

# Классическое машинное обучение

Минкин Даниэль

9 февраля 2025 г.

## Содержание

<b>1</b>	<b>Линейные модели</b>	<b>1</b>
1.1	Регрессия . . . . .	1
1.2	Классификация . . . . .	1
1.3	Общее . . . . .	2
1.4	Оценка по МНК . . . . .	2
1.4.1	Общее . . . . .	2
1.4.2	Невырожденность матрицы $X^T X$ . . . . .	3
1.4.3	Переход к новому скалярному произведению . . . . .	9

## 1 Линейные модели

Линейные модели — класс моделей которые используют линейное преобразование для вектора входных фичей.

### 1.1 Регрессия

Формализуем задачу регрессии, пусть у нас есть вектор  $\bar{x} \in \mathbb{R}^n$ . Тогда предсказание может быть сделано с помощью такой формулы:

$$y_{pred} = \bar{x} \cdot \bar{w} + w_0 \quad (1)$$

Т.е мы ищем такой вектор  $\bar{w} \in \mathbb{R}^{n+1}$ , который будет выдавать наиболее близкие  $y_{pred}$  к  $y_{true}$ .

### 1.2 Классификация

В случае решения задачи классификации через линейные модели сначала делается предсказание как в случае регрессии, а после к результату предсказания применяется разделяющее правило (например в случае бинарной классификации правило может задаваться так: если  $y_{pred} > 0$  мы относим объект к положительному классу — если нет, то к отрицательному)

В случае бинарной классификации с разделяющим правилом из примера уравнение регрессии задает гиперплоскость, которая разделяет исходное

пространство: i.e  $\sum_{i=1}^n w_i \cdot x_i + c > 0$  — плоскость в  $n$ -мерном пространстве, которая делит пространство на положительный и отрицательные классы

### 1.3 Общее

Несколько фактов:

- **При использовании OneHot Encoding-а мы можем избавиться от одной encoded фичи.** Все просто: пусть у нас есть веса для каждой закодированной фичи  $w_1, w_2 \dots w_n$ , также у нас добавляется константа  $c$  к предсказанию по фичам. Давайте удалим последнюю фичу, тогда в случае если  $x_1 = 0 \dots x_{n-1} = 0$  нам нужно добавить к константе еще и  $w_n$ , иначе результат изменится, следовательно константа в новой модели должна быть равна  $w_n + c$ . Однако тогда нам нужно внести поправку в веса в случае если один из  $x_i \mid i < n$  равен 1, чтобы результат остался таким же. Мы можем просто вычесть из старых весов  $w_n$ , за счет того, что мы добавили его к const ничего не изменится. Таким образом мы успешно исключили одну encoded фичу, оставив результаты предсказаний без изменений. **Важно заметить, что если мы работаем в модели без константы, то тогда данный подход не сработает**
- **Для более сложных зависимостей необходимо использовать новые фичи которые являются функциями от старых.** Т.е мы включаем в модель фичи задаваемые как  $f(x_1, x_2 \dots x_n)$
- Если между признаками есть приближённая линейная зависимость, коэффициенты в линейной модели могут совершенно потерять физический смысл

### 1.4 Оценка по МНК

#### 1.4.1 Общее

В первую очередь опустим свободный член, так как можно считать, что у нас просто есть доа признак, который всегда равен 1. Пусть функция потерь задается как Евклидова норма между предсказанными и истинными значениями. Т.е мы решаем следующую задачу:

$$\|Xw - b\|_{euclidean} \rightarrow \min_w \quad (2)$$

где  $X$  — матрица размера  $(N, k)$ ,  $N$  — размер выборки, а  $k$  — кол-во фичей, т.е это матрица где в строки записаны вектора, по которым нужно сделать предсказания.

Однако нам также нужно сделать поправку на размер выборки, чтобы значения функции потерь можно было сравнивать между собой для разных выборок. Получается задача выглядит так:

$$\frac{\|Xw - b\|_{euclidean}}{N} \rightarrow \min_w \quad (3)$$

**Функция потерь является функционалом, так как принимает на вход три значения — матрицу наблюдений, true значения предсказываемой переменной и функцию, которая возвращает некое значение по вектору наблюдений**

Значение коэффициентов может быть получено через псевдообратную матрицу, по ее свойству:

$$\|XX^+b - b\|_{euclidean} \leq \|Xw - b\|_{euclidean} \quad (4)$$

для любого  $w$

Так как деление на константу — монотонное преобразование, данное решение минимизирует нашу функцию потерь

**Важно заметить:** псевдообратная матрица может быть использована и для других норм, однако тогда нам нужно перейти в евклидово пространство  $A$  с новым скалярным произведением, тогда псевдообратная матрица будет минимизировать норму задаваемую как  $\sqrt{\langle u, u \rangle_A}$ . Об этом будет рассказано позже

Таким образом решением данного СЛАУ будет

$$w_* = A^+b \quad (5)$$

В случае если  $N \geq k$

$$w_* = (X^T X)^{-1} X^T b \quad (6)$$

А в противном случае, когда  $N < k$

$$w_* = X^T (X X^T)^{-1} b \quad (7)$$

Мы можем считать, что  $X$  — матрица полного столбцового или строчного ранга, зачастую у нас не будет полностью ЛЗ столбцов или строк, а даже если они и есть их можно исключить из-за бессмысленности

#### 1.4.2 Невырожденность матрицы $X^T X$

В реальных задачах матрицах  $X^T X$  или  $X X^T$  являются невырожденными, однако нам нужно оценить насколько они “невырождены”, так как у нас есть погрешность при вычислении детерминанта, например мы могли получить неотрицательное значение из-за логики работы чисел с плавающей точкой или же в изначальных данных содержится погрешность.

#### Число обусловленности

Число обусловленности — максимальное значение отношения относительного изменения функции и относительного изменения аргумента

$$\mu(f, x) = \max_{\Delta x} \frac{\frac{\|f(x+\Delta x) - f(x)\|}{\|f(x)\|}}{\frac{\|\Delta x\|}{\|x\|}} \quad (8)$$

в малой окрестности  $\Delta x$

Легко заметить, что чем больше значение  $\mu(f, x)$ , тем хуже так как наше решение может очень сильно измениться от малого приращения аргумента, а мы бы хотели иметь “стабильное” решение

Покажем на зависимость числа обусловленности задачи  $f(A) = (A^T A)^{-1}$  от матрицы корреляции. Мы воспользуемся определением числа обусловленности:

$$\mu(f, A) = \max_{\Delta A} \frac{\frac{\|f(A+\Delta A) - (A^T A)^{-1}\|}{\|(A^T A)^{-1}\|}}{\frac{\|\Delta A\|}{\|A\|}} \quad (9)$$

Тут нам понадобится матричное дифференцирование, нам нужно найти дифференциал  $(A^T A)^{-1}$ , заметим, что это композиция двух функций. Во-первых, вспомним определение дифференциала, пусть задана функция  $f : R^n \rightarrow R^m$  Тогда ее дифференциал может быть найден так:

$$f(x_0 + \Delta x) - f(x_0) = [D_{x_0} f](\Delta x) + o(\|\Delta x\|) \quad (10)$$

При этом дифференциал  $[D_{x_0} f]$  — линейное отображение из  $R^n$  в  $R^m$ , т.е. мы находим линейную часть приращения отображения

Найдем дифференциал  $f(X) = X^T X$  где  $X$  - матрица  $(n, m)$  и  $n > m$ . Тогда:

$$\begin{aligned} f(X + \Delta X) - f(X) &= (X + \Delta X)^T (X + \Delta X) - X^T X \\ &= (X^T + (\Delta X)^T)(X + \Delta X) - X^T X \\ &= X^T X + X^T \Delta X + (\Delta X)^T X + (\Delta X)^T \Delta X - X^T X \\ &= X^T \Delta X + (\Delta X)^T X + (\Delta X)^T \Delta X \end{aligned}$$

Можно легко показать, что  $[D_X f](\Delta X) = X^T \Delta X + (\Delta X)^T X$  является линейным оператором по  $\Delta X$ . При этом  $(\Delta X)^T \Delta X$  и есть  $o(\|\Delta X\|)$ .

Найдем дифференциал  $g(X) = X^{-1}$ . Мы будем использовать равенство  $E = X^{-1} X$ . Продифференцируем выражение с обеих сторон.

$$0 = [D_{X_0}(X^{-1} X)](H) \quad (11)$$

По правилам дифференцирования умножения мы получаем

$$0 = [D_{X_0}(X^{-1})](H) \cdot X_0 + X_0^{-1} \cdot [D_{X_0}(X)](H) \quad (12)$$

Таким образом

$$[D_{X_0}(X^{-1})](H) \cdot X_0 = -X_0^{-1} \cdot [D_{X_0}(X)](H) \quad (13)$$

То все сводится к формуле

$$[D_{X_0}(X^{-1})](H) = -X_0^{-1} \cdot [D_{X_0}(X)](H) \cdot X_0^{-1} \quad (14)$$

При этом  $[D_{X_0}(X)](H) = H$ , таким образом итоговая формула имеет вид

$$[D_{X_0}(X^{-1})](H) = -X_0^{-1} \cdot H \cdot X_0^{-1} \quad (15)$$

Мы имеем дело с функцией  $g(f(X))$ , дифференциал такой функции представим как

$$[D_{f(X_0)}(g)]([D_{X_0}(f)](\Delta X)) \quad (16)$$

Следовательно  $[D_{X_0}((X^T X)^{-1})](\Delta X)$ , может быть найдено так

$$(X_0^T X_0)^{-1} \cdot (X_0^T \Delta X + (\Delta X)^T X_0) \cdot (X_0^T X_0)^{-1} \quad (17)$$

Таким образом в изначальной формуле числа обусловленности мы можем записать приращение функции как

$$\mu(A) = \max_{\Delta A} \frac{\frac{\|(A^T A)^{-1} \cdot (A^T \Delta A + (\Delta A)^T A) \cdot (A^T A)^{-1}\|}{\|(A^T A)^{-1}\|}}{\frac{\|\Delta A\|}{\|A\|}} \quad (18)$$

Оценим наше выражение сверху, сделаем несколько предположений на будущее, **пусть мы используем одно и тоже семейство норм для входного и выходного пространства и для этого семейства выполняется свойство субмультипликативности**, тогда для  $(A^T A)^{-1}$  и  $(A^T \Delta A + (\Delta A)^T A)$  можно будет использовать следующее мажорирование

$$\|(A^T A)^{-1}(A^T \Delta A + (\Delta A)^T A)(A^T A)^{-1}\| \leq \|(A^T A)^{-1}\|^2 \|A^T \Delta A + (\Delta A)^T A\| \quad (19)$$

В итоге мы можем мажорировать наше выражение так:

$$\begin{aligned} & \frac{\frac{\|(A^T A)^{-1} \cdot (A^T \Delta A + (\Delta A)^T A) \cdot (A^T A)^{-1}\|}{\|(A^T A)^{-1}\|}}{\frac{\|\Delta A\|}{\|A\|}} \leq \\ & \frac{\|(A^T A)^{-1}\| \cdot \|A^T \Delta A + (\Delta A)^T A\|}{\frac{\|\Delta A\|}{\|A\|}} \end{aligned}$$

Продолжая мажорирование мы получим следующую ситуацию

$$\frac{\|(A^T A)^{-1}\| \cdot \|A^T \Delta A + (\Delta A)^T A\|}{\frac{\|\Delta A\|}{\|A\|}} \leq \frac{\|(A^T A)^{-1}\| \cdot (\|A^T \Delta A\| + \|(\Delta A)^T A\|)}{\frac{\|\Delta A\|}{\|A\|}}$$

Предположим, что норма выходного пространства устойчива к транспонированию, также вспомним условие про субмультипликативность, тогда наше выражение примет вид

$$\begin{aligned} & \frac{\|(A^T A)^{-1}\| \cdot (\|A^T \Delta A\| + \|(\Delta A)^T A\|)}{\frac{\|\Delta A\|}{\|A\|}} = \\ & 2 \cdot \frac{\|(A^T A)^{-1}\| \cdot \|A^T \Delta A\|}{\frac{\|\Delta A\|}{\|A\|}} \leq \\ & 2 \cdot \frac{\|(A^T A)^{-1}\| \cdot \|A^T\| \cdot \|\Delta A\|}{\frac{\|\Delta A\|}{\|A\|}} = \\ & 2 \cdot \frac{\|(A^T A)^{-1}\| \cdot \|A^T\|}{\frac{1}{\|A\|}} = \\ & 2 \cdot \|(A^T A)^{-1}\| \cdot \|A\|^2 \end{aligned}$$

Таким образом, при определенных требованиях к семейству норм мы получаем, что

$$\max_{\Delta A} \frac{\frac{\|(A^T A)^{-1} \cdot (A^T \Delta A + (\Delta A)^T A) \cdot (A^T A)^{-1}\|}{\frac{\|\Delta A\|}{\|A\|}}}{2 \cdot \|(A^T A)^{-1}\| \cdot \|A\|^2} \leq$$

Вместо точного максимума мы можем использовать оценку сверху для всех значений.

Покажем влияние на данный результат близости матрицы к “вырожденной”. Применим сингулярное разложение к матрице  $A$ , если она имеет размеры  $(n, m)$ , где  $n > m$

$$A = U_{(n,n)} \Sigma_{(n,m)} V_{(m,m)}^T \quad (20)$$

$\Sigma$  — диагональная матрица где на диагонали находятся сингулярные числа матрицы, а  $U$  и  $V$  — унитарные матрицы, тогда  $A^T A$  равно  $V \Sigma^T U^T U \Sigma V^T = V \Sigma^T E \Sigma V^T$ . Тогда:

$$(A^T A)^{-1} = (V \Sigma^T E \Sigma V^T)^{-1} \quad (21)$$

При этом  $\Sigma^T E \Sigma$  — матрица  $(m, m)$ , если матрица  $A$  имеет размер  $(n, m)$ , где  $n > m$ . Так как матрица  $A$  имеет  $m$  сингулярных значений и является

матрицей полного столбцового ранга — матрица  $\Sigma^T \Sigma$  является обратимой. Продолжим раскрывать наше выражение

$$(A^T A)^{-1} = V(\Sigma^T \Sigma)^{-1} V^T \quad (22)$$

Мажорируем наше старое выражение, оно будет равно

$$2 \cdot \|(A^T A)^{-1}\| \cdot \|A\|^2 \leq 2 \cdot \|(\Sigma^T \Sigma)^{-1}\| \cdot \|V\|^2 \cdot \|A\|^2 \quad (23)$$

$\Sigma^T \Sigma$  — квадратная матрица на диагонали у которой находятся квадраты сингулярных значений матрицы, следовательно, обратная матрица к ней будет иметь вид  $\text{diag}(\frac{1}{\sigma_1^2}, \frac{1}{\sigma_2^2} \dots \frac{1}{\sigma_m^2})$

$$2 \cdot \|(\Sigma^T \Sigma)^{-1}\| \cdot \|V\|^2 \cdot \|A\|^2 = 2 \cdot \|\text{diag}(\frac{1}{\sigma_1^2}, \frac{1}{\sigma_2^2} \dots \frac{1}{\sigma_m^2})\| \cdot \|V\|^2 \cdot \|A\|^2 \quad (24)$$

При этом сингулярные числа показывают близость матрицы к невырожденной, или же близость к матрице полного столбцового/строкового ранга, если матрица  $A$  имеет сильную зависимость между столбцами то  $\min\{\sigma_1, \sigma_2, \dots\} \rightarrow 0$ , что ведет к увеличению числа обусловленности (сингулярные числа рассмотрю позже)

### Второй подход к числу обусловленности

Мы можем сделать проще и рассмотреть число обусловленности матрицы  $A$  как линейного оператора. Рассмотрим уравнение:

$$Ax = b \quad (25)$$

Тогда решением в общем виде является

$$x = A^+ b \quad (26)$$

Рассмотрим число обусловленности

$$\mu(A, b) = \max_{\Delta b} \frac{\|A^+(b+\Delta b) - A^+ b\|}{\frac{\|b+\Delta b - b\|}{\|b\|}} = \max_{\Delta b} \frac{\frac{\|A^+ \Delta b\|}{\|A^+ b\|}}{\frac{\|\Delta b\|}{\|b\|}} \quad (27)$$

Перенеся деление

$$\max_{\Delta b} \frac{\|A^+ \Delta b\| \cdot \|b\|}{\|\Delta b\| \cdot \|A^+ b\|} = \max_{\Delta b} \left( \frac{\|A^+ \Delta b\|}{\|\Delta b\|} \right) \cdot \left( \frac{\|b\|}{\|A^+ b\|} \right) \quad (28)$$

**Аналогично прошлому случаю потребуем свойство субмультипликативности от нормы**, тогда мы можем мажорировать первый множитель с помощью  $\|A^+\|$ . Следовательно:

$$\mu(A, b) = \|A^+\| \frac{\|b\|}{\|A^+ b\|} = \|A^+\| \frac{\|Ax + \epsilon\|}{\|x\|} \quad (29)$$

где  $\epsilon$  - вектор отклонений

$$\|A^+\| \frac{\|Ax + \epsilon\|}{\|x\|} \leq \|A^+\| \frac{\|Ax\| + \|\epsilon\|}{\|x\|} \quad (30)$$

Таким образом мы получаем

$$\|A^+\| \frac{\|Ax\|}{\|x\|} + \|A^+\| \frac{\|\epsilon\|}{\|x\|} \leq \|A^+\| \frac{\|Ax\|}{\|x\|} + \|A^+\| \frac{\|\epsilon\|}{\|x\|} \quad (31)$$

Финальное выражение выглядит так:

$$\|A^+\| \|A\| + \|A^+\| \frac{\|\epsilon\|}{\|x\|} \quad (32)$$

В случае если используется сингулярная норма: Сингулярная норма матрицы — ее максимальное сингулярное значение, а псевдообратная матрица имеет обратные сингулярные значения к исходной матрице (легко показать), таким образом:

$$\|A\|_{spec} = \sigma_{max}(A) \quad (33)$$

и

$$\|A^+\|_{spec} = \frac{1}{\sigma_{min}(A)} \quad (34)$$

Таким образом мы получаем:

$$\mu(A, b) = \frac{\sigma_{max}(A)}{\sigma_{min}(A)} + \|A^+\| \frac{\|\epsilon\|}{\|x\|} \quad (35)$$

Если пренебречь вторым слагаемым, то мы получим то, что описано в учебнике ШАДА: “Пожертвовав математической строгостью, мы можем считать, что число обусловленности матрицы  $X$  — это корень из отношения наибольшего и наименьшего из собственных чисел матрицы  $X^T X$ ”.

**Пояснение:** Сингулярные значения матрицы  $A$  — это квадратные корни из собственных значений матрицы  $A^T A$

### Intuition проблемы числа обусловленности

Матрица  $A$  — линейное отображение из одного пространства в другое, если в новом пространстве есть вектора с сильной линейной связью, то оно становится более сжатым, т.е. непохожие вектора в исходном пространстве могут стать сильно более похожими в новом пространстве за счет того, что оно сжато в размерах, таким образом если мы немного изменим целевой вектор (из нового пространства), то вектор который его образовал может сильно отличаться от того, что мы получили в прошлый раз, так как за счет сжатия они лежат рядом в новом пространстве, однако не в новом, таким образом: число обусловленности — лишь следствие того, что новое пространство сжато



Для наглядности можно рассмотреть, двумерный случай, рассмотрим матрицу перехода в новое пространство

$$\begin{pmatrix} 1 & 1 \\ 1 & 1 + \epsilon \end{pmatrix}$$

где  $\epsilon > 0$

Мы видим, что  $\begin{bmatrix} 1 \\ 0 \end{bmatrix}$  отобразится в  $\begin{bmatrix} 1 \\ 1 \end{bmatrix}$ , а  $\begin{bmatrix} 0 \\ 1 \end{bmatrix}$  в  $\begin{bmatrix} 1 \\ 1 + \epsilon \end{bmatrix}$ . При маленьких  $\epsilon$  матрица остается невырожденной, однако пространство “сжимается” делая  $\begin{bmatrix} 1 \\ 0 \end{bmatrix}$  и  $\begin{bmatrix} 0 \\ 1 \end{bmatrix}$  все менее отличимыми в новом пространстве

### 1.4.3 Переход к новому скалярному произведению

Мы хотим перейти в пространство с новым скалярным произведением, однако большинство компьютерных систем заточены под обычное Евклидово скалярное произведение, поэтому нам бы хотелось перейти в такой базис, чтобы операции в этом базисе с использованием “обычного” скалярного произведения были равносильны операциям с использованием “нового” скалярного произведения в исходном пространстве

Заметим, что любое произвольное скалярное произведение может быть задано как квадратичная форма, исходя из свойства линейности.

$$\langle u, v \rangle_G = u^T G v \quad (36)$$

Заметим, что  $G$  — положительно определена и симметрична, по свойству скалярного произведения. Мы можем найти матрицу перехода  $P$  таким образом, через разложение Холецкого.

$$G = P^T P \quad (37)$$

Тогда вектор  $k$  и матрица  $M$  в новом пространстве превращаются в  $Pk$  и  $PMP^T$  соответственно

#### Пример

Пусть мы работаем в новом пространстве и хотим умножить матрицу на вектор в новом пространстве. У нас есть  $k' = Pk$  и  $M' = PMP^T$ , мы хотим, чтобы итоговый вектор имел вид:

$$Mk = \begin{bmatrix} \langle M_{1,*}, k \rangle_G \\ \langle M_{2,*}, k \rangle_G \\ \vdots \\ \langle M_{n,*}, k \rangle_G \end{bmatrix}$$

это равносильно:

$$Mk = \begin{bmatrix} M_{1,*}^T Gk \\ M_{2,*}^T Gk \\ \vdots \\ M_{3,*}^T Gk \end{bmatrix}$$

следовательно

$$[Mk]_{new \ dot \ product} = MGk \quad (38)$$

Рассмотрим линейный оператор  $MG$ , тогда в новом базисе данный оператор может быть представлен как  $PMGP^{-1}$

$$PMGP^{-1} = PMP^T PP^{-1} = PMP^T \quad (39)$$

Проверим, что это работает

$$P^{-1}(PMP^T Pk) = P^{-1}(PMGk) = MGk \quad (40)$$

Таким образом

$$P^{-1}M'k' = MGk \quad (41)$$

Покажем, что это сработает и для умножения на вектор справа

$$k'^T M' P^{-1T} = k^T P^T PMP^T P^{-1T} = k^T GM \quad (42)$$

Таким образом мы выполнили поставленную перед нами задачу