

# Some Math for ML

Шмидт Ян  
[Мой телеграм](#)  
[Лендинг](#)  
[RoadMap](#)

12 февраля 2025 г.



# Оглавление

<b>1</b>	<b>Дисклеймер</b>	<b>5</b>
<b>2</b>	<b>Алгебра</b>	<b>7</b>
2.1	Нотация суммирования и произведения . . . . .	7
2.2	Функции . . . . .	7
2.3	Линейная комбинация . . . . .	8
2.4	Линейная зависимость и независимость . . . . .	8
2.5	Скалярное произведение векторов . . . . .	8
2.6	Что такое полином, количество корней полинома над полем $\mathbb{C}$ . . . . .	9
2.7	Матрицы . . . . .	9
2.7.1	Ранг матрицы . . . . .	9
2.7.2	Элементарные преобразования над матрицами . . . . .	9
2.7.3	Виды матриц . . . . .	10
2.7.4	Решение систем линейных уравнений (СЛУ) через матрицы . . . . .	10
2.7.5	Определитель матрицы . . . . .	10
2.8	Собственные числа, собственные векторы . . . . .	11
<b>3</b>	<b>Математический анализ I</b>	<b>13</b>
3.1	Кванторы . . . . .	13
3.2	Интуитивное понятие предела $\lim$ . . . . .	13
3.3	Производные функций (определение, свойства) . . . . .	13
3.4	Интеграла Римана . . . . .	14
3.5	Экстремум функции . . . . .	14
3.6	Разложение функций в ряд Тейлора . . . . .	15
3.7	Ряды Фурье . . . . .	15
3.8	Частные производные . . . . .	15
3.9	Градиент . . . . .	16
3.10	Лапласиан . . . . .	16
3.11	Матрица Якоби . . . . .	17
3.12	Гессиан . . . . .	17
3.13	Метрика (расстояние) . . . . .	18
3.14	Норма . . . . .	18
<b>4</b>	<b>Теория вероятностей и математическая статистика</b>	<b>21</b>
4.1	Основные понятия . . . . .	21
4.2	Типы случайных величин: дискретные и непрерывные . . . . .	21
4.3	Плотность распределения случайной величины . . . . .	21
4.4	Некоторые распределения . . . . .	22
4.4.1	Нормальное (гауссовское) . . . . .	22

4.4.2	Экспоненциальное . . . . .	22
4.4.3	Бернулли . . . . .	22
4.4.4	Пуассона . . . . .	22
4.5	Математическое ожидание . . . . .	23
4.6	Дисперсия . . . . .	23
4.7	Корреляция . . . . .	23
4.8	Выборка, генеральная совокупность . . . . .	23
<b>5</b>	<b>Численные методы. Методы оптимизации</b>	<b>25</b>
5.1	Определение функционала . . . . .	25
5.2	Что такое минимизация функционала . . . . .	25
5.3	Градиент. Градиентный спуск . . . . .	25
<b>6</b>	<b>Теоретическая информатика</b>	<b>27</b>
6.1	Энтропия Шеннона . . . . .	27

# Глава 1

## Дисклеймер

**Обозначения.** Красным текстом обозначены штуки, которые можно почитать по желанию, они не обязательны.

**Параграфы «Дополнительно».** В параграфах «Дополнительно» набросал всякое для прошаренных и тех, кто вдруг захочет когда-нибудь с головой окунуться в матешу. Некоторое из перечисленного в этих параграфах 99% не пригодится на твоей будущей работе в качестве МЛщика, это может пригодиться только лютым ризерчерам, которые с утра до вечера пилят кастомные модели, разрабатывают качественно новые алгоритмы — да, там может пригодиться понимание подкапотных тонкостей. Этими темами можно покрыть хорошую часть курса крепких технических вузов. При продвижении по грейдам и вкате в МЛ это можно спокойно скипать

**Литература.** То, что я указываю в «Литературе», — это тоже больше относится к «Дополнительному». Можно спокойно ограничиться видосами на ютубе (ODS, 3b1b, etc.), википедией (желательно, англоязычной), курсами на Stepik, промптингом нейронки типа чатгпт, deepseek — интернетом, короче.



