А Rn×n . Разработав этот механизм, мы завершаем наше первоначальное обсуждение числовой линейной алгебры выводом и использованием одной окончательной матричной факторизации, которая существует для любой матрицы А Rm×n : разложения по сингулярным числам (SVD).

6.1 Получение SVD

Для А Rm×n , мы можем думать о функции х Ах как об отображении, переводящем точки в Rn в точки в Rm. С этой точки зрения мы могли бы задаться вопросом, что происходит с геометрией Rn в процессе, и, в частности, с влиянием А на длины и углы между векторами.

Применяя нашу обычную отправную точку для проблем с собственными значениями, мы можем задать вопрос о влиянии A на длины векторов путем изучения критических точек отношения

при различных значениях х. Масштабирование х не имеет значения, так как

 $|\alpha|$ Таким образом, мы можем ограничить наш поиск до x с x = 1. Более того, поскольку R(x) 0, вместо этого мы можем рассматривать [R(x)]2 = Ax = x A Ax. Однако, как мы показали в предыдущих главах, критические точки xAx при условии x = 1 — это в точности собственные векторы xi , удовлетворяющие условию A Axi = λ ixi ; заметим, что λ i 0 и xi · xj = 0, когда i = j, так как AA симметрична и положительно полуопределена.

Основываясь на нашем использовании функции R, базис {xi} является разумным для изучения геометрических эффектов A. Возвращаясь к этой первоначальной цели, определим yi — Axi . Мы можем сделать дополнительное наблюдение о yi , раскрывающее еще более сильную структуру собственных значений:

$$\lambda$$
iyi = λ i · Axi по определению yi

- $= A(\lambda ixi)$
- = A(A Axi), так как хі является собственным вектором АА
- = (AA)(Axi) по ассоциативности
- = (AA) yи

Таким образом, имеем два случая:

1. Когда $\lambda i = 0$, то yi = 0. В этом случае xi является собственным вектором AA и yi = Axi является $AAxi = \frac{2}{x}$ соответствующий собственный вектор AA с yi = Axi = Axi λixi .

2. При λі = 0 уі =0.

Идентичное доказательство показывает, что если у — собственный вектор АА, то х — А у — либо нуль, либо собственный вектор АА с тем же собственным значением.

Возьмем за k число строго положительных собственных значений $\lambda i > 0$, обсуждавшихся выше. По нашей конструкции выше мы можем взять $x1, \ldots, xk$ Rn — собственные векторы AA и соответствующие им собственные векторы $y1, \ldots, yk$ Rm группы AA такие, что

$$Aкси = \lambda ixi$$

$$AAyi = \lambda iyi$$

для собственных значений $\lambda i > 0$; здесь мы нормализуем так, что xi = yi = 1 для всех i. Следуя традиционным обозначениям, мы можем определить матрицы V^- Rn×k и U^- Rm×k , столбцами которых являются xi и yi соответственно.

Мы можем изучить влияние этих новых базисных матриц на А. Возьмем еі за і-й стандартный базисный вектор. Затем,

$$U^-AV^-ei = U^-Axi$$
 по определению V^-

"=" $\frac{1}{\lambda i}$ $U^-A(\lambda ixi)$, так как мы предполагали $\lambda i > 0$

"=" $\frac{1}{\lambda i}$ $U^-A(A Axi)$, так как хі является собственным вектором AA

"=" $\frac{1}{U^-}$ $U^-(AA)Axi$ по ассоциативности λi 1

"=" $\frac{1}{U^-}$ $U^-(AA)yi$, поскольку мы перемасштабировали так, что $yi = 0$

1 $\lambda i = \lambda i U^-yi$, поскольку AAyi = $\lambda i yi$

= $\lambda i ei$

Возьмем Σ^- = diag(λ 1, ..., λ k). Тогда приведенный выше вывод показывает, что $U^-AV^- = \Sigma^-$

Дополните столбцы U $^-$ и V $^-$ до U $^-$ Rm $^+$ m и V $^-$ Rn $^+$ n , добавив ортонормированные векторы xi и yi с Axi =0 и AAyi =0 соответственно. В этом случае легко показать UAVei =0 и/или UAV =0 . Таким образом, если мы возьмем $^{\rm e}_{\rm a}$

тогда мы можем расширить наше предыдущее соотношение, чтобы показать UAV = Σ или, что то же самое,

$$A = U\Sigma V$$
.

SVD обеспечивает полную геометрическую характеристику действия А. Поскольку U и V ортогональны, их можно рассматривать как матрицы вращения; как диагональная матрица Σ просто масштабирует отдельные координаты. Таким образом, все матрицы A Rm×n являются композицией поворота, масштаба и второго поворота.

6.1.1 Вычисление SVD

Напомним, что столбцы V просто являются собственными векторами AA, поэтому их можно вычислить с помощью методов, обсуждавшихся в предыдущей главе. Поскольку $A = U\Sigma V$, мы знаем, что $AV = U\Sigma .$ Таким образом, столбцы U, соответствующие ненулевым сингулярным числам в Σ , просто являются нормированными столбцами AV; остальные столбцы удовлетворяют AAui = 0, что можно решить с помощью факторизации LU.

Эта стратегия ни в коем случае не является самым эффективным или стабильным подходом к вычислению SVD, но она достаточно хорошо работает для многих приложений. Мы опустим более специализированные подходы к нахождению SVD, но отметим, что многие из них являются простыми расширениями степенной итерации и других уже рассмотренных нами стратегий, которые работают без явного формирования АА или АА.

6.2 Применение СВД

Оставшуюся часть этой главы мы посвящаем знакомству со многими приложениями SVD. SVD появляется бесчисленное количество раз как в теории, так и в практике численной линейной алгебры, и ее важность трудно переоценить.

6.2.1 Решение линейных систем и псевдообратных систем

В частном случае, когда A Rn×n является квадратным и обратимым, важно отметить, что SVD можно использовать для решения линейной задачи Ax =b. В частности, имеем UΣ V x =b, или

$$x = V\Sigma$$
 1U 6.

В этом случае Σ является квадратной диагональной матрицей, поэтому 1 просто матрица, диагональные элементы которой Σ 1/σi.

Чтобы охватить все три случая, мы можем решить задачу оптимизации следующего вида:

минимизировать x^2 такое, что A Ax = A b

Другими словами, эта оптимизация требует, чтобы x удовлетворял нормальным уравнениям с наименьшей возможной нормой. Теперь давайте напишем A = UΣ V . Затем,

AA = (UΣ V) (UΣ V)
$$= VΣ \quad U \quad UΣ \quad V, \text{ так как (AB)} = BA$$
$$= VΣ \quad Σ \quad V, \text{ так как U ортогонален}$$

Таким образом, спрашивать, что A Ax = A b, то же самое, что спрашивать

$$V\Sigma \Sigma V x = V\Sigma U G$$

Или, что то же самое, Σ y = d

если мы возьмем d Ub и y V x. Обратите внимание, что у = x, поскольку U ортогонален, поэтому наша оптимизация принимает вид:

минимизировать 2 y так, чтобы Σ y = d

Однако, поскольку Σ является диагональным, условие Σ у = d просто утверждает оіуі = di ; поэтому всякий раз, когда оі = 0, мы должны иметь уі = di/оі . Когда оі = 0, ограничений на уі нет , поэтому, поскольку мы минимизируем у, мы могли бы также взять уі = 0. Другими словами, решением этой оптимизации является у = Σ + d, где Σ + Rn×m имеет следующий вид:

$$\Sigma^{+}$$
 1/оі і = j, оі = 0 и і k 0 в противном случае

Эта форма, в свою очередь, дает $x = Vy = V\Sigma + d = V\Sigma + Ub$.

С этой мотивацией мы делаем следующее определение:

Определение 6.1 (псевдообратное). Псевдообратным к $A = U\Sigma V Rm \times n$ является $A + V\Sigma + U Rn \times m$.

Наш вывод выше показывает, что псевдоинверсия А обладает следующими свойствами:

- 1 Когда A является квадратным и обратимым, A [∓] A
- Когда А переопределено, A +b дает решение методом наименьших квадратов для Ах b.
- Когда A недоопределенно, A +b дает решение методом наименьших квадратов для Ax b с минимальным (евклидова) норма.

Таким образом, мы, наконец, можем объединить недоопределенный, полностью определенный и переопределенный случаи Ax b.

6.2.2 Разложение на внешние произведения и аппроксимации низкого ранга

Если мы разложим произведение $A = U\Sigma V$, легко показать, что из этого соотношения следует:

где $\min\{m, n\}$, a ui и vi — i-e столбцы U и V соответственно. Наша сумма достигает только $\min\{m, n\}$, так как мы знаем, что оставшиеся столбцы U или V будут обнулены с помощью Σ .

Это выражение показывает, что любую матрицу можно разложить как сумму внешних произведений векторы:

Определение 6.2 (Внешнее произведение). Внешнее произведение $u=Rm\ u\ v=Rm\ v=Rm\$

Предположим, мы хотим написать произведение Ах. Тогда вместо этого мы могли бы написать:

$$Ax = \prod_{g=1}^g \sigma i u i v i x$$

$$= \sigma i u i (v \prod_{g \in N(K)}^g N(K))$$

$$= \prod_{g \in G} i (v i \cdot x) u i, поскольку x \cdot y = xy$$

Таким образом, применение A к x равносильно линейному объединению векторов ui с весами σi(vi · x). Эта стратегия вычисления Ax может обеспечить значительную экономию, когда количество ненулевых значений σi относительно невелико. Кроме того, мы можем игнорировать малые значения σi , эффективно усекая эту сумму для аппроксимации Ax с меньшими затратами.

Точно так же из §6.2.1 мы можем записать псевдоинверсию А как:

$$A + = \frac{\text{Buy s}}{\sigma i = 0}.$$

Очевидно, мы можем применить тот же прием для вычисления A +x, и фактически мы можем аппроксимировать A +x, оценивая только те члены в сумме, для которых оі относительно мало. На практике мы вычисляем сингулярные значения оі как квадратные корни из собственных значений AA или AA, и такие методы, как степенная итерация, могут использоваться для выявления частичного, а не полного набора собственных значений. Таким образом, если нам нужно будет решить ряд задач наименьших квадратов Axi bi для различных bi и нас удовлетворит аппроксимация xi , может быть полезно сначала вычислить наименьшие значения оі и использовать приведенную выше аппроксимацию. Эта стратегия также избавляет от необходимости вычислять или сохранять полную матрицу A + и может быть точной, когда A имеет широкий диапазон сингулярных значений.

Возвращаясь к нашему исходному обозначению $A = U\Sigma V$, наши рассуждения выше эффективно показывают, что потенциально полезным приближением A является $A^{\sim} U\Sigma^{\sim} V$, где Σ^{\sim} округляет малые значения Σ до нуля. Легко проверить, что пространство-столбец A^{\sim} имеет размерность, равную количеству ненулевых значений на диагонали $\Sigma^{\sim} \Sigma$. На самом деле, это приближение не является специальной оценкой, а скорее решает сложную задачу оптимизации с помощью следующей известной теорема (изложена без доказательства):

Теорема 6.1 (Экарт-Янг, 1936). Предположим, что А получается из A = UΣ V путем усечения всех, кроме A A Fro и A A 2, значений оі множества A, до нуля. Тогда A минимизирует оба ограничения, зависящих от k наибольших сингулярных заключающиеся в том, что пространство столбцов A имеет размерность не более k.

6.2.3 Нормы матрицы

Построение SVD также позволяет нам вернуться к нашему обсуждению матричных норм из §3.3.1. Например, вспомните, что мы определили норму Фробениуса для А как

Если мы напишем A = U Σ V , мы можем упростить это выражение:

А
$$\frac{2}{\Phi_{\text{PO}}}$$
 $\frac{2}{\Theta_{\text{PO}}}$ $\frac{2}{\Theta_$

Таким образом, норма Фробениуса множества А Rm×n есть сумма квадратов его сингулярных значений.

Этот результат представляет теоретический интерес, но на практике основное определение нормы Фробениуса уже несложно оценить. Что еще интереснее, напомним, что индуцированная двойная норма А задается выражением

A
$$\frac{2}{2} = \max\{\lambda : \text{существует } x \in \mathbb{R}$$
 C A Ax = λx }.

Теперь, когда мы изучили проблемы с собственными значениями, мы понимаем, что это значение представляет собой квадратный корень из наибольшего собственного значения AA или, что то же самое,

A2 =
$$max{\sigma i}$$
.

Другими словами, мы можем прочитать двойную норму А непосредственно из его собственных значений.

Точно так же напомним, что число обусловленности А определяется выражением cond A = A2A 12. По нашему вывод A +, сингулярные значения A должны быть величины, обратные сингулярным числам A. В сочетании с нашим упрощением A2 получается:

усл A =
$$\frac{\sigma max}{max}$$
.

Это выражение дает стратегию для оценки обусловленности А. Конечно, вычисление от требует решения систем Ax = b, процесс, который сам по себе может страдать от плохой обусловленности A; если это проблема, то обусловливание можно ограничить и аппроксимировать, используя различные аппроксимации сингулярных значений A.

6.2.4 Проблема Прокруста и выравнивание

Многие методы компьютерного зрения включают выравнивание трехмерных фигур. Например, предположим, что у нас есть трехмерный сканер, который собирает два облака точек одного и того же жесткого объекта с разных точек зрения. Типичной задачей может быть совмещение этих двух облаков точек в единую систему координат.

Поскольку объект является жестким, мы ожидаем, что существуют некоторая матрица поворота R и перемещение t R3 , такие что вращение первого облака точек на R и последующее перемещение на байт выравнивает два набора данных. Наша работа состоит в том, чтобы оценить эти два объекта.

Если два сканирования перекрываются, пользователь или автоматизированная система могут отметить п соответствующих точек, которые соответствуют между двумя сканированиями; мы можем хранить их в двух матрицах X1, X2 R3×n. Тогда для каждого столбца x1i из X1 и x2i из X2 мы ожидаем Rx1i +t = x2i. Мы можем написать функцию энергии, измеряющую, насколько верно это соотношение:

Если мы зафиксируем R и минимизируем по tot, оптимизация E, очевидно, станет задачей наименьших квадратов. Теперь предположим, что мы оптимизируем для R с фиксированным значением. Это то же самое, что минимизировать это X2, переведенные на байт, при условии, что R является матрицей вращения 3 × 3, RX1 X, где столбцы X, то есть RR = I3×3. Это известно как ортогональная проблема Прокруста.

Для решения этой задачи введем след квадратной матрицы следующим образом:

Определение 6.3 (След). След A Rn×n есть сумма его диагоналей:

Таким образом, мы хотим максимизи ровать $tr(X RX1) c RR = I3 \times 3$. В упражнениях вы докажете, что tr(AB) = tr(BA). Таким образом, наша цель может быть немного упрощена до tr(RC) c C X1X 2. Применяя СВД, если разложить $C = U\Sigma V$, то можно упростить еще больше:

```
tr(RC) = tr(RU\Sigma\ V\ ) по определению = tr((V\ RU)\Sigma\ ),\ \tauак как tr(AB) = tr(BA) = tr(R^{\sim}\ \Sigma\ ) , ecли определить R^{\sim} = V\ RU, что также является ортогональным = \sigmair^{\sim} ii , так как \Sigma диагональна
```

Поскольку R^{\sim} ортогонален, все его столбцы имеют единичную длину. Отсюда следует, что r^{\sim} іі 1, иначе норма столбца і была бы слишком велика. Поскольку σ і 0 для всех і, этот аргумент показывает, что мы можем максимизировать tr(RC), взяв $R^{\sim}=I3\times3$. Отменив наши замены, мы получим $R=VRU^{\sim}=VU$.

В более общем виде мы показали следующее:

Теорема 6.2 (ортогональный прокруст). Ортогональная матрица R минимизирует RX Y VU, где 2 дан кем-то SVD применяется к фактору XY = UΣ V .

Возвращаясь к проблеме выравнивания, одной из типичных стратегий является альтернативный подход:

- 1. Зафиксировать R и минимизировать E относительно tot.
- 2. Зафиксировать полученное t и минимизировать E по R при условии, что RR = I3×3.
- 3. Вернитесь к шагу 1.

Энергия E уменьшается с каждым шагом и, таким образом, сходится к локальному минимуму. Поскольку мы никогда не оптимизируем и R одновременно, мы не можем гарантировать, что результатом будет минимально возможное значение E, но на практике этот метод работает хорошо.

6.2.5 Анализ основных компонентов (АПК)

Вспомните настройку из §5.1.1: мы хотим найти низкоразмерную аппроксимацию набора точек данных, которые мы можем сохранить в матрице X Rn×k для k наблюдений в n измерениях. Ранее мы показали, что если нам разрешено только одно измерение, наилучшее возможное направление задается доминирующим собственным вектором XX.

Вместо этого предположим, что нам разрешено проектировать на диапазон d векторов c d min{k, n}, и мы хотим выбрать эти векторы оптимально. Мы могли бы записать их в матрице C размера n × d; поскольку мы можем применить Грамма-Шмидта к любому набору векторов, мы можем предположить, что столбцы C ортонормированы, показывая CC = Id×d . Поскольку C имеет ортонормированные столбцы, по нормальным уравнениям проекция X на пространство столбцов C задается CCX.

$$X-CCX$$
 $\frac{2}{\Phi_{po}} = tr((X CCX) (X CCX))$ поскольку A $\frac{2}{\Phi_{po}} = tp(AA)$ $= tr(X X 2X CCX + X CCCCX) = const.$ $tr(X CCX)$, поскольку CC = Id×d $= -CX$

Таким образом, мы можем максимизировать
фро; для статистиков это показывает, когда строки X имеют среднее значение СX, равное нулю, что мы хотим максимизировать дисперсию проекции СX.

Теперь предположим, что мы факторизуем $X = U\Sigma V$. Тогда мы хотим максимизировать $C = C^{\sim} \Sigma$ Fro = $U\Sigma V$ Fro Σ^{\sim} CFro по ортогональности V, если мы возьмем $C^{\sim} = CU$. Если элементами C^{\sim} являются C^{\sim} іј, то разложение этой нормы дает

$$\Sigma$$
 C[~] 2_{\phipo} = 2 g 2 g c⁻ ij.

По ортогональности столбцов C^{\sim} мы знаем, что $i c^{\sim} = 1$ дл $_{\rm p}^2$ всех j, u, поскольку C^{\sim} может иметь менее n столбцов, $j c^{\sim}$ не боль $j e^2$ 1 в 1 приведенный выше функцине образом, что $j e^2$ не боль $j e^2$ путем сортировки столбцов $j e^2$ таким образом, что $j e^2$ $j e^2$ не образом, что $j e^2$ $j e^2$

Мы показали, что SVD для X можно использовать для решения такой задачи анализа основных компонентов (PCA). На практике строки X обычно сдвигаются так, чтобы иметь среднее значение, равное нулю, перед выполнением SVD; как показано на рис. HOMEP, это центрирует набор данных относительно начала координат, предоставляя более значимые векторы PCA ui.

6.3 Проблемы

Machine Translated by Google

Часть 3 Нелинейные методы



Глава 7

Нелинейные системы

Как бы мы ни старались, просто невозможно выразить все системы уравнений в линейной структуре, которую мы разработали в последних нескольких главах. Едва ли нужно мотивировать использование логарифмов, экспонент, тригонометрических функций, модулей, многочленов и т. д. в практических задачах, но, за исключением нескольких частных случаев, ни одна из этих функций не является линейной. Когда эти функции появляются, мы должны использовать более общий, хотя и менее эффективный набор механизмов.

7.1 Задачи с одной переменной

Мы начнем наше обсуждение с рассмотрения задач с одной скалярной переменной. В частности, для заданной функции f(x): R R мы хотим разработать стратегии для нахождения точек x \mathring{R} , таких что f(x) = 0; мы называем x корнем f . Задачи c одной переменной в линейной алгебре, в частности, не = b/a. Однако интересный; в конце концов, мы можем решить уравнение ax b = 0 в закрытой форме, поскольку x b = 0 в закрытой форме, посколь

7.1.1 Характеристика проблем

Мы больше не можем предполагать, что функция f линейна, но без каких-либо предположений о ее структуре мы вряд ли продвинемся вперед в решении систем с одной переменной. Например, решатель гарантированно не найдет нули f (x), заданного выражением

$$e(x) = \begin{cases} 1 & x & 0 & 1 \\ x > 0 & \end{cases}$$

Или хуже:

$$e(x) = \begin{cases} 1 & x & Q & 1 \\ uhave \end{cases}$$

Эти примеры тривиальны в том смысле, что разумный клиент программы для поиска корней вряд ли будет ожидать успеха в этом случае, но гораздо менее очевидные случаи не намного сложнее. строить.

По этой причине мы должны добавить некоторые «упорядочивающие» предположения о том, что f обеспечивает точку опоры для возможности разработки методов поиска корней. Ниже приведены типичные такие предположения, перечисленные в порядке возрастания силы:

- Непрерывность: функция f является непрерывной, если ее можно нарисовать, не поднимая пера; более формально f непрерывна, если разность f(x) f(y) обращается в нуль при x y.
- Липшиц: функция f является липшицевой, если существует константа C такая, что | e (x) e (y) | C|x y|; Липшицевы функции не обязательно должны быть дифференцируемыми, но их скорость изменения ограничена.
- Дифференцируемость: функция f является дифференцируемой, если ее производная f существует для всех х.

прои**хк**о**дниве онб дуиффаруниц**ируемо k раз и каждая из этих k производных непрерывна; • C : Функция C означает, что все С и непрерывны.

По мере того как мы будем добавлять все более и более сильные предположения о f, мы сможем разработать более эффективные алгоритмы для решения f(x) = 0. Мы проиллюстрируем этот эффект, рассмотрев несколько алгоритмов ниже.

7.1.2 Непрерывность и бисекция

Предположим, что все, что мы знаем о f, это то, что оно непрерывно. В этом случае мы можем сформулировать интуитивную теорему из стандартного исчисления с одной переменной:

Теорема 7.1 (теорема о промежуточном значении). Предположим, что f:[a,b] R непрерывна. Предположим, что f(x) < u < f(y). Тогда существует z между x и y такое, что f(z) = u.

Другими словами, функция f должна достигать всех значений между f(x) и f(y).

Предположим, что нам дана на вход функция f, а также два значения u r такие, что $f() \cdot f(r) < 0$; обратите внимание, что это означает, что f() u f(r) имеют противоположный знак. Тогда по теореме о промежуточном значении мы знаем, что где-то между u r есть корень f() Это обеспечивает очевидную стратегию деления пополам для нахождения x

Эта стратегия просто итеративно делит интервал [,r] пополам, каждый раз сохраняя сторону, в которой, как известно, существует корень. Ясно, что по теореме о промежуточном значении он сходится безоговорочно, в том смысле, что до тех пор, пока f() · f(r) < 0, в конечном счете и r, и r гарантированно сходятся к допустимому корню х

7.1.3 Анализ поиска корней

Разделение пополам — самый простой, но не обязательно самый эффективный метод поиска корней. Как и в большинстве методов собственных значений, деление пополам по своей сути является итеративным и может никогда не дать точного решения. Однако *** мы можем спросить, насколько близко значение ск функции с в k-й итерации к корню х , который мы надеемся вычислить. Этот анализ обеспечит основу для сравнения с другими методами.

В общем, предположим, что мы можем установить границу ошибки Ek, такую, что оценка xk корня во время k-й итерации метода меснахождения корня удовлетворяет условию | xk | x | < Эк . Очевидно, что любой алгоритм c Ek | 0 представляет собой сходящуюся схему; скорость сходимости, однако, может быть охарактеризована скоростью, с которой Ek приближается к 0. находятся в интервале [k ,rk], верхняя граница Например, в делении пополам, поскольку ошибка как

делим интервал пополам на каждой итерации, мы знаем, что Ek+1 Čk, так и х определяется выражением Ek | рк - к |. Поскольку мы = 1/2Ek . Поскольку Ek+1 линейно в Ek , мы говорим, что деление пополам обладает линейной сходимостью.

7.1.4 Итерация с фиксированной точкой

Разделение пополам гарантированно сходится к корню для любой непрерывной функции f , но если мы узнаем больше о f, мы сможем сформулировать алгоритмы, которые сходятся быстрее.

Например, предположим, что мы хотим найти х такое, что g(x) = x *; конечно, эта установка эквивалентна задаче поиска корня, так как решение f(x) = 0 совпадает с решением f(x) + x = x. Однако в качестве дополнительной информации мы также можем знать, что g(x) + x = x.

Система g(x) = x предлагает потенциальную стратегию, которую мы могли бы предположить:

- 1. Возьмем х0 за начальное предположение корня.
- 2. Повторить xk = g(xk 1).

Если эта стратегия сходится, ясно, что результатом является неподвижная точка q, удовлетворяющая вышеуказанным критериям.

К счастью, свойство Липшица гарантирует, что эта стратегия сходится к корню, если он существует.

Если взять Ek = |xk - x|, то мы имеем следующее свойство:

Индуктивное применение этого утверждения показывает, что Ek C ^K |E0 | 0 при k . Таким образом, итерация с фиксированной точкой сходится к желаемому х *!

В самом деле, если g липшицева с константой C < 1 в окрестности [х δ , х $*+\delta$], то до тех пор, пока x0 выбран в этом интервале, итерация с фиксированной точкой будет сходиться. Это верно, поскольку приведенное выше выражение для Ek показывает, что оно сжимается с каждой итерацией. и [g (х)]

Один важный случай имеет место, когда g есть C, 1 < 1. В силу непрерывности g в этом случае мы + δ] , в котором |g (x)| < 1 ϵ известно, что существует некоторая окрестность N = [x δ , x для * для любого x N, некоторого выбора достаточно малого $\epsilon > 0.1$. Возьмем любые x, y N. Тогда имеем

$$|\Gamma(x) - \Gamma(y)| = |\Gamma(\theta)| \cdot |x - y|$$
 по теореме о среднем значении основного исчисления для некоторого θ [x, y] < (1 ϵ) |x y|

Это показывает, что g является липшицевым c константой 1 ϵ < 1 в N. Таким образом, когда g непрерывно дифференцируема и g (х) < 1, итерация c фиксированной точкой будет сходиться к х , когда начальное приближениех0 близко.

¹Это утверждение трудно разобрать: убедитесь, что вы его понимаете!

Пока что у нас мало причин использовать итерацию с фиксированной точкой: мы показали, что она гарантированно сходится только тогда, когда g является липшицевой, а наши рассуждения о Ek показывают линейную сходимость, подобную делению пополам. Однако есть один случай, когда итерация с фиксированной точкой дает преимущество.

Предположим, что g дифференцируема c g (x) = 0. Тогда член первого порядка исчезает в ряду Тейлора для g, оставляя после себя:

$$1 g(xk) = g(x) + g(x^{-})(xk x) 2 + O(x\kappa - x^{*})$$

Таким образом, в данном случае имеем:

Таким образом, в этом случае Ek является квадратичным по Ek 1, поэтому мы говорим, что итерация с фиксированной точкой может иметь квадратичную сходимость; обратите внимание, что это доказательство квадратичной сходимости верно только потому, что мы уже знаем Ek 0 из нашего более общего доказательства сходимости. Это означает, что Ek 0 намного быстрее, поэтому нам потребуется меньше итераций, чтобы получить разумный корень.

Пример 7.1 (Сходимость итерации по фиксированной точке).

7.1.5 Метод Ньютона

Мы еще раз сузим наш класс функций, чтобы получить метод, который имеет более последовательную квадратичную сходимость. Теперь предположим, что мы снова хотим решить f(x) = 0, но теперь мы предполагаем, что f является 1, а условием C, несколько более жестким, чем условие Липшица.

В точке xk R, поскольку f теперь дифференцируема, мы можем аппроксимировать ее касательной: f(x)

$$f(xk) + f(xk)(x xk)$$

Решение этого приближения для f(x) 0 дает корень f(xk) f (xk)

Повторение этой формулы известно как метод Ньютона для поиска корней, и он сводится к итеративному решению линейной аппроксимации нелинейной задачи.

Обратите внимание, что если мы определим

$$(x) = x - f(x) \qquad \frac{f(x) \, \Gamma}{},$$

тогда метод Ньютона сводится к итерации с фиксированной точкой на g. Дифференцируя, находим:

$$(x) = 1 - \frac{x(x) r^2 f(x) f(x) по}{2}$$
 фактор-правилу $f(x) f(x) f(x)$

итерация с вывода предположим, что х фиксированной точкой выше, мы знаем, что метод Ньютона сходится квадратично к х * для достаточно близкого первоначальная догадка. Таким образом, когда f дифференцируемо с простым корнем, метод Ньютона обеспечивает формулу итерации с фиксированной точкой, которая гарантированно сходится квадратично; однако, когда х не простое, сходимость может быть линейной или хуже.

Вывод метода Ньютона предполагает другие методы, полученные с использованием большего количества членов в ряду Тейлора. Например, «метод Галлея» добавляет члены, включающие f, в итерации, а класс «методов Хаусхолдера» принимает произвольное число производных. Эти методы обеспечивают сходимость еще более высокого порядка за счет необходимости оценивать более сложные итерации и возможность более экзотических режимов отказа. Другие методы заменяют ряд Тейлора другими основными формами; например, линейно-дробная интерполяция использует рациональные функции для лучшей аппроксимации функций с асимптотной структурой.

7.1.6 Метод секущих

Одна проблема эффективности, которую мы еще не рассмотрели, — это стоимость оценки f и ее производных.

Если f — очень сложная функция, мы можем захотеть свести к минимуму количество вычислений f или, что еще хуже, f. Более высокие порядки сходимости помогают решить эту проблему, но мы также можем разработать численные методы, которые позволяют избежать оценки дорогостоящих производных.

Пример 7.2 (Дизайн). Предположим, мы проектируем ракету и хотим знать, сколько топлива добавить в двигатель. Для данного количества галлонов х мы можем написать функцию f (x), определяющую максимальную высоту ракеты; наши инженеры уточнили, что мы хотим, чтобы ракета достигла высоты h, поэтому нам нужно решить f(x) = h. Вычисление f(x) включает в себя моделирование взлета ракеты и отслеживание расхода топлива, что является дорогостоящим мероприятием, и хотя мы можем подозревать, что f дифференцируема, мы не сможем оценить f за какое-то практическое время.

Одной из стратегий разработки методов с меньшим воздействием является максимально возможное повторное использование данных. Для например, мы легко могли бы аппроксимировать:

$$\frac{f(xk) - f(xk - 1) f(xk)}{xk - xk-1}.$$

То есть, поскольку нам нужно было вычислить f(xk 1) на предыдущей итерации, мы просто используем наклон f(xk) для аппроксимации производной. Конечно, это приближение работает хорошо, особенно когда xk близки к сходимости.

Подключив наше приближение к методу Ньютона, мы получим новую итеративную схему:

$$x\kappa+1 = x\kappa - \frac{f(xk)(xk xk 1)}{f(xk) f(xk 1)}$$

Обратите внимание, что пользователь должен будет предоставить два начальных предположения x0 и x-1 , чтобы запустить эту схему, или может запустить одну итерацию Newton, чтобы запустить ее.

Анализ метода секущих несколько сложнее, чем другие методы, которые мы рассматриваем, поскольку он использует как f(xk), так и f(xk 1); доказательство его сходимости выходит за рамки нашего обсуждения. Интересно, что анализ ошибок показывает, что ошибка уменьшается со скоростью 1+ 5/2 («золотое сечение») между линейной и квадратичной; поскольку сходимость близка к сходимости метода Ньютона без необходимости вычисления f, метод секущих может обеспечить сильную альтернативу.

7.1.7 Гибридные методы

Можно выполнить дополнительную разработку, чтобы попытаться объединить преимущества различных алгоритмов поиска корней. Например, мы могли бы сделать следующие замечания о двух рассмотренных нами методах:

- Гарантируется безусловная сходимость пополам, но только с линейной скоростью.
- Метод секущих сходится быстрее, когда достигает корня, но в некоторых случаях может не сходится. сходятся.

Предположим, мы заключили в скобки корень f(x) в интервале [k, rk], как при делении пополам. Мы можем сказать, Если наша текущая оценка x * определяется выражением xk = -K что когда $|-\varphi(\kappa)| < |-f(rk)|$ и xk = rk в противном случае. мы отслеживаем xk и xk = 1, то мы могли бы взять xk+1 как следующую оценку корня, полученную методом секущих. Однако если xk находится вне интервала [k, rk], мы можем заменить его на xk+rk/2. Это исправление гарантирует, что xk+1 = -k0, и независимо от выбора мы можем обновить до допустимой скобки [k+1, rk+1]1, как при делении пополам, проверив знак xk+11. Этот алгоритм известен как «метод Деккера».

Стратегия, описанная выше, пытается объединить безусловную сходимость деления пополам с более сильными корневыми оценками метода секущих. Во многих случаях он успешен, но скорость его сходимости довольно трудно анализировать; специализированные режимы отказа могут свести этот метод к линейной сходимости или к худшему — фактически в некоторых случаях деление пополам неожиданно может сходиться быстрее! Другие методы, например, «метод Брента», делают шаги пополам чаще, чтобы избежать этого случая, и могут демонстрировать гарантированное поведение за счет несколько более сложной реализации.

7.1.8 Случай с одной переменной: сводка

Теперь мы представили и проанализировали ряд методов решения f(x) = 0 в случае одной переменной. На данный момент, вероятно, очевидно, что мы коснулись только поверхности таких методов; существует много итерационных схем поиска корней, все с разными гарантиями, скоростями сходимости и предостережениями. Несмотря на это, на основе нашего опыта мы можем сделать ряд наблюдений:

- Из-за возможной общей формы f маловероятно мы сможем точно найти корни х , и вместо того довольствуемся итерационными схемами.
- Мы хотим, чтобы последовательность xk корневых оценок достигла x* как можно быстрее. Если Ek является границей ошибки, то мы можем охарактеризовать ряд ситуаций сходимости, предполагая, что Ek 0 при k . Полный список условий, которые должны выполняться, когда k достаточно велико, приведен ниже:
 - 1. Линейная сходимость: Ek+1 СЕк для некоторого С < 1 2.
 - Суперлинейная сходимость: Ek+1 СЕг для r > 1 (теперь мы не требуем C < 1, поскольку, если Ek достаточно мало, степень r может сократить последствия C)
 - 3. Квадратичная сходимость: Ek+1 CE2 4. $_{\rm K}$ Кубическая сходимость: Ek+1 CE3 $_{\rm K}$ (и так далее)
- Метод может сходиться быстрее, но на каждой отдельной итерации требует дополнительных вычислений; по этой причине может быть предпочтительнее сделать больше итераций более простого метода, чем меньше итераций более сложного.

7.2 Многовариантные задачи

В некоторых приложениях может потребоваться решение более общей задачи f(x) =0 для функции f : Rn Rm. Мы уже видели один пример этой проблемы при решении Ax = b, что эквивалентно нахождению корней f (x) Ax b, но общий случай значительно сложнее. В частности, такие стратегии, как деление пополам, трудно расширять, поскольку теперь мы в значительной степени гарантируем, что все m различных значений равны нулю одновременно.

7.2.1 Метод Ньютона

К счастью, одна из наших стратегий расширяется очень просто. Напомним, что для f : Rn — Rm мы можем записать матрицу Якоби, которая дает производную каждой компоненты f в каждом из координатных направлений:

(Df)ij dxj
$$\frac{A \varphi^{\mu}}{}$$

Мы можем использовать якобиан f, чтобы расширить наш вывод метода Ньютона на несколько измерений. В частности, приближение первого порядка f определяется выражением:

$$f(x)$$
 $f(xk) + Df(xk) \cdot (x xk)$.

Замена желаемого f (x) = 0 дает следующую линейную систему для следующей итерации xk+1:

$$D f(xk) \cdot (xk+1 \quad xk) = f(xk)$$

Это уравнение можно решить, используя псевдообратное, когда m < n; когда m > n, можно попытаться применить метод наименьших квадратов, но существование корня и сходимость этого метода маловероятны. Однако, когда D f является квадратным, что соответствует методу f: Rn Rn: , мы получаем типичную итерацию для Ньютона

$$xk+1 = xk$$
 [D $f(xk)$] 1 $f(xk)$, где, как

всегда, мы не вычисляем явно матрицу [D f(xk)] 1, а используем ее для обозначения решения линейной системы.

Сходимость методов с фиксированной точкой, таких как метод Ньютона, который итерирует xk+1 = g(xk), требует, чтобы собственное значение максимальной величины якобиана Dg было меньше 1. После проверки этого предположения рассуждения, аналогичные одномерному случаю, показывают что метод Ньютона может , для квадратичную сходимость вблизи корней х которого D f(x) невырождена. имеют

7.2.2 Делаем Ньютона быстрее: Квазиньютон и Бройен

По мере увеличения m и n метод Ньютона становится очень дорогим. Для каждой итерации необходимо инвертировать другую матрицу D f(xk); поскольку oн так часто меняется, предварительная факторизация D f(xk) = LkUk не помогает.

Некоторые квазиньютоновские стратегии пытаются применять разные стратегии аппроксимации для упрощения отдельных итераций. Например, один простой подход может повторно использовать D f из предыдущих итераций при повторном вычислении f(xk) в предположении, что производная не меняется очень быстро. Мы вернемся к этим стратегиям, когда будем обсуждать применение метода Ньютона к оптимизации.

Другой вариант — попытаться провести параллель с нашим выводом метода секущих. Так как метод секущих по-прежнему содержит деление, такие приближения не обязательно избавят от необходимости инвертировать матрицу, но они позволяют проводить оптимизацию без явного вычисления якобиана D f . Такие расширения не совсем очевидны, поскольку разделенные разности не дают полной приближенной матрицы Якоби.

Напомним, однако, что производная по направлению f в направлении v задается как Dv f = D $f \cdot v$. Как и в случае с методом секущих, мы можем использовать это наблюдение в своих интересах, попросив, чтобы наша аппроксимация J якобиана удовлетворяла

$$J \cdot (xk \quad xk \quad 1) \quad f(xk) \quad f(xk \quad 1).$$

Метод Бройдена является одним из таких расширений метода секущих, который отслеживает не только оценку xk для x , но и матрицу Jk , оценивающую якобиан; должны быть предоставлены начальные оценки J0 и x0 . Предположим, у нас есть предыдущая оценка Jk 1 якобиана из предыдущей итерации. Теперь у нас есть новая точка данных xk , в которой мы оценили f(xk), поэтому мы хотели бы обновить Jk 1 до нового якобиана Jk с учетом этого нового наблюдения. Одна разумная модель состоит в том, чтобы попросить, чтобы новое приближение было похоже на старое приближение, за исключением направления xk xk 1:

минимизировать Jk Jk
$$\frac{2}{\Phi_{po}}$$
 Jk 1 такое, что Jk · (xk xk 1) = f(xk) f(xk 1)

Чтобы решить эту задачу, определим J Jk Jk 11 μ d f(xk) f(xk 1) Jk 1 · x. x xk Выполнение этих замен дает следующую форму:

минимизировать
$$\frac{2}{\Phi_{po}}$$

 Ј Јтак, что Ј $x = d$

Если мы возьмем λ как множитель Лагранжа, эта минимизация эквивалентна нахождению критических точек лагранжиана Λ :

$$\Lambda = J$$
 $\frac{2}{\Phi_{po}} + \lambda (J \cdot x d)$

Дифференцирование по (Ј)іј показывает:

$$0 = \frac{\Lambda}{1 = 2(J)ij + \lambda i(x)j} = J = \lambda(x) (J)ij 2$$

Подстановка в $J \cdot x = d$ показывает $\lambda(x)(x) = 2d$ или, что то же самое, $\lambda = 2d/x$ 2 . Наконец, мы можем заменить, чтобы найти:

$$J = 2 - \lambda(x) = x_2 \frac{d(x)}{d(x)}$$

Расширение нашей замены показывает:

$$Jk = Jk \quad 1 + J$$

$$= Jk \quad 1 + \frac{A(-x)}{x \cdot 2}$$

$$= Jk \quad 1 + \frac{(f(xk) - f(xk - 1) - Jk - 1 \cdot x)}{xk \cdot xk \cdot 1} (xk - xk - 1)$$

Таким образом, метод Бройдена просто чередует это обновление с соответствующим ньютоновским f(xk). Дополнительную шаг xk+1 = xk \int_{κ}^{1} эффективность в некоторых случаях можно получить, отслеживая матрица \int_{κ}^{1} явно, а не матрица Jk , которую можно обновить по аналогичной формуле.

7.3 Кондиционирование

Мы уже показали в примере 1.7, что число условий нахождения корня в одной переменной:

Как показано на рисунке HOMEP, этот номер условия показывает, что наилучшая возможная ситуация для нахождения корня возникает, когда f быстро меняется вблизи x, поскольку в этом случåе возмущение x заставит f принимать значения, далекие от 0.

Применение идентичного аргумента, когда f многомерно, показывает число обусловленности D f (x $\,$) $\,$ 1 . Обратите внимание, что когда D f необратима, число обусловленности бесконечно. Эта странность происходит из-за того, что возмущение x первого $\,$ * сохраняет f(x) = 0, и такое условие действительно может быть выполнено. порядка создает сложные случаи поиска корня, подобные показанным на рисунке HOMEP.

7.4 Проблемы

Множество возможностей, в том числе:

- Множество возможных схем итерации с фиксированной точкой для данной задачи поиска корня, графическая версия итерации с фиксированной точкой
- Средняя итерация поля в ML
- Метод Мюллера комплексные корни
- Итерационные методы более высокого порядка методы Хаусхолдера
- Интерпретация собственного материала как поиск корня
- Сходимость метода секущих
- Корни многочленов
- Метод Ньютона-Фурье (!)
- «Модифицированный метод Ньютона в случае неквадратичной сходимости»
- Конвергенция ¿ спектральный радиус для многомерного Ньютона; квадратичная сходимость
- Обновление Шермана-Моррисона для Broyden

Machine Translated by Google

Глава 8

Неограниченная оптимизация

В предыдущих главах мы выбрали в значительной степени вариационный подход к построению стандартных алгоритмов вычислительной линейной алгебры. То есть мы определяем целевую функцию, возможно, с ограничениями, и ставим наши алгоритмы как задачу минимизации или максимизации. Пример из нашего предыдущего обсуждения приведен ниже:

Проблема	Цель	Ограничения
Наименьших квадратов	E(x) = Ax -б ²	Никто
Проецироватьb наа	E(c) = ca -b	Никто
Собственные векторы симметричной м	атрицы E(x) = x Ax	x = 1
Псевдоинверсия	E(x) = x	Ах = А б
Анализ основных компонентов	$E(C) = X$ $CCXFro CC = Id \times d$	
шаг Бройдена	$E(Jk) = Jk Jk 1 \frac{2}{\Phi_{po}}$	$\int Jk \cdot (xk + xk + 1) = f(xk) + f(xk + 1)$

Очевидно, что такая постановка задач является мощным и общим подходом. По этой причине важно разработать алгоритмы, которые функционируют в отсутствие специальной формы для энергии Е, точно так же, как мы разработали стратегии для нахождения корней f, не зная формы априори.

8.1 Неограниченная оптимизация: мотивация

В этой главе мы будем рассматривать задачи без ограничений, то есть задачи, которые могут быть поставлены как минимизация или максимизация функции f: Rn R без каких-либо требований к входным данным. На практике с такими проблемами столкнуться нетрудно; мы перечисляем несколько примеров ниже.

Пример 8.1 (Нелинейный метод наименьших квадратов). Предположим, нам дано некоторое количество пар (xi , yi) таких, что f(xi) уi , и мы хотим найти наилучшую аппроксимацию f в определенном классе. Например, мы можем ожидать, что f экспоненциальна, и в этом случае мы должны иметь возможность написать f(x) = сеах для некоторого с и некоторого а; наша задача найти эти параметры. Одной из простых стратегий может быть попытка минимизировать следующую энергию:

$$E(a, c) = (yi ceaxi)^{2}$$

Эта форма для Е не является квадратичной по а, поэтому наши линейные методы наименьших квадратов неприменимы.

Пример 8.2 (оценка максимального правдоподобия). В машинном обучении проблема оценки параметров включает в себя изучение результатов рандомизированного эксперимента и попытку обобщить их с использованием распределения вероятностей определенной формы. Например, мы можем измерить рост каждого ученика в классе, получив список значений роста hi для каждого ученика i. Если у нас много учеников, мы можем смоделировать распределение роста учеников, используя нормальное распределение:

g(h;
$$\mu$$
, σ) = σ $\frac{1}{2\pi} e^{(h \mu) \frac{2}{2\sigma^2}}$,

где μ — среднее значение распределения, а σ — стандартное отклонение.

При таком нормальном распределении вероятность того, что мы наблюдаем рост hi учащегося i, определяется как g(hi; μ, σ), и при (разумном) предположении, что рост учащегося i вероятностно не зависит от роста учащегося j, вероятность наблюдения всего набора наблюдаемых высот определяется произведением

$$P({h1, ..., hn}; \mu, \sigma) = g(hi; \mu, \sigma).$$

Обычный метод оценки параметров μ и σ функции g состоит в том, чтобы максимизировать P как функцию μ и σ при фиксированном {hi}; это называется оценкой максимального правдоподобия μ и σ. На практике мы обычно оптимизируем логарифмическое правдоподобие (μ, σ) log P({h1, ..., hn}; μ, σ); эта функция имеет те же максимумы, но обладает лучшими числовыми и математическими свойствами.

Пример 8.3 (Геометрические задачи). Многие геометрические задачи, встречающиеся в графике и зрении, не сводятся к энергии наименьших квадратов. Например, предположим, что у нас есть несколько точек x1, . . . ,xk R3. Если мы хотим сгруппировать эти точки, мы могли бы суммировать их с помощью минимизации одного x:

Точка х R3 , минимизирующая E, называется геометрической медианой {x1, ..., xк}. Обратите внимание, что норма разности х xi в E не возведена в квадрат, поэтому энергия больше не является квадратичной по компонентам x.

Пример 8.4 (Физические равновесия, адаптировано из СІТЕ). Предположим, мы прикрепляем объект к набору пружин; каждая пружина закреплена в точке хі R3 и имеет натуральную длину Li и постоянную ki . В отсутствие гравитации, если наш объект находится в точке р R3еть источников обладает потенциальной энергией

E (p) = 2
$$\frac{1}{\sum_{\substack{Ku \ (p \ xi2 \ Li)}}}$$

Равновесия этой системы задаются минимумами E и отражают точки p, в которых все силы пружины уравновешены. Такие системы уравнений используются для визуализации графов G = (V, E) путем присоединения вершин в V пружинами для каждой пары в E.

8.2 Оптимальность

Прежде чем обсуждать, как минимизировать или максимизировать функцию, мы должны четко понимать, что мы ищем; обратите внимание, что максимизация f — это то же самое, что минимизация f, поэтому для нашего рассмотрения достаточно проблемы минимизации. Для конкретного f: Rn R и х *Rn нам нужно вывести). условия оптимальности, имеет наименьшее возможное значение f(х подтверждающие, что х Конечно, в идеале мы хотели бы найти глобальные оптимумы f:

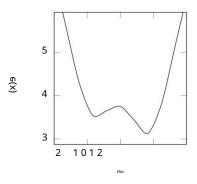


Рисунок 8.1: Функция f(x) с несколькими оптимумами.

Определение 8.1 (Глобальный минимум). Точка x f(x) f(x) * Rn является глобальным минимумом функции f : Rn R, если для всех x Rn

Поиск глобального минимума f без какой-либо информации о структуре f фактически требует поиска в темноте. Например, предположим, что алгоритм оптимизации идентифицирует локальный минимум около x = 1 в функции на рис. 8.1. Почти невозможно понять, что существует второй, более низкий минимум около x = 1, просто угадывая значения x — насколько нам известно, может существовать третий, даже более низкий минимум f при x = 1000!

Таким образом, во многих случаях мы удовлетворяемся тем, что находим локальный минимум:

выражением Это определение требует, * достигает наименьшего значения в некоторой окрестности, определяемой чтобы х был радиусом є. Обратите внимание, что у алгоритмов локальной оптимизации есть серьезное ограничение: они не могут гарантировать, что они дадут наименьшее возможное значение f, как на рис. 8.1, если будет достигнут левый локальный минимум; многие стратегии, эвристические и другие, применяются для изучения ландшафта возможных значений x, чтобы помочь обрести уверенность в том, что локальный минимум имеет наилучшее возможное значение.

8.2.1 Дифференциальная оптимальность

Знакомая история из исчисления с одной и несколькими переменными состоит в том, что поиск потенциальных минимумов и максимумов функции f: Rn R более прост, когда f дифференцируема. Напомним, что вектор градиента f = (f/ x1, ..., f/ xn) указывает направление, в котором f увеличивается больше всего; вектор f указывает в направлении наибольшего убывания. Один из способов убедиться в этом — вспомнить, что вблизи а f выглядит как точка линейной функции x0 Rn

$$f(x)$$
 $f(x0) + f(x0) \cdot (x x0)$

Если мы возьмем х $x0 = \alpha$ f(x0), то найдем:

$$f(x0 + \alpha \quad f(x0)) \quad f(x0) + \alpha \quad f(x0)^{2}$$

Когда f(x0) > 0, знак α определяет, увеличивается или уменьшается f.

Нетрудно формализовать приведенное выше рассуждение, чтобы показать, что если x0 является локальным минимумом, то мы должны иметь f(x0) = 0. Обратите внимание, что это условие необходимо, но недостаточно: максимумы и седло

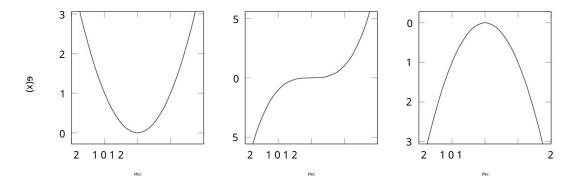


Рисунок 8.2: Критические точки могут принимать разные формы; здесь мы показываем локальный минимум, седловую точку и локальный максимум.

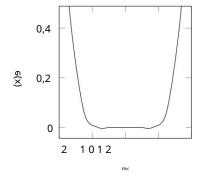


Рисунок 8.3: Функция со многими стационарными точками.

точки также имеют f(x0) = 0, как показано на рис. 8.2. Тем не менее, это наблюдение о минимумах дифференцируемых функций дает общую стратегию поиска корней:

- 1. Найдите точки xi, удовлетворяющие условию f(xi) =0.
- 2. Проверьте, какая из этих точек является локальным минимумом, а не максимумом или седловой точкой.

Учитывая их важную роль в этой стратегии, мы даем точкам, которые мы ищем, особое имя:

Определение 8.3 (Стационарная точка). Стационарной точкой функции f: Rn R называется точка x Rn , удовлетворяющая условию f(x) = 0.

То есть наша стратегия минимизации может состоять в том, чтобы найти стационарные точки f и затем исключить те, которые не являются минимумами.

Важно помнить, когда мы можем ожидать успеха наших стратегий минимизации.

В большинстве случаев, таких как показанные на рис. 8.2, стационарные точки f изолированы, что означает, что мы можем записать их в дискретный список {x0,x1,...}. Однако вырожденный случай показан на рис. 8.3; здесь весь интервал [1/2, 1/2] состоит из стационарных точек, что делает невозможным рассмотрение их по одной. По большей части мы будем игнорировать такие вопросы, как вырожденные случаи, но вернемся к ним, когда будем рассматривать обусловленность задачи минимизации.

Предположим, мы идентифицируем точку х R как стационарную точку f и теперь хотим проверить, является ли она локальным минимумом. Если f дважды дифференцируема, одна из стратегий, которую мы можем использовать, состоит в том, чтобы записать его гессиан

матрица:

Мы можем добавить еще один член в наше разложение Тейлора f, чтобы увидеть роль Hf:

$$f(x) = f(x0) + f(x0) \cdot (x = x0) + 2$$
 $\frac{1}{-\frac{(x-x0) + f(x-x0)}{2}}$

Если мы подставим стационарную точку х *, то по определению мы знаем:

1
$$f(x) f(x) + (x x) + (x x) + (x x) = 1$$

Если Hf положительно определен, то это выражение показывает, что f(x) f(x), и, таким образом, x^* ^{является локальным минимумом.} В более общем случае может возникнуть одна из нескольких ситуаций:

- Если Hf положительно определена, то χ^* является локальным минимумом f .
- Если Hf отрицательно определен, то х * является локальным максимумом f .
- Если Hf неопределенно, то х * является седловой точкой f .
- Если Нf необратима, то могут возникнуть странности, такие как функция на рис. 8.3.

Проверить, является ли матрица положительно определенной, можно, проверив, существует ли ее факторизация Холецкого, или, что медленнее, проверив, что все ее собственные значения положительны. Таким образом, когда гессиан функции f известен, мы можем проверить стационарные точки на оптимальность, используя приведенный выше список; многие алгоритмы оптимизации, включая те, которые мы обсудим, просто игнорируют последний случай и уведомляют пользователя, поскольку это относительно маловероятно.

8.2.2 Оптимальность через функциональные свойства

Иногда, если мы знаем больше информации о f : Rn R , мы можем предоставить условия оптимальности, которые сильнее или легче проверяются, чем приведенные выше.

Одним из свойств f, которое имеет большое значение для оптимизации, является выпуклость, показанная на рис. НОМЕР:

Определение 8.4 (выпуклое). Функция f: Rn R называется выпуклой, если для всех x,y Rn и α (0, 1) выполняется следующее соотношение:

$$f((1 \quad \alpha)x + \alpha y) \quad (1 \quad \alpha)f(x) + \alpha f(y).$$

Когда неравенство строгое, функция строго выпуклая.

Выпуклость означает, что если вы соедините в Rn две точки линией, значения f вдоль линии будут меньше или равны тем, которые вы получили бы с помощью линейной интерполяции.

Выпуклые функции обладают многими сильными свойствами, самые основные из которых следующие:

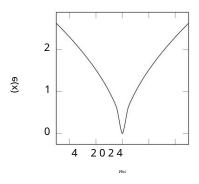


Рис. 8.4: Квазивыпуклая функция.

Предложение 8.1. Локальный минимум выпуклой функции f: Rn R обязательно является глобальным минимумом.

Доказательство. Возьмем x за такой локальный минимум и предположим, что x * = x c f(x) < f(x). существует. Тогда для α (0, 1)

Но переход α 0 показывает, что x не может быть локальным минимумом.

Это предложение и связанные с ним наблюдения показывают, что можно проверить, достигли ли вы глобального минимума выпуклой функции, просто применяя оптимальность первого порядка. Таким образом, полезно проверить вручную, не окажется ли оптимизируемая функция выпуклой, что на удивление часто происходит в научных вычислениях; одно достаточное условие, которое легче проверить, когда f дважды дифференцируема, состоит в том, что Hf везде положительно определена.

Другие методы оптимизации имеют гарантии при других предположениях о f. Для экзамена ple, одной из более слабых версий выпуклости является квазивыпуклость, достигаемая, когда

$$f((1 \quad \alpha)x + \alpha y) \quad max(f(x), f(y)).$$

Пример квазивыпуклой функции показан на рис. 8.4; хотя она не имеет характерной формы «чаши» выпуклой функции, она имеет уникальный оптимум.

8.3 Одномерные стратегии

Как и в предыдущей главе, мы начнем с одномерной оптимизации f : R R , а затем расширим наши стратегии до более общих функций f : Rn R.

8.3.1 Метод Ньютона

Наша основная стратегия минимизации дифференцируемых функций f: Rn R будет состоять в том, чтобы найти sta стационарные , удовлетворяющую f(x) = 0. Предположим, что мы можем проверить, являются ли точки условными точками x максимумов, минимумов или седловых точек в качестве пост- шаге обработки, сосредоточимся на задаче нахождени $\dot{\pi}$ стационарных точек x

1
$$f(x)$$
 $f(xk) + f(xk)(x xk) + f(xk)(x xk) 2^{2}$

Аппроксимация в правой части представляет собой параболу, вершина которой расположена в точке xk — f (xk)/f (xk). Конечно, в действительности f не обязательно является параболой, поэтому метод Ньютона просто итерирует формулу

$$xk+1 = xk$$
 $f \frac{f(xk)}{(xk)}$.

Эту технику легко проанализировать, учитывая ту работу, которую мы уже проделали для понимания метода Ньютона для поиска корней в предыдущей главе. В частности, альтернативный способ вывести приведенную выше формулу исходит из нахождения корня f(x), поскольку стационарные точки удовлетворяют f(x) = 0.

Таким образом, в большинстве случаев метод оптимизации Ньютона демонстрирует квадратичную сходимость при условии, что начальное приближение x0 достаточно близко к x^* .

Возникает естественный вопрос: можно ли аналогичным образом применить метод секущих?

Наш вывод метода Ньютона, приведенный выше, находит корни f , поэтому метод секущих можно использовать для устранения вычисления f, но не f ; ситуации, в которых мы знаем f, но не f, относительно редки. Более подходящей параллелью является замена отрезков, используемых для аппроксимации f в методе секущих, параболами. Эта стратегия, известная как последовательная параболическая интерполяция, также минимизирует квадратичное приближение f на каждой итерации, но вместо использования f(xk), f (xk) и f (xk) для построения приближения, которое использует f(xk), f(xk 1) и f(xk 2). Вывод этого метода относительно прост, и он сходится сверхлинейно.

8.3.2 Поиск золотого сечения

Мы пропустили деление пополам в нашем параллельном методе поиска корней с одной переменной. Причин этого упущения много. Наша мотивация для деления пополам заключалась в том, что оно использовало только самое слабое предположение относительно f, необходимое для нахождения корней: непрерывность. Однако теорема о промежуточном значении не применима к минимумам каким-либо интуитивным образом, поэтому кажется, что такого прямого подхода не существует.

Однако полезно иметь по крайней мере одну доступную стратегию минимизации, которая не требует дифференцируемости f в качестве базового предположения; ведь существуют недифференцируемые функции, имеющие четкие минимумы, например f(x) | | | при x = 0. С этой целью одним из альтернативных предположений может быть то, что f унимодулярна:

Другими словами, унимодулярная функция некоторое время убывает, а затем начинает возрастать; локализованные минимумы не допускаются. Обратите внимание, что такие функции, как |x| не дифференцируемы, но все же унимодулярны.

Предположим, у нас есть два значения x0 и x1, такие что a < x0 < x1 < b. Мы можем сделать два обсерва ции, которые помогут нам сформулировать методику оптимизации:

• Если f(x0) f(x1), то мы знаем, что f(x) f(x1) для всех x [a, x0]. Таким образом, интервал [a, x0] можно отбросить при поиске минимума f.

• Если f(x1) f(x0), то мы знаем, что f(x) f(x0) для всех x [x1, b], и поэтому можем отбросить [x1, 6].

Эта структура предлагает потенциальную стратегию минимизации, начиная с интервала [a, b] и итеративно удаляя части в соответствии с приведенными выше правилами.

Однако остается одна важная деталь. Наша гарантия сходимости для алгоритма деления пополам основывалась на том факте, что мы могли удалить половину рассматриваемого интервала на каждой итерации.

Мы могли бы действовать аналогичным образом, каждый раз удаляя треть интервала; это требует двух оценок f во время каждой итерации в новых местах x0 и x1. Однако, если оценка f требует больших затрат, мы можем захотеть повторно использовать информацию из предыдущих итераций, чтобы избежать по крайней мере одной из этих двух оценок.

На данный момент a=0 и b=1; стратегии, которые мы выведем ниже, будут работать в более общем плане за счет сдвига и масштабирования. В отсутствие дополнительной информации о f мы могли бы также сделать симметричный выбор $x0=\alpha$ и x1=1 α для некоторого α (0, 1/2). Предположим, наша итерация удаляет крайний правый интервал [x1, b]. Тогда интервал поиска становится равным [0, 1 α], и мы знаем $f(\alpha)$ из предыдущей итерации. Следующая итерация разделит [0, 1 α] так, что $x0=\alpha(1 \alpha)$ и $x1=(1\alpha)$ хотим повторно использовать $f(\alpha)$ из предыдущей итерации, мы могли бы установить $f(\alpha)$ $f(\alpha)$ 0 что дает:

$$\alpha = \frac{1}{2}(3 \quad 5)$$

$$1 - a = 2 \quad -(5 - 1)$$

Значение 1 α т, указанное выше, является золотым сечением! Это позволяет повторно использовать одну из оценок функции из предыдущих итераций; симметричный аргумент показывает, что тот же выбор α работает, если мы удалили левый интервал вместо правого.

Алгоритм поиска золотого сечения использует эту конструкцию (CITE):

- 1. Возьмем $\tau = \frac{1}{12}(5 1)$., и инициализировать а и b так, чтобы f была унимодулярной на [a, b].
- 2. Сделать начальное подразделение x0 = a + (1 t)(b a) и x1 = a + t(b a).
- 3. Инициализировать f0 = f(x0) и f1 = f(x1).
- 4. Повторяйте до тех пор, пока b а не станет достаточно малым:
 - (a) Если f0 f1, то удалить интервал [a, x0] следующим образом:
 - Переместить влево: а х0
 - Повторно использовать предыдущую итерацию: х0

x1, f0 f1 • Создать новую выборку: x1
$$a + \tau(b - a)$$
, f1 $f(x1)$

- (b) Если f1 > f0, то удалить интервал [x1, b] следующим образом:
 - Переместить вправо: b x1
 - Повторно использовать предыдущую итерацию: x1

$$x0$$
, f1 f0 • Создать новую выборку: $x0$ a + (1 τ)(b a), f0 f($x0$)

Очевидно, этот алгоритм сходится безусловно и линейно. Когда f не является глобально унимодальным, может быть трудно найти [a, b] такое, что f является унимодальным на этом интервале, что несколько ограничивает применение этого метода; обычно [a, b] угадывается при попытке заключить в скобки локальный минимум f.

8.4 Многовариантные стратегии

Мы продолжаем параллель с нашим обсуждением поиска корней, расширяя наше обсуждение до многовариантных проблем. Как и в случае поиска корня, задачи с несколькими переменными значительно сложнее, чем задачи с одной переменной, но на практике они возникают так часто, что заслуживают внимательного рассмотрения.

Здесь мы будем рассматривать только случай, когда f : Rn R дифференцируема. Методы оптимизации, более похожие на поиск золотого сечения для недифференцируемых функций, имеют ограниченное применение и их трудно сформулировать.

8.4.1 Градиентный спуск

Напомним из нашего предыдущего обсуждения, что f(x) указывает в направлении «самого крутого подъема» f в точке x; аналогично вектор f(x) является направлением «самого крутого спуска». По крайней мере, это определение гарантирует, что при f(x) = 0 для малых $\alpha > 0$ мы должны иметь

$$f(x \quad \alpha \quad f(x)) \quad f(x)$$
.

Предположим, что наша текущая оценка местоположения минимума f равна xk. Затем мы можем выбрать xk+1 так, чтобы f(xk+1) < f(xk) для итеративной стратегии минимизации. Одним из способов упростить поиск xk+1 было бы использование одного из наших одномерных алгоритмов из \$8.3 для решения более простой задачи. В частности, рассмотрим функцию gk(t) f(xk) f(xk), ограничивающую f на прямую, проходящую через xk, параллельную f(xk). Благодаря нашему обсуждению градиента мы знаем, что малое f приведет f(xk) уменьшению f(xk).

Алгоритм градиентного спуска итеративно решает эти одномерные задачи, чтобы улучшить нашу оценку хк:

- 1. Выбрать начальную оценку х0
- 2. Итерировать до сходимости xk :

```
(a) Возьмем gk(t) f(xk 	 t 	 f(xk)) (b) 
Используем одномерный алгоритм для нахождения ^* минимизация gk по всем t 	 0 («линейный поиск») t (c) Возьмем xk+1 	 xk 	 t 	 f(xk)
```

Каждая итерация градиентного спуска уменьшает f(xk), поэтому целевые значения сходятся. Алгоритм завершается только тогда, когда f(xk) 0, показывая, что градиентный спуск должен как минимум достигать локального минимума; Однако сходимость для большинства функций f медленная. Процесс линейного поиска можно заменить методом, который просто уменьшает цель на незначительную, если не оптимальную величину, хотя в этом случае труднее гарантировать сходимость.

8.4.2 Метод Ньютона

Параллельно с нашим выводом случая с одной переменной, мы можем написать аппроксимацию ряда Тейлора для f: Rn R, используя его гессиан Hf:

$$f(x) \quad f(xk) + \quad f(xk) \cdot (x \quad xk) + 2 \qquad \qquad 1 \\ - \quad (x \quad xk) \cdot Hf(xk) \cdot (x \quad xk)$$

Дифференцирование по х и установка результата равным нулю дает следующую итеративную схему:

$$xk+1 = xk$$
 [Hf(xk)] 1 f(xk)

Легко проверить еще раз, что это выражение является обобщением выражения из § 8.3.1, и еще раз оно сходится квадратично, когда х0 близко к минимуму.

Метод Ньютона может быть более эффективным, чем градиентный спуск, в зависимости от цели оптимизации f. Напомним, что каждая итерация градиентного спуска потенциально требует множества вычислений f во время процедуры линейного поиска. С другой стороны, мы должны оценивать и инвертировать гессиан Hf во время каждой итерации метода Ньютона. Обратите внимание, что эти факторы не влияют на количество итераций, но влияют на время выполнения: это компромисс, который может быть неочевидным при традиционном анализе.

Интуитивно понятно, почему метод Ньютона быстро сходится, когда он близок к оптимуму. В частности, градиентный спуск не знает Hf; это происходит аналогично спуску с горы, когда вы смотрите только на свои ноги. Используя Hf, метод Ньютона имеет большую картину формы f поблизости.

Однако, когда Hf не является положительно определенным, объектив локально может выглядеть как седло или пик, а не как чаша. В этом случае переход к приблизительно стационарной точке может не иметь смысла.

Таким образом, адаптивные методы могут проверять, является ли Hf положительно определенным, прежде чем применять шаг Ньютона; если он не является положительно определенным, методы могут вернуться к градиентному спуску, чтобы найти лучшее приближение к минимуму. В качестве альтернативы они могут модифицировать Hf, например, путем проецирования на ближайшую положительно определенную матрицу.

8.4.3 Оптимизация без производных: BFGS

Метод Ньютона может быть сложно применить к сложным функциям f:Rn R. Вторая производная от f может быть значительно более сложной, чем форма f, и Hf меняется c каждой итерацией, что затрудняет повторное использование результатов предыдущих итераций. Кроме того, Hf имеет размер $n \times n$, поэтому для хранения Hf требуется O(n-2) пространство, что может быть неприемлемо.

Как и в нашем обсуждении поиска корней, методы минимизации, имитирующие метод Ньютона, но использующие приближенные производные, называются квазиньютоновскими методами. Часто они могут иметь столь же сильные свойства сходимости без необходимости явной переоценки и даже факторизации гессиана на каждой итерации. В нашем обсуждении мы будем следить за развитием (CITE NO CEDAL AND WRIGHT).

Предположим, мы хотим минимизировать f: Rn R, используя итеративную схему. Рядом с текущей оценкой xk корня, мы могли бы оценить f с помощью квадратичной модели:

1
$$f(xk + \delta x)$$
 $f(xk) + f(xk) \cdot \delta x + \frac{1}{2} (\delta x) Bk(\delta x)$.

Обратите внимание, что мы попросили, чтобы наше приближение согласовывалось с f до первого порядка в точке xk; однако, как и в методе Бройдена для нахождения корня, мы позволим нашей оценке гессиана Вk варьироваться.

Эта квадратичная модель минимизируется, если взять $\delta x = B$ $\frac{1}{K}$ f(xk). В случае, когда $\delta x 2$ велико и мы не хотим делать такой значительный шаг, мы позволим себе масштабировать эту разницу на величину шага αk , что даст

$$x\kappa+1 = x\kappa - a\kappa B$$
 κ $f(xk)$.

Наша цель — найти разумную оценку Bk+1, обновив Bk, чтобы мы могли повторить этот процесс.

Гессе функции f есть не что иное, как производная от f, поэтому мы можем записать условие в стиле секущей на Bk+1:

$$Bk+1(xk+1 xk) = f(xk+1) f(xk).$$

Подставим sk xk+1 xk и yk f(xk+1) f(xk), что приведет к эквивалентному условию Bk+1sk = yk.

Учитывая проводимую оптимизацию, мы хотим, чтобы Вk обладал двумя свойствами:

- Вк должна быть симметричной матрицей, такой как гессиан Hf.
- Вк должно быть положительно (полу)определенным, так что мы ищем минимумы, а не максимумы или седловые точки.

Условия симметрии достаточно, чтобы исключить возможность использования оценки Бройдена, развитой нами в предыдущей главе.

Ограничение положительной определенности неявно накладывает условие на связь между sk и B частности, предварительно умножая связь Bk+1sk = yk на s Bk+1sk = s k yk . Для юк . $_{\rm K}$ $^{\rm показывает \, c}_{\rm K}$ Чтобы Bk+1 было положительно определенным, мы должны иметь sk · yk > 0. Это наблюдение может направить нас при выборе α k ; легко видеть, что оно выполняется при достаточно малых α k > 0.

Предположим, что sk и yk удовлетворяют нашему условию совместимости. Имея это на месте, мы можем написать оптимизация в стиле Бройдена, ведущая к возможному приближению Bk+1 :

минимизировать Bk+1 Bk+1 Bk

такое, что B k+1 =
$$E_{K+1}$$

Вместо того, чтобы вдаваться в детали схемы DFP, мы перейдем к более популярному методу, известному как формула BFGS (Бройден-Флетчер-Гольдфарб-Шанно), которая используется во многих современных системах. Заметьте, что — игнорируя наш выбор αk на данный момент — наше приближение второго порядка было минимизировано, приняв δх = В

1 f(xk). Таким образом, в итоге поведение нашей итерационной схемы диктуется обратной матрицей В k . Запрос о том, что Вk+1 Вk мал, все же может подразумевать относительно плохие различия между действием Вк

1 и что из Б! к+1

Учитывая это наблюдение, схема BFGS вносит небольшое изменение в приведенный выше вывод. Вместо того, чтобы вычислять Вk на каждой итерации, мы можем напрямую вычислить обратное значение Hk В. 1

Теперь наше условие Bk+1sk = yk меняется на sk = Hk+1yk; условие симметричности Bk равносильно требованию симметричности Hk. Мы решаем оптимизацию

```
минимизировать Hk+1 Hk+1 Hk
    _{\text{такое, что H}} = Hk+1
    _{\text{k+1 sk}} = Hk+1yk
```

У этой конструкции есть приятное дополнительное преимущество, заключающееся в том, что для вычисления δx = Hk f(xk) не требуется обращения матриц.

Чтобы вывести формулу для Hk+1, мы должны определить матричную норму . . . Как и в нашем предыдущем обсуждении, норма Фробениуса выглядит наиболее близкой к оптимизации методом наименьших квадратов, что делает более вероятным, что мы можем сгенерировать выражение в закрытой форме для Hk+1, а не решать приведенную выше минимизацию как подпрограмму оптимизации BFGS.

Однако у нормы Фробениуса есть один серьезный недостаток для матриц Гессе. Напомним, что матрица Гессе имеет элементы (Hf)ij = fi/ xj . Часто величины xi для разных i могут иметь разные единицы измерения; например, рассмотрите максимизацию прибыли (в долларах), полученной от продажи чизбургера радиуса r (в дюймах) и цены р (в долларах), что приводит к f : (дюймы, доллары) доллары. Возводить в квадрат эти различные количества и складывать их не имеет смысла.

Предположим, мы нашли симметричную положительно определенную матрицу W такую, что Wsk = yk; в упражнениях проверим, что такая матрица существует. Такая матрица переводит единицы sk = xk+1 xk в единицы yk = f(xk+1) f(xk). Вдохновившись нашим выражением A = $Tr(A \cap A)$, мы можем определить взвешенную норму Фробениуса матрицы A как

A
$$^2_{BT}$$
 Tp(A BAB)

Несложно проверить, что это выражение имеет согласованные единицы применительно к нашей оптимизации для Hk+1 . Когда и W, и A симметричны со столбцами w i andai соответственно, расширение приведенного выше выражения показывает:

$$\begin{array}{ccc} A & \begin{array}{ccc} 2 \\ B\tau & = & (w \ i \cdot aj)(w \ j \cdot ai). \end{array} \\ & & ij \end{array}$$

Этот выбор нормы в сочетании с выбором W дает особенно чистую формулу для Hk+1 и ук: при заданных Hk, sk,

$$Hk+1 = (In \times n \quad \rho ksky k \quad)Hk(In \times n \quad \rho kyks k) + \rho ksks k,$$

где pk 1/y·s. В приложении к этой главе мы покажем, как вывести эту формулу.

Алгоритм BFGS избавляет от необходимости вычислять и инвертировать матрицу Гессе для f, но по-прежнему требует O(n). Полерживарицы длязыестный как L-BFGS («BFGS с ограниченной памятью»), позволяет избежать этой проблемы, сохраняя ограниченную историю векторов ук andsk и применение Нk путем рекурсивного расширения его формулы.Этот подход на самом деле может иметь лучшие числовые свойства, несмотря на его компактное использование пространства; в частности, старые векторы ук иsk могут больше не иметь значения и должны быть проигнорированы.

8.5 Проблемы

Список идей:

- Вычислить Гаусс-Ньютон
- Стохастические методы, АдаГрад
- Алгоритм VSCG
- условия Вульфа для градиентного спуска; подключиться к БФГС
- Формула Шермана-Моррисона-Вудбери для Bk для BFGS
- Докажите, что BFGS сходится; показать существование матрицы W
- (Обобщенный) алгоритм уменьшенного градиента
- Номер условия для оптимизации

Приложение: Получение BFGS Update1

$$\Lambda$$
 (w i · (hj h * _))(w j · (hi h * _)) α ij(Hij Hji) λ (Hyk sk) $_{\text{s} < j}$ = (w i · (hj h * _))(w j · (hi h * _)) α ijHij λ (Hyk sk) , если предположить, что α ij = α ji ij

Взятие производных для нахождения критических точек показывает (для у уk, s sk):

$$0 = \frac{\Lambda}{M} = 2$$
wi(w j · (h h)) аіј λ iyj

Hij = 2 wi(W(H H))j аіј λ iyj

= 2 (W(H H))jwi аіј λ iyj в силу симметрии W

= 2(W(H H)W)ji аіј λ iyj

= 2(W(H H)W)ij аіј λ iyj в силу симметрии W и H

Итак, в матричной форме мы имеем следующий список фактов:

$$0 = 2W(H \quad H \quad)W \quad A \quad \lambda y$$
 , где $Aij = \alpha ij$ $A = \quad A$, $W = W$, $H = H$, H ,

Мы можем получить пару отношений, используя транспозицию в сочетании с симметрией Н и W и асимметрией А:

$$0 = 2W(H \quad H \quad)W \quad A \quad \lambda y$$

$$0 = 2W(H \quad H \quad)W + A \quad y\lambda = \quad 0 =$$

$$4W(H \quad H \quad)W \quad \lambda y \quad y\lambda$$

Пост-умножение этого отношения на показывает:

$$0 = 4(y \quad WH \quad y) \quad \lambda(y \cdot s) \quad y(\lambda \cdot s)$$

Теперь возьмем скалярный продукт с:

$$0 = 4(y \cdot s)$$
 $4(y \cdot H \cdot y)$ $2(y \cdot s)(\lambda \cdot s)$

Это показывает:

$$\lambda \cdot s = 2\rho y$$
 (s H y), при ρ 1/y·s

¹Особая благодарность Тао Ду за отладку нескольких частей этого вывода.

Теперь подставим это в наше векторное равенство:

$$0 = 4$$
(у WH y) λ (у·s) у(λ ·s) из предыдущего
$$= 4$$
(у WH y) λ (у·s) у[2ру (s H y)] из нашего упрощения
$$= \lambda = 4$$
р(у WH y) $2\rho^{-2}$ у (s H y)у

Пост-умножение на у показывает:

$$\lambda y = 4\rho(y \text{ WH } y)y \text{ } 2\rho \text{ }^2y \text{ } (s \text{ H } y)yy$$

Выполняя транспонирование,

$$y\lambda = 4\rho y(y \quad yH \quad W) \quad 2\rho \quad ^2 y (s \quad H \quad y)yy$$

Объединение этих результатов и деление на четыре показывает:

$$\frac{1}{4}(\lambda y + y\lambda) = \rho(2yy \quad WH \quad yy \quad yy H \quad W) \quad \rho \qquad \qquad ^2y (s \quad H \quad y)yy$$

Теперь мы предварительно и после умножим на W 1. Поскольку Ws = y, мы можем эквивалентно писать = W 1y; кроме того, в силу симметрии W мы знаем, что y W 1 = s. Применение этих тождеств к приведенному выше выражению показывает:

$$\frac{1}{4}$$
W 1 ($\lambda y + y\lambda$)W 1 = 2 ρ ss ρ H *ys - ρ sy H * ρ^2 (ys)ss + ρ 2(y H y)ss = 2 ρ ss ρ H *ys - ρ sy H * ρ ss + s ρ 2(y H y)s по определению ρ = ρ ss ρ H *ys - ρ sy H * +s ρ 2(y H y)s

Наконец, мы можем завершить вывод шага BFGS следующим образом:

Это окончательное выражение и есть шаг BFGS, представленный в этой главе.

Глава 9

Ограниченная оптимизация

Мы продолжаем рассмотрение задач оптимизации, изучая случай с ограничениями. Эти проблемы принимают следующий общий вид:

минимизировать

f(x) так, чтобы g(x) = 0

h(x) 0

Здесь f: Rn R, g: Rn Rm u h: Rn Rp. Очевидно, что эта форма является чрезвычайно общей, поэтому нетрудно предсказать, что алгоритмы решения таких задач при отсутствии дополнительных предположений относительно f, g или h могут быть трудно сформулировать и будут подвержены вырождениям, таким как локальные минимумы и отсутствие конвергенция. На самом деле эта оптимизация кодирует другие проблемы, которые мы уже рассмотрели; если мы возьмем f(x) = 0, то эта оптимизация с ограничениями станет нахождением корней на g, a если мы возьмем g(x) = h(x) = 0, то она сведется g(x) = h(x) = h(x)0, то она сведется g(x) = h(x)1.

Несмотря на эту несколько мрачную перспективу, оптимизация для общего случая с ограничениями может оказаться полезной, когда f, g и h не имеют полезной структуры или слишком специализированы, чтобы заслуживать специального обращения. Кроме того, когда f в любом случае является эвристическим, простое нахождение возможного x, для которого f (x) < f (x0) для начального предположения x0 , ценно. Одним из простых приложений в этой области может быть экономическая система, в которой f измеряет затраты; очевидно, мы хотим минимизировать затраты, но если x0 представляет текущую конфигурацию, любое уменьшение f на x является ценным выходом.

9.1 Мотивация

На практике нетрудно столкнуться с задачами условной оптимизации. На самом деле, мы уже перечислили множество приложений этих задач, когда обсуждали собственные векторы и собственные значения, поскольку эту задачу можно поставить как нахождение критических точек х Ах при х2 = 1; конечно, частный случай вычисления собственных значений допускает специальные алгоритмы, упрощающие эту задачу.

Здесь мы перечисляем другие оптимизации, которые не пользуются структурой задач на собственные значения:

Пример 9.1 (геометрическая проекция). Многие поверхности S в R3 можно неявно записать в виде g(x) = 0 для некоторого g. Например, единичная сфера получается, если взять g(x) = x 1, тогда как куб может $\frac{2}{2}$

можно построить, взяв g(x) = x1 1. Фактически, некоторые среды 3D-моделирования позволяют пользователям задавать «кляксы»-объекты, как на рисунке HOMEP, как суммы

$$r(x)$$
 c + aie bix xi $\frac{2}{2}$.

Предположим, что нам дана точка у R3 , и мы хотим найти точку на S, ближайшую к у. Эта проблема решается с помощью следующей условной минимизации:

минимизировать
$$x x$$

y2 такое, что $g(x) = 0$

Пример 9.2 (Производство). Предположим, у вас есть m различных материалов; у вас есть si единиц каждого материала i на складе. Вы можете производить k различных продуктов; продукт j дает вам прибыль pj и использует сij материала i для его производства. Чтобы максимизировать прибыль, вы можете решить следующую оптимизацию для общего количества xj , которое вы должны произвести каждого изделия j:

Первое ограничение гарантирует, что вы не сделаете отрицательные числа для любого продукта, а второе гарантирует, что вы не используете больше, чем ваш запас каждого материала.

Пример 9.3 (неотрицательный метод наименьших квадратов). Мы уже видели множество примеров задач наименьших квадратов, но иногда отрицательные значения в векторе решения могут не иметь смысла. Например, в компьютерной графике анимированная модель может быть представлена как деформирующаяся костная структура плюс сетчатая «кожа»; для каждой точки на коже можно рассчитать список весов, чтобы аппроксимировать влияние положения суставов костей на положение вершин кожи (СІТЕ). Такие веса должны быть неотрицательными, чтобы избежать вырожденного поведения при деформации поверхности. В таком случае мы можем решить задачу «неотрицательных наименьших квадратов»:

Недавние исследования включают характеристику разреженности неотрицательных решений методом наименьших квадратов, которые часто имеют несколько значений xi, точно удовлетворяющих xi = 0 (CITE).

Пример 9.4 (Корректировка связки). Предположим, что в компьютерном зрении мы делаем снимок объекта с нескольких ракурсов. Естественной задачей является воссоздание трехмерной формы объекта. Для этого мы могли бы отметить соответствующий набор точек на каждом изображении; в частности, мы можем взять хіј R2 в качестве положения характерной точки ј на изображении і. На самом деле у каждой характерной точки есть положение уј R3 в пространстве, которое мы хотели бы вычислить. Кроме того, мы должны найти положения самих камер, которые мы можем представить как

неизвестные проекционные матрицы Pi . К этой проблеме, известной как корректировка пакетов, можно подойти с помощью стратегии оптимизации:

Ограничение ортогональности гарантирует разумность преобразований камеры.

9.2 Теория оптимизации с ограничениями

В нашем обсуждении мы будем предполагать, что f, g и h дифференцируемы. Существуют некоторые методы, которые делают только предположения о слабой непрерывности или липшицевости, но эти методы весьма специализированы и требуют расширенного аналитического рассмотрения.

Хотя мы еще не разработали алгоритмы общей оптимизации с ограничениями, мы неявно использовали теорию таких задач при рассмотрении методов собственных значений. В частности, вспомним метод множителей Лагранжа, введенный в теореме 0.1. В этом методе критические точки f(x), зависящие от g(x), характеризуются как критические точки неограниченной функции множителей Лагранжа Λ(x,λ) f(x) λ · g(x) относительно обоих λ и x одновременно.

Эта теорема позволила нам дать вариационную интерпретацию задач на собственные значения; в более общем смысле, он дает альтернативный (необходимый, но недостаточный) критерий того, что х является критической точкой оптимизации, ограниченной равенством.

Однако простое нахождение x, удовлетворяющего ограничениям, может стать серьезной проблемой. Мы можем разделить эти проблемы, сделав несколько определений:

Определение 9.1 (допустимая точка и допустимое множество). Допустимой точкой задачи оптимизации с ограничениями является любая точка x, удовлетворяющая условиям g(x) =0 и h(x) 0. Допустимое множество — это множество всех точек x, удовлетворяющих этим ограничениям.

Определение 9.2 (Критическая точка оптимизации с ограничениями). Критическая точка оптимизации с ограничениями — это точка, удовлетворяющая ограничениям, которая также является локальным максимумом, минимумом или седловой точкой f в допустимом множестве.

Ограниченная оптимизация сложна, потому что она одновременно решает проблемы нахождения корня (ограничение g(x) =0), проблемы выполнимости (ограничение h(x) 0) и минимизацию (функция f). Помимо этого, чтобы довести наши дифференциальные методы до полной общности, мы должны найти способ добавить ограничения неравенства в систему множителей Лагранжа. Предположим, мы нашли минимум оптимизации, обозначенный x. Для каждого ограничения неравенства hi(x) 0 у нас есть два варианта:

- hi(x) = 0: такое ограничение активно, что, вероятно, указывает на то, что если ограничение будет снято, оптимум может измениться.
- hi(x) > 0: такое ограничение неактивно, то есть в окрестности x , если бы мы сняли это ограничение, мы все равно достигли бы того же минимума.

Конечно, мы не знаем, какие ограничения будут активны или неактивны в точке х , пока она не будет вычислена.

Если бы все наши ограничения были активны, то мы могли бы изменить наше ограничение h(x) 0 на равенство, не влияя на минимум. Это может мотивировать изучение следующей системы множителей Лагранжа:

$$\Lambda(x,\lambda,\mu)$$
 $f(x)$ $\lambda \cdot g(x)$ $\mu \cdot h(x)$

Однако мы больше не можем сказать, ^{*}что х является критической точкой Л, потому что неактивные ограничения удалили бы приведенные выше члены. Игнорируя пока этот (важный!) вопрос, мы могли бы действовать вслепую и запросить критические точки этого нового Л по х, которые удовлетворяют следующему:

$$0 = f(x) \qquad \underset{s}{\lambda i} gi(x) \qquad \mu j \quad hj(x)$$

Здесь мы выделили отдельные компоненты д и h и обработали их как скалярные функции, чтобы избежать сложных обозначений.

Хитрый прием может распространить это условие оптимальности на системы, ограниченные неравенством. Обратите внимание, что если бы мы взяли µj = 0 всякий раз, когда hj неактивен, то это удаляет нерелевантные члены из условий оптимальности. Другими словами, мы можем добавить ограничение на множители Лагранжа:

$$\mu$$
jhj(x) = 0.

При наличии этого ограничения мы знаем, что по крайней мере одно из значений µj и hj(x) должно быть равно нулю, и, следовательно, наше условие оптимальности первого порядка по-прежнему выполняется!

До сих пор наша конструкция не различала ограничение hj(x) 0 и ограничение hj(x) 0. Если ограничение неактивно, его можно было бы снять, не влияя локально на результат оптимизации, поэтому мы рассмотрим случай, когда ограничение активно. Интуитивно 1 в этом случае мы ожидаем, что будет способ уменьшить f, нарушив ограничение. Локализованно, направление убывания f равно

f(x), а направление убывания hj есть hj(x). Таким образом, начиная с x, мы можем еще больше уменьшить f, нарушив dграничение hj(x) 0, когда f(x) 0.

Конечно, с произведениями градиентов f и hj работать сложно. Однако напомним, что при нашем условии оптимальности * x первого порядка мы получаем:

$$f(x \quad) = \quad \begin{array}{ccc} \lambda & \star & & \star \\ & i & gi(x \quad) + & & \star \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & \\ & & & \\ &$$

Неактивные значения μ равны нулю и могут быть удалены. На самом деле, мы можем снять ограничения g(x) = 0, добавив к h ограничения-неравенства g(x) = 0 и g(x) = 0; это математическое удобство для написания доказательства, а не численный маневр. Затем, скалярное произведение с hk для любого фиксированного k показывает:

$$\int_{M \times A \times A \times A \times M \times M \times M}^{*} hj(x) \cdot hk(x) = f(x) \cdot hk(x)$$

Векторизация этого выражения показывает, что Dh(x) Dh(x) μ * 0. Так как Dh(x) Dh(x) положительно полуопределенно инициальна, из этого * 0. Так им образом, наблюдение f(x) • hj(x) 0 проявляется просто как следует μ тот факт, что μ 0.

Наши наблюдения можно формализовать, чтобы доказать условие оптимальности первого порядка для оптимизаций с ограничениями неравенства:

¹ Не следует считать наше обсуждение формальным доказательством, так как мы не рассматриваем многих граничных случаев.

Теорема 9.1 (условия Каруша-Куна-Таккера (ККТ)). Вектор х Rn является критической точ \dot{k} ой для минимизации f при g(x) = 0 и h(x) 0, когда существуют λ Rm и μ Rp такие, что:

• 0 = f(x) i
$$\lambda$$
i gi(x) j μ j hj(x) («стационарность»)
• g(x) = 0 и h(x) 0 («первичная допустимость») • μ jhj(x) = 0 для всех ј («дополнительная нежесткость») • μ j

0 для всех j («двойственная допустимость »)

Обратите внимание, что при удалении h эта теорема сводится к критерию множителя Лагранжа.

Пример 9.5 (Простая оптимизация2). Предположим, мы хотим решить

максимизировать

В этом случае у нас не будет λ и будет три μ . Возьмем f(x, y) = xy, h1(x, y) = x y и h3(x, y) = y. Условия $\frac{2}{x}$, h2(x, y) = x, kkt

Стационарность:
$$0 = y + \mu 1 \quad \mu 2 \ 0$$

$$= x + 2\mu 1y \quad \mu 3$$
Первичная выполнимость: $x^2 = 2$

$$+ yx, y = 0$$
Дополнительная нежесткость: $\mu 1(2 - x - y \mu 2x^{-2}) = 0$

$$= 0 \ \mu 3y = 0$$

Двойная выполнимость: µ1, µ2, µ3 0

Пример 9.6 (Линейное программирование). Рассмотрим оптимизацию:

минимизекс
$$6 \cdot x$$
 такое, что $Ax \cdot c$

Обратите внимание, что пример 9.2 можно записать таким образом. Условия ККТ для этой задачи:

Стационарность: А
$$\mu$$
 = b Первичная выполнимость: А x с Дополнительная нежесткость: μ i(ai · x ci) = 0 i, где μ это строка i из А Двойная осуществимость: μ 0

Как и в случае множителей Лагранжа, мы не можем предполагать, что любой х , удовлетворяющий условиям ККТ, автоматически минимизирует f с учетом ограничений, даже локально. Один из способов проверить локальную оптимальность — исследовать гессиан функции f, ограниченный подпространством Rn , в котором х может перемещаться, не нарушая ограничений; если этот «приведенный» гессиан положительно определен, то оптимизация достигла локального минимума.

²Из http://www.math.ubc.ca/~israel/m340/kkt2.pdf

9.3 Алгоритмы оптимизации

Тщательное рассмотрение алгоритмов оптимизации с ограничениями выходит за рамки нашего обсуждения; к счастью, существует много стабильных реализаций этих методов, и многое можно сделать как «клиент» этого программного обеспечения, а не переписывать его с нуля. Тем не менее, полезно набросать некоторые возможные подходы, чтобы получить представление о том, как работают эти библиотеки.

9.3.1 Последовательное квадратичное программирование (SQP)

Подобно BFGS и другим методам, которые мы рассматривали при обсуждении безусловной оптимизации, одной из типичных стратегий условной оптимизации является аппроксимация f, g и h более простыми функциями, решение аппроксимированной оптимизации и итерация.

Предположим, у нас есть предположение xk решения задачи условной оптимизации. Мы могли бы применить разложение Тейлора второго порядка к f и приближение первого порядка к g и h, чтобы определить следующую итерацию следующим образом:

$$xk+1$$
 $xk+$ аргумент мин. $\frac{1}{2}$ $dHf(xk)d+ f(xk)\cdot d+f(xk)$ $f(xk)\cdot d+f(xk)$ $f(xk)\cdot d=0$ $f(xk)+ f(xk)\cdot d=0$

Оптимизация для нахождения d имеет квадратичную цель с линейными ограничениями, для которых оптимизация может быть значительно проще с использованием одной из многих стратегий. Она известна как квадратичная программа.

Конечно, это приближение Тейлора работает только в окрестности оптимальной точки. Когда хорошее начальное предположение x0 недоступно, эти стратегии, скорее всего, потерпят неудачу.

Ограничения равенства Когда единственными ограничениями являются равенства и h удалено, квадратичная программа для d имеет условия оптимальности множителя Лагранжа, полученные следующим образом:

$$1 \Lambda(d,\lambda) \qquad \overline{2} dHf(xk)d + f(xk) \cdot d + f(xk) + \lambda (g(xk) + Dg(xk)d)$$

$$= 0 = d\Lambda = Hf(xk)d + f(xk) + [Dg(xk)]\lambda$$

В сочетании с условием равенства получается линейная система:

$$\begin{array}{cccc} Hf(xk) \left[Dg(xk)\right] & & \Gamma & & f(xk) \\ \mathcal{J}r(x\kappa) \ 0 & & \lambda & & "=" & g(xk) \end{array}$$

Таким образом, каждую итерацию последовательного квадратичного программирования при наличии только ограниченийравенств можно выполнить, решая эту линейную систему на каждой итерации, чтобы получить xk+1 xk + d. Важно отметить, что приведенная выше линейная система не является положительно определенной, поэтому в больших масштабах ее может быть трудно решить.

Расширения этой стратегии работают как BFGS, и аналогичные аппроксимации работают для оптимизации без ограничений, вводя аппроксимации гессиана Hf. Стабильность также может быть введена путем ограничения расстояния, пройденного за одну итерацию.

Ограничения неравенства Существуют специализированные алгоритмы для решения квадратичных программ, а не общих нелинейных программ, и их можно использовать для генерации шагов SQP. Одна примечательная стратегия состоит в том, чтобы поддерживать «активный набор» ограничений, которые активны как минимум по отношению к d; тогда можно применить описанные выше методы с ограничениями равенства, игнорируя неактивные ограничения. Итерации оптимизации активного набора обновляют активный набор ограничений, добавляя нарушенные ограничения к активному набору и удаляя те ограничения-неравенства hj , для которых f · hj 0, как в нашем обсуждении условий ККТ.

9.3.2 Барьерные методы

Другой вариант минимизации при наличии ограничений состоит в том, чтобы изменить ограничения на энергетические условия. Например, в случае ограничения равенством мы могли бы минимизировать «расширенную» цель следующим образом:

$$f\rho(x) = f(x) + \rho g(x)^{\frac{2}{2}}$$

Обратите внимание, что принятие р заставит g(x) быть как можно меньше, поэтому в конечном итоге мы достигнем g(x) 0. Таким образом, барьерный метод оптимизации с ограничениями применяет итерационные методы оптимизации без ограничений к fp и проверяет, насколько хорошо выполняются ограничения; если они не находятся в пределах заданного допуска, р увеличивается, и оптимизация продолжается с использованием предыдущей итерации в качестве отправной точки.

Барьерные методы просты в реализации и использовании, но они могут демонстрировать некоторые пагубные виды отказов. В частности, с ростом р влияние f на целевую функцию уменьшается, и гессиан функции fр становится все менее обусловленным.

Барьерные методы также могут применяться к ограничениям неравенства. Здесь мы должны убедиться, что hi(x) 0 для всех i; типичный выбор барьерных функций может включать 1/hi(x) («обратный барьер») или log hi(x) («логарифмический барьер»).

9.4 Выпуклое программирование

Вообще говоря, методы, подобные тем, которые мы описали для оптимизации с ограничениями, практически не дают гарантий качества выходных данных. Конечно, эти методы не могут получить глобальные минимумы без хорошего начального приближения x0, а в некоторых случаях, например, когда гессиан

рядом с х* не является положительно определенной, они могут вообще не сходиться.

Есть одно заметное исключение из этого правила, которое встречается во многих важных оптимизациях: выпуклое программирование. Идея здесь состоит в том, что когда f — выпуклая функция, а само допустимое множество выпукло, то оптимизация имеет единственный минимум. Мы уже определили выпуклую функцию, но нам нужно понять, что означает, что набор ограничений

выпуклый:

Определение 9.3 (выпуклое множество). Множество S Rn называется выпуклым, если для любых x,y S точка tx + (1 t)y также принадлежит S для любого t [0, 1].

Как показано на рисунке HOMEP, интуитивно множество является выпуклым, если его граничная форма не может изгибаться ни внутрь, ни наружу.

Пример 9.7 (Круги). Диск $\{x Rn : x2 1\}$ выпукл, а единичная окружность $\{x Rn : x2 = 1\}$ — нет.

Легко видеть, что выпуклая функция имеет единственный минимум, даже если эта функция ограничена выпуклой областью. В частности, если функция имеет два локальных минимума, то линия точек между этими минимумами должна давать значения f не выше, чем в конечных точках.

Сильные гарантии сходимости доступны для выпуклых оптимизаций, которые гарантируют нахождение глобального минимума, пока f выпукло, а ограничения на g и h создают выпуклое допустимое множество. Таким образом, ценным упражнением почти для любой задачи оптимизации является проверка ее выпуклости, поскольку такое наблюдение может во много раз повысить уверенность в качестве решения и шансы на успех.

Новая область, называемая дисциплинированным выпуклым программированием, пытается объединить простые правила о выпуклости для создания выпуклых оптимизаций (СІТЕ CVX), позволяя конечному пользователю комбинировать простые выпуклые энергетические условия и ограничения, если они удовлетворяют критериям, делающим окончательную оптимизацию выпуклой. Полезные утверждения о выпуклости в этой области включают следующее:

- Пересечение выпуклых множеств выпукло; таким образом, добавление нескольких выпуклых ограничений является допустимой операцией.
- Сумма выпуклых функций выпукла.
- Если f и g выпуклы, то h(x) max{ f(x), g(x)}.
- Если f выпуклая функция, множество {x : f(x) с} выпукло.

Такие инструменты, как библиотека CVX, помогают отделить реализацию различных выпуклых задач от их минимизации.

Пример 9.8 (выпуклое программирование).

- Неотрицательная задача наименьших квадратов в примере 9.3 является выпуклой, поскольку Ах b2 является выпуклой функция от x, a множество x 0 выпукло.
- Задача линейного программирования в примере 9.6 является выпуклой, поскольку имеет линейную цель и линейную ограничения.
- Мы можем включить х1 в задачу выпуклой оптимизации, введя переменную у. Для этого мы добавляем ограничения уі хі и уі хі для каждого і и цели і уі . Члены этой суммы не меньше |хі | и что энергия и ограничения выпуклы. Как минимум, мы должны иметь уі = |хі | поскольку мы наложили ограничение на уі |хі | и мы хотим минимизировать энергию. «Дисциплинированные» выпуклые библиотеки могут выполнять такие операции за кулисами, не раскрывая такие подстановки конечному пользователю.

Особенно важным примером выпуклой оптимизации является линейное программирование из примера 9.6. Знаменитый симплекс-алгоритм отслеживает активные ограничения, вычисляет результирующий х с помощью линейной системы и проверяет, нужно ли обновлять активное множество; приближения Тейлора не нужны, потому что целевое и допустимое множество задаются линейным механизмом. Стратегии линейного программирования внутренних точек, такие как барьерный метод, также успешно решают эти задачи. По этой причине линейные программы можно решать в огромных масштабах — до миллионов или миллиардов переменных! — и они часто появляются в таких задачах, как планирование или ценообразование.

9.5 Проблемы

- Получить симплекс?
- Двойственность линейного программирования

Machine Translated by Google

Глава 10

Итерационные линейные решатели

В предыдущих двух главах мы разработали стратегии для решения нового класса задач, связанных с минимизацией функции f(x) с ограничениями на x или без них. При этом мы отказались от численной линейной алгебры и, в частности, метода исключения Гаусса, согласно которому мы должны найти точное решение системы уравнений, и вместо этого обратились к итерационным схемам, которые гарантированно приближают минимум функции все лучше и лучше по мере того, как они повторять все больше и больше. Даже если мы никогда точно не найдем минимум, мы знаем, что в конечном итоге мы найдем x0 с f (x0) 0 с произвольным качеством уровней, в зависимости от количества выполненных итераций.

Теперь у нас есть повторение нашей любимой задачи из числовой линейной алгебры, решение Ax = b относительно x, но мы применяем итеративный подход, а не ожидаем найти решение в закрытой форме. Эта стратегия открывает новый класс решателей линейных систем, которые могут находить надежные приближения x за невероятно малое число итераций. Мы уже предложили, как подойти к этому в нашем обсуждении линейной алгебры, предполагая, что решения линейных систем являются среди прочего минимумами энергии Ax b.

Зачем создавать еще один класс решателей линейных систем? До сих пор большинство наших прямых подходов требовали от нас представления А в виде полной матрицы размера n × n, а такие алгоритмы, как LU, QR или факторизация Холецкого, приводили около тысячи ³) время. Есть два случая, о которых следует помнить для роten причин для использования итерационных схем:

- 1. Когда А является разреженным, такие методы, как исключение Гаусса, имеют тенденцию вызывать заполнение, а это означает, что даже) если А содержит O(n) ненулевых значений, промежуточные шаги исключения могут ввести O(n ненулевых значений. Это свойство может быстро привести к нехватке памяти в системах линейной алгебры. Напротив, алгоритмы в этой главе требуют только того, что вы можете применить А для векторов, что может быть выполнено за время, пропорциональное количеству ненулевых значений в матрице.
- 2. Мы можем захотеть победить O(n ³) время выполнения стандартных методов матричной факторизации. В частности, если итерационная схема может найти довольно точное решение Ax = b за несколько итераций, время выполнения можно значительно сократить.

Кроме того, обратите внимание, что многие из рассмотренных нами нелинейных методов оптимизации, в частности, зависящие от ньютоновского шага, требуют решения линейной системы на каждой итерации! Таким образом, разработка максимально быстрого решателя может иметь большое значение при реализации крупномасштабных методов оптимизации, требующих одного или нескольких линейных решений на итерацию. На самом деле, в этом случае неточное, но быстрое решение линейной системы может быть приемлемым, поскольку оно в любом случае используется в более крупной итерационной методике.

Обратите внимание, что большая часть нашего обсуждения связана с CITE, хотя наше развитие может быть несколько короче, учитывая развитие предыдущих глав.

10.1 Градиентный спуск

Мы сосредоточим наше обсуждение на решении Ах = b, где A имеет три свойства:

- 1. A Rn×n является квадратным
- 2. А симметрична, т. е. А = А
- 3. А положительно определена, т. е. для всех x = 0, x Ax > 0

Ближе к концу этой главы мы ослабим эти предположения. А пока заметьте, что мы можем заменить Ax = b нормальными уравнениями A Ax = A b, чтобы удовлетворить этим критериям, хотя, как мы обсуждали, эта замена может создать проблемы с числовыми условиями.

10.1.1 Получение итерационной схемы

В этом случае легко проверить, что решения Ax = b являются минимумами функции f(x), заданной квадратичной формой

1
$$f(x)$$
 x Ax bx + c 2 для

любого с R. В частности, взятие производной от f показывает

$$f(x) = Ax$$
 b,

и установка f(x) = 0 дает желаемый результат.

Вместо того, чтобы решать f(x) = 0 напрямую, как мы это делали в прошлом, предположим, что мы применяем стратегия градиентного спуска к этой минимизации. Вспомним базовый алгоритм градиентного спуска:

- 1. Вычислить направление поиска d k f(xk-1) = b Axk 1.
- 2. Определить $xk = xk + 1 + \alpha kd k$, где αk выбрано так, что f(xk) < f(xk + 1)

Для общей функции f определение значения α k может быть сложной одномерной задачей «линейного поиска», сводящейся к минимизации f(xk 1 + α kd k) как функции одной переменной α k 0. Для наш конкретный выбор квадратичной формы f(x) = x Ax bx + c, однако мы можем выполнять лине α t поиск в закрытой форме. В частности, определить

Таким образом, если мы хотим минимизировать g по α , мы просто выбираем

$$\alpha = \frac{d(b - Ax) dAd}{dAd}$$

В частности, для градиентного спуска мы выбрали d K = b Axk, поэтому на самом деле αk принимает красивый вид:

В итоге наша формула для линейного поиска дает следующую схему итеративного градиентного спуска для решения Ax =b в симметричном положительно определенном случае:

$$A_{K} = b - Axk-1$$

$$aK = \frac{AA_{K}}{A_{K}}$$

$$xk = xk + 1 + \alpha kd k$$

10.1.2 Конвергенция

По построению наша стратегия градиентного спуска уменьшает f(xk) при k . Тем не менее, мы не показали, что алгоритм достигает минимально возможного значения f, и мы не смогли описать, сколько итераций мы должны выполнить, чтобы достичь разумного уровня уверенности в том, что Axk b.

Одна простая стратегия для понимания сходимости алгоритма градиентного спуска для нашего выбора f состоит в том, чтобы исследовать изменение обратной ошибки от итерации к итерации.1 Предположим, что это решение, которое мы * х ищем, то есть Ax = b. Затем мы можем изучить соотношение ошибка от итерации к итерации:

Rk
$$\frac{f(xk)}{f(xk-1)}$$
 $\frac{f(x-1)}{f(x-1)}$

Очевидно, ограничение Rk < β < 1 для некоторого β показывает, что градиентный спуск сходится.

Для удобства мы можем разложить f(xk):

¹ Этот аргумент представлен, например, в http://www-personal.umich.edu/~mepelman/teaching/IOE511/Handouts/511notes07-7.pdf.

$$= f(xk - 1) \qquad d \qquad \frac{AA \times K}{K} \qquad AA \times K \times 2 \qquad \frac{1}{\Gamma_{K}} \qquad \frac{AA \times K}{\Gamma_{K} \qquad \text{объявление}} \qquad \Gamma \qquad Ad \times K \text{ по определению } ak \times K$$

$$= f(xk - 1) \qquad \frac{\left(d \times A \times\right)^{2}}{\Gamma_{M} \qquad \text{объявление}} \qquad 1 + \frac{\left(A \times A \times\right)^{2}}{A_{M} \qquad \text{объявление}} = f(xk - 1) \qquad \frac{\left(d \times A \times\right)^{2}}{2d} \qquad \frac{2}{2d} \qquad \frac{d}{d} = \frac{d}$$

Таким образом, мы можем вернуться к нашей дроби:

$$PK = rac{e (x \kappa - 1) - rac{(A_K A_K)^2}{2A_K} f(x)}{f(x k 1) f(x)}$$
 по нашей формуле для $f(x k)$

$$= 1 - rac{(A_K A_K)^2}{2d_K}$$
 Ad $k(f(x k 1) f(x))$

Обратите внимание, что Ax = b, поэтому мы можем написать:

Таким образом,

$$Rk = 1 - 2d k$$
 $\frac{\left(\text{Д к } \text{Д к} \right)^2}{\text{Ad k}(f(\text{xk } 1) \quad f(\text{x} \quad)) d}$ $= 1 - \text{Д} \quad \frac{\left(\text{Д к } \text{k} \right)}{\text{Ad k} \cdot \text{d k }_{\text{K}} \quad \text{A} \quad 1d \text{k}}$ нашим последним упрощением $= 1 - \text{Д} \quad \frac{\Gamma_{\text{K}} \quad \text{Д k}}{\frac{\Gamma_{\text{K}} \quad \text{L}}{\frac{\Gamma_{\text{K}} \quad \text{L}}}{\frac{\Gamma_{\text{K}} \quad \text{L}}{\frac{\Gamma_{\text{K}} \quad \text{L}}{\frac{\Gamma_{\text{K}}$

Потребовалось значительное количество алгебры, но мы доказали важный факт: Сходимость градиентного спуска по f зависит от обусловленности A. То есть, чем лучше обусловлен A, тем быстрее будет сходиться градиентный спуск. Кроме того, поскольку cond A

1, мы знаем, что наша вышеприведенная стратегия градиентного спуска безоговорочно сходится к x, хотя

*, сходимость может быть медленной, когда A плохо обусловлен.

Рисунок НОМЕР иллюстрирует поведение градиентного спуска для хорошо и плохо обусловленных матриц. Как вы можете видеть, градиентный спуск может изо всех сил пытаться найти минимум нашей квадратичной функции f, когда собственные значения A имеют большой разброс.

10.2 Сопряженные градиенты

Напомним, что для решения Ax = b для A Rn×n использовала срвемя. Пересмотр градиентного спуска стратегия O(n, описанная выше, мы видим, что каждая 2) времени, так как мы должны вычислить матрицу-вектор итерация требует O(n произведений между A, xk 1 и dk . Таким образом, если градиентный спуск занимает более n итераций, мы с тем же успехом можно было бы применить исключение Гаусса, которое восстановит точное решение за то же время. К сожалению, мы не можем показать, что градиентный спуск должен выполнять конечное число итераций, и на самом деле в плохо обусловленных случаях это может потребоваться огромное количество итераций, чтобы найти минимум.

По этой причине мы разработаем алгоритм, который гарантированно сходится не более чем за n шагов () в сохраняя O(n), что ³ наихудшем случае для решения линейных систем. По пути мы найдем этот алгоритм на самом деле демонстрирует лучшие свойства сходимости в целом, что делает его разумным выбором, даже если мы не запустим его до конца.

10.2.1 Мотивация

Наш вывод алгоритма сопряженных градиентов мотивирован довольно простым наблюдением. Предположим, мы знали решение х по-другому:

* форму № ВхТогда мы можем записать нашу квадратичную

$$f(x) = x Ax - bx + c$$
 по определению 2 1 (x x $x = x - b$) $A(x + c) + x Ax = x + b$ $A(x + c) + x + c$ $A(x + c) + x$

Таким образом, с точностью до постоянного сдвига f то же самое, $\frac{1}{12}$ (x x) A(x x). Конечно, мы делаем что и произведение, не знающее x, но это наблюдение показывает нам природу f; это просто измерение от x к x относительно «А-нормы» v На самом деле, $\frac{2}{\Delta}$ расстояния v Av.

поскольку А симметрична и положительно определена, даже если это может быть медленным для выполнения на практике, мы знаем, что ее можно разложить на множители с использованием стратегии Холецкого как A = ЛЛ. Имея в руках эту факторизацию, f принимает еще более приятный вид:

$$1 \varphi(x) = {}_{2} - L(x - x *)$$
 2 + KOHCT. 2

Поскольку L является обратимой матрицей, эта норма действительно является мерой расстояния между x и x * - Определите y L x и y * L X * Затем с этой новой точки зрения мы минимизируем 2 . Конечно, если f(y) = . y - y бы мы действительно могли добраться до этого с помощью факторизации Холецкого, оптимизируя f было бы чрезвычайно легко, но чтобы вывести схему этой минимизации без L, мы рассмотрим возможность минимизации f, используя только линейный поиск, полученный в § 10.1.1.

Мы делаем простое наблюдение о минимизации нашей упрощенной функции f с помощью такой стратегии. egy, показанный на PИСУНКЕ HOMEP:

Предложение 10.1. Предположим, что {w 1, ..., wn } ортогональны в Rn . Затемf минимизируется не более чем за n шагов путем поиска строки в направлении w 1, затем в направлении w 2 и так далее.

Доказательство. Возьмем столбцы матрицы Q Rn×n в качестве векторов w i ; Q — ортогональная матрица. Поскольку Q f(y) = у мы можем написать 2 ; другими словами**у юртооввраем**ваем так далее. Если мы напишем z Qy и после второй итерации

*г Qy и *, то, очевидно, после первой итерации должно быть z1 = z
$$\mathring{1}$$
, z2 = z2 , * так далее. После п шагов мы достигаем zn = z \mathring{n} , что дает желаемый результат.

Таким образом, оптимизацию f всегда можно выполнить с помощью n строковых поисков, если эти поиски осуществляются в ортогональных направлениях.

Все, что мы сделали для перехода от f к f, — это повернули координаты, используя L . Такое линейное преобразование переводит прямые линии в прямые линии, поэтому выполнение поиска строки f вдоль некоторого вектора w эквивалентно равносильно поиску строки вдоль (L) 1w нашей исходной квадратичной функции f. Наоборот, если мы делаем n строковых поисков по f в направлениях vi , таких что L vi w i ортогональны, то по предложению 10.1 мы должны найти х

 * . Обратите внимание, что запрос w i \cdot w = 0 аналогичен заданию j

$$0 = w \ i \cdot w \ j \qquad = (L \ vi) \ (L \ vj) = v \ i \ (LL) vj = v \ i \ Avj \ .$$

Мы только что доказали важное следствие предложения 10.1. Определим сопряженные векторы следующим образом:

Определение 10.1 (А-сопряженные векторы). Два вектора v, w являются A -сопряженными, если v Aw = 0.

Затем, основываясь на нашем обсуждении, мы показали:

Предложение 10.2. Предположим, {v1, ..., vn} A-сопряжены. Затем f минимизируется не более чем за n шагов путем поиска строки в направлении v1, затем в направлении v2 и так далее.

На высоком уровне алгоритм сопряженных градиентов просто применяет это предложение, генерируя и ища по Асопряженным направлениям, а не перемещаясь по f. Обратите внимание, что этот результат может показаться несколько нелогичным: мы не обязательно движемся по направлению наискорейшего спуска, а скорее просим, чтобы наш набор направлений поиска удовлетворял глобальному критерию, чтобы убедиться, что мы не повторяем работу. Эта установка гарантирует сходимость за конечное число итераций и признает структуру f в терминах f, обсуждавшуюся выше.

Напомним, что мы мотивировали А-сопряженные направления тем, что заметили, что они ортогональны после применения L из факторизации A = LL. С этой точки зрения мы имеем дело с двумя скалярными произведениями: $xi \cdot xj$ и $yi \cdot yj$ (L xi) · (L xj) = xi i LLxj = xi Axj . Эти два продукта будут фигурировать в нашем последующем обсуждении в равных количествах, поэтому мы обозначаем «А-внутренний продукт» как

$$u,vA$$
 (L u) · (L v) = u Av.

10.2.2 Субоптимальность градиентного спуска

На данный момент мы знаем, что если мы сможем найти п А-сопряженных направлений поиска, мы сможем решить Ах =b за п шагов с помощью линейного поиска вдоль этих направлений. Остается найти стратегию, позволяющую максимально эффективно находить эти направления. Для этого мы рассмотрим еще одно свойство алгоритма градиентного спуска, которое вдохновит на более совершенный подход.

Предположим, что мы находимся в точке xk во время метода итеративного поиска строки на f(x); мы будем называть направление наискорейшего спуска функции f в точке xk невязкой rk b Axk. Мы не можем решить выполнить линейный поиск вдоль rk, как при градиентном спуске, поскольку направления градиента не обязательно являются A-сопряженными. Итак, слегка обобщая, мы найдем xk+1 поиском по строке по еще не определенному направлению vk+1.

Из нашего вывода градиентного спуска мы должны выбрать $xk+1 = xk + \alpha k + 1 vk + 1$, где $\alpha k + 1$ определяется выражением

$$a\kappa+1 = \frac{{}^{B}k+1 rk}{v k+1 A v k+1}$$

Применяя это расширение xk+1, мы можем написать альтернативную формулу обновления для остатка:

Эта формула верна независимо от нашего выбора vk+1 и может быть применена к любому итеративному методу поиска строк.

Однако в случае градиентного спуска мы выбрали vk+1 rk. Этот выбор дает рекуррентное соотношение:

$$rk+1 = rk$$
 $\alpha k+1 A r k$.

Эта простая формула приводит к поучительному утверждению:

Предложение 10.3. При выполнении градиентного спуска на f, $span\{r0, ..., rk\} = span\{r0, Ar0, ..., A кр0<math>\}$.

Доказательство. Это утверждение следует по индукции из нашей формулы для rk+1 выше.

Раскрываемая нами структура начинает очень походить на методы Крылова, упомянутые в главе 5: Это не ошибка!

Градиентный спуск достигает xk, двигаясь по r0, затем по r1 и так далее по rk. Таким образом, в итоге мы знаем, что итерация xk градиентного спуска на f лежит где-то в плоскости x0 + A k 1r0} по предложению 10.3. К сожалению, что если мы запустим градиентный спуск, итерация xk это span {r0,r1,...,rk 1} = x0 + span {r0, Ar0,..., неверно, будет оптимальной в этом подпространстве. Иными словами, в общем случае может быть так, что

В идеале замена этого неравенства на равенство гарантирует, что генерация xk+1 из xk не «отменит» никакой работы, проделанной во время итераций с 1 по k 1.

Если мы пересмотрим наше доказательство предложения 10.1 с учетом этого факта, мы сможем сделать наблюдение, подсказывающее, как мы могли бы использовать сопряжение для улучшения градиентного спуска. В частности, когда zi переключается * а z, оно никогда не меняет значение в будущей итерации. Другими словами, после поворота от z к i , х выполняется следующее предложение:

Предложение 10.4. Возьмем xk за k-ю итерацию процесса из предложения 10.1 после поиска вдоль vk . Затем,

$$xk$$
 $x0 = arg min v span e (x0 + v)$
 $\{v1,...,vk\}$

Таким образом, в лучшем из возможных миров, пытаясь превзойти градиентный спуск, мы могли бы надеяться, что А найдет А-сопряженные направления {v1, . . . ,vn} такие, что охватывают {v1, . . . ,vk} = span {r0, Ar0, . . . , для каждогок 1r0} k; тогда наша итерационная схема гарантированно будет работать не хуже, чем градиентный спуск на любой заданной итерации. Но мы жадно хотим сделать это без ортогонализации или хранения более чем конечного числа векторов за раз.

10.2.3 Генерация А-сопряженных направлений

Конечно, для любого набора направлений мы можем сделать их А-ортогональными, используя такой метод, как ортогонализация Грама Шмидта. К сожалению, ортогонализация {r0, Ar0, . . .} найти набор направлений поиска дорого и потребовало бы от нас ведения полного списка направлений vk; эта конструкция, вероятно, превысила бы требования по времени и памяти даже при исключении Гаусса.

Мы покажем одно заключительное наблюдение о Граме-Шмидте, которое делает сопряженные градиенты управляемыми, создавая сопряженные направления без дорогостоящего процесса ортогонализации.

Игнорируя эти проблемы, мы могли бы написать «метод сопряженных направлений» следующим образом:

Обновить направление поиска (неверный шаг Грама-Шмидта):
$$vk = A$$

$$R = \frac{A k + ro, via ro}{Bu , BuA}$$

Ви , виА

Поиск строки: $\alpha k = \frac{B k rk + 1}{B k ABK}$

Обновление оценки: $xk = xk + 1 + \alpha k vk$

Обновить остаток: $rk = rk + 1 + \alpha k A vk$

Здесь мы вычисляем k-е направление поиска vk, просто проецируя v1, ..., vk 1 из вектора A k 1r0. Этот алгоритм, очевидно, обладает свойством span $\{v1, ..., vk\}$ = span $\{r0, Ar0, ..., A k 1r0\}$ предложено в $\{s10.2.2, ho umeet две проблемы:$

- 1. Подобно итерации мощности для собственных векторов, мощность A k 1r0 , вероятно, будет выглядеть в основном как первый собственный вектор A, делая проекцию все более и более плохо обусловленной.
- 2. Мы должны сохранить v1, . . . ,vk 1 для вычисления vk; таким образом, каждая итерация этого алгоритма требует больше памяти и времени, чем предыдущая.

Мы можем решить первую проблему относительно простым способом. В частности, прямо сейчас мы проецируем предыдущие направления поиска из A k 1r0, но на самом деле мы можем проецировать предыдущие направления из любого вектора w при условии, что

k 1 w span {r0, Ar0, ..., A r0}\span {r0, Ar0, ..., A
k
 $^{2}\Gamma$ 0},

то есть до тех пор, пока w имеет некоторый компонент в новой части пространства.

Альтернативным выбором w с этим свойством является невязка rk 1. Это свойство следует из остаточного обновления rk = rk 1 αkAvk; в этом выражении мы умножаем vk на A, вводя нужную нам новую мощность A. Этот выбор также более точно имитирует алгоритм градиентного спуска, который взял vk =rk 1. Таким образом, мы можем немного обновить наш алгоритм:

Теперь мы не выполняем арифметические действия с плохо обусловленным вектором A k 1r0 , но все еще имеем описанную выше проблему «памяти».

На самом деле удивительное наблюдение по поводу шага ортогонализации, описанного выше, состоит в том, что большинство слагаемых в сумме равны нулю! Это удивительное наблюдение позволяет выполнять каждую итерацию сопряженных градиентов без увеличения использования памяти. Мы зафиксируем этот результат в предложении:

Предложение 10.5. В приведенном выше методе «сопряженного направления» rk ,vA = 0 для всех < k.

Доказательство. Поступим индуктивно. Для базового случая k = 1 доказывать нечего , поэтому предположим, что k > 1 и что результат верен для всех k < k. По формуле остаточного обновления мы знаем:

$$rk$$
, $vA = rk$ 1, vA $\alpha kAvk$, $vA = rk$ 1, vA αkvk , AvA ,

где второе равенство следует из симметрии А.

Во-первых, предположим, что < k 1. Тогда первый член указанной выше разности равен нулю по индукции. Кроме того, по построению Av span {v1, . . . ,v+1}, поэтому, поскольку мы построили наши направления поиска как Асопряженные, мы знаем, что второй член также должен быть равен нулю.

Чтобы завершить доказательство, рассмотрим случай = k 1. Используя формулу обновления невязки, мы знаем:

Avk
$$1 = \frac{1}{\alpha k + 1}$$
 (rk 2 rk 1)

Предварительное умножение byrk показывает:

rk,vk 1A =
$$\frac{1}{a\kappa-1}p_{\kappa} (rk 2 rk 1)$$

 Таким образом, наше доказательство выше показывает, что мы можем найти новое направление vk следующим образом:

$$vk = rk - 1$$
 $\frac{rk - 1, viA}{s_{1} + s_{2} + s_{3}} vi$ по формуле Грама-Шмидта $s_{2} + s_{3} + s_{4} + s_{5} + s_$

Поскольку суммирование по і исчезает, стоимость вычисления vk не зависит от k.

10.2.4 Формулировка алгоритма сопряженных градиентов

Теперь, когда у нас есть стратегия, которая дает А-сопряженные направления поиска с относительно небольшими вычислительными затратами, мы просто применяем эту стратегию, чтобы сформулировать алгоритм сопряженных градиентов. В частности, предположим, что х0 является начальным предположением решения задачи Ax = b, и возьмем r0 b Ax0. Для удобства возьмем v0 0. Затем мы итеративно обновим xk 1 до xk, используя серию шагов для k = 1, 2, . . .:

Обновить направление поиска:
$$vk = rk$$
 1 — $\frac{rk + 1, vk + 1A}{vk + 1, vk + 1A} \frac{vk}{vk}$ 1

Поиск строки: $\alpha k = \frac{\frac{B}{k} \frac{k}{rk} \frac{1}{k}}{\frac{B}{k} \frac{ABK}{k}}$

Обновление оценки: $xk = xk + 1 + \alpha kvk$

Обновить остаток: $rk = rk + 1 = \alpha kAvk$

Эта итеративная схема является лишь незначительной корректировкой алгоритма градиентного спуска, но имеет много желаемых свойств по конструкции:

- f(xk) ограничено сверху значением k-й итерации градиентного спуска.
- Алгоритм сходится к х * за n шагов
- На каждом шаге итерация хk является оптимальной в подпространстве, натянутом на первые k направлений поиска.

В интересах выжать максимальное численное качество из сопряженных градиентов, мы можем попытаться упростить численные выражения приведенных выше выражений. Например, если мы подставим обновление направления поиска в формулу для αk , по ортогональности мы можем написать:

$$a\kappa = \frac{p_{k} + 1 rk + 1}{B_{K} + ABK}$$

Теперь числитель этой дроби гарантированно неотрицательный без числовой точности. проблемы.

Точно так же мы можем определить константу βk , чтобы разделить обновление направления поиска на два этапа:

$$\beta \kappa - \frac{\text{rk 1,vk 1A}}{\text{BK-1,BK-1A}}$$

$$vk = \text{rk 1 + }\beta kvk 1$$

Мы можем упростить нашу формулу для βk:

Это выражение показывает, что βk 0, свойство, которое могло бы не сохраняться после проблем с численной точностью. У нас есть одно оставшееся вычисление ниже: (rk 2 αk 1Avk 1) по

С этими упрощениями у нас есть альтернативная версия сопряженных градиентов:

Обновить направление поиска:
$$\beta k = \frac{\frac{p_{k-1} \text{ rk} - 1}{p_{k-2} \text{ rk} - 2}}{\frac{p_{k-2} \text{ rk} - 2}{p_{k-2} \text{ rk} - 2}}$$

$$vk = rk - 1 + \beta k v k - 1$$
Поиск строки: $\alpha k = \frac{\frac{p_{k-1} \text{ rk} - 1}{p_{k-2} \text{ rk} - 2}}{\frac{p_{k-1} \text{ rk} - 1}{p_{k-2} \text{ rk} - 2}}$
Обновление оценки: $xk = xk - 1 + \alpha k v k$

Обновить остаток: rk = rk 1 $\alpha kAvk$

По числовым причинам иногда вместо использования формулы обновления forrk рекомендуется использовать формулу остатка rk =b Ахk . Эта формула требует дополнительного умножения матрицы на вектор, но исправляет числовой «дрейф», вызванный округлением с конечной точностью. Обратите внимание, что нет необходимости хранить длинный список предыдущих остатков или направлений поиска: сопряженные градиенты занимают постоянный объем памяти от итерации к итерации.

10.2.5 Условия сходимости и остановки

По построению алгоритм сопряженных градиентов (CG) гарантированно сходится не медленнее, чем градиентный спуск по f, при этом не сложнее в реализации и имеет ряд других особенностей.

положительные свойства. Подробное обсуждение сходимости компьютерной графики выходит за рамки нашего обсуждения, но в целом алгоритм лучше всего ведет себя на матрицах с равномерно распределенными собственными значениями в небольшом диапазоне. Одна грубая оценка, параллельная нашей оценке в §10.1.2, показывает, что алгоритм СG удовлетворяет:

$$\frac{f(xk) \quad f(x)}{f(x0) \quad f(x)} \quad 2 \quad -\frac{\overline{\kappa}-1}{\overline{\kappa}+1}$$

где к сond A. В более общем смысле количество итераций, необходимых для сопряженного градиента для достижения заданного значения ошибки, обычно может быть ограничено функцией к, тогда как границы сходимости градиентного спуска пропорциональны к.

Мы знаем, что сопряженные градиенты гарантированно сходятся к х ровно за п шагов, но когда п велико, может быть предпочтительнее остановиться раньше. Фактически, формула для βk будет делить на ноль, когда невязка становится очень короткой, что может вызвать проблемы с численной точностью вблизи минимума f. Таким образом, на практике ЦГ обычно останавливается, когда отношение rk/r0 достаточно мало.

10.3 Предварительная подготовка

Теперь у нас есть две мощные итерационные схемы для нахождения решений Ax = b, когда A симметрично и положительно определено: градиентный спуск и сопряженные градиенты. Обе стратегии сходятся безоговорочно, а это означает, что независимо от начального предположения x0 при достаточном количестве итераций мы получим сколь угодно близкое к истинному решению x ровно за конечное число *; на самом деле, сопряженные градиенты гарантируют, что мы достигнем x* итераций. Конечно, время, необходимое для достижения решения Ax = b для обоих этих методов, прямо пропорционально количеству итераций, необходимых для достижения x в пределах приемлемого допуска. Таким образом, имеет смысл * максимально настроить стратегию, чтобы свести к минимуму количество итераций для сходимости.

С этой целью мы замечаем, что мы можем охарактеризовать скорость сходимости обоих алгоритмов и многих других связанных итерационных методов с точки зрения числа обусловленности cond A. То есть, чем меньше значение cond A, тем меньше времени это займет. решить Ax = b. Обратите внимание, что эта ситуация несколько отличается для исключения Гаусса, которое требует одинакового количества шагов независимо от A; другими словами, обусловленность A влияет не только на качество результатов итерационных методов, но и на скорость, с которой х

* приближается.

Конечно, для любой обратимой матрицы P решение PAx = Pb эквивалентно решению Ax = b. Хитрость, однако, заключается в том, что число обусловленности PA не обязательно должно быть таким же, как у A; в (недостижимой) крайности, конечно, если бы мы взяли P = A, мы полностью устранили бы проблемы с обусловливанием! В более общем случае предположим, что P A cond A, и поэтому может быть целесообразно применить P перед решением

1. Тогда мы ожидаем cond PA линейной системы. В этом случае мы будем называть P предобуславливателем.

Хотя идея предварительной обработки кажется привлекательной, остаются два вопроса:

- 1. Хотя А может быть симметричным и положительно определенным, произведение РА в общем случае не будет обладать эти свойства.
 - 1 2. Нужно найти Р А это легче вычислить, чем А 1 себя.

Мы рассматриваем эти вопросы в следующих разделах.

10.3.1 CG с предварительным кондиционированием

Мы сосредоточим наше обсуждение на сопряженных градиентах, поскольку они обладают лучшими свойствами сходимости, хотя большинство наших конструкций довольно легко применимы и к градиентному спуску. С этой точки зрения, если мы посмотрим на наши конструкции в § 10.2.1, станет ясно, что наша конструкция СG сильно зависит как от симметрии, так и от положительной определенности A, поэтому запуск CG на PA обычно не будет сходиться из коробки.

Предположим, однако, что предобуславливатель Р сам по себе симметричен и положительно определен. Это разумное предположение, поскольку А должен удовлетворять этим свойствам. Тогда мы снова можем написать а = EE. Делаем следующее Факторизация Холецкого обратного Р Предложение

1 наблюдение:

10.6. Номер состояния РА такой же, как у Е-1АЕ-.

Доказательство. Мы показываем, что РА и Е 1AE имеют одинаковые сингулярные значения; число обусловленности — это отношение максимального к минимальному сингулярному значению, так что этого утверждения более чем достаточно. В частности, ясно, что Е 1AE симметричен и положительно определен, поэтому его собственные векторы являются его сингулярными значениями. Итак, предположим, что Е 1AE x = λx . Мы знаем Р. Таким образом, если мъргарим ельно умножим обе части выражения нашего собственного вектора на Е , мы найдем РАЕ x = λx .

Определение у E х показывает, что PAy = λy . Таким образом, PA и E -1AE- имеют полные собственные пространства и идентичные собственные значения.

Это предложение подразумевает, что если мы выполним СG на симметричной положительно определенной матрице E

1AE , мы получим те же преимущества кондиционирования, которые были бы у нас, если бы мы могли повторять PA. Как и при
доказательстве предложения 10.6, мы могли бы выполнить наше новое решение для у = E x в два этапа:

- 1. Решите E 1AE y = E 1b.
- 2. Решите x = E y.

Нахождение E было бы неотъемлемой частью этой стратегии, но, вероятно, это сложно, но мы вскоре докажем, что в этом нет необходимости.

Игнорируя вычисление Е, мы могли бы выполнить шаг 1, используя СG, следующим образом:

Обновить направление поиска:
$$\beta k = \frac{p_k - 1 \text{ rk} - 1}{p_k - 2 \text{ rk} - 2}$$

$$vk = rk \quad 1 + \beta kvk \quad 1$$
 Поиск строки: $\alpha k = \frac{p_k \quad 1 \ rk \quad 1}{B_K \quad 3 - 1A3 - BK}$

Обновить оценку: yk = yk 1 + αkvk

Обновить невязку: rk =rk 1 αkE 1AE vk

Эта итерационная схема будет сходиться в соответствии с условиями нашей матрицы Е 1АЕ .

Определим r^{\sim} k Erk, v^{\sim} k E vk и xk Eyk. Если мы вспомним соотношение P = E E, 1, мы можем перепишем нашу итерацию предварительно обусловленных сопряженных градиентов, используя эти новые переменные:

Обновить направление поиска:
$$\beta k = \frac{{^{\sim} k \ 1P^{\sim} k \ 1}}{{^{\sim} k \ 2P^{\sim} k \ 2}}$$

$$v^{\sim} \ k = Pr^{\sim} \ k \ 1 + \beta k v^{\sim} \ k \ 1$$
 Поиск строки: $\alpha k = \frac{\frac{r}{k} \ 1_{\prod K-1}}{\frac{v^{\sim}}{k} \ A v^{\sim} \ k}$ Обновить оценку: $xk = xk \ 1 + \alpha k v^{\sim} \ k$ Обновить остаток: $r^{\sim} \ k = r^{\sim} \ k \ 1 \ \alpha k A v^{\sim} \ k$

Эта итерация не зависит от факторизации Холецкого Р, выполненной с использованием

1 вообще, а вместо этого может быть исключительно приложений Р и А. Легко видеть, что схема хк х пользуется преимуществами

*, так на самом деле это предобусловливания без необходимости факторизации предобуславливателя.

В качестве примечания, еще более эффективное предварительное обусловливание можно выполнить, заменив А на РАQ для второй матрицы Q, хотя для применения этой второй матрицы потребуются дополнительные вычисления. Этот пример представляет собой распространенный компромисс: если само предварительное обусловливание требует слишком много времени для применения в одной итерации CG или другого метода, оно может не стоить уменьшенного количества итераций.

10.3.2 Общие предварительные условия

Поиск хороших предобуславливателей на практике — это не только наука, но и искусство. Поиск наилучшего зависит от аппроксимация P of A и т.д. ¹ структуры A, конкретного приложения и Однако даже грубые приближения могут значительно помочь сходимости, поэтому редко появляются приложения компьютерной графики, в которых не используется предобуславливатель.

Наилучшая стратегия для формулирования Р часто зависит от приложения, и интересная проблема инженерной аппроксимации включает в себя разработку и тестирование различных Р для наилучшего предобуславливателя. Две общие стратегии приведены ниже:

- Диагональный (или «Якоби») преобуславливатель просто принимает P за матрицу, полученную инвертированием диагональных элементов A; то есть P диагональная матрица с элементами 1/аii. Эта стратегия может смягчить неравномерное масштабирование от строки к строке, которое является частой причиной плохой обработки.
- Разреженный аппроксимативный обратный предобуславливатель формулируется путем решения подзадачи minP S AP IFro, где P ограничен набором S матриц, по которым оптимизация такой задачи является менее сложной. Например, общим ограничением является предписание шаблона разреженности для P, например, только ненулевые значения по диагонали или там, где A имеет ненулевые значения.
- Факторы неполного предобуславливателя Холецкого A L L _{*} а затем приближается к A путем решения соответствующих задач прямой и обратной замены. Например, популярная стратегия включает в себя выполнение шагов факторизации Холецкого, но сохранение результата только в позициях (i, j), где aij = 0.
- Ненулевые значения в А можно рассматривать как граф, а удаление ребер в графе или группировка узлов может привести к разъединению различных компонентов; результирующая система является блочно-диагональной после перестановки строк и столбцов и, таким образом, может быть решена с использованием последовательности меньших решений.

 Такая стратегия декомпозиции области может быть эффективна для линейных систем, возникающих из дифференциальных уравнений, подобных рассмотренным в главе НОМЕР.

Некоторые предобуславливатели поставляются с ограничениями, описывающими изменения в обусловливании А после замены его на РА, но по большей части это эвристические стратегии, которые следует протестировать и уточнить.

10.4 Другие итерационные схемы

Алгоритмы, которые мы подробно разработали в этой главе, применяются для решения Ах = b, когда A является квадратным, симметричным и положительно определенным. Мы сосредоточились на этом случае, потому что он очень часто встречается на практике, но бывают случаи, когда A асимметрично, неопределенно или даже прямоугольно. Выведение итерационных алгоритмов в каждом случае выходит за рамки нашего обсуждения, поскольку многие из них требуют некоторого специализированного анализа или расширенной разработки, но мы суммируем здесь некоторые методы на высоком уровне (СІТЕ ЕАСН):

- Методы расщепления разлагают A = M N и отмечают, что Ax = b эквивалентно Mx = Nx +b. Если M легко инвертировать, то схему с фиксированной точкой можно вывести, написав Mxk = Nxk 1 +b (СІТЕ); эти методы просты в реализации, но имеют сходимость в зависимости от спектра матрицы G = M 1N и, в частности, могут расходиться, когда спектральный радиус G больше единицы. Одним из популярных вариантов M является диагональ A. Такие методы, как последовательная чрезмерная релаксация (SOR), взвешивают эти два термина для лучшей сходимости.
- Метод остатка уравнения нормального сопряженного градиента (CGNR) просто применяет алгоритм CG к нормальным уравнениям A Ax = A b. Этот метод прост в реализации и гарантирует сходимость до тех пор, пока A имеет полный ранг, но сходимость может быть медленной из-за плохой обработки AA, как обсуждалось в главе NUMBER.
- Метод ошибки уравнения нормали сопряженного градиента (CGNE) аналогичным образом решает AAy =b; затем решение Ax =b есть просто A y.
- Такие методы, как MINRES и SYMMLQ, применяются к симметричным, но не обязательно положительно определенным матрицам A путем замены нашей квадратичной формы f(x) на g(x) b Ax2; эта функция g минимизируется на решениях Ax =b независимо от определенности A.
- Учитывая плохую обусловленность CGNR и CGNE, алгоритмы LSQR и LSMR также минимизируют g(x) с меньшим количеством допущений относительно A, в частности позволяя решать системы наименьших квадратов.
- Обобщенные методы, включая GMRES, QMR, BiCG, CGS и BiCGStab, решают Ax =b с единственной оговоркой, что А является квадратным и обратимым. Они оптимизируют аналогичные энергии, но часто должны хранить больше информации о предыдущих итерациях и, возможно, должны учитывать промежуточные матрицы, чтобы гарантировать сходимость с такой общностью.
- Наконец, методы Флетчера-Ривза, Полака-Рибьера и другие возвращаются к более общей задаче минимизации неквадратичной функции f, применяя шаги сопряженного градиента для поиска новых направлений поиска линии. Функции f, которые хорошо аппроксимируются квадратичными уравнениями, могут быть очень эффективно минимизированы с использованием этих стратегий, хотя они не обязательно используют гессиан; например, метод Флетчера-Ривза просто заменяет невязку в итерациях компьютерной графики отрицательным градиентом f. Можно охарактеризовать сходимость этих методов, когда они сопровождаются достаточно эффективными стратегиями поиска линий.

Многие из этих алгоритмов почти так же легко реализовать, как компьютерную графику или градиентный спуск, и существует множество реализаций, которые просто требуют ввода A и b. Многие из перечисленных выше алгоритмов требуют применения как A, так и A, что в некоторых случаях может быть технической проблемой. Как правило, чем более обобщенным является метод, то есть чем меньше допущений метод делает в отношении структуры матрицы A, тем больше итераций может потребоваться, чтобы компенсировать это отсутствие допущений. При этом не существует жестких и быстрых правил, просто взглянув на наиболее успешную итеративную схему, хотя существует ограниченное теоретическое обсуждение, сравнивающее преимущества и недостатки каждого из этих методов (СІТЕ).

10.5 Проблемы

- Получить CGNR и/или CGNE
- Получение МИНРЕС
- Получить Флетчер-Ривз
- Слайд 13 из http://math.ntnu.edu.tw/~min/matrix_computation/Ch4_Slide4_CG_2011. пдф

Часть IV

Функции, производные и интегралы



Глава 11

Интерполяция

До сих пор мы вывели методы для анализа функций f , например, для нахождения их минимумов и корней. Вычисление f(x) при конкретном x Rn может быть дорогостоящим, но фундаментальное допущение методов, разработанных нами в предыдущих главах, состоит в том, что мы можем получить f(x) , когда захотим, независимо от x.

Есть много контекстов, когда это предположение нереалистично. Например, если мы делаем снимок с помощью цифровой камеры, мы получаем сетку n × m значений цвета пикселей, отображающую континуум света, попадающего в объектив камеры. Мы можем думать о фотографии как о непрерывной функции от положения изображения (x, y) до цвета (r, g, b), но на самом деле мы знаем значение изображения только в точках, разделенных нм на плоскости изображения. Точно так же в машинном обучении и статистике часто нам даются образцы функции только в точках, где мы собирали данные, и мы должны интерполировать их, чтобы получить значения в другом месте; в медицинских учреждениях мы можем наблюдать за реакцией пациента на разные дозы лекарства, но можем только предсказать, что произойдет при дозировке, которую мы явно не пробовали.

В этих случаях, прежде чем мы сможем минимизировать функцию, найти ее корни или даже вычислить значения f(x) в произвольных точках x, нам нужна модель для интерполяции f(x) ко всему Rn (или некоторому его подмножеству) при заданном набор выборок f(xi). Конечно, методы, решающие эту проблему интерполяции, по своей сути являются приближенными, поскольку мы не знаем истинных значений f, поэтому вместо этого мы стремимся к тому, чтобы интерполируемая функция была гладкой и служила «разумным» предсказанием значений функции.

В этой главе мы будем предполагать, что значения f(xi) известны с полной уверенностью; в этом случае мы могли бы также думать о проблеме как о расширении f на оставшуюся часть домена, не нарушая значения ни в одном из входных местоположений. В главе HOMEP (НАПИШИТЕ МНЕ В 2014 ГОДУ) мы рассмотрим задачу регрессии, в которой значение f(xi) известно с некоторой неопределенностью, и в этом случае мы можем полностью отказаться от сопоставления f(xi) в пользу того, чтобы сделать f более гладким. .

11.1 Интерполяция в одной переменной

Прежде чем рассматривать наиболее общий случай, мы разработаем методы интерполяции функций одной переменной f: R R. В качестве входных данных мы возьмем набор из k пар (xi , yi) с предположением f(xi) = yi; наша задача — найти f(x) для x {x1, ..., xк}.

Наша стратегия в этом и других разделах будет черпать вдохновение из линейной алгебры, записывая f(x) в базисе. То есть набор всех возможных функций f: R — R слишком велик, чтобы с ним можно было работать, и включает в себя множество функций, которые непрактичны в вычислительной среде. Таким образом, мы упрощаем пространство поиска, заставляя f быть записанной как линейная комбинация более простого стандартного блока

Базисные функции. Эта стратегия уже знакома из базового исчисления: разложение Тейлора записывает функции в виде многочленов, тогда как ряды Фурье используют синус и косинус.

11.1.1 Полиномиальная интерполяция

Возможно, самый простой интерполянт состоит в том, чтобы предположить, что f(x) находится в R[x], множестве многочленов. Многочлены гладкие, и найти полином степени k - 1 через k точек выборки несложно .

На самом деле в примере 3.3 уже проработаны детали такого метода интерполяции. Как напомним, предположим, что мы хотим найти f(x) а0 + a1x + a2x ; здесь нашиминаkзвеtxными являются значения а0, . . . , ак 1 . Подстановка выражения yi = f(xi) для каждого і показывает, что вектор а удовлетворяет системе Вандермонда k × k:

Таким образом, полиномиальную интерполяцию степени k можно выполнить с помощью линейного решения аk × k, применяя наши общие стратегии из предыдущих глав, но на самом деле мы можем добиться большего.

Один из способов думать о нашей форме для f(x) состоит в том, что она написана в базисе. Так же, как базой для Rn является набор из n линейно независимых векторов v1, . . . , vn, здесь пространство полиномов степени x k 1}. Это может в диапазоне мономов {1, x, x R[x], но наш текущий выбор 2,..., быть наиболее очевидный базис для k 1, записанный имеет несколько свойств, которые делают его полезным для задачи интерполяции. Один из способов увидеть эту проблему — изобразить последовательность функций 1, 2 , 3 x, . . . для х [0, 1]; в этом выглядеть одинаково. интервале, легко видеть, что при увеличении k фун**кциk** хкоторые начинают

Продолжая применять нашу интуицию из линейной алгебры, мы можем решить записать наш многочлен в базисе, который больше подходит для рассматриваемой проблемы. На этот раз напомним, что нам дано k пар (х1, у1), . . . ,(хк , ук). Мы будем использовать эти (неподвижные) точки для определения интерполяционного базиса Лагранжа ф1,..., фk, написав:

$$\varphi i(x)$$
 $\frac{j=i(x-xj)}{j=i(xi-xj)}$

 ϕ i(x) $\frac{j=i \ (x \ xj)}{j=i \ (xi \ xj)}$ рм x^{k-1} , легко видеть, что каждый ϕ i по-прежнему полигон. Хотя это и не записано в базисе 1, х, х номиналом степени к 1. Кроме того, базис Лагранжа обладает следующим желательным свойством:

$$\phi i(x) =$$
 1, когда = і 0 в противном случае.

Таким образом, найти единственный полином степени k 1, соответствующий нашим парам (xi, yi), легко в базисе Лагранжа:

В частности, если мы подставим x = xj, то получим: f(xj)

= уј в соответствии с нашим выражением для фі(х) выше.

Таким образом, в базисе Лагранжа можно записать замкнутую формулу для f(x), не требующую решения системы Вандермонда. Недостаток, однако, заключается в том, что каждый фі(x) требует O(k) времени для оценки с использованием приведенной выше формулы для данного x, поэтому время; если мы найдем коэффициенты нахождение f(x) берет O(n ai из системы Вандермонда в явном виде, однако оценка время можно сократить до O(n).

Помимо времени вычислений, у базиса Лагранжа есть дополнительный численный недостаток. Обратите внимание, что знаменатель является произведением ряда членов. Если хі близки друг к другу, то произведение может включать много членов, близких к нулю, поэтому мы делим на потенциально небольшое число. Как мы видели, эта операция может создавать числовые проблемы, которых мы хотим избежать.

Одной из баз для полиномов степени k - 1, которая пытается найти компромисс между числовым качеством мономов и эффективностью базиса Лагранжа, является базис Ньютона, определяемый следующим образом:

$$y=1$$
 $\psi_i(x) = (x \quad x_i)$ $y=1$

Определим ψ 1(x) 1. Заметим, что ψ i(x) — многочлен степени i 1. По определению ψ i ясно, что ψ i(x) = 0 для всех < i. Если мы захотим написать f(x) = i ci ψ i(x) и выписать это наблюдение более явно, мы найдем:

$$f(x1) = c1\psi1(x1) f(x2) =$$
 $c1\psi1(x2) + c2\psi2(x2) f(x3) = c1\psi1(x3) +$
 $c2\psi2(x3) + c3\psi3(x3)$
 \vdots \vdots

Другими словами, мы можем решить следующую систему сил нижнего треугольника:

Эту систему можно решить за O(n, время с использованием прямой замены, а не O(n время необходимое для решения системы Вандермонда.

Теперь у нас есть три стратегии интерполяции k точек данных с использованием полинома степени k - 1 путем записи его в мономиальном, лагранжевом и ньютоновском базисах. Все три представляют различные компромиссы между числовым качеством и скоростью. Однако важным свойством является то, что результирующая интерполированная функция f(x) одинакова в каждом случае. Точнее говоря, существует ровно один многочлен степени k - 1, проходящий через набор из k точек, поэтому, поскольку все наши интерполянты имеют степень k - 1, они должны иметь одинаковый выход.

11.1.2 Альтернативные базы

Хотя полиномиальные функции особенно хорошо поддаются математическому анализу, нет фундаментальной причины, по которой наш интерполяционный базис не может состоять из различных типов функций. Например, венчающий результат анализа Фурье подразумевает, что большой класс функций хорошо аппроксимируется суммами тригонометрических функций соs(kx) и sin(kx) при k N. Конструкция

как система Вандермонда все еще применима в этом случае, и на самом деле алгоритм быстрого преобразования Фурье (который заслуживает более широкого обсуждения) показывает, как выполнить такую интерполяцию еще быстрее.

Меньшее расширение разработки в §11.1.1 относится к рациональным функциям вида:

$$= 2 \frac{p0 + p1x + p2x f(x)^{2} + \dots + mx}{q0 + q1x + q2x + \dots + qnx} + \frac{M}{q0}$$

Обратите внимание, что если нам дано k пар (xi , yi), то нам потребуется m + n + 1 = k, чтобы эта функция была корректно определена. Необходимо установить одну дополнительную степень свободы, чтобы учесть тот факт, что одна и та же рациональная функция может быть выражена несколькими способами путем одинакового масштабирования числителя и знаменателя.

Рациональные функции могут иметь асимптоты и другие закономерности, которые невозможно получить, используя только полиномы, поэтому они могут быть желательными интерполянтами для функций, которые быстро изменяются или имеют полюса. На самом деле, когда m и n фиксированы, коэффициенты pi и qi по-прежнему можно найти с помощью линейных методов, умножив обе части на знаменатель:

$$yi(q0 + q1xi + q2x + 2x + qnx + qnx + m) = p0 + p1xi + p2x + 2x + qnx + mx$$

Опять же, неизвестными в этом выражении являются р и q.

Однако гибкость рациональных функций может вызвать некоторые проблемы. Например, рассмотрим следующий пример: Пример 11.1

(Ошибка рациональной интерполяции, Булирш-Стёр, $\S 2.2$). Предположим, мы хотим найти f(x) со следующими точками данных: (0, 1), (1, 2), (2, 2). Мы могли бы выбрать m = n = 1. Тогда наши линейные условия станут такими:

$$q0 = p0$$

 $2(q0 + q1) = p0 + p1 \ 2(q0 + 2q1) = p0 + 2p1$

Одно нетривиальное решение этой системы:

$$p0 = 0$$
 $p1 = 2$
 $q0 = 0$
 $q1 = 1$

Отсюда следует следующая форма для f(x):

Эта функция имеет вырождение при x = 0, и на самом деле сокращение x в числителе и знаменателе не дает f(0) = 1, как хотелось бы.

Этот пример иллюстрирует более широкое явление. Наша линейная система для нахождения р и q может столкнуться с проблемами, когда результирующий знаменатель рх имеет корень при любом из фиксированных хі . Можно показать, что в этом случае не существует рациональной функции с фиксированным выбором m и n, интерполирующих заданные значения. Типичное частичное разрешение в этом случае представлено в (СІТЕ), которое попеременно увеличивает m и n до тех пор, пока не будет найдено нетривиальное решение. Однако с практической точки зрения специализированный характер этих методов является хорошим индикатором того, что альтернативные стратегии интерполяции могут быть предпочтительнее, когда основные рациональные методы не работают.

11.1.3 Кусочная интерполяция

До сих пор мы строили наши стратегии интерполяции, комбинируя простые функции на всех R. Однако, когда количество точек данных k становится большим, становятся очевидными многие вырождения. Например, на рисунке NUMBER показаны примеры, в которых подгонка полиномов высокой степени к входным данным может привести к неожиданным результатам. Кроме того, на рисунке NUMBER показано, что эти стратегии нелокальны, а это означает, что изменение любого отдельного значения уі во входных данных может изменить поведение f для всех x, даже тех, которые далеки от соответствующего хі . Почему-то это свойство нереалистично: мы ожидаем, что только входные данные вблизи заданного x повлияют на значение f(x), особенно когда имеется большое облако входных точек.

По этим причинам, когда мы проектируем набор базисных функций φ1, . . . , φk , желательно не только то, что с ними легко работать, но и то, что они имеют компактный носитель:

Определение 11.1 (Компактная опора). Функция g(x) имеет компактный носитель, если существует C = 0 для любого = 0 для

То есть функции с компактным носителем имеют только конечный диапазон точек, в которых они могут принимать ненулевые значения.

Обычная стратегия построения интерполяционных базисов с компактным носителем состоит в том, чтобы делать это кусочно. В частности, большая часть литературы по компьютерной графике основана на построении кусочных полиномов, которые определяются путем разбиения R на набор интервалов и записи разных полиномов в каждом интервале. Для этого мы упорядочим наши точки данных так, чтобы x1 < x2 < · · · < xk . Затем два простых примера кусочных интерполянтов следующие:

- Кусочная постоянная (НОМЕР РИСУНКА): для заданного х найдите точку данных хі , минимизирующую |x| хи |x| и определить |x| уі .
- Кусочно-линейный (НОМЕР РИСУНКА): Если x < x1, возьмите f(x) = y1, а если x > xk, возьмите f(x) = yk. В противном случае найдите интервал с x = [xi, xi+1] и определите

$$f(x) = yi+1 \cdot + \underbrace{yi \cdot x \cdot xu}_{Xi+1} \quad xj+1 \quad xi \quad \frac{x \cdot xu}{Xi+1}$$

В более общем случае мы можем записать разные многочлены в каждом интервале [xi , xi+1]. Обратите внимание на нашу схему: кусочно-постоянные полиномы разрывны, а кусочно-линейные функции непрерывны. Легко видеть, что кусочно-линейные функции непрерывны. Легко видеть, что кусочно-хевадратичные могут быть С и т.д. Эта повышенная непрерывность и дифференцируемносты и ученосты и дифференцируемносты и ученосты и дифференцируемносты и ученость и дифференцируемносты и ученость и дифференцируемность и диференцируемность и диференцируемн

Однако такая повышенная преемственность имеет свои недостатки. С каждой дополнительной степенью дифференцируемости мы делаем более сильное предположение о гладкости f . Это допущение может быть нереалистичным: многие физические явления действительно зашумлены или прерывисты, и эта повышенная гладкость может негативно повлиять на интерполяционные результаты. Одна область, в которой этот эффект особенно заметен, — это когда интерполяция используется в сочетании с инструментами физического моделирования. Моделирование турбулентных потоков жидкости с помощью сверхсглаженных функций может устранить прерывистые явления, такие как ударные волны, которые желательны в качестве выходных данных.

Помимо этих проблем, кусочные полиномы все еще могут быть записаны как линейные комбинации базисных функций. Например, следующие функции служат основой для кусочно-постоянных функций:

$$\phi$$
i(x) = $\begin{pmatrix} 1, \text{ когда xi} & \frac{1+xi}{2} & x < & \frac{xu+xu+1}{2} \\ 0 \text{ иначе} & 2 \end{pmatrix}$

Этот базис просто помещает константу 1 рядом с хі и 0 в другом месте; кусочно-постоянная интерполяция множества точек (хі , уі) записывается как f(x) = і уіфі(x). Точно так же так называемый базис «шляпы», показанный на рисунке HOMEP, охватывает набор кусочно-линейных функций с острыми краями в наших точках данных. хи:

Опять же, по построению кусочно-линейная интерполяция заданных точек данных есть $f(x) = i yi \psi i(x)$.

11.1.4 Гауссовские процессы и кригинг

Не включено в CS 205A, осень 2013 г.

11.2 Многопараметрическая интерполяция

Существует множество расширений описанных выше стратегий для интерполяции функции с заданными точками данных (xi , yi) , где xi Rn теперь может быть многомерным. Стратегии интерполяции в этом случае, однако, не столь ясны, поскольку менее очевидно разбить Rn на небольшое число областей вокруг xi . По этой причине общепринятой схемой является интерполяция с использованием функций относительно низкого порядка, то есть предпочтение упрощенных и эффективных стратегий интерполяции тем, которые выводят функции С .

Если все, что нам дано, это набор входных и выходных данных (хі , уі), то одна кусочно-постоянная стратегия интерполяции состоит в использовании интерполяции ближайшего соседа. В этом случае f(x) просто принимает значение уі , соответствующее хі, сводящее к минимуму х хі2; простые реализации перебирают все і, чтобы найти это значение, хотя структуры данных, такие как деревья kd, могут быстрее находить ближайших соседей. Точно так же, как кусочно-постоянные интерполяции делят R на интервалы вокруг точек данных хі , стратегия ближайшего соседа делит Rn на набор ячеек Вороного:

Определение 11.2 (ячейка Вороного). Дан набор точек $S = \{x1, x2, \dots, xk\}$ Rn , соответствующее , камера Вороного конкретному xi, есть множество Vi $\{x: x \mid xi2 < x \mid xj2 \text{ для всех } j=i\}$. То есть это множество точек ближе к xi , чем к любому другому xj в S.

На рисунке НОМЕР показан пример ячеек Вороного о наборе точек данных в ячейках R2, обладающих многими благоприятными свойствами; например, они являются выпуклыми многоугольниками и локализованы относительно каждарсамом деле связность ячеек Вороного — это хорошо изученная проблема вычислительной геометрии, которая привела к построению знаменитой триангуляции Делоне.

Существует много вариантов непрерывной интерполяции функций на Rn , каждый со своими преимуществами и недостатками. Например, если мы хотим расширить нашу стратегию ближайшего соседа, описанную выше, мы могли бы вычислить несколько ближайших соседей x и интерполировать f (x) на основе x —

хі2 для каждого ближайшего соседа хі . Определенные структуры данных «k-ближайших соседей» могут ускорить запросы, когда вы хотите найти несколько точек в наборе данных, ближайших к заданному х.

Другой стратегией, часто встречающейся в литературе по компьютерной графике, является барицентрическая интерполяция. Предположим, у нас есть ровно n + 1 выборочная точка (x1, y1), . . . ,(xn+1, yn+1), где xi Rn , и, как всегда, мы хотим интерполировать значения у на все Rn ; например, на плоскости нам дали бы три значения, связанные с вершинами треугольника. Любую точку x Rn можно однозначно записать в виде линейной комбинации x = i=1 aixi с дополнительным ограничением, что i ai = 1; другими словами, мы записываем x как средневзвешенное значение точек xi . Барицентрическая интерполяция в этом случае просто пишет f(x) = i ai(x)yi .

На плоскости R2 барицентрическая интерполяция представляет собой прямую геометрическую интерполяцию в областях с вращающимся треугольником, как показано на рисунке HOMEP. Кроме того, легко проверить, что полученная интерполированная функция f(x) аффинна, т. е. ее можно записать как f(x) = c + d · x для некоторых с R и d Rn

В общем, система уравнений, которую мы хотим решить для барицентрической интерполяции при некотором x = Rn:

При отсутствии вырождений эта система fora обратима при наличии n + 1 точек xi . Однако при наличии большего количества xi система для а становится недоопределенной. Это означает, что существует несколько способов записать данный x как средневзвешенное значение xi .

Одним из решений этой проблемы является добавление дополнительных условий на вектор усреднения весова. Результатом этой стратегии являются обобщенные барицентрические координаты, что является предметом исследований в области современной математики и инженерии. Типичные ограничения на а требуют, чтобы она была гладкой как функция на Rn и неотрицательной внутри множества xi , когда эти точки определяют многоугольник или многогранник. На рисунке NUMBER показан пример обобщенных барицентрических координат, вычисленных по точкам данных на многоугольнике с более чем n + 1 точкой.

Альтернативное решение недоопределенной проблемы для барицентрических координат связано с идеей использования кусочных функций для интерполяции; мы ограничим наше обсуждение здесь хі R2 для простоты, хотя расширения на более высокие измерения относительно очевидны.

Много раз нам давали не только множество точек xi, но и разложение интересующей нас области (в данном случае некоторое подмножество R2) на n + 1-мерные объекты, использующие эти точки в качестве вершин. Например, на рисунке НОМЕР показано такое разбиение части R2 на треугольники.

Интерполяция в этом случае проста: внутренняя часть каждого треугольника интерполируется с использованием барицентрических координат.

Пример 11.2 (Затенение). В компьютерной графике одним из наиболее распространенных представлений формы является набор треугольников в сетке. В повершинной модели затенения один цвет вычисляется для каждой вершины сетки. Затем, чтобы отобразить изображение на экране, эти значения для каждой вершины интерполируются с использованием барицентрической интерполяции к внутренней части треугольников. Подобные стратегии используются для текстурирования и других общих задач. На рисунке NUMBER показан пример этой простой модели затенения. Кроме того, одна проблема, специфичная для компьютерной графики, — это взаимодействие между перспективными преобразованиями и стратегиями интерполяции. Барицентрическая интерполяция цвета на трехмерной поверхности с последующим проецированием этого цвета на плоскость изображения отличается от проецирования треугольников на плоскость изображения с последующей интерполяцией цветов внутрь треугольника; таким образом, алгоритмы в этой области должны применять перспективную коррекцию, чтобы учесть эту ошибку.

Учитывая набор точек в R2, проблема триангуляции далеко не тривиальна, и алгоритмы для выполнения такого рода вычислений часто плохо распространяются на Rn. Таким образом, в более высоких измерениях предпочтительными становятся стратегии ближайшего соседа или регрессии (см. НОМЕР Главы).

Барицентрическая интерполяция приводит к обобщению кусочно-линейных шляпных функций из §11.1.3, показанных на рисунке НОМЕР. Напомним, что наш результат интерполяции полностью определяется значениями уі в вершинах треугольников. На самом деле, мы можем думать о f(x) как о линейной комбинации і уіфі(x), где каждая фі(x) — это кусочно-барицентрическая функция, полученная путем помещения 1 на хі и 0 во всех остальных местах, как на рис. . Эти треугольные функции шляпы составляют основу «метода конечных элементов первого порядка», который мы рассмотрим в следующих главах; специальные конструкции, использующие полиномы более высокого порядка, известные как «элементы более высокого порядка», могут использоваться для обеспечения дифференцируемости вдоль ребер треугольника.

Альтернативное и не менее важное разложение области определения f происходит, когда точки хі встречаются на регулярной сетке в Rn . Следующие примеры иллюстрируют ситуации, когда это так:

Пример 11.3 (Обработка изображения). Как упоминалось во введении, типичная цифровая фотография представлена в виде сетки m × n интенсивностей красного, зеленого и синего цветов. Мы можем думать об этих значениях как о живущих на решетке в Z × Z. Предположим, однако, что мы хотим повернуть изображение на угол, не кратный 90 . Затем, как показано на рисунке NUMBER, мы должны искать значения изображения в потенциально нецелочисленных позициях, требуя интерполяции цветов до R × R.

Пример 11.4 (Медицинская визуализация). Типичный вывод устройства магнитно-резонансной томографии (МРТ) представляет собой сетку значений ат х р, представляющих плотность ткани в разных точках; теоретически типичная модель этой функции такова: f:R3 R. Мы можем извлечь внешнюю поверхность конкретного органа, показанного на рисунке НОМЕР, найдя набор уровней {x:f(x) = c} для некоторого с. Чтобы найти этот набор уровней, нам нужно распространить f на всю сетку вокселей, чтобы точно найти, где он пересекает с.

Стратегии интерполяции на основе сетки обычно применяют одномерные формулы из §11.1.3 по одному измерению за раз. Например, схемы билинейной интерполяции в R2 линейно интерполируют одно измерение за раз, чтобы получить выходное значение:

Пример 11.5 (Билинейная интерполяция). Предположим, f принимает следующие значения:

- f(0, 0) = 1
- f(0, 1) = 3
- f(1, 0) = 5
- f(1, 1) = 11

а промежуточный f получается билинейной интерполяцией. Чтобы найти $\phi(\frac{1}{14}, \frac{7}{2})$, мы сначала интерполируем по х1 , чтобы найти:

$$\phi = \frac{1}{4}$$
, $31f(0, 0) + f$
-(1, 0) = 20 = 44

$$\frac{1}{4}$$
, $\frac{3 \ 1 \ f(0, 1) + f}{(1, 1) = 5 \ T = 4 \ 4}$

Далее интерполируем в х2:

$$\phi$$
 $\frac{1}{4}$, $\frac{1}{2}$ "=" $\frac{1}{2}$ ϕ $\frac{1}{4}$, ϕ 0 + $_2$ 4 $\frac{1}{4}$, 31 = -2

Важным свойством билинейной интерполяции является то, что мы получаем один и тот же результат, интерполируя сначала по x2 , а затем по x1.

Методы более высокого порядка, такие как бикубическая интерполяция и интерполяция Ланцоша, снова используют больше полиномиальных членов, но их вычисления выполняются медленнее. В частности, в случае интерполяции изображений бикубические стратегии требуют больше точек данных, чем квадрат ближайших к точке х значений функции; эти дополнительные затраты могут замедлить работу графических инструментов, для которых каждый поиск в памяти требует дополнительного времени вычислений.

11.3 Теория интерполяции

До сих пор наше отношение к интерполяции было довольно эвристическим. Полагаясь на нашу интуицию в отношении того, что «разумная» интерполяция для набора значений функции по большей части является приемлемой стратегией, могут возникнуть тонкие проблемы с различными методами интерполяции, которые важно признать.

11.3.1 Линейная алгебра функций

Мы начали наше обсуждение с того, что представили различные стратегии интерполяции в качестве различных базисов для набора функций f: R R. Эта аналогия с векторными пространствами распространяется на полную геометрическую теорию функций и, по существу, на ранние работы в области функционального анализа. расширяет геометрию Rn до наборов функций. Здесь мы обсудим функции одной переменной, хотя многие аспекты расширения до более общих функций легко осуществить.

Точно так же, как мы можем определить понятия размаха и линейной комбинации для функций, для фиксированных a, b R мы можем определить скалярное произведение функций f(x) и g(x) следующим образом:

ж, г
$$f(x)g(x) dx$$
.

Подобно тому, как А-внутреннее произведение векторов помогло нам вывести алгоритм сопряженных градиентов и имело много общего с скалярным произведением, функциональное скалярное произведение можно использовать для определения методов линейной алгебры для работы с пространствами функций и понимания их диапазона. Мы также определяем норму функции как f f, f.

Пример 11.6. Внутренний продукт функции Возьмем $pn(x) = {}^{H}$ быть n-м мономом. Тогда при a = 0 и xb = 1, мы имеем:

Обратите внимание, что это показывает:

Это значение приблизительно равно 1, когда n m, но n = m, что подтверждает наше предыдущее утверждение о том, что мономы значительно «перекрываются» на [0, 1].

Учитывая этот внутренний продукт, мы можем применить алгоритм Грама-Шмидта, чтобы найти ортонормированный базис для набора многочленов. Если мы возьмем а = 1 и b = 1, мы получим полиномы Лежандра, изображенные на рисунке HOMEP:

P0(x) = 1
P1(x) = x
P2(x) =
$$\frac{1}{2}(3x^2 - 1)$$

P3(x) = $(5x-21^3 - 3x)$
P4(x) = $\frac{1}{8}(35x^4 - 30x^2 + 3)$
 \vdots \vdots

Эти полиномы обладают многими полезными свойствами благодаря своей ортогональности. Например, предположим, что мы хотим аппроксимировать f(x) суммой i aiPi(x). Если мы хотим минимизировать f i aiPi по функциональной норме, это задача наименьших квадратов! В силу ортогональности базиса Лежандра для R[x] простое расширение наших методов проецирования показывает:

Таким образом, аппроксимация f c помощью полиномов может быть достигнута простым интегрированием f по элементам базиса Лежандра; в следующей главе мы узнаем, как этот интеграл может быть выполнен приближенно.

Для положительной функции w(x) мы можем определить более общий скалярный продукт $\cdot, \cdot w$, написав

f,
$$rB = w(x)f(x)g(x) dx$$
.

Если мы возьмем $w(x) = \frac{1}{x}$ с a = 1 и b = 1, то применение Грамма-Шмидта дает чебышевский

1 х 2 полиномы:

T0(x) = 1
T1(x) = x
T2(x) = 2x
2
 1
T3(x) = 4x 3 3x
T4(x) = 8x 4 8x 2 + 1
: :

На самом деле для этих многочленов выполняется удивительное тождество:

$$Tk(x) = cos(k arccos(x)).$$

Эту формулу можно проверить, явно проверив ее для Т0 и Т1, а затем индуктивно применив наблюдение:

Tk+1(x) =
$$cos((k + 1) arccos(x)) = 2x$$

 $cos(k arccos(x))$ $cos((k - 1) arccos(x))$ по тождеству $cos((k + 1)\theta)$
= $2 cos(k\theta) cos(\theta)$ $cos((k - 1)\theta) = 2xTk(x)$ Tk 1(x)

Эта формула «рекуррентности с тремя членами» также позволяет легко генерировать полиномы Чебышева.

Как показано на рисунке HOMEP, благодаря тригонометрической формуле для полиномов Чебышева легко увидеть, что минимумы и максимумы Tk колеблются между +1 и -1.

Кроме того, эти экстремумы расположены в точках соs(iπ/k) (так называемые «точки Чебышева») для і от 0 до k; это хорошее распределение экстремумов позволяет избежать колебательных явлений, подобных показанному на рис. НОМЕР, при использовании конечного числа полиномиальных членов для аппроксимации функции. На самом деле, более технический подход к полиномиальной интерполяции рекомендует размещать хі для интерполяции рядом с точками Чебышева для получения гладкого вывода.

11.3.2 Аппроксимация кусочными полиномами

Предположим, мы хотим аппроксимировать функцию f(x) полиномом степени n на отрезке [a, b]. Определим х как расстояние b а. Одна мера ошибки аппроксимации зависит от х, которая должна обращаться в нуль при х 0. Тогда, если мы аппроксимируем f кусочными полиномами, этот тип анализа говорит нам, как далеко друг от друга мы должны разнести полиномы, чтобы получить желаемый уровень приближения.

Например, предположим, что мы аппроксимируем f константой интерполяции $\frac{a+6}{2}$), как в кусочно-постоянной c = f(. Ecли предположить, что | <math>f(x)| < M для всех x [a, b], мы имеем:

Таким образом, мы ожидаем ошибки О(х) при использовании кусочно-постоянной интерполяции.

Предположим вместо этого, что мы аппроксимируем f, используя кусочно-линейную интерполяцию, то есть взяв

$$(x) = f(a) + \frac{6 - x f}{6 - a}$$
 $\frac{x - a}{6 - a}$ $\phi(6)$.

По теореме о среднем значении мы знаем, что $f(x) = f(\theta)$ для некоторого θ [a, b]. Разложение Тейлора относительно θ показывает, что $f(x) = f(\theta) + f(\theta)(x - \theta) + O(-x - x)$) на [a, b], а нашу линейку $f(x) = f(\theta) + f(\theta)(x - \theta)$ аппроксимация по можно переписать . Таким образом, вычитание этих двух выражений показывает, что). Нетрудно мере того, как ошибка аппроксимации f уменьшается до O^{-2} предсказать эту ошибку аппроксимации , хотя на практике (-х с полиномом степени n, делает O(-x) квадратичная сходимость линейных аппроксимаций достаточной для большинства приложений.

11.4 Проблемы

Идеи:

- Метод Хорнера для вычисления многочленов
- Рекурсивная стратегия для полиномиальных коэффициентов Ньютона.
- Сплайны, де Кастельжо
- Проверить интерполяцию площади треугольника барицентрической интерполяции

Глава 12

Численное интегрирование и Дифференциация

В предыдущей главе мы разработали инструменты для заполнения разумных значений функции f(x) с учетом выборки значений (xi , f(xi)) в области определения f . Очевидно, что эта интерполяционная проблема полезна сама по себе для завершения функций, о которых известно, что они непрерывны или дифференцируемы, но значения которых известны только в наборе изолированных точек, но в некоторых случаях мы затем хотим изучить свойства этих функций. В частности, если мы хотим применить инструменты исчисления к f , иметь возможность аппроксимировать его интегралы и мы должны производные.

На самом деле существует множество приложений, в которых ключевую роль в вычислениях играют численное интегрирование и дифференцирование. В самом простом случае некоторые известные функции определяются как интегралы. Например, «функция ошибки», используемая в качестве кумулятивного распределения кривой Гаусса или колоколообразной кривой, записывается:

erf(x)
$$\frac{2}{\pi}$$
 τ^2 ΔT

Аппроксимации erf (x) необходимы во многих статистических контекстах, и один из разумных подходов к нахождению этих значений состоит в том, чтобы выполнить приведенный выше интеграл численно.

В других случаях численные аппроксимации производных и интегралов являются частью более крупной системы. Например, методы, которые мы разработаем в будущих главах для аппроксимации решений дифференциальных уравнений, будут сильно зависеть от этих аппроксимаций. Точно так же в вычислительной электродинамике интегральные уравнения, решающие для неизвестной функции ф при заданном ядре К и выходе f, появляются в соотношении:

e (x) =
$$K(x,y)\phi(y)$$
 dy.

Эти типы уравнений должны быть решены для оценки электрических и магнитных полей, но если ϕ и K не являются очень особыми, мы не можем надеяться найти такой интеграл в закрытой форме, но решить это уравнение самостоятельно для неизвестной функции ϕ .

В этой главе мы разработаем различные методы численного интегрирования и дифференцирования на основе выборки значений функции. Эти алгоритмы обычно представляют собой довольно простые аппроксимации, поэтому для их сравнения мы также разработаем несколько стратегий, которые оценивают ожидаемую эффективность различных методов.

12.1 Мотивация

Нетрудно сформулировать простые приложения численного интегрирования и дифференцирования, учитывая, как часто инструменты исчисления появляются в основных формулах и методах физики, статистики и других областей. Здесь мы предлагаем несколько менее очевидных мест, где проявляются интеграция и дифференциация.

Пример 12.1 (выборка из дистрибутива). Предположим, что нам дано распределение вероятностей p(t) на интервале [0, 1]; то есть, если мы случайным образом выбираем значения в соответствии с этим распределением, мы ожидаем, что p(t) будет пропорциональна количеству раз, когда мы рисуем значение рядом с t. Распространенной задачей является генерация случайных чисел, распределенных как p(t).

Вместо того, чтобы разрабатывать специальный метод, чтобы делать это каждый раз, когда мы получаем новое значение p(t), можно сделать полезное наблюдение. Мы определяем кумулятивную функцию распределения р как

$$\Phi(\tau) = \int_0^{\tau} p(x) \, dx.$$

Тогда, если X — случайное число, равномерно распределенное в [0, 1], можно показать, что F 1 (X) распределено как р, где F 1 — обратное число F. Таким образом, если мы можем аппроксимировать F или F 1 мы можем генерировать случайные числа в соответствии с произвольным распределением p; это приближение сводится к интегрированию p, которое, возможно, придется выполнять численно, когда интегралы не известны в точной форме.

Пример 12.2 (Оптимизация). Напомним, что большинство наших методов минимизации и нахождения корней функции f зависело от наличия не только значений f(x), но и ее градиента f(x) и даже гессиана Hf. Мы видели, что такие алгоритмы, как BFGS и метод Бройдена, строят грубые приближения производных f в процессе оптимизации. Однако, когда f имеет высокие частоты, может быть лучше аппроксимировать f вблизи текущей итерации xk, а не использовать значения из потенциально далеких точек x для < k.

Пример 12.3 (рендеринг). Уравнение рендеринга из трассировки лучей и других алгоритмов для высококачественного рендеринга представляет собой интеграл, утверждающий, что свет, покидающий поверхность, равен интегралу света, попадающего на поверхность по всем возможным входящим направлениям после его отражения и рассеяния; по сути, в нем говорится, что световая энергия должна сохраняться до и после взаимодействия света с объектом. Алгоритмы рендеринга должны аппроксимировать этот интеграл, чтобы вычислить количество света, излучаемого поверхностью, отражающей свет в сцене.

Пример 12.4 (обработка изображения). Предположим, мы думаем об изображении как о функции двух переменных I(x, y). Многие фильтры, в том числе размытие по Гауссу, можно рассматривать как свертки, заданные формулой

$$(I g)(x, y) = I(u, v)g(x u, y v) du dv.$$

Например, чтобы размыть изображение, мы могли бы принять g за гауссиан; в этом случае (I g)(x, y) можно рассматривать как средневзвешенное значение цветов I вблизи точки (x, y). На практике изображения представляют собой дискретные сетки пикселей, поэтому этот интеграл необходимо аппроксимировать.

Пример 12.5 (правило Байеса). Предположим, что X и Y — непрерывнозначные случайные величины; мы можем использовать P(X) и P(Y) для выражения вероятности того, что X и Y принимают определенные значения. Иногда знание X может повлиять на наше знание Y. Например, если X — артериальное давление пациента, а Y — вес пациента,

тогда, зная, что у пациента большой вес, можно предположить, что у него также высокое кровяное давление. Таким образом, мы также можем записать условные распределения вероятностей P(X|Y) (читай «вероятность X при заданном Y»), выражающие такие отношения.

Основа современной теории вероятностей гласит, что P(X | Y) и P(Y | X) связаны следующим образом: P(Y | X)P(X)

$$P(X|Y) = \frac{}{P(Y|X)P(X) dY}$$

Оценка интеграла в знаменателе может быть серьезной проблемой в алгоритмах машинного обучения, где распределения вероятностей принимают сложные формы. Таким образом, приблизительные и часто рандомизированные схемы интеграции необходимы для алгоритмов выбора параметров, которые используют это значение как часть более крупного метода оптимизации.

12.2 Квадратура

Мы начнем с рассмотрения задачи численного интегрирования, или квадратуры. Эту проблему — с одной переменной — можно выразить так: «Дана выборка из n точек из некоторой функции f(x), f(x) dx». В предыдущем разделе мы найти приближение

представили несколько ситуаций, которые

уваривания именно к этой методике.

Есть несколько вариантов проблемы, которые требуют немного другого лечения или адаптации. ция:

- Конечные точки а и b могут быть фиксированными, или мы можем захотеть найти квадратурную схему, которая может эффективно аппроксимировать интегралы для многих пар (a, b) .
- Мы можем запросить f(x) при любом x, но хотим аппроксимировать интеграл, используя относительно небольшое количество выборок, или нам может быть предоставлен список предварительно вычисленных пар (xi , f(xi)) и мы вынуждены использовать эти данные. точек в нашем приближении.

Эти соображения следует иметь в виду, когда мы разрабатываем различные алгоритмы для квадратурной задачи.

12.2.1 Интерполяционная квадратура

Многие стратегии интерполяции, разработанные в предыдущей главе, могут быть расширены до методов квадратур с использованием очень простого наблюдения. Предположим, мы запишем функцию f(x) в терминах набора базисных функций фі(x):

$$f(x) = ai\phi i(x)$$
.

Тогда мы можем найти интеграл от f следующим образом:

$$f(x) dx = \begin{cases} 6 & \text{ if } f(x) dx = \\ a & \text{ if } f(x) dx = \\ a & \text{ if } f(x) dx = \\ &$$

Другими словами, интегрирование f просто включает в себя линейное комбинирование интегралов базисных функций, составляющих f.

Пример 12.6 (одночлены). Предположим, мы пишем f(x) = k akx K . Мы знаем

$$_{0}^{1}$$
 к дх = $_{---}^{1}$,

поэтому, применяя приведенный выше вывод, мы знаем

$$\int_{0}^{1} f(x) dx = \frac{a\kappa}{\kappa + 1}$$

Другими словами, в наших обозначениях выше мы определили $ck = \frac{1}{1K+1}$.

Схемы, в которых мы интегрируем функцию путем интерполяции отсчетов и интегрирования интерполированной функции, известны как интерполяционные квадратурные правила; почти все методы, которые мы представим ниже, могут быть написаны таким образом. Конечно, мы можем столкнуться с проблемой курицы и яйца, если сам интеграл фі(х) dx не известен в точной форме. Некоторые методы конечных элементов более высокого порядка решают эту проблему, затрачивая дополнительное вычислительное время на выполнение высококачественной численной аппроксимации интеграла одного фі , а затем, поскольку все ф имеют одинаковую форму, применяют формулы замены координат к напишите интегралы для остальных базисных функций. Этот канонический интеграл можно аппроксимировать в автономном режиме с помощью высокоточной схемы, а затем использовать пов

12.2.2 Правила квадратуры

Если нам дан набор пар (xi , f(xi)) , наше обсуждение выше предлагает следующую форму квадратурного правила для аппроксимации интеграла от f на некотором интервале:

Различные веса wi дают разные аппроксимации интеграла, которые, как мы надеемся, становятся все более похожими по мере того, как мы собираем более плотную выборку xi.

На самом деле, даже классическая теория интегрирования предполагает, что эта формула является разумной отправной точкой. Например, интеграл Римана, представленный во многих вводных классах исчисления, принимает форму:

a
$$f(x) = \lim_{\kappa \to 0} xk = 0$$
 $f(x^{\kappa} k)(xk+1) = xk$

Здесь интервал [a, b] разбит на куски $a = x1 < x2 < \cdots < xn = b$, где xk = xk+1 xk, $x^k = xk+1$ любая точка из [xk , xk+1]. Для фиксированного набора xk до предела этот интеграл явно может быть записан в форме Q[f] , приведенной выше.

С этой точки зрения выбор {xi} и {wi} полностью определяет стратегию квадратуры. Есть много способов определить эти значения, как мы увидим в следующем разделе и как мы уже видели для интерполяционной квадратуры.

Пример 12.7 (Метод неопределенных коэффициентов). Предположим, мы фиксируем x1, . . . , xn и хотят найти разумный набор сопутствующих весов wi , чтобы i wi f(xi) была подходящей аппроксимацией интеграла

выключенный . Альтернативой описанной выше стратегии базисной функции является использование метода неопределенных коэффициентов. В этой стратегии мы выбираем n функций f1(x), . . . , fn(x) , интегралы которых известны, и попросим, чтобы наше квадратурное правило точно восстановило интегралы этих функций:

Это создает линейную систему уравнений размера $n \times n$ для wi .

Один из распространенных вариантов - взять k 1, то есть убедиться, что квадратурная схема восстанавливается (x) = x интегралы от полиномов низкого порядка. Мы знаем

$$\kappa \times \Delta X = \frac{6 \kappa + 1 - a^{\kappa + 1}}{\kappa + 1}.$$

Таким образом, мы получаем следующую линейную систему уравнений для wi:

Эта система в точности является системой Вандермонда, обсуждавшейся в § 11.1.1.

12.2.3 Квадратура Ньютона-Котеса

Квадратура правит, когда правит х. _яs равномерно распределены в [a, b], известны как квадратуры Ньютона-Котеса Как показано на рис. НОМЕР, есть два разумных выбора равномерно расположенных образцов:

• Замкнутая квадратура Ньютона-Котеса помещает хі в точки а и b. В частности, при k {1, . . . , п} мы брать

xк a +
$$\frac{(\kappa - 1)(6 - a)}{n - 1}$$
.

• Открытая квадратура Ньютона-Котеса не помещает хі в точку а или b:

$$x\kappa = a + \frac{\kappa (6 - a)}{\pi + 1}$$

После этого выбора формулы Ньютона-Котеса просто применяют полиномиальную интерполяцию для аппроксимации интеграла от а до b; степень многочлена, очевидно, должна быть n - 1, чтобы квадратурное правило оставалось корректным.

В общем, мы будем держать n относительно небольшим. Таким образом, мы избегаем колебательных и шумовых явлений, возникающих при подгонке полиномов высокой степени к набору точек данных. Как и при кусочно-полиномиальной интерполяции, при интегрировании по большому интервалу [a, b] мы объединяем небольшие фрагменты в составные правила .

Закрытые правила. Замкнутые квадратурные стратегии Ньютона-Котеса требуют n 2, чтобы избежать деления на ноль. На практике часто встречаются две стратегии:

• Правило трапеций получается для n = 2 (так что x1 = a и x2 = b) путем линейной интерполяции от f(a) до f(b). В нем говорится, что

$$\int_{a}^{6} dx$$
 (b a) 2 $\frac{f(a) + f(b) f(x)}{a}$

• Правило Симпсона исходит из принятия n = 3, так что теперь мы имеем

$$x1 = a$$

$$x2 = \frac{a+6}{2}$$

$$x3 = 6$$

Интегрирование параболы, проходящей через эти три точки, дает

$$\int_{a}^{6} f(x) dx \qquad \frac{6 - a}{6} \qquad f(a) + 4 f \qquad \frac{a + 6}{2} \qquad + f(6) .$$

Открытые правила. Открытые правила квадратур допускают возможность n = 1, что дает упрощенное правило средней точки:

$$\int_{a}^{6} f(x) dx \quad (b \quad a)f \qquad \frac{a+6}{2} \quad \cdot$$

Большие значения п дают правила, подобные правилу Симпсона и правилу трапеций.

Композитная интеграция. Обычно нам может понадобиться интегрировать f(x) с более чем одним, двумя или тремя значениями xi. Очевидно, как построить составную линейку из средней точки или трапециевидной линейки, приведенной выше, как показано на рисунке HOMEP; просто суммируйте значения по каждому интервалу. Например, если мы разделим [a, b] на k интервалов, то мы можем взять x и xi $\frac{6-k}{a}$ + i x. Тогда составное правило средней точки:

Точно так же правило составных трапеций:

$$\int_{a}^{6} f(x) dx = \int_{\pi=1}^{\kappa} \frac{f(xi) + f(xi+1)}{2} \times \int_{\pi=1}^{\pi} \int_{\pi=1}^$$

путем разделения двух усредненных значений f в первой строке и переиндексации

Альтернативный подход к составному правилу средней точки состоит в применении интерполяционной квадратурной формулы из § 12.2.1 к кусочно-линейной интерполяции; точно так же составная версия правила трапеций получается из кусочно-линейной интерполяции.

Составная версия правила Симпсона, показанная на рисунке HOMEP, объединяет три точки одновременно, чтобы получить параболические приближения. Смежные параболы встречаются в точках хі с четными индексами и могут не иметь общих касательных. Это суммирование, которое существует только при четном п, принимает вид:

Точность. До сих пор мы разработали ряд квадратурных правил, которые эффективно комбинируют один и тот же набор f(xi) различными способами для получения различных приближений интеграла от f. Каждое приближение основано на другом инженерном предположении, поэтому неясно, лучше ли какое-либо из этих правил, чем какоелибо другое. Таким образом, нам необходимо разработать оценки ошибок, характеризующие их соответствующее поведение. Мы будем использовать наши интеграторы Newton-Cotes выше, чтобы показать, как можно проводить такие сравнения, как это представлено в СІТЕ.

Во-первых, рассмотрим квадратурное правило средней точки на одном интервале [a, b]. Определите с $\frac{1}{12}$ (a + б). Ряд Тейлора f о с равен:

1
$$f(x) = f(c) + f(c)(x - c) + f(c)(x - c) + 24$$
 1 2 $f(c)(x - c) + 6$ 1 3 $f(c)(x - c) + 24$ 4 + · · ·

Таким образом, в силу симметрии относительно с нечетные члены выпадают:

Обратите внимание, что первый член этой суммы в точности соответствует $\int_{a}^{6} f(x) dx$, обеспечиваемый средней точкой оценке правила, так что это правило точно до O($\frac{3}{2}$).

Теперь, подставив а и b в наш ряд Тейлора для f о c, мы получим:

$$f(a) = f(c) + f(c)(a - c) + 1 - f(c)(a - c) + 2 + (c)(a - c) + 6 + 1 + (c)(a - c) + (c)(a - c)(a - c) + (c)(a - c)(a - c) + (c)(a - c)(a - c)$$

Если сложить их вместе и умножить обе части на b а/2, получится:

Член f (c) обращается в нуль по определению с. Обратите внимание, что левая часть представляет собой интегральную оценку по правилу трапеций, а правая часть согласуется с нашим рядом Тейлора для f(x) dx с точностью до кубического члена. Другими словами, правило трапеций также O(x

Здесь мы остановимся, чтобы отметить поначалу неожиданный результат: правила трапеций и средней точки имеют один и тот же порядок точности! На самом деле, изучение члена третьего порядка показывает, что правило средней точки примерно в два раза точнее, чем правило трапеций. Этот результат кажется противоречащим здравому смыслу, поскольку правило трапеций использует линейную аппроксимацию, а правило средней точки является постоянным. Однако, как показано на рисунке НОМЕР, правило средней точки фактически восстанавливает интеграл линейных функций, что объясняет его дополнительную степень точности.

Аналогичный аргумент применим к нахождению оценки ошибки для правила Симпсона. [НАПИСАТЬ ОБЪЯСНЕНИЕ). 5 НАРОД ЗДЕСЬ; ИСКЛЮЧИТЬ ИЗ 205А]. В конце концов мы обнаруживаем, что правило Симпсона имеет ошибку типа O(x К такого рода анализу относится важное предостережение. В общем, теорема Тейлора применима только тогда, когда х достаточно мало. Если выборки находятся далеко друг от друга, тогда проявляются недостатки полиномиальной интерполяции, а колебательные явления, как обсуждалось в разделе НОМЕР, могут привести к нестабильным результатам для схем интегрирования высокого порядка.

Таким образом, возвращаясь к случаю, когда а и b далеко друг от друга, мы теперь разделим [а, b] на интервалы ширины х и применим любое из наших квадратурных правил внутри этих интервалов. Обратите внимание, что наше общее количество интервалов равно b a/ x, поэтому в этом случае мы должны умножить наши оценки ошибок на 1/ x . В частности, выполняются следующие порядки точности:

- Составная средняя точка: O(x
 Составная трапеция: O(x
- Композитный Симпсон: O(x ⁴)

12.2.4 Квадратура Гаусса

В некоторых приложениях мы можем выбирать местоположения xi, в которых f выбирается. В этом случае мы можем оптимизировать не только веса квадратурного правила, но и местоположения xi, чтобы получить максимальное качество. Это наблюдение приводит к сложным, но теоретически привлекательным квадратурным правилам.

Детали этой техники выходят за рамки нашего обсуждения, но мы предлагаем один простой способ ее получения. В частности, как и в примере 12.7, предположим, что мы хотим оптимизировать x1, ..., xn и w1, ..., wn одновременно для повышения порядка схемы интегрирования. Теперь у нас есть 2n вместо n известных, поэтому мы можем обеспечить равенство для 2n примеров:

```
6

f1(x) dx = w1 f1(x1) + w2 f1(x2) + ··· + wn f1(xn)

a

f2(x) dx = w1 f2(x1) + w2 f2(x2) + ··· + wn f2(xn)

\vdots \qquad \vdots

6

f2n(x) dx = w1 fn(x1) + w2 fn(x2) + ··· + wn fn(xn)
```

Теперь и хі , и wi неизвестны, поэтому эта система уравнений больше не является линейной. Например, если мы хотим оптимизировать эти значения для многочленов на интервале [1, 1] , мы должны

необходимо решить следующую систему полиномов (CITE):

Может случиться так, что системы, подобные этой, имеют несколько корней и другие вырождения, которые зависят не только от выбора fi (обычно многочленов), но и от интервала, на котором мы аппроксимируем интеграл. Кроме того, эти правила не являются прогрессивными в том смысле, что набор хi для n точек данных не имеет ничего общего с набором для k точек данных, когда k = n, поэтому повторно использовать данные для получения более точной оценки сложно. С другой стороны, когда они применимы, квадратура Гаусса имеет наивысшую возможную степень при фиксированном n. Квадратурные правила Кронрода пытаются избежать этой проблемы, оптимизируя квадратуру с 2n + 1 точками при повторном использовании точек Гаусса.

12.2.5 Адаптивная квадратура

Как мы уже показали, существуют некоторые функции f, интегралы которых лучше аппроксимируются данным квадратурным правилом, чем другие; например, правила средней точки и правила трапеций интегрируют линейные функции с полной точностью, в то время как при быстрых колебаниях f могут возникнуть проблемы с дискретизацией и другие проблемы.

Напомним, что квадратурное правило Гаусса предполагает, что размещение хі может влиять на качество квадратурной схемы. Однако осталась одна часть информации, которую мы не использовали: значения f(xi). Ведь именно они определяют качество нашей квадратурной схемы.

Имея это в виду, адаптивные квадратурные стратегии исследуют текущую оценку и генерируют новые хі , где подынтегральное выражение более сложное. Стратегии адаптивной интеграции часто сравнивают результаты нескольких квадратурных методов, например трапеции и средней точки, с предположением, что они согласуются там, где достаточно выборки f (см. НОМЕР РИСУНКА). Если они не согласуются с какимлибо допуском на заданном интервале, генерируется дополнительная точка выборки и обновляются интегральные оценки.

ДОБАВИТЬ ПОДРОБНЕЕ ИЛИ ПРИМЕР; ОБСУЖДЕНИЕ РЕКУРСИВНОГО АЛГОРИТМА; ГАНДЕР И ГАУЧИ

12.2.6 Множественные

переменные Много раз мы хотели бы интегрировать функции f(x), где x Rn. Например, когда n = 2, мы можем интегрировать по прямоугольнику, вычисляя

В более общем смысле, как показано на рис. ЧИСЛО $_{\rm i}$, мы можем захотеть найти интеграл, где $_{_{\rm om}}$ f(x) dx, Ω — некоторое подмножество Rn $_{\rm om}$

«Проклятие размерности» экспоненциально усложняет интеграцию по мере увеличения размерности. В частности, количество выборок f, необходимых для достижения сравнимой квадратурной точности для интеграла в Rk, увеличивается как O(n). Это наблюдение может быть обескураживающим, но в некоторой степени разумным: чем больше входных измерений для f, тем больше выборок требуется, чтобы понять его поведение во всех измерениях.

Простейшей стратегией интегрирования в Rk является интегрированный интеграл. Например, если f — функция Для двух переменных предположим, что мы хотим найти f(x,y) dx dy. Для фиксированного у мы можем аппроксимировать внутренний интеграл по x, используя одномерное квадратурное правило; затем мы интегрируем эти значения по y, используя другое квадратурное правило. Очевидно, что обе схемы интегрирования вызывают некоторую ошибку, поэтому нам может потребоваться выборка x более плотно, чем в одном измерении, для достижения желаемого качества вывода.

В качестве альтернативы, точно так же, как мы разделили [a, b] на интервалы, мы можем подразделить Ω на треугольники и прямоугольники в 2D, многогранники или прямоугольники в 3D и т. д. и использовать простые интерполяционные квадратурные правила в каждой части. Например, одним из популярных вариантов является интегрирование результатов барицентрической интерполяции внутри многогранников, поскольку этот интеграл известен в закрытой форме.

Однако, когда n велико, нецелесообразно делить домен, как предлагается. В этом случае мы можем использовать рандомизированный метод Монте-Карло. В этом случае мы просто генерируем k случайных точек xi Ω , например, c равномерной вероятностью. Усреднение значений f(xi) дает аппроксимацию f(x) dx, которая сходится как 1/ k – $^{\circ}$ независимо от размерности Ω ! Так, в больших размерах

оценка Монте-Карло предпочтительнее описанных выше детерминированных квадратурных методов.

ПОДРОБНЕЕ О КОНВЕРГЕНЦИИ МОНТЕ-КАРЛО И ВЫБОРЕ РАСПРЕДЕЛЕНИЙ

БОЛЕЕ Ом

12.2.7 Кондиционирование

До сих пор мы рассматривали качество квадратурного метода с использованием значений точности O(х,); очевидно, что по этой метрике предпочтительнее набор квадратурных весов с большими k.

Однако другая мера уравновешивает измерения точности, полученные с использованием аргументов Тейлора. В частности, напомним, что мы записали наше квадратурное правило как Q[f] $\hat{f}(x)$ $\hat{f}(x)$. Предположим, мы возмущаем f некоторым другим f . Определить f $\hat{f}(x)$ \hat{f}

Таким образом, устойчивость или обусловленность квадратурного правила зависит от нормы набора весов ш .

В общем, легко проверить, что по мере увеличения порядка квадратурной точности условие w ухудшается, поскольку wi принимает большие отрицательные значения; это контрастирует со случаем всех положительных значений, где обусловленность ограничена b - a, потому что i wi = b - a для полиномиальных интерполяционных схем, а большинство методов низкого порядка имеют только положительные коэффициенты (СНЕСК). Этот факт является отражением той же интуиции, что мы не должны интерполировать функции, используя многочлены высокого порядка. Таким образом, на практике мы обычно предпочитаем составные квадратуры методам высокого порядка, которые могут давать более точные оценки, но могут быть неустойчивыми при численном возмущении.

12.3 Дифференциация

Численное интегрирование является относительно стабильной задачей. в том, что влияние любого отдельного значения $f(x_3^6)$ на f(x) dx уменьшается до нуля, когда а и b становятся далеко друг от друга. С другой стороны, аппроксимация производной функции f(x) таким свойством устойчивости не обладает. С точки зрения анализа Фурье можно показать, что интеграл f(x) обычно имеет более низкие частоты, чем f, в то время как дифференцирование для получения f усиливает высокие частоты f, что делает ограничения выборки, условия и стабильность особенно сложными для аппроксимации f.

Несмотря на сложные обстоятельства, аппроксимации производных обычно относительно легко вычислить, и они могут быть устойчивыми в зависимости от рассматриваемой функции. На самом деле, при разработке правила секущих, метода Бройдена и т. д. мы использовали простые аппроксимации производных и градиентов, чтобы помочь в процедурах оптимизации.

Здесь мы сосредоточимся на аппроксимации f для f: R R. Нахождение градиентов и якобианов часто достигается путем дифференцирования в одном измерении за раз, что эффективно сводится к одномерной задаче, которую мы здесь рассматриваем.

12.3.1 Дифференциация базовых функций

Простейший случай дифференцирования подходит для функций, построенных с использованием подпрограмм интерполяции. Как и в § 12.2.1, если мы можем записать f(x) = i aiфi(x), то по линейности мы знаем

Другими словами, мы можем думать о функциях ϕ і как о базисе для производных функций, записанных в базисе ϕ і!

Пример этой процедуры показан на рисунке НОМЕР. Это явление часто связывает разные интерполяционные схемы. Например, кусочно-линейные функции имеют кусочно-постоянные производные, полиномиальные функции имеют полиномиальные производные меньшей степени и т. д.; мы вернемся к этой структуре, когда будем рассматривать дискретизацию уравнений в частных производных. Между тем, в этом случае полезно знать, что f известно с полной уверенностью, хотя, как и на рисунке ЧИСЛО, его производные могут иметь нежелательные разрывы.

12.3.2 Конечные разности

Более распространен случай, когда у нас есть функция f(x), которую мы можем запросить, но производные которой неизвестны. Это часто происходит, когда f принимает сложную форму или когда пользователь предоставляет f(x) как подпрограмму без аналитической информации о ее структуре.

Определение производной предлагает разумный подход:

$$\lim h \qquad \frac{f(x+h)-f(x)f(x)}{}$$

Как и следовало ожидать, для конечного h > 0 с малым |h| выражение в пределе обеспечивает возможное значение, приближающееся к f(x).

Чтобы подтвердить эту интуицию, мы можем использовать ряд Тейлора, чтобы написать:

1
$$f(x + h) = f(x) + f(x)h + f(x)h 2^{-1}$$
 2 + · · ·

Преобразование этого выражения показывает:

$$f(x) = \frac{f(x + h) - f(x) f}{f(x + h) - f(x) f} + O(4) 4$$

Таким образом, следующая прямая разностная аппроксимация f имеет линейную сходимость:

$$f(x+h) \qquad f(x) f(x)$$

Точно так же изменение знака h показывает, что обратные различия также имеют линейную сходимость:

$$f(x) = \begin{cases} \frac{f(x)}{h} & f(x - h) \end{cases}$$

На самом деле мы можем улучшить сходимость нашего приближения, используя хитрость. Тейлор теорему мы можем написать:

$$f(x+h) = f(x) + f(x)h + \frac{1}{2} \phi(x) + \frac{1}{6} \phi(x) + \frac{3}{6} \phi(x) + \frac{3}{6} \phi(x) + \frac{1}{6} \phi($$

Таким образом, эта центрированная разность дает приближение f (x) с квадратичной сходимостью; это высший порядок сходимости, которого мы можем ожидать с разделенной разностью. Однако мы можем добиться большей точности, оценив f в других точках, например, x + 2h, хотя на практике это приближение редко используется в пользу простого уменьшения h.

Построение оценок производных высших порядков может происходить по аналогичным построениям. Для Например, если мы сложим разложения $Teйлopa\ f(x+h)$ и f(x-h), мы увидим

$$f(x + h) + f(x - h) = 2 f(x) + f(x)h f(x^{2} + O(4^{3}))$$

$$= \frac{h + h - 2 f(x) + f(x - h) = f(x)}{h + O(4 + 2)} + O(4 + 2^{2})$$

Чтобы предсказать подобные комбинации для высших производных, нужно заметить, что наша вторая формулу производной можно разложить по-разному:

$$\frac{f(x+h)-2\ f(x)+f(x-h)\ h\ 2}{(x+h)-2\ f(x)+f(x-h)\ h} = \frac{f(x+h)-f(x)\ h}{(x+h)-f(x)+f(x-h)\ h}$$

То есть наша аппроксимация второй производной представляет собой «конечную разность конечных разностей». Один из способов интерпретации этой формулы показан на рисунке HOMEP. Когда мы вычисляем аппроксимацию прямой разности f между x и x + h, мы можем думать об этом наклоне как о живущем при x + h/2; мы аналогичным образом можем использовать обратные разности, чтобы поместить наклон в x - h/2. Нахождение наклона между этими значениями возвращает приближение к x.

Одной из стратегий, которая может улучшить сходимость приведенных выше приближений, является экстраполяция Ричардсона. В качестве примера более общего шаблона предположим, что мы хотим использовать прямые разности для аппроксимации f. Определять

D (4)
$$\frac{e(x + 4) - e(x)}{4}$$
.

Очевидно, что D(h) приближается к f (x) при h $\,$ 0. Однако, более конкретно, из нашего обсуждения в § 12.3.2 мы знаем, что D(h) принимает форму:

$$D(h) = f(x) + \frac{1}{2}f(x)h + O(h^{-2})$$

Предположим, что мы знаем D(h) и $D(\alpha h)$ для некоторого $0 < \alpha < 1$. Мы знаем:

$$D(\alpha h) = f(x) + \frac{1}{2}f(x)\alpha h + O(h^{-2})$$

Мы можем записать эти два отношения в матрицу:

1
$$\frac{1}{2}$$
 $\frac{1}{2}$ \times (x) \times (y) \times (x) \times (x) (x) \times (x) (x) \times (x) (x) \times (x) \times (x) (x) \times (x) (x) \times (x) (x)

Или, что то же самое,

$$x(x)$$
 $x = 1$ $x =$

То есть мы взяли аппроксимацию f(x) за O(h) с помощью D(h) и превратили ее в (x) ок аппроксимацию за O(h)! аппроксимация D(h). Тот же прием можно применить и ко многим другим задачам, записав аппроксимацию D(h) = a + bhn + O(hm), где m > n, где a - m величина, которую мы надеемся оценить, и b - следующий член в разложении Тейлора. На самом деле, даже экстраполяцию Ричардсона можно применять рекурсивно для получения приближений все более и более высокого порядка.

12.3.3 Выбор размера шага

В отличие от квадратуры, численное дифференцирование обладает любопытным свойством. Оказывается, что любой выбранный нами метод может быть сколь угодно точным, просто выбрав достаточно малое значение h. Это наблюдение привлекательно с точки зрения того, что мы можем добиться более качественных приближений без дополнительного времени вычислений. Загвоздка, однако, в том, что мы должны делить на h и сравнивать все больше и больше похожих значений f(x) и f(x + h); в арифметике с конечной точностью сложение и/или деление на почти нулевые значения вызывает числовые проблемы и нестабильность. Таким образом, существует диапазон значений h, которые недостаточно велики, чтобы вызвать значительную ошибку дискретизации, и недостаточно малы, чтобы создавать численные проблемы; На рисунке NUMBER показан пример дифференциации простой функции в арифметике с плавающей запятой IEEE.

12.3.4 Интегрированные количества

Не включено в CS 205A, осень 2013 г.

12.4 Проблемы

- Квадратура Гаусса всегда содержит средние точки, стратегия с использованием ортогональных полиномов
- Адаптивная квадратура
- Применение экстраполяции Ричардсона в других местах

Глава 13

Обыкновенные дифференциальные уравнения

Мы мотивировали проблему интерполяции в главе 11, перейдя от анализа к поиску функций. То есть в таких задачах, как интерполяция и регрессия, неизвестное — это функция f, а задача алгоритма — заполнить недостающие данные.

Мы продолжим это обсуждение, рассмотрев аналогичные задачи, связанные с заполнением значений функций. Здесь наше неизвестное по-прежнему является функцией f, но вместо того, чтобы просто угадывать пропущенные значения, мы хотели бы решить более сложные задачи проектирования. Например, рассмотрим следующие проблемы:

- Найти f, аппроксимирующую некоторую другую функцию f0 , но удовлетворяющую дополнительным критериям (гладкость, непрерывность, низкая частота и др.).
- ullet Смоделируйте некоторую динамическую или физическую взаимосвязь как f(t), где t время.
- Найдите f со значениями, аналогичными f0 , но с некоторыми общими свойствами с другой функцией. гO.

В каждом из этих случаев наша неизвестная является функцией f, но наши критерии успеха более сложны, чем «соответствие заданному набору точек данных».

Теории обыкновенных дифференциальных уравнений (ОДУ) и уравнений в частных производных (УЧП) изучают случай, когда мы хотим найти функцию f (x) на основе информации о ее производных или взаимосвязях между ними. Обратите внимание, что мы уже решили простую версию этой проблемы при обсуждении квадратуры: при заданном f (t) методы квадратуры обеспечивают способы аппроксимации f (t) с помощью интегрирования.

В этой главе мы рассмотрим случай обыкновенных дифференциальных уравнений и, в частности, начальные задачи. Здесь неизвестной является функция f(t): R Rn, и мы дали уравнение, которому удовлетворяют f,и ее производные, а также f(0); наша цель — предсказать f(t) для t > 0. Мы приведем несколько примеров ОДУ, встречающихся в литературе по информатике, а затем перейдем к описанию общих методов решения.

В качестве примечания мы будем использовать обозначение f для обозначения производной d f/dt от f : [0,) Rn .Наш цель будет состоять в том, чтобы найти f(t) с заданными отношениями между t, f(t), f (t), f (t) и так далее.

13.1 Мотивация

ОДУ появляются практически в любой части научного примера, и нетрудно встретить практические ситуации, требующие их решения. Например, основные законы физического движения задаются ОДУ:

Пример 13.1 (второй закон Ньютона). Продолжая §5.1.2, вспомним, что Второй закон Ньютона гласит, что F = ma, то есть общая сила, действующая на объект, равна произведению его массы на его ускорение.

Если мы моделируем n частиц одновременно, то можно представить себе объединение всех их положений в вектор х R3n Точно так же мы можем написать функцию F(t,x,x) R3n, которая принимает время, положения частиц и их скорости и возвращает общую силу, действующую на каждую частицу. Эта функция может учитывать взаимосвязь между частицами (например, гравитационные силы или пружины), внешние эффекты, такие как сопротивление ветра (которое зависит от x), внешние силы, изменяющиеся во времени t, и так далее.

Затем, чтобы найти положение всех частиц в зависимости от времени, мы хотим решить уравнение x = F(t,x,x)/m. Обычно нам даются положения и скорости всех частиц в момент времени t = 0 в качестве начального условия.

Пример 13.2 (Сворачивание белков). В меньшем масштабе уравнения, управляющие движением молекул, также являются обыкновенными дифференциальными уравнениями. Одним из особенно сложных случаев является укладка белка, в которой геометрическая структура белка предсказывается путем моделирования межмолекулярных сил с течением времени. Эти силы часто принимают нелинейные формы, которые продолжают бросать вызов исследователям в области вычислительной биологии.

Пример 13.3 (Градиентный спуск). Предположим, мы хотим минимизировать функцию энергии Е(x) по всем x. В главе 8 мы узнали, что Е(x) указывает в направлении, в котором E больше всего уменьшается при заданном x, поэтому мы провели линейный поиск вдоль этого направления от x, чтобы минимизировать E локально. Альтернативный вариант, популярный в некоторых теориях, состоит в том, чтобы решить ОДУ формы x = - E(x); другими словами, думайте о x как о функции времени x(t), которая пытается уменьшить E, спускаясь с горы.

Например, предположим, что мы хотим решить Ax =b для симметричного положительно определенного A. Мы знаем из что это эквивалентно минимизации E(x) $x = f(x) - \frac{5}{12}$ 10.1.1 x Ax bx + c. Таким образом, мы могли бы попытаться решить ОДУ =b Ax. При t мы ожидаем, что x(t) будет все лучше и лучше удовлетворять линейной системе.

Пример 13.4 (моделирование толпы). Предположим, мы пишем программное обеспечение для видеоигр, требующее реалистичного моделирования виртуальных скоплений людей, животных, космических кораблей и т.п. Одна из стратегий создания правдоподобного движения, показанная на рисунке НОМЕР, заключается в использовании дифференциальных уравнений. Здесь скорость члена толпы определяется как функция его окружения; например, в толпе людей близость других людей, расстояние до препятствий и т. д. могут влиять на направление движения данного агента. Эти правила могут быть простыми, но в совокупности их взаимодействие сложно. В основе этого механизма лежат стабильные интеграторы для дифференциальных уравнений, поскольку мы не хотим иметь заметно нереалистичное или нефизическое поведение.

13.2 Теория ОДУ

Полное рассмотрение теории обыкновенных дифференциальных уравнений выходит за рамки нашего обсуждения, и мы отсылаем читателя к СІТЕ за более подробной информацией. Помимо этого, мы упомянем здесь некоторые основные моменты, которые будут иметь отношение к нашей разработке в следующих разделах.

13.2.1 Основные понятия

Наиболее общая проблема начального значения ОДУ принимает следующую форму:

Здесь F — некоторая связь между f и всеми ее производными; мы используем f () для обозначения -th производной от f . Мы можем думать об ОДУ как об определяющих эволюцию f во времени t; мы знаем f и ее производные в нулевое время и хотим предсказать их в будущем.

ОДУ принимают разные формы даже в одной переменной. Например, обозначим у = f(t) и предположим, что у R1 . Затем примеры ОДУ включают:

- у = 1 + cos t: Это ОДУ можно решить путем интегрирования обеих частей, например, с помощью квадратурного уравнения. МЕТОДЫ
- у = ау: это ОДУ линейно по у
- у = ау + е т: ОДУ зависит от времени и положения.
- у + 3у у = t: это ОДУ включает несколько производных от у
- у грех у = е ты : Это ОДУ нелинейно по у и t.

Очевидно, что решить самые общие ОДУ может быть непросто. Мы ограничим большую часть нашего обсуждения случаем явных ОДУ, в которых можно выделить производную высшего порядка:

Определение 13.1 (явное ОДУ). ОДУ является явным, если его можно записать в виде

$$f(k)(t) = F[t, f(t), f(t), f(t), ..., f(k 1)(t)].$$

Например, явная форма второго закона Ньютона: $x(t) = \frac{1}{M} a(T,X(T),X(T))$.

Удивительно, но, обобщая прием из \$5.1.2, на самом деле любое явное ОДУ можно преобразовать в уравнение первого порядка f (t) = F[t, f(t)], где f имеет многомерный выход. Это наблюдение подразумевает, что нам нужно не более одной производной при рассмотрении алгоритмов ОДУ. Чтобы увидеть эту связь, мы просто вспомним, что d 2y/dt2 = d/dt(dy/dt). Таким образом, мы можем определить промежуточную переменную z dy/dt и понимать d 2y/dt2 как dz/dt с ограничением z = dy/dt. В более общем случае, если мы хотим решить явную задачу

$$f(k)(t) = F[t, f(t), f(t), f(t), ..., f(k 1)(t)],$$

где f:R Rn , то мы определяем g(t):R Rkn с помощью ОДУ первого порядка:

Здесь мы обозначаем, что gi(t): R Rn содержит n компонент g. Тогда g1(t) удовлетворяет исходному ОДУ. Чтобы убедиться в этом, мы просто проверим, что из приведенного выше уравнения следует, что g2(t) = g1(t), g3(t) = g2(t) = g1(t) и так далее. Таким образом, выполнение этих замен показывает, что последняя строка кодирует исходный ODE.

Приведенный выше прием упростит наши обозначения, но следует соблюдать некоторую осторожность, чтобы понять, что этот подход не упрощает вычисления. В частности, во многих случаях наша функция f(t) будет иметь только один выход, а ОДУ будет состоять из нескольких производных. Заменим этот случай одной производной и несколькими выходами.

Пример 13.5 (расширение ОДУ). Предположим, мы хотим решить y = 3y 2y + y, где y(t) : R R. Это уравнение эквивалентно:

Подобно тому, как наш вышеприведенный трюк позволяет нам рассматривать только ОДУ первого порядка, мы можем еще больше ограничить наши обозначения автономными ОДУ. Эти уравнения имеют вид f(t) = F[f(t)], т. е. F(t) = F[f(t)] больше не зависит от t. Для этого мы могли бы определить

$$\Gamma (T)$$

$$\frac{f(t)}{g^{-}(t)}$$

Тогда вместо этого мы можем решить следующее ОДУ для g:

$$\Gamma (\tau) = \begin{cases} f(\tau) & \text{F[} f(t), g^{-}(t)] \\ g^{-}(\tau) & 1 \end{cases}$$

В частности, $g^{-}(t) = t$, если принять $g^{-}(0) = 0$.

Поведение ОДУ можно визуализировать разными способами, как показано на рисунке НОМЕР. Например, если неизвестная функция f(t) является функцией одной переменной, то мы можем думать о F[f(t)] как о характеристике наклона f(t), как показано на рисунке НОМЕР. В качестве альтернативы, если f(t) имеет выход в R2 , , мы больше не можем визуализировать зависимость от времени t, но можем нарисовать фазовое пространство, которое показывает тангенс f(t) при каждом (x, y) R2

13.2.2 Существование и уникальность

Прежде чем мы сможем перейти к дискретизации задачи о начальных значениях, мы должны кратко признать, что не все задачи дифференциальных уравнений разрешимы. Кроме того, некоторые дифференциальные уравнения допускают несколько решений.

Пример 13.6 (Неразрешимое ОДУ). Рассмотрим уравнение y = 2y/t с заданным y(0) = 0; обратите внимание, что мы не делим на ноль, потому что y(0) предписывается. Переписать как

$$\frac{1}{v} \frac{dy}{dt} = \frac{2}{\tau}$$

и интегрирование по t с обеих сторон показывает:

$$пер |y| = 2 ln t + c Или,$$

что то же самое, у = Ct2 для некоторого С R. Обратите внимание, что у(0) = 0 в этом выражении, что противоречит нашим начальным условиям. Таким образом, это ОДУ не имеет решения при заданных начальных условиях.

Пример 13.7 (неединственные решения). Теперь рассмотрим то же ОДУ с y(0) = 0. Рассмотрим y(t), заданное равенством y(t) = Ct2 для любого С R. Тогда y(t) = 2Ct. Таким образом,

$$\frac{2\Gamma}{T}$$
 "=" $\frac{2KT2}{T}$ = 2Ct = y (t),

показывая, что ОДУ решается этой функцией независимо от С. Таким образом, решения этой задачи неединственны.

К счастью, существует богатая теория, характеризующая поведение и устойчивость решений дифференциальных уравнений. Наше развитие в следующей главе будет иметь более сильный набор условий, необходимых для существования решения, но на самом деле при слабых условиях на f можно показать, что ОДУ f (t) = F[f (t)] имеет решение. Например, одна из таких теорем гарантирует локальное существование решения:

Теорема 13.1 (Локальное существование и единственность). Предположим, что F непрерывна и липшицева, т. е. F[y] F[x]2 Ly x2 для некоторого L. Тогда ОДУ f(t) = F[f(t)] допускает ровно одно решение для всех t = 0 независимо от начальных условий.

В нашем последующем развитии мы будем предполагать, что ОДУ, которое мы пытаемся решить, удовлетворяет условиям такой теоремы; это предположение довольно реалистично в том смысле, что, по крайней мере локально, должно быть достаточно вырожденное поведение, чтобы разрушить такие слабые предположения.

13.2.3 Уравнения модели

Один из способов понять поведение ОДУ состоит в том, чтобы изучить поведение решений некоторых простых модельных уравнений, которые можно решить в замкнутой форме. Эти уравнения представляют собой линеаризацию более практичных уравнений и, таким образом, локально моделируют тип поведения, который мы можем ожидать.

Начнем с ОДУ в одной переменной. Учитывая наши упрощения в § 13.2.1, простейшее уравнение, с которым мы могли бы работать, было бы y = F[y], где y(t) : R R. Использование линейного приближения дало бы уравнения типа y = ay + b. Подстановка y = y + b показывает: y = y = ay + b = a(y = b/a) + b = ay = ay. Таким образом, в наших модельных уравнениях константа b = 0.

Согласно приведенному выше аргументу, мы локально можем понять поведение y = F[y] , изучая линейную уравнение y = ay. На самом деле, применение стандартных рассуждений исчисления показывает, что

$$y(t) = Ceat$$
.

Очевидно, что есть три случая, показанные на рис. НОМЕР:

- 1. a > 0: в этом случае решения становятся все больше и больше; на самом деле, если y(t) и $y^{\hat{}}(t)$ оба удовлетворяют ОДУ с немного разными начальными условиями, при t они расходятся.
- 2. а = 0: система в этом случае решается константами; решения с разными начальными точками остаются на одном и том же расстоянии друг от друга.
- 3. a < 0: Тогда все решения ОДУ стремятся к 0 при t

Мы говорим, что случаи 2 и 3 стабильны в том смысле, что незначительное возмущение у(0) дает решения, которые становятся все ближе и ближе с течением времени; случай 1 нестабилен, так как небольшая ошибка в задании входного параметра у(0) будет усиливаться с течением времени t. Нестабильные ОДУ порождают некорректные вычислительные проблемы; без тщательного рассмотрения мы не можем ожидать, что численные методы генерируют полезные решения в этом случае, поскольку даже теоретические выходные данные очень чувствительны к возмущениям входных данных. С другой стороны, устойчивые задачи хорошо поставлены, поскольку небольшие ошибки в у(0) со временем уменьшаются.

Переходя к нескольким измерениям, мы могли бы изучить линеаризованное уравнение

$$y = Ay$$
.

Как объяснялось в §5.1.2, если у1, \cdots , уk — собственные векторы оператора A с собственными значениями λ 1, \ldots , λ k и у(0) тогда c1у1 + \cdots + ckyk, $\stackrel{=}{}$

$$y(\tau) = c1e$$
 $\lambda 1t \lambda k t y 1 + \cdots + cke$

Другими словами, собственные значения A занимают место а в нашем одномерном примере. Из этого результата нетрудно интуитивно понять, что многомерная система устойчива именно тогда, когда ее спектральный радиус меньше единицы.

В действительности мы хотим решить у = F[y] для общих функций F. Предполагая, что F дифференцируема, мы можем написать F[y] F[y0] + JF(y0)(y y0), что дает приведенное выше модельное уравнение после смены. Таким образом, на коротких промежутках времени мы ожидаем поведение, подобное уравнению модели. Кроме того, условия в теореме 13.1 можно рассматривать как ограничение поведения JF, обеспечивающее связь с менее локализованными теориями ОДУ.

13.3 Схемы временного шага

Ключевое значение для нашего рассмотрения имеет идея стабильности. Точно так же, как ОДУ могут быть стабильными или нестабильными, так же могут быть и дискретизации. Например, если h слишком велико, некоторые схемы будут накапливать ошибки с экспоненциальной скоростью; Напротив, другие методы устойчивы в том смысле, что даже если h велико, решения останутся ограниченными. Однако стабильность может конкурировать с точностью; часто схемы, устойчивые во времени, являются плохим приближением у(t), даже если они гарантированно не ведут себя дико.

13.3.1 Форвард Эйлер

Наша первая стратегия ODE исходит из нашего построения схемы прямого дифференцирования в §12.3.2:

F[yk] = y(t) =
$$\frac{y_K+1 - y_K}{}$$
 + O(4) 4

Решение этой зависимости для показов yk+1

$$yk+1 = yk + hF[yk] + O(h^{2})$$
 $yk + hF[yk].$

Таким образом, прямая схема Эйлера применяет формулу справа для оценки yk+1. Это одна из наиболее эффективных стратегий временного шага, поскольку она просто вычисляет F и прибавляет результат к yk. По этой причине мы называем его явным методом, т. е. существует явная формула для yk+1 через yk и F.

Проверить точность этого метода довольно просто. Обратите внимание, что наше приближение), поэтому каждый ук +1 равно O(h ² шаг вызывает квадратичную ошибку. Мы называем эту ошибку локализованной ошибкой усечения, потому что это ошибка, вызванная одним шагом; слово «усечение» относится к тому факту, что мы усекли ряд Тейлора, чтобы получить эту формулу. Конечно, наш итератейк уже может быть неточным из-за накопленных ошибок усечения от предыдущих итераций. Если мы проинтегрируем от t0 до t с шагами O(1/h), то наша общая ошибка будет выглядеть как O(h); эта оценка представляет собой глобальную ошибку усечения, и поэтому мы обычно пишем, что прямая схема Эйлера является «точной первого порядка».

Стабильность этого метода требует несколько большего рассмотрения. В нашем обсуждении мы будем работать над устойчивостью методов в случае одной переменной у = ау, имея в виду, что аналогичные утверждения переносятся на многомерные уравнения путем замены а на спектральный радиус. В этом случае мы знаем

$$yk+1 = yk + ahyk = (1 + ah)yk$$
.

Другими словами, yk = (1 + ah) ky0. Таким образом, интегратор устойчив при |1 + ah| 1, так как иначе |yk | экспоненциально. Предполагая, что a < 0 (иначе задача некорректна), мы можем упростить:

Таким образом, прямой Эйлер допускает ограничение шага по времени для устойчивости, определяемое нашим окончательным условием на h. Другими словами, выход прямого метода Эйлера может резко увеличиться, даже если у = ау стабильно, если h недостаточно мало. На рисунке НОМЕР показано, что происходит, когда это условие соблюдается или нарушается. В нескольких измерениях мы можем заменить это ограничение аналогичным ограничением, используя спектральный радиус А. Для нелинейных ОДУ эта формула дает руководство по стабильности, по крайней мере, локально во времени; глобально h, возможно, придется скорректировать, если якобиан F станет хуже обусловленным.

13.3.2 Обратный Эйлер

Точно так же мы могли бы применить схему обратного дифференцирования в yk+1 для разработки интегратора ОДУ:

$$F[yk+1] = y(t) = h$$
 $\frac{yk+1 - yk}{} + O(u)$

Таким образом, мы решаем следующую потенциально нелинейную систему уравнений относительно yk+1:

$$yk = yk+1$$
 $hF[yk+1].$

Поскольку мы должны решить это уравнение для yk+1, обратный Эйлер является неявным интегратором.

Этот метод обладает точностью первого порядка, как и форвард Эйлер, по идентичному доказательству. Однако стабильность этого метода значительно отличается от метода Эйлера. Снова рассматривая модельное уравнение у = ау, пишем:

$$yk \ yk = yk+1$$
 hayk+1 = $yk+1 = \frac{1}{1}$ ha

Параллельно с нашим предыдущим аргументом обратный Эйлер стабилен при следующем условии:

Очевидно, мы всегда берем h 0, поэтому обратный Эйлер безусловно устойчив.

Конечно, даже если обратный Эйлер устойчив, он не обязательно точен. Если h слишком велико, уk слишком быстро приблизится к нулю. При моделировании ткани и других физических материалов, для реалистичности которых требуется большое количество высокочастотных деталей, обратный метод Эйлера может оказаться неэффективным выбором. Кроме того, мы должны инвертировать F[·], чтобы найти уk+1.

Пример 13.8 (обратный Эйлер). Предположим, мы хотим решить у = Ау для А Rn×n . Затем, чтобы найти уk+1, решаем следующую систему:

$$yk = yk+1$$
 $hAyk+1 = yk+1 = (In \times n + hA)$ $_{1 r \kappa}$.

13.3.3 Трапециевидный метод

Предположим, что ук известно в момент времени tk, а yk+1 представляет значение в момент времени tk+1 = tk + h. Предположим, что мы также знаем yk+1/2 на полпути между этими двумя шагами. Тогда из нашего вывода центрированного дифференцирования мы знаем:

$$yk+1 = yk + hF[yk+1/2] + O(h^{-3})$$

Из нашего вывода правила трапеций:

$$\frac{F[y_{K}+1] + F[y_{K}]}{2} = F[y_{K}+1/2] + O(h^{2})$$

Подстановка этого соотношения дает нашу первую схему интегрирования второго порядка, метод трапеций для интегрирования ОДУ:

$$yk+1 = yk + h2 \frac{F[yk+1] + F[yk]}{}$$

Подобно обратному Эйлеру, этот метод является неявным, поскольку мы должны решить это уравнение относительно уk+1 .

Снова проводя анализ устойчивости по у = ау, находим в этом случае временные шаги метода трапеций решения

$$y_{K+1} = y_K + \frac{1}{2} - ra(y_{K+1} + y_K)$$

Другими словами,

$$y_K = \frac{1 + \frac{1}{2}ra}{1 - \frac{1}{2}xa_K}y_0.$$

Таким образом, метод устойчив, когда

$$\frac{1 + 1 \, \text{ra}}{1 - \frac{2}{3} \, 1 \, \text{ra}} < 1.$$

Легко видеть, что это неравенство выполняется при a < 0 и h > 0, что показывает безусловную устойчивость метода трапеций.

Однако, несмотря на более высокий порядок точности при сохранении стабильности, метод трапеций имеет некоторые недостатки, которые делают его менее популярным, чем обратный метод Эйлера. В частности, рассмотрим

p
$$\frac{y_{K+1}}{\ddot{y}_{K}} = \frac{1 + \frac{1}{2}ra}{2 \cdot 1 \cdot 1 \cdot \frac{1}{2}ra}$$

Когда a < 0, для достаточно больших h это отношение в конце концов становится отрицательным; на самом деле при h имеем R 1. Таким образом, как показано на рисунке ЧИСЛО, если временные шаги слишком велики, трапециевидный метод интегрирования имеет тенденцию демонстрировать нежелательное колебательное поведение, которое совсем не похоже на то, что мы могли бы ожидать для решений у = ау .

13.3.4 Методы Рунге-Кутты

Класс интеграторов можно вывести, сделав следующее наблюдение:

$$y_{K}+1=y_{K}+$$
 y (t) dt по основной теореме исчисления
$$=y_{K}+$$
 t_{K} $t_{K}+t_{K}$ $t_{K}+$

Конечно, прямое использование этой формулы не подходит для формулировки метода временного шага, поскольку мы не знаем y(t), но тщательное применение наших квадратурных формул из предыдущей главы может привести к выработке возможных стратегий.

Например, предположим, что мы применяем метод трапеций для интегрирования. Затем находим:

$$4y_{K+1} = y_{K} + \frac{1}{2} (F[y_{K}] + F[y_{K+1}]) + O(h)$$

Это формула, которую мы написали для метода трапеций в § 13.3.3.

Однако , если мы не хотим вычислять yk+1 в неявном виде, мы должны найти выражение для аппроксимации). помощник F[yk+1]. Однако, используя метод Эйлера, мы знаем, что yk+1 = yk + hF[yk] + O(h) эта 2 Изготовление замена yk+1 не влияет на порядок аппроксимации трапецеидального временного шага, описанного выше, поэтому мы можем написать:

=
$$yk + 2$$
 $-(F[yk] + F[yk + hF[yk]]) + O(hyk+1)^{3}$

Игнорирование $O(h^{-3})$ дает новую стратегию интеграции, известную как метод Хойна, которая второго порядка точного и явного.

Если мы изучим поведение устойчивости метода Хойна для у = ау для а < 0, мы знаем:

Таким образом, метод устойчив, когда

Неравенство справа показывает, что h a < $\sqrt{2}$, а тот, что слева, всегда верен при h > 0 и |a|, 0, поэтому условием устойчивости является $\frac{2}{|a|}$.

h Метод Хойна является примером метода Рунге-Кутты, полученного путем применения квадратурных методов к приведенному выше интегралу и подстановки шагов Эйлера в F[·]. Прямой Эйлер — это точный метод Рунге-Кутты первого порядка, а метод Хойна — метод второго порядка. Популярный метод Рунге-Кутты четвертого порядка (сокращенно «RK4») определяется следующим образом:

= ук + 6
$$\overline{\Psi}$$
 (к1 + 2к2 + 2к3 + к4) ук+1 где k1 = F [уk]
$$K2 = FyK + _{-} \frac{1}{2}_{\Psi ac1}$$

$$K3 = FyK + _{-} \frac{1}{2}_{\Psi ac2}$$

$$k4 = Fyk + hk3 _{-}$$

Этот метод можно получить, применив квадратурное правило Симпсона.

Методы Рунге-Кутты популярны, потому что они являются явными и поэтому их легко оценить, обеспечивая при этом высокую степень точности. Цена такой точности, однако, состоит в том, что F[·] необходимо оценивать большее количество раз. Кроме того, стратегии Рунге-Кутты могут быть расширены до неявных методов, которые могут решать жесткие уравнения.

13.3.5 Экспоненциальные интеграторы

Один класс интеграторов, который обеспечивает высокую точность, когда F[·] приблизительно линейна, заключается в явном использовании нашего решения модельного уравнения. В частности, если бы мы решали ОДУ у = Ау, используя собственные векторы А (или любой другой метод), мы могли бы найти явное решение y(t), как объяснено в § 13.2.3. Обычно мы пишем yk+1 = е Ahyk, кодирует наше возведение в степень собственных значений (на самом деле мы можем найти матрицу е Теперь, если мы из этого выражения, которое решает ОДУ к времени h). напишем

$$y = Ay + G[y],$$

где G — нелинейная, но малая функция, мы можем достичь довольно высокой точности, интегрируя часть A в явном виде, а затем отдельно аппроксимируя нелинейную часть G. Например, экспоненциальный интегратор первого порядка применяет прямой Эйлер к нелинейному члену G:

$$y\kappa+1 = e$$
 Ax $y\kappa - A$ 1 $(1 e Ah)G[yk]$

Анализ, выявляющий преимущества этого метода, более сложен, чем то, что мы написали, но интуитивно ясно, что эти методы будут вести себя особенно хорошо, когда G мало.

13.4 Многозначные методы

Преобразования в §13.2.1 позволили нам значительно упростить обозначения в предыдущем разделе, приведя все явные ОДУ к форме у = F[y]. На самом деле, хотя все явные ОДУ могут быть записаны таким образом, неясно, всегда ли они должны это делать.

В частности, когда мы сводили ОДУ k-го порядка к ОДУ первого порядка, мы вводили ряд переменных, представляющих производные от первого до k - 1-го желаемого результата. Фактически, в нашем окончательном решении нас интересует только нулевая производная, то есть сама функция, поэтому порядки точности временных переменных менее важны.

С этой точки зрения рассмотрим ряд Тейлора.

$$h$$
2 y(tk + h) = y(tk) + hy (tk) + y (tk) + O(h 2—

Если мы знаем только у с точностью до O(h²), это не влияет на нашу аппроксимацию, так как у умножается на h. Точно так же, если мы знаем только у до O (h), это приближение не повлияет на приведенные выше члены ряда Тейлора, потому что оно будет умножено на h 2/2. Таким образом, теперь мы рассматриваем «многозначный» meth (k) (t) = F[t,y (t),y (t), . . . ,y различных производных функции (k 1) (t)] с разным порядком точности для од, предназначенных для интегрирования у у.

Учитывая важность второго закона Ньютона F = ma, мы ограничимся случаем y = F[t,y,y]; существует множество расширений для менее распространенного случая k-го порядка. Введем вектор «скорости» v(t) = y (t) и вектор «ускорения» а. Используя наше предыдущее сокращение, мы хотим решить следующую систему первого порядка:

$$y(t) = v(t) v(t)$$

= $a(t) a(t) =$
 $F[t,y(t),v(t)]$

Наша цель — разработать интегратор, специально адаптированный для этой системы.

13.4.1 Схемы Ньюмарка

Начнем с вывода знаменитого класса интеграторов Ньюмарка.1 Обозначим yk , vk и ak как векторы положения, скорости и ускорения в момент времени tk ; наша цель — перейти ко времени tk+1 tk + h.

¹Мы следим за развитием в http://www.stanford.edu/group/frg/course_work/AA242B/CA-AA242B-Ch7. пдф.

Используйте y(t), v(t) и a(t) для обозначения функций времени, предполагая, что мы начинаем c tk . Тогда, очевидно мы можем написать

$$BK+1 = BK + {TK+1 \atop TK} a(T) дT$$

Мы также можем записать yk+1 как интеграл, включающий a(t), выполнив несколько шагов:

$$y_{K}+1 = y_{K} + \sum_{TK}^{TK+1} v(t) dt$$
 $= y_{K} + [tv(t)]tk_{TK}^{+}1 \qquad tk+1 ta(t) dt$ после интегрирования по частям

 $= y_{K} + T_{K}+1 B_{K}+1 - T_{K}B_{K} - \sum_{TK}^{K} tk+1 ta(t) dt$ путем разложения разностного члена

 $= y_{K} + hv_{K} + tk+1 v_{K}+1 - tk+1 v_{K} - tk+1 ta(t) dt$ путем сложения и вычитания hvk

 $= y_{K} + hv_{K} + tk+1 (v_{K}+1 - v_{K}) \qquad ta(t) dt$ после факторинга

 $= y_{K} + x_{BK} + x_{K} + x_{K}$

Предположим, мы выбрали т [tk , tk+1]. Тогда мы можем написать выражения для ak и ak+1 , используя ряд Тейлора относительно т:

$$ak = a(\tau) + a(\tau)(tk - \tau) + O(h$$
 2)
 $ak+1 = a(\tau) + a(\tau)(tk+1 - \tau) + O(h$ 2)

Для любой константы ү R, если мы масштабируем первое уравнение на 1 ү, а второе на ү и суммируем результаты, мы находим:

$$a(\tau) = (1 \quad y)ak + yak + 1 + a (\tau)((y \quad 1)(tk \quad \tau) \quad y(tk + 1 \quad \tau)) + O(h$$
 $= (1 \quad y)ak + yak + 1 + a (\tau)(\tau \quad hy \quad tk) + O(h$ $^2)$ после замены $tk + 1 = tk + h$

В конце концов, мы хотим проинтегрировать от tk до tk+1 , чтобы получить изменение скорости. Таким образом, мы вычисляем:

$$a(\tau) d\tau = (1 \quad y) hak + y hak + 1 +$$

$$= (1 \quad y) hak + y hak + 1 + O(h^{2}),$$
 $a(\tau) d\tau = (1 \quad y) hak + y hak + 1 + O(h^{2}),$

где выполняется второй шаг, поскольку подынтегральная функция равна O(h) , а интервал интегрирования имеет ширину h. Другими словами, теперь мы знаем:

$$vk+1 = vk + (1 y)hak + yhak+1 + O(h^2)$$

Чтобы сделать аналогичное приближение для yk+1 , мы можем написать

Таким образом, мы можем использовать наши более ранние отношения, чтобы показать:

$$y_{K}+1=y_{K}+x_{BK}+$$
 $(t_{K}+1)$ $(t_$

Здесь мы используем β вместо γ (и поглощаем коэффициент два в процессе), потому что γ , который мы выбираем для аппроксимации γ не обязательно должен быть таким же, как тот, который мы выбираем для аппроксимации γ не обязательно должен быть таким же, как тот, который мы выбираем для аппроксимации γ не обязательно должен быть таким же, как тот, который мы выбираем для аппроксимации γ не обязательно должен быть таким же, как тот, который мы выбираем для аппроксимации γ не обязательно должен быть таким же, как тот, который мы выбираем для аппроксимации γ не обязательно должен быть таким же, как тот, который мы выбираем для аппроксимации γ не обязательно должен быть таким же, как тот, который мы выбираем для аппроксимации γ не обязательно должен быть таким же, как тот, который мы выбираем для аппроксимации γ не обязательно должен быть таким же, как тот, который мы выбираем для аппроксимации γ не обязательно должен быть таким же, как тот, который мы выбираем для аппроксимации γ не обязательно должен быть таким же, как тот, который мы выбираем для аппроксимации γ не обязательно должен быть таким же, как тот, который мы выбираем для аппроксимации γ не обязательно должен быть γ не обязательно должен быть γ не обязательно должен γ не обязател

После всего этого интегрирования мы получили класс схем Ньюмарка с двумя входными параметры у и β, который имеет точность первого порядка по приведенному выше доказательству:

$$y_{K+1} = y_{K} + x_{BK} + 4 - \beta + 2^{2} a_{K} + \beta h$$
 $_{2 a_{K}+1}$
 $v_{K+1} = v_{K} + (1 \quad y)_{hak} + y_{hak+1} a_{K} =$
 $F[t_{K}, y_{K}, v_{K}]$

Разный выбор β и у приводит к разным схемам. Например, рассмотрим следующие примеры:

• β = γ = 0 дает интегратор постоянного ускорения:

$$yk+1 = yk + hvk + h ak 2 - \frac{12}{2}$$

 $BK+1 = BK + xaK$

Этот интегратор является явным и выполняется точно, когда ускорение является постоянной функцией.

• β = 1/2, γ = 1 дает постоянный неявный интегратор ускорения:

$$yk+1 = yk + hvk + h ak+1 + \frac{1}{2}$$
 $bk+1 = bk + xak+1$

Здесь скорость неявно ступенчата с использованием обратного Эйлера, что дает точность первого порядка. Обновление у, однако, может быть записано

$$y_{K}+1 = y_{K} + 2$$
 $+ y_{K} + y_{$

показывая, что он локально обновляется по правилу средней точки; это наш первый пример схемы, в которой обновления скорости и положения имеют разные порядки точности. Но даже в этом случае можно показать, что этот метод, однако, имеет глобальную точность первого порядка по у.

• β = 1/4, γ = 1/2 после некоторой алгебры дает следующую схему трапеций второго порядка:

$$xk+1 = xk + 2 \cdot 1 + h(vk + vk+1)$$

$$BK+1 = BK + 2 - 4(aK + aK+1)$$

• β = 0, γ = 1/2 дает схему центральных разностей второго порядка точности. В каноническом форма, у нас есть

$$xk+1 = xk + hvk + h ak 2 - \frac{1}{2}$$

$$1$$
 $BK+1 = BK + 2 - 4(aK + aK + 1)$

Метод получил свое название, потому что упрощение приведенных выше уравнений приводит к альтернативной форме:

$$vk+1 = \frac{yk+2 \quad yk}{2h \ yk+1}$$
h 2 \frac{2yk+1 + yk \ ak+1 = 1}{2k}

Можно показать, что методы Ньюмарка безусловно устойчивы при $4\beta > 2\gamma > 1$, а точность второго порядка возникает именно при $\gamma = 1/2$ (ПРОВЕРКА).

13.4.2 Ступенчатая сетка

Другой способ достижения точности второго порядка iny заключается в использовании центрированных разностей относительно времени tk+1/2 tk + h/2:

того, чтобы использовать аргументы Тейлора, чтобы попытаться переместить wk+1/2, мы можем просто сохранить скорости v в половине точек на сетке временных шагов.

Затем мы можем использовать аналогичное обновление, чтобы увеличить скорость:

$$BK+3/2 = BK+1/2 + XaK+1$$
.

Обратите внимание, что это обновление на самом деле имеет второй порядок точности и для x, поскольку, если мы подставим наши выражения для vk+1/2 и vk+3/2, мы можем написать:

$$aK+1 = \frac{1}{42} (yK+2 - 2yK+1 + yK)$$

Наконец, для члена ускорения достаточно простого приближения, поскольку это член более высокого порядка:

$$ak+1 = F tk+1, xk+1, 2$$
 $+(vk+1/2 + vk+3/2)$

Это выражение можно подставить в выражение для vk+3/2.

Когда $F[\cdot]$ не зависит от v, метод полностью явный:

Это известно как метод интегрирования чехарды благодаря шахматной временной сетке и тому факту, что каждая средняя точка используется для обновления следующей скорости или положения.

В противном случае, если обновление скорости зависит от v, метод становится неявным.

Часто зависимость от скорости симметрична; например, сопротивление ветра просто меняет знак, если вы движетесь в обратном направлении. Это свойство может привести к симметричным матрицам на неявном этапе обновления скоростей, что позволяет использовать сопряженные градиенты и связанные с ними быстрые итерационные методы для решения.

13.5 Сделать

- Определить жесткое ОДУ
- Приведите таблицу методов временного шага для F[t;y]
- Используйте обозначение Y более последовательно.

13.6 Проблемы

- ТВД РК
- Многошаговые/многозначные методы а-ля Хит
- Интеграция Verlet
- Симплектические интеграторы

Machine Translated by Google

Глава 14

Уравнения с частными производными

Наше интуитивное представление об обыкновенных дифференциальных уравнениях обычно основано на эволюции физических систем во времени. Уравнения, подобные второму закону Ньютона, определяющему движение физических объектов во времени, доминируют в литературе по таким проблемам с начальными значениями; дополнительные примеры приходят из химических концентраций, реагирующих с течением времени, популяций хищников и жертв, взаимодействующих от сезона к сезону, и так далее. В каждом случае исходная конфигурация — например, положения и скорости частиц в системе в нулевое время — известна, и задача состоит в том, чтобы предсказать поведение с течением времени.

Однако в этой главе мы рассматриваем возможность связывания отношений между различными производными функции. Нетрудно найти примеры, когда такая связь необходима. Например, при моделировании количества дыма или газов, таких как «градиенты давления», производные от давления газа в пространстве, учитывайте, как газ движется во времени; такая структура является разумной, поскольку газ естественным образом диффундирует из областей высокого давления в области низкого давления. При обработке изображений производные связываются еще более естественно, поскольку измерения изображений, как правило, происходят одновременно в направлениях х и у.

Уравнения, связывающие вместе производные функций, известны как уравнения в частных производных. Они являются предметом богатой, но очень тонкой теории, достойной более широкого рассмотрения, поэтому наша цель здесь будет состоять в том, чтобы обобщить ключевые идеи и предоставить достаточный материал для решения проблем, обычно возникающих на практике.

14.1 Мотивация

Уравнения в частных производных (УЧП) возникают, когда неизвестной является некоторая функция f: Rn Rm. Нам задано одно или несколько соотношений между частными производными f, и цель состоит в том, чтобы найти f, удовлетворяющее критериям. PDE появляются практически в любой области прикладной математики, и мы перечислим лишь некоторые из них ниже.

Кроме того, прежде чем вводить конкретные УЧП, мы должны ввести некоторые обозначения. В частности, есть несколько комбинаций частных производных, которые часто встречаются в мире частных производных. Если f: R3 R и v: R3 R3 , то стоит запомнить следующие операторы:

Дивергенция:
$$\cdot$$
 v ++ $\frac{v_1}{x_1} \frac{v_2}{x_2} \frac{v_3}{x_3} \frac{v_3}{v_3} \frac{}{2} \frac{}{2}$

Завиток: \times v $\frac{2 \text{ fff}}{x_1} \frac{}{x_2} \frac{v_1}{x_3} \frac{v_2}{x_1} \frac{v_1}{x_1} \frac{}{x_2}$

Лапласиан: $2f + + x\overline{2} x 2 \overline{x} 2\overline{1} 23$

Пример 14.1 (Моделирование жидкости). Поток жидкостей, таких как вода и дым, регулируется уравнениями Навье-Стокса, системой УЧП со многими переменными. В частности, предположим, что в некоторой области Ω R3 движется жидкость . Мы определяем следующие переменные, показанные на рисунке NUMBER:

- t [0,): время
- v(t) : Ω R3 : скорость жидкости
- $\rho(t)$: Ω R: плотность жидкости
- p(t): Ω R: давление жидкости
- f(t) : Ω R3 : Внешние силы, такие как гравитация, действующие на жидкость.

Если жидкость имеет вязкость µ, то, если предположить, что она несжимаема, уравнения Навье-Стокса утверждают:

$$p = \frac{v}{t} + v \cdot v = p + \mu + 2v + f$$

Здесь 2v = 2v1/ x 2 + 2v2/ x 2 † 2v3/ x 2; мы думаем о градиенте 3 как о градиенте в пространстве, а не как о f f x 1 , x 2 ,

времени, т. е. $f = \frac{1}{f-x^3}$. Эта система уравнений определяет временную динамику движения жидкости и (на самом деле можно построить, применяя второй закон Ньютона к отслеживанию «частиц» жидкости. Однако его формулировка включает не только производные по времени — гно также и производные в пространстве , что делает его УЧП.

Пример 14.2 (уравнения Максвелла). Уравнения Максвелла определяют взаимодействие электрических полей Е и магнитных полей В во времени. Как и в случае с уравнениями Навье-Стокса, мы думаем о градиенте, дивергенции и завитке как о частных производных по пространству (а не по времени t). Тогда систему Максвелла (в «сильной» форме) можно записать:

Закон Гаусса для электрических полей:
$$\cdot$$
 E = $\frac{p}{\epsilon^0}$

Закон Гаусса для магнетизма: $\cdot B = 0$

$$B$$
 Закон Фарадея: $\times E = t$

$$= \frac{1}{2}$$
 Закон Ампера: × B = μ0 J + ε0 t

Здесь ε0 и μ0 — физические константы, а J кодирует плотность электрического тока. Как и уравнения Навье-Стокса, уравнения Максвелла связывают производные физических величин по времени t с их производными в пространстве (задаваемыми терминами ротора и дивергенции).

Пример 14.3 (уравнение Лапласа). Предположим, Ω — область в R2 с границей Ω и что нам дана функция g : Ω R, показанная на рис. НОМЕР. Мы можем захотеть интерполировать g внутрь Ω . Однако, когда Ω имеет неправильную форму, наши стратегии интерполяции из главы 11 могут не сработать.

Предположим, мы определяем f(x): Ω R как интерполирующую функцию. Затем одна стратегия, вдохновленная нашим подходом к методу наименьших квадратов, заключается в определении функционала энергии:

$$E[x] = \phi(x)^{\frac{2}{2}}A^{x}$$

$$f E[f]$$
 такое, что $f(x) = g(x)$ x Ω

Эта установка выглядит как оптимизация, которую мы решили в предыдущих примерах, но теперь наша неизвестная — это функция f, a не точка в Rn!

Если f минимизирует E, то E[f + h] E[f] для всех функций h(x). Это утверждение верно даже для малых возмущений E[f + ɛh] при ϵ 0. Деля на ϵ и переходя к пределу при ϵ 0, мы должны иметь E[f + ɛh] ϵ 0; это то же самое, что приравнять производные по направлению к нулю, чтобы найти ее минимумы. Мы можем упростить:

Дифференциальные шоу:

Эта производная должна равняться нулю для всех h, так что, в частности, мы можем выбрать h(x) = 0 для всех x применяя интегрирование по частям, имеем:

$$\frac{\Gamma}{--} E[f + \varepsilon h] | \varepsilon = 0 = 2 d\varepsilon \int_{0}^{\infty} h(x) 2 f(x) dx$$

Это выражение должно быть равно нулю для всех (всех!) возмущений h, поэтому мы должны иметь 2 f(x) = 0 для всех $x \Omega \Omega$ (формальное доказательство выходит за рамки нашего обсуждения). То есть проблема интерполяции выше

¹ Используемое здесь обозначение E[·] не означает «ожидание», как это могло бы быть в теории вероятностей, а скорее просто функционал «энергия»; это стандартное обозначение в областях функционального анализа.

можно решить с помощью следующего уравнения в

частных производных:

 $2 f(x) = 0 f(x) = g(x) x \Omega$

Это уравнение известно как уравнение Лапласа, и его можно решить, используя разреженные положительно определенные линейные методы, подобные тому, что мы рассмотрели в главе 10. Как мы видели, его можно применять к задачам интерполяции для нерегулярных областей Ω; кроме того, E[f] может быть увеличено для измерения других свойств f , например, насколько хорошо f аппроксимирует некоторую зашумленную функцию f0, чтобы получить связанные УЧП путем параллельного рассуждения выше.

Пример 14.4 (уравнение Эйконала). Предположим, что Ω Rn — некоторая замкнутая область пространства. Тогда мы могли бы взять d(x) как функцию, измеряющую расстояние от некоторой точки x0 до x полностью в пределах Ω . Когда Ω выпукло, мы можем записать d в закрытой форме:

$$d(x) = x x02.$$

Однако, как показано на рисунке НОМЕР, если Ω невыпуклая или представляет собой более сложную область, такую как поверхность, вычисление расстояний усложняется. В этом случае функции расстояния d удовлетворяют локализованному условию, известному как уравнение эйконала:

d2 = 1.

Если мы сможем его вычислить, d можно будет использовать для таких задач, как планирование пути роботов путем минимизации расстояния, которое они должны пройти, с ограничением, что они могут двигаться только в Ω.

Специализированные алгоритмы, известные как методы быстрого марша, используются для нахождения оценок d при данных x0 и Ω путем интегрирования уравнения эйконала. Это уравнение нелинейно относительно производной d, поэтому методы интегрирования для этого уравнения несколько специализированы, а доказательство их эффективности сложно. Интересно, но неудивительно, что многие алгоритмы решения уравнения эйконала имеют структуру, аналогичную алгоритму Дейкстры для вычисления кратчайших путей на графах.

Пример 14.5 (Гармонический анализ). Различные объекты по-разному реагируют на вибрации, и в значительной степени эти реакции зависят от геометрии объектов. Например, виолончели и фортепиано могут играть одну и ту же ноту, но даже неопытный музыкант легко различает издаваемые ими звуки. С математической точки зрения мы можем принять Ω R3 как форму, представленную либо в виде поверхности, либо в виде объема.

Если зажать ребра фигуры, то ее частотный спектр задается решениями следующей дифференциальной задачи на собственные значения:

$$2 f = \lambda f f(x)$$

= 0 x Ω , где 2 —

лапласиан области Ω , а Ω — граница области Ω . На рисунке HOMEP показаны примеры этих функций в разных областях Ω .

$$\sin kx = k \cos kx \quad x \cdot 2 \quad x$$

$$2 = k \sin kx$$

$$\sin k \cdot 0 = 0$$

$$\sin k \cdot 2\pi = 0$$

Таким образом, собственными функциями являются sin kx с собственными значениями k 2 ·

14.2 Основные определения

Используя обозначения СІТЕ, мы будем предполагать, что наша неизвестная есть некоторая функция f : Rn R. Для уравнений, содержащих до трех переменных, мы будем использовать индексы для обозначения частных производных:

и так далее.

Частные производные обычно указываются как отношения между двумя или более производными f в следующем:

• Линейная, однородная: fxx + fxy fy = 0

• Линейная: fxx y fyy + f = xy2

• Нелинейный: $f = \frac{2}{xx} = \phi y$

Как правило, мы действительно хотим найти $f:\Omega$ R для некоторого Ω Rn . Точно так же, как ОДУ были сформулированы как задачи с начальным значением, мы будем формулировать большинство УЧП как задачи с граничными значениями. То есть наша задача будет заключаться в том, чтобы заполнить f внутри Ω заданными значениями на его границе. На самом деле, мы можем думать о проблеме начального значения ОДУ следующим образом: областью является Ω = $[0, \)$ с границей Ω = $\{0\}$, где мы предоставляем входные данные. Рисунок NUMBER иллюстрирует более сложные примеры. Граничные условия для этих задач принимают разные формы:

- Условия Дирихле просто определяют значение f(x) на Ω
- Условия Неймана задают производные f(x) на Ω
- Смешанные условия или условия Робина объединяют эти два

14.3 Уравнения модели

Вспомним из предыдущей главы, что мы смогли понять многие свойства ОДУ, исследуя модельное уравнение у = ау. Мы можем попытаться использовать аналогичный подход для PDE, хотя мы обнаружим, что история становится более тонкой, когда производные связаны друг с другом.

Как и в случае уравнения модели для ОДУ, мы будем изучать линейный случай с одной переменной. Мы также ограничимся системами второго порядка, т. е. системами, содержащими не более чем вторую производную от u; модельное ОДУ было первого порядка, но здесь нам нужно по крайней мере два порядка, чтобы изучить, как производные взаимодействуют нетривиальным образом.

Линейное УЧП второго порядка имеет следующий общий вид:

$$c = 0 \text{ aij} \frac{f}{xi} \frac{f}{xj} + \frac{bi +}{xi}$$

$$ij$$

Формально мы можем определить «оператор градиента» как:

Вы должны проверить, что этот оператор является разумной записью в том смысле, что такие выражения, как f, · v и × v, обеспечивают правильные выражения. В этих обозначениях УЧП можно рассматривать как принимающую матричную форму:

$$(A + \cdot b + c)f = 0.$$

Эта форма имеет много общего с нашим изучением квадратичных форм в сопряженных градиентах, и на самом деле мы обычно характеризуем УЧП структурой A:

- Если А положительно или отрицательно определено, система является эллиптической.
- Если А положительно или отрицательно полуопределенно, система является параболической.
- Если А имеет только одно собственное значение отличного от остальных знака, система является гиперболической.
- Если А не удовлетворяет ни одному из критериев, система является ультрагиперболической.

Эти критерии перечислены примерно в порядке сложности решения каждого типа уравнений. Мы рассмотрим первые три случая ниже и приведем примеры соответствующего поведения; ультрагиперболические уравнения встречаются на практике реже и требуют для своего решения узкоспециализированных методов.

TODO: Приведение к канонической форме с помощью собственного материала A (не в 205A)

14.3.1 Эллиптические УЧП

Точно так же, как положительно определенные матрицы позволяют использовать специализированные алгоритмы, такие как разложение Холецкого и сопряженные градиенты, которые упрощают их обращение, эллиптические УЧП имеют особенно сильную структуру, которая приводит к эффективным методам решения.

Модельное эллиптическое УЧП представляет собой уравнение Лапласа, задаваемое формулой 2 f = g для некоторой заданной функции g в виде в примере 14.3. Например, с двумя переменными уравнение Лапласа принимает вид

$$fxx + fyy = r$$
.

На рисунке NUMBER показаны некоторые решения уравнения Лапласа для различных вариантов выбора и и f в двумерной области.

Эллиптические уравнения обладают многими важными свойствами. Особое теоретическое и практическое значение имеет идея эллиптической регулярности, согласно которой решения эллиптических УЧП автоматически являются гладкими функциями в С (Ω). Это свойство не сразу очевидно: УЧП второго порядка в f требует только, чтобы f было дважды дифференцируемо, чтобы иметь смысл, но на самом деле при слабых условиях они автоматически бесконечно дифференцируемы. Это свойство подтверждает физическую интуицию, согласно которой эллиптические уравнения представляют собой устойчивые физические равновесия, подобные позе покоя растянутого резинового листа.

Эллиптические уравнения второго порядка в приведенной выше форме также гарантированно допускают решения, в отличие от УЧП в некоторых других формах.

Пример 14.6 (Пуассон по одной переменной). Уравнение Лапласа с g = 0 получило специальное название уравнения Пуассона. В одной переменной можно записать f(x) = 0, что тривиально решается как $f(x) = \alpha x + \beta$. Этого уравнения достаточно, чтобы исследовать возможные граничные условия на [a, b]:

- Условия Дирихле для этого уравнения просто определяют f(a) и f(b); очевидно, существует единственная линия, проходящая через (a, f(a)) и (b, f(b)), которая обеспечивает решение уравнения.
- Условия Неймана определяют f (a) и f (b). Но f (a) = f (b) = α для f(x) = α x + β . Таким образом, граничные значения для задач Неймана могут подчиняться условиям совместимости, необходимым для принятия решения. Кроме того, выбор β не влияет на граничные условия, поэтому при их выполнении решение не единственно.

14.3.2 Параболические УЧП

Продолжая параллель со структурой линейной алгебры, с положительно-полуопределенными системами уравнений иметь дело лишь немного сложнее, чем с положительно-определенными. В частности, положительно полуопределенные матрицы допускают нулевое пространство, с которым нужно обращаться осторожно, но в остальных направлениях матрицы ведут себя так же, как и определенный случай.

Уравнение теплопроводности представляет собой модель параболического УЧП. Предположим, что f(0; x, y) — распределение тепла в некоторой области Ω R2 в момент времени t = 0. Тогда уравнение теплопроводности определяет, как тепло распространяется с течением времени t как функция f(t; x, y):

$$\frac{f}{---}=\alpha$$
 2f, t

где $\alpha > 0$, и мы снова считаем 2 лапласианом в пространственных переменных x и y, то есть 2 = 2/x 2 + 2/y 2. Это уравнение должно быть параболическим, так как перед fxx и fyy стоит один и тот же коэффициент α , но ftt в уравнении не фигурирует.

Рисунок НОМЕР иллюстрирует феноменологическую интерпретацию уравнения теплопроводности. Мы можем думать о 2 f как об измерении выпуклости f , как на рисунке ЧИСЛО (а). Таким образом, уравнение теплопроводности увеличивает u со временем, когда его значение «закрыто» вверх, и уменьшает f в противном случае. Эта отрицательная обратная связь устойчива и приводит к равновесию при t .

Для уравнения теплопроводности необходимы два граничных условия, оба из которых прямые физические интерпретации:

- Распределение тепла f(0; x, y) в момент времени t = 0 во всех точках (x, y) Ω
- Поведение f при t > 0 в точках (x, y) Ω . Эти граничные условия описывают поведение на границе области. Здесь условия Дирихле дают f(t; x, y) для всех t 0 и (x, y) Ω , что соответствует ситуации, когда посторонний агент фиксирует температуры на границе области. Эти условия могут возникнуть, если Ω представляет собой кусок фольги, расположенный рядом с источником тепла, на температуру которого фольга не оказывает существенного влияния, например, большой холодильник или духовка. Напротив, условия Неймана задают производную f в направлении, нормальном к границе Ω , как на рисунке HOMEP; они соответствуют фиксации потока тепла из Ω , вызванного различными типами изоляции.

14.3.3 Гиперболические УЧП

Окончательное уравнение модели представляет собой волновое уравнение, соответствующее случаю неопределенной матрицы:

$$-\frac{2}{t^2}f \quad c \quad 2 \quad 2f = 0$$

Волновое уравнение является гиперболическим, потому что вторая производная по времени имеет противоположный знак от двух пространственных производных. Это уравнение определяет движение волн через упругую среду, подобную резиновому листу; например, его можно получить, применив второй закон Ньютона к точкам на отрезке резинки, где х и у - положения на листе, а f (t; x, y) - высота отрезка резинки в момент времени t.

Рисунок НОМЕР иллюстрирует одномерное решение волнового уравнения. Поведение волны значительно отличается от диффузии тепла в том смысле, что при t энергия может не распространяться. В частности, волны могут бесконечно отражаться взад и вперед по домену. По этой причине мы увидим, что стратегии неявного интегрирования могут не подходить для интегрирования гиперболических УЧП, поскольку они имеют тенденцию гасить движение.

Граничные условия для волнового уравнения аналогичны граничным условиям для уравнения теплопроводности, но теперь мы должны указать как f(0; x, y) , так и ft(0; x, y) в нулевое время:

- Условия при t = 0 определяют положение и скорость волны в начальный момент времени.
- Граничные условия на Ω определяют, что происходит на концах материала. Условия Дирихле соответствуют фиксации сторон волны, например, защипывание струны виолончели, которая удерживается двумя концами на инструменте. Условия Неймана соответствуют тому, что концы волны остаются нетронутыми, как конец кнута.

14.4 Производные как операторы

В PDE и в других местах мы можем думать о производных как об операторах, действующих на функции так же, как матрицы действуют на векторы. Наш выбор обозначений часто отражает эту параллель: производная d f/dx выглядит как произведение оператора d/dx и функции f . На самом деле дифференцирование является линейным оператором, как и умножение матриц, поскольку для всех f , g : R и a, b R

$$\frac{dd (a f (x) + bg (x)) = a f (x) + b g}{(x). дx дx дx}$$

На самом деле, когда мы дискретизируем УЧП для численного решения, мы можем провести эту аналогию полностью. Например, рассмотрим функцию f на [0, 1], дискретизированную с использованием n + 1 равноотстоящих отсчетов, как показано на рисунке HOMEP. Напомним, что расстояние между двумя образцами равно h = 1/n. В главе 12 мы разработали приближение для второй производной f (x):

h) f (x) =
$$\frac{f(x + h)}{2}$$
 2 f(x) + f(x

Предположим, что наши n выборок f(x) на [0, 1] таковы: y0 f(0), y1 f(h), y2 f(2h), . . . , yn = f(nh). Затем применение нашей формулы выше дает стратегию для аппроксимации f в каждой точке сетки:

То есть вторую производную функции на сетке точек можно вычислить с помощью шаблона 1 - 2 - 1, показанного на рисунке HOMEP (a).

Одна тонкость, которую мы не рассмотрели, заключается в том, что происходи π ри у 0 и у n, поскольку приведенная выше формула требует у 1 и уn+1. Фактически это решение кодирует граничные условия, введенные в §14.2. Принимая во внимание, что у0 = f(0) и уn = f(1), примеры возможных граничных условий для f включают:

- Граничные условия Дирихле: у 1 = yn+1 = 0, то есть просто зафиксировать значение у за пределами конечные точки
- Граничные условия Неймана: у 1 = y0 и yn+1 = yn, кодирующие граничное условие f(0) = f(1) = 0.
- Периодические граничные условия: у 1 = yn и yn+1 = y0, что делает отождествление $f(0) = \phi(1)$

Предположим, мы складываем выборки yk в вектор у Rn+1 , а выборки yk — вектор w Rn+1 . ^{В секунду} Тогда из нашей вышеприведенной конструкции легко увидеть, что h 2w = L1y, где L1 — один из следующих выборов:



То есть матрицу L можно рассматривать как дискретизированную версию оператора у Rn+1 , а $\frac{A}{2 \text{ д} \text{ X} 2}$ действующий на не функции f : [0, 1] R.

Мы можем написать аналогичную аппроксимацию для 2 f, когда мы выбираем $f:[0,1] \times [0,1]$ R с сеткой значений, как на рисунке HOMEP. В частности, напомним, что в этом случае 2 f = fxx + fyy, так что, в частности, мы можем просуммировать вторые производные x и y, как мы это делали в приведенном выше одномерном примере. Это приводит к двойному трафарету 1 - 2 - 1, как на рисунке HOMEP. Если мы пронумеруем наши выборки как yk, f(kh, h), то наша формула для лапласиана f в этом случае f(kh, h) выбираем f(kh, h) го наша формула для лапласиана f в этом случае f(kh, h) го наша f(kh, h) го наша

$$(2y)k$$
, $\frac{1}{h-2}(y(k-1), +yk$, $(-1)+y(k+1), +yk$, $(+1)$ $4yk$,

Если мы еще раз объединим наши выборки у и у в у и w, то, используя аналогичную конструкцию и выбор граничных условий, мы снова можем записать h 2w = L2y. Этот двумерный сеточный лапласиан L2 используется во многих приложениях для обработки изображений, где (k,) используется для индексации пикселей на изображении.

Естественный вопрос, который следует задать после обсуждения выше, заключается в том, почему мы перешли к лапласиану второй производной в нашем обсуждении выше, а не к дискретизации первой производной f (x). В принципе, нет никаких причин, по которым мы не могли бы сделать аналогичные матрицы D, реализующие прямую, обратную или симметричную разностную аппроксимацинано несколько технических моментов несколько усложняют эту задачу, как подробно описано ниже.

Самое главное, мы должны решить, какое приближение первой производной использовать. Если мы напишем ,

1 как прямая разница h — например, у k (уk+1 — уk), то мы окажемся в неестественно асимметричном
положении, когда нам понадобится граничное условие в уп , но не в у0. Напротив, мы могли бы использовать
симметричную разность (уk+1 — уk — 1), но эта дискретизация страдает от более тонкого ограждения
проблема, показанная на рисунке НОМЕР. В частности, эта версия у k игнорирует значение уk и смотрит только на
своих соседей уk — 1 и уk+1 , что может создавать артефакты, поскольку каждая строка D включает уk только для
четного или нечетного k, но не для обоих.

Если мы используем прямые или обратные производные, чтобы избежать проблем со столбами забора, мы теряем порядок точности, а также страдаем от описанной выше асимметрии. Как и в случае интеграции чехарды

Алгоритм в §13.4.2, один из способов избежать этих проблем состоит в том, чтобы думать о производных как о живущих в половинных узлах сетки, как показано на рисунке HOMEP. В одномерном случае это изменение соответствует (уk+1 разницы в вершинах, ребрах и — yk) как уk+1/2 . Этот метод размещения различных производных для маркировки центрах ячеек сетки особенно распространен в моделировании жидкости, которое поддерживает давление, скорость жидкости и т. д. в местах, которые упрощают вычисления.

Помимо этих тонкостей, наш главный вывод из этого обсуждения состоит в том, что если мы дискретизируем функцию f(x), отслеживая выборки (xi , yi) , то наиболее разумные приближения производных f будут вычислимы как произведение Lx для некоторой матрицы L. Это наблюдение завершает аналогию: «Производные действуют на функции, как матрицы действуют на векторы». Или в стандартизированной записи экзамена: Производные : Функции :: Матрицы : Векторы

14.5 Численное решение УЧП

Многое еще предстоит сказать о теории УЧП. Вопросы существования и уникальности, а также возможность характеристики решений различных УЧП приводят к детальным обсуждениям с использованием передовых аспектов реального анализа. Хотя полное понимание этих свойств необходимо для строгого доказательства эффективности дискретизации PDE, у нас уже есть достаточно, чтобы предложить несколько методов, которые используются на практике.

14.5.1 Решение эллиптических уравнений

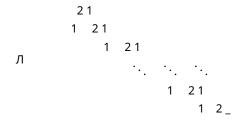
Мы уже сделали большую часть работы по решению эллиптических УЧП в § 14.4. В частности, предположим, что мы хотим решить линейное эллиптическое УЧП вида Lf = g. Здесь L — дифференциальный оператор; например, чтобы решить уравнение Лапласа, мы возьмем L 2 лапласиана. Затем в §14.4 мы показали, что если мы дискретизируем f, взяв набор отсчетов в векторе у с yi = f(xi), то соответствующее приближение Lf можно записать как Ly для некоторой матрицы L. Если мы также дискретизируем g используя отсчеты в векторе b, затем решение эллиптического УЧП Lf = g аппроксимируется решением линейной системы Ly =b.

Пример 14.7 (Эллиптическая дискретизация УЧП). Предположим, мы хотим аппроксимировать решения f(x) = g(x) на [0, 1] с граничными условиями f(0) = f(1) = 0. Будем аппроксимировать f(x) вектором у Rn выборка f(x) следующим образом:

y1	ф(ч)
y2	ф(2ч)
÷	÷
П	ф (нч)

где h = 1/n+1. Мы не добавляем образцы в x = 0 или x = 1, так как граничные условия определяют значения там. Мы будем использовать b для хранения аналогичного набора значений для q(x).

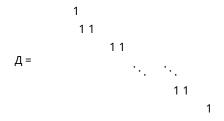
Учитывая наши граничные условия, мы дискретизируем f(x) как $\frac{1}{h}$ 2 Ly, где L определяется как:



Таким образом, наше приближенное решение УЧП определяется выражением y = h 2L 1b.

Точно так же, как эллиптические УЧП являются наиболее простыми для решения, их дискретизация с использованием матриц, как в приведенном выше примере, является наиболее простой для решения. На самом деле дискретизация эллиптического оператора L обычно представляет собой положительно определенную и разреженную матрицу, идеально подходящую для методов решения, описанных в главе 10.

Пример 14.8 (Эллиптические операторы как матрицы). Рассмотрим матрицу L из примера 14.7. Мы можем показать, что L отрицательно определена (и, следовательно, положительно определенная система Ly = h 2b может быть решена с использованием сопряженных градиентов), заметив, что L = DD для матрицы D R(n+1)×n, заданной формулой:



Эта матрица есть не что иное, как первая производная с конечной разностью, поэтому это наблюдение соответствует тому факту, что d 2 f/dx2 = d/dx(d f/dx). Таким образом, x Lx = x DDx = Dx 0, что показывает, что L является отрицательно полуопределенным. Легко видеть, что Dx = 0 именно тогда, когда x = 0, что завершает доказательство того, что L на самом деле отрицательно определена.

14.5.2 Решение параболических и гиперболических уравнений

Параболические и гиперболические уравнения обычно вводят в формулировку временную переменную, которая также является дифференцированной, но потенциально более низкого порядка. Поскольку решения параболических уравнений обладают многими свойствами устойчивости, численные методы для работы с этой переменной времени часто устойчивы и хорошо обусловлены; напротив, необходимо уделять больше внимания гиперболическому поведению и предотвращать затухание движения с течением времени.

Полудискретные методы Вероятно, наиболее простым подходом к решению простых уравнений, зависящих от времени, является использование полудискретного метода. Здесь мы дискретизируем домен, но не переменную времени, что приводит к ОДУ, которое можно решить с помощью методов из главы 13.

Пример 14.9 (полудискретное уравнение теплопроводности). Рассмотрим уравнение теплопроводности с одной переменной, заданное формулой ft = fxx, где f(t; x) представляет теплоту в точке x и в момент времени t. В качестве граничных данных пользователь предоставляет

функция f0(x) такая, что f(0;x) f0(x); мы также присоединяем границу x $\{0,1\}$ к холодильнику и тем самым обеспечиваем f(t;0) = f(t;1) = 0.

Предположим, мы дискретизируем переменную х, определив:

где, как и в примере 14.7, мы берем h = 1/n+1 и опускаем выборки при х {0, 1} , так как они обеспечиваются граничными условиями.

Комбинируя эти fi, мы можем определить f(t): R Rn как полудискретную версию f, в которой мы выбрали пространство, но не время. По нашей конструкции приближение полудискретного УЧП представляет собой ОДУ, заданное формулой f(t) = L f(t).

В предыдущем примере показан пример очень общего шаблона для параболических уравнений.

Когда мы моделируем непрерывные явления, такие как тепло, проходящее через область, или химические вещества, диффундирующие через мембрану, обычно имеется одна временная переменная, а затем несколько пространственных переменных, которые дифференцируются эллиптическим образом. Когда мы полудискретно дискретизируем эту систему, мы можем использовать стратегии интеграции ОДУ для их решения. Фактически, точно так же, как матрица, используемая для решения линейного эллиптического уравнения, как в § 14.5.1, обычно является положительно или отрицательно определенной, когда мы записываем полудискретное параболическое УЧП f = L f, матрица L обычно является отрицательно определенной. Это наблюдение означает, что f, решающая это непрерывное ОДУ, безусловно устойчива, поскольку отрицательные собственные значения со временем затухают.

Как отмечалось в предыдущей главе, у нас есть много вариантов решения ОДУ во времени, возникающих в результате пространственной дискретизации. Если временные шаги малы и ограничены, могут быть приемлемы явные методы. Неявные решатели часто применяются для решения параболических уравнений в частных производных; диффузионное поведение неявного Эйлера может привести к неточности, но с точки зрения поведения похоже на диффузию, обеспечиваемую уравнением теплопроводности, и может быть приемлемо даже при довольно больших временных шагах. Гиперболические УЧП могут потребовать неявных шагов для стабильности, но расширенные интеграторы, такие как «симплектические интеграторы», могут предотвратить чрезмерное сглаживание, вызванное этими типами шагов.

Один противоположный подход состоит в том, чтобы записывать решения полудискретных систем f = L f через собственные векторы L. Предположим, что $v1, \ldots, vn$ — собственные векторы оператора L c собственными значениями $\lambda 1, \ldots, \lambda n$ и что мы знаем $f(0) = c1v1 + \cdots$ + cnvn. Затем cnvn вспомните, что решение f = L f определяется выражением:

$$\lambda i t$$
 cie vi

В этой формуле нет ничего нового, кроме §5.1.2, который мы ввели во время обсуждения собственных векторов и собственных значений. Однако собственные векторы оператора L могут иметь физический смысл в случае полудискретного УЧП, как в примере 14.5, где показано, что собственные векторы лапласиана L соответствуют различным резонансным колебаниям области. Таким образом, этот подход на основе собственных векторов можно применять, например, для разработки «низкочастотных аппроксимаций» исходных данных путем усечения приведенной выше суммы по i, с тем преимуществом, что зависимость t известна точно без временного шага.

Пример 14.10 (собственные функции лапласиана). На рисунке NUMBER показаны собственные векторы матрицы L из примера 14.7. Собственные векторы с низкими собственными значениями соответствуют низкочастотным функциям на [0, 1] со значениями, фиксированными на концах, и могут быть хорошими приближениями f (x), когда он относительно гладкий.

Полностью дискретные методы В качестве альтернативы мы могли бы обращаться с переменными пространства и времени более демократично и дискретизировать их обе одновременно. Эта стратегия дает систему уравнений для решения, больше похожую на §14.5.1. Этот метод легко сформулировать, распараллелив эллиптический случай, но результирующие линейные системы уравнений могут быть большими, если зависимость между временными шагами имеет глобальный охват.

Пример 14.11 (Полностью дискретная диффузия тепла). Явный, неявный, Крэнк-Николсон. Не входит в СS 205А.

Важно отметить, что, в конце концов, даже полудискретные методы можно считать полностью дискретными в том смысле, что метод ОДУ с временным шагом по-прежнему дискретизирует переменную t; разница в основном заключается в классификации того, как были получены методы. Однако одним из преимуществ полудискретных методов является то, что они могут корректировать временной шаг для t в зависимости от текущей итерации, например, если объекты движутся быстро в физическом моделировании, может иметь смысл сделать больше временных шагов и разрешить это движение. Некоторые методы даже регулируют дискретизацию области значений x в случае, если требуется большее разрешение вблизи локальных разрывов или других артефактов.

14.6 Метод конечных элементов

Не входит в 205А.

14.7 Практические примеры

Вместо тщательного рассмотрения всех часто используемых методов PDE в этом разделе мы приводим примеры того, где они применяются на практике в компьютерных науках.

- 14.7.1 Обработка изображений в области градиента
- 14.7.2 Фильтрация с сохранением границ
- 14.7.3 Жидкости на основе сетки

14.8 Сделать

- Подробнее о существовании/уникальности
- Условия КЛЛ
- Теорема Лакса об эквивалентности
- Последовательность, стабильность и друзья

14.9 Проблемы

- Показать, что 1d лапласиан может быть факторизован как DD для первой производной матрицы D
- Решить УЧП первого порядка

Благодарности

[Обсуждение идет здесь]

Я очень благодарен студентам Стэнфордского университета CS 205A, осенью 2013 г., за обнаружение многочисленных опечаток и ошибок при написании этой книги. Ниже приводится, несомненно, неполный список студентов, внесших свой вклад в эту работу: Тао Ду, Леннарт Янссон, Майлз Джонсон, Люк Неппер, Минджэ Ли, Ниша Машарани, Джон Рейна, Уильям Сонг, Бен-Хан Сун, Мартина Троеш, Ожан. Тургут, Патрик Уорд, Чжунъёб Ё и Ян Чжао.

Особая благодарность Яну Хейланду и Тао Ду за помощь в разъяснении вывода алгоритма BFGS.

[Подробнее обсуждение здесь]