# Глава 2

## Алгебра

Литература<sup>1 2 3 4 5</sup>

### 2.1 Нотация суммирования и произведения

$$\begin{aligned} a_1 + a_2 + \dots + a_n &= \sum_{i=1}^n a_i \\ a_1 \cdot a_2 \cdot \dots \cdot a_n &= \prod_{i=1}^n a_i \end{aligned} \tag{2.1}$$

### 2.2 Функции

Функцией  $f$  называется такое правило, по которому мы задаем соответствие рассматриваемому множеству  $X$ :

$$\begin{aligned} f(X) &= Y \\ X &\mapsto Y \end{aligned} \tag{2.2}$$

$X$  называют *прообразом*, а  $Y$  — *образом*. Например, функция  $f(x) = x^2$  каждому числу  $x$  сопоставляет  $x^2$ :

$$x \mapsto x^2 \tag{2.3}$$

Числу 2 будет соответствовать  $2^2 = 4$ , то есть

$$2 \mapsto 4, \tag{2.4}$$

пятерке —  $5^2 = 25$ , то есть

$$5 \mapsto 25, \tag{2.5}$$

5 — прообраз, 25 — образ для функции  $f(x) = x^2$ , и так далее.

---

<sup>1</sup>Задачник по высшей алгебре Фаддеева

<sup>2</sup>Конспекты по алгебре Вавилова ([Alg\\_lectures\\_II](#))

<sup>3</sup>Вообще конспекты Вавилова и лабы им. Чебышева

<sup>4</sup>[Сайт](#) с литературой по алгебре

<sup>5</sup>[Огромный список литературы](#) и конспектов по различным курсам вышмата, который читается на матмехе и в Чебышевке (hard)

## 2.3 Линейная комбинация

**Линейная комбинация** — выражение, построенное на множестве элементов путём умножения каждого элемента на коэффициенты с последующим сложением результатов (например, линейной комбинацией  $x$  и  $y$  будет выражение вида

$$ax + by, \quad (2.6)$$

где  $a$  и  $b$  — коэффициенты).

Если  $v_1, \dots, v_n$  — векторы, а  $a_1, \dots, a_n$  — скаляры, то линейная комбинация этих векторов со скалярами в качестве коэффициентов — это:

$$a_1v_1 + a_2v_2 + a_3v_3 + \dots + a_nv_n = \sum_{i=1}^n a_iv_i \quad (2.7)$$

## 2.4 Линейная зависимость и независимость

При **линейной зависимости** существует нетривиальная линейная комбинация элементов этого множества, равная нулевому элементу: для векторов  $v_1, \dots, v_n$  найдутся такие числа  $a_1, \dots, a_n$ , не равные одновременно нулю, что

$$\sum_{i=1}^n a_iv_i = 0 \quad (2.8)$$

При отсутствии такой комбинации, то есть, когда коэффициенты линейной комбинации равны нулю, векторы называются **линейно независимым**. Иными словами, если

$$\sum_{i=1}^n a_iv_i = 0 \quad (2.9)$$

выполняется только при  $a_i = 0$ ,  $i = \overline{1, n}$ .

## 2.5 Скалярное произведение векторов

Скалярное произведение — полезная штука в МЛ, потому что может встречаться, например, в определении косинусного расстояния.

$$\langle x, y \rangle = |x||y| \cos(\widehat{x, y}) = \sum_{i=1}^n x_i y_i, \quad (2.10)$$

$$x, y \in \mathbb{R}^n,$$

$\widehat{x, y}$  — угол между векторами  $x, y$ .

Свойства:

- $\mathbf{a} \cdot \mathbf{b} = \mathbf{b} \cdot \mathbf{a}$  — коммутативность
- $(\alpha \mathbf{a} + \beta \mathbf{b}) \cdot (\gamma \mathbf{c} + \delta \mathbf{d}) = \alpha \gamma (\mathbf{a} \cdot \mathbf{c}) + \alpha \delta (\mathbf{a} \cdot \mathbf{d}) + \beta \gamma (\mathbf{b} \cdot \mathbf{c}) + \beta \delta (\mathbf{b} \cdot \mathbf{d})$  — билинейность



## 2.6 Что такое полином, количество корней полинома над полем $\mathbb{C}$

У полинома

$$f(x) = a_0x^n + a_1x^{n-1} + \dots + a_{n-1}x + a_n \quad (2.11)$$

всегда ровно  $n$  корней с учетом кратностей (какие-то корни могут иметь кратность выше 1). Корни могут быть как  $\mathbb{R}$ , так и  $\mathbb{C}$ .

## 2.7 Матрицы

$$A = \begin{pmatrix} a_{11} & \dots & a_{1n} \\ \vdots & \ddots & \vdots \\ a_{n1} & \dots & a_{nn} \end{pmatrix} \quad (2.12)$$

Матрицы можно складывать-вычитать, умножать, обращать (= делить), транспонировать. Пример транспонирования:

$$A = \begin{pmatrix} a_{11} & a_{12} & a_{13} \\ a_{21} & a_{22} & a_{23} \\ a_{31} & a_{32} & a_{33} \end{pmatrix}^T = \begin{pmatrix} a_{11} & a_{21} & a_{31} \\ a_{12} & a_{22} & a_{32} \\ a_{13} & a_{23} & a_{33} \end{pmatrix} \quad (2.13)$$

### 2.7.1 Ранг матрицы

Строки и столбцы матрицы являются элементами соответствующих векторных пространств:

- Столбцы матрицы  $A$  составляют элементы пространства размерности  $m$ ;
- Строки матрицы  $A$  составляют элементы пространства размерности  $n$ .

**Рангом матрицы** называют количество линейно независимых столбцов матрицы (столбцовый ранг матрицы) или количество линейно независимых строк матрицы (строчный ранг матрицы). Этому определению эквивалентно определение ранга матрицы как порядка максимального отличного от нуля минора матрицы.

При элементарных преобразованиях ранг матрицы не меняется.

### 2.7.2 Элементарные преобразования над матрицами

Элементарными преобразованиями строк называют:

- Перестановку местами любых двух строк матрицы;
- Умножение любой строки матрицы на обратимую ненулевую константу  $k$ ;
- Прибавление к любой строке матрицы другой строки, умноженной на некоторую константу.

Аналогично определяются элементарные преобразования столбцов. Элементарные преобразования обратимы.

Обозначение:  $A \sim B$  указывает на то, что матрица  $A$  может быть получена из  $B$  путём элементарных преобразований (или наоборот).

### 2.7.3 Виды матриц

Диагональные, квадратные, прямоугольные, треугольные, симметричные, ортогональные, блочные и многие другие. Список можно посмотреть на англоязычной странице в [вики](#).

### 2.7.4 Решение систем линейных уравнений (СЛУ) через матрицы

$$\begin{cases} a_{11}x_1 + a_{12}x_2 + \dots + a_{1n}x_n = b_1 \\ a_{21}x_1 + a_{22}x_2 + \dots + a_{2n}x_n = b_2 \\ \dots \\ a_{n1}x_1 + a_{n2}x_2 + \dots + a_{nn}x_n = b_n \end{cases} \quad (2.14)$$

Решение:

$$Ax = b, \quad (2.15)$$

$$A = \begin{pmatrix} a_{11} & \dots & a_{1n} \\ \vdots & \ddots & \vdots \\ a_{n1} & \dots & a_{nn} \end{pmatrix}_{n \times n} \quad (2.16)$$

— матрица коэффициентов,

$$b = \begin{pmatrix} b_1 \\ \vdots \\ b_n \end{pmatrix}_{n \times 1} \quad (2.17)$$

— вектор свободных коэффициентов,

$$x = A^{-1}b \quad (2.18)$$

— решение системы через обратную матрицу

### 2.7.5 Определитель матрицы

$$A = \begin{vmatrix} a_{11} & \dots & a_{1n} \\ \vdots & \ddots & \vdots \\ a_{n1} & \dots & a_{nn} \end{vmatrix} = \det(A) \quad (2.19)$$

**Определитель** матрицы полезен тем, что от его значения зависит, с какой матрицей мы имеем дело. Если  $\det A = 0$  (такая матрица называется **вырожденной**), то результат применения данной матрицы в качестве линейного оператора приведет к вырождению рассматриваемого пространства. Так, если применить вырожденную матрицу  $2 \times 2$  к плоскости  $\mathbb{R}^2$ , то она выродится в одну-единственную прямую.

**Интерпретация.** Значение определителя показывает, по сколько раз растягивается (масштабируется) преобразуемое пространство.  $\det A = 3$  будет означать, что все единичные кубы в  $\mathbb{R}^n$  будут растянуты так, что их объем увеличится в 3 раза. Если  $\det A = k < 0$ , то пространство масштабируется в  $k$  раз, но поменяет свою ориентацию. Например, если рассмотреть снова пример с плоскостью, то при отрицательном

определителе плоскость растянется или сожмется (если  $|\det A| < 1$ ) и при этом зеркально отразится: оси  $x$  и  $y$  поменяются местами. В  $\mathbb{R}^3$  чуть сложнее: в таком случае перестанет работать правило буравчика: пары осей поменяют свое положение.

Определитель не равен нулю тогда и только тогда, когда матрица обратима и соответствующее линейное отображение является изоморфизмом<sup>6</sup>:

$$\det A_{n \times n} \neq 0 \Leftrightarrow \begin{cases} \exists A^{-1} \forall X, Y \subset \mathbb{R}^n: \\ X \xrightarrow{A} Y \\ Y \xrightarrow{A^{-1}} X \end{cases} \quad (2.20)$$

## 2.8 Собственные числа, собственные векторы

Формально, пусть  $A, v$  – матрица линейного оператора (отображения) и произвольный вектор соответственно. Если после применения матрицы к вектору (отображая как-то линейно вектор), выходит, что он всего лишь масштабировался на множитель  $\lambda$ , то говорят, что  $v$  – собственный вектор, а  $\lambda$  – собственное число.

$$Av = \lambda v \quad (2.21)$$

**Как вычислять?** Вычислять с. ч. можно так:

$$\begin{aligned} \det(A - \lambda E) &= \begin{vmatrix} a_{11} - \lambda & \dots & a_{1n} \\ \vdots & \ddots & \vdots \\ a_{n1} & \dots & a_{nn} - \lambda \end{vmatrix} = \\ &= b_1 \lambda^n + b_2 \lambda^{n-1} + \dots b_{n-1} \lambda + b_n = 0, \end{aligned} \quad (2.22)$$

$b_1 \lambda^n + b_2 \lambda^{n-1} + \dots b_{n-1} \lambda + b_n = 0$  – характеристический полином,  $\lambda_i$  – собственные числа, это корни характеристического полинома,  $E$  – единичная диагональная матрица

$$E_{n \times n} = \begin{pmatrix} 1 & 0 & \dots & 0 \\ 0 & 1 & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & 1 \end{pmatrix}, \quad (2.23)$$

С.ч. и векторы можно вычислять численными методами, например, [степенным методом](#).

**Дополнительно** Если интересно, можно почитать еще про миноры матрицы, вычисление определителя через миноры, теорему Сильвестра, квадратичные и эрмитовы формы\*\*, про след матрицы, Жордановы нормальные формы\*\*, линейные и векторные пространства, прямую сумму подпространств\*\*, подпространства, многообразия\*\*, функциональные пространства\*\*, теоремы о разложении пространств в прямые суммы\*\*, корневые пространства\*\*.

<sup>6</sup>Изоморфизм – обратимое отображение между двумя структурами. Можно всегда перейти от одной структуры к другой и обратно. Удобно, когда в исходном пространстве задачу решить сложно, но мы можем перейти в другое пространство, где задача решается значительно проще, а затем вернуться обратно и получить ответ в исходных терминах. Пример – SVM (метод опорных векторов), где скалярное произведение заменяется на произвольное ядро и происходит переход в пространство с новым скалярным произведением, где разделяющая гиперплоскость линейная, а в исходном имеет весьма сложную и неочевидную структуру



# Глава 3

## Математический анализ I

Литература<sup>1 2 3</sup>

### 3.1 Кванторы

Квантор всеобщности —  $\forall$ , квантор существования —  $\exists$ ; отрицание будет такое:  $\nexists$ .  $\exists!$  — существует и единственно.

### 3.2 Интуитивное понятие предела $\lim$

Если коротко, то  $\lim$  — это предельное состояние чего-то, как следует из названия. Можно рассматривать числовой предел последовательности точек, предел последовательности функций, предел значений функции и так далее. Функция  $f(x) = x^2$  своим пределом будет иметь бесконечность:

$$\lim_{x \rightarrow \infty} f(x) = \lim_{x \rightarrow \infty} x^2 = \infty \quad (3.1)$$

Предел последовательности  $\left\{\frac{1}{n^2}\right\}$  при  $n \rightarrow \infty$  равен 0:  $1, \frac{1}{4}, \frac{1}{9}, \frac{1}{16}, \frac{1}{25}, \dots$  — в пределе получим нереально мелкое число, сколько угодно близкое к 0.

### 3.3 Производные функций (определение, свойства)

$$f'(x) = \lim_{\Delta x \rightarrow 0} \frac{f(x + \Delta x) - f(x)}{\Delta x}, \quad (3.2)$$

$\Delta x$  — маленькое приращение  $x$ , то есть  $x + \Delta x$  и  $x$  очень мало отличаются друг от друга (на величину  $\Delta$ ).

---

<sup>1</sup>Фихтенгольц, 3 тома, но он бывает очень душным, классический учебник, но довольно старый

<sup>2</sup>Зорич, 2 тома, менее душный, но написан чуть более тяжелым языком, чуть другой подход изложения материала, больше тем и большая глубина повествования

<sup>3</sup>Сборник задач по матану, Демидович Б.П. — классика классик, куча задач, есть решебник в виде китайского антидемировича

### 3.4 Интеграла Римана

$$\int_a^b f(x)dx = \lim_{\max \Delta x_i \rightarrow 0} \sum_{i=1}^n f(x_i) \Delta x_i \quad (3.3)$$

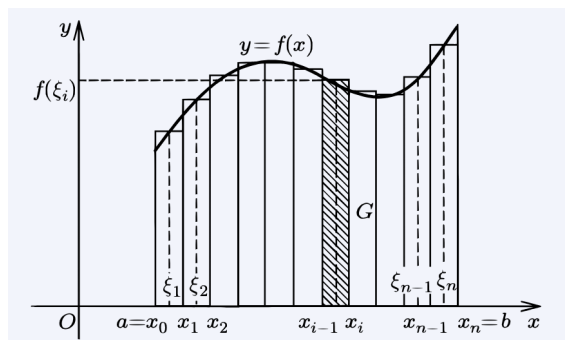
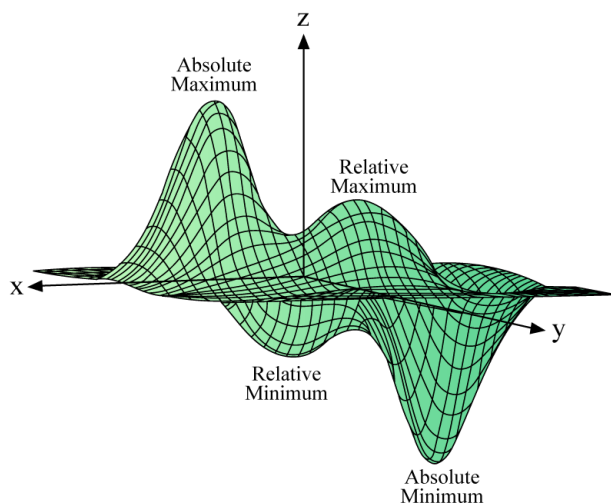


Иллюстрация интеграла. Делим область под графиком на много вертикальных столбиков, вычисляем их площади и складываем. При стремлении ширины этих прямоугольников к 0 получим интеграл. Значит, интеграл — это площадь под графиком, потому что в пределе ( $\max \Delta x \rightarrow 0$ ) мы суммируем бесконечно тонкие прямоугольники под графиком

### 3.5 Экстремум функции

**Экстремум** — максимальное или минимальное значение функции на заданном множестве. Точка, в которой достигается экстремум, называется *точкой экстремума*. Соответственно, если достигается минимум — точка экстремума называется *точкой минимума*, а если максимум — *точкой максимума*. В математическом анализе выделяют также понятие *локальный экстремум* (соответственно минимум или максимум).



Calcworkshop.com

Иллюстрация экстремума функции двух переменных

### 3.6 Разложение функций в ряд Тейлора

Это нужно для того, чтобы сложные хитрые функции представить в виде суммы более простых, потому что их можно будет проще дифференцировать и интегрировать, если их представить в виде суммы простых слагаемых.

Ряд Тейлора для функции  $f(x)$  в точке  $x = a$  можно записать следующим образом:

$$f(x) = f(a) + \frac{f'(a)(x-a)}{1!} + \dots + \frac{f^{(n)}(a)(x-a)^n}{n!} + \dots \quad (3.4)$$

$$f(x) = f(a) + \sum_{i=1}^{\infty} \frac{f^{(i)}(a)(x-a)^i}{i!} \quad (3.5)$$

### 3.7 Ряды Фурье

Та же идея, но тут функции представляются через суммы синусов и косинусов с различными коэффициентами.

$$f(x) = \frac{a_0}{2} + \sum_{n=1}^{\infty} (a_n \cos nx + b_n \sin nx), \quad (3.6)$$

$$\begin{cases} a_0 = \frac{1}{\pi} \int_{-\pi}^{\pi} f(x) dx, \\ a_n = \frac{1}{\pi} \int_{-\pi}^{\pi} f(x) \cos(nx) dx, \\ b_n = \frac{1}{\pi} \int_{-\pi}^{\pi} f(x) \sin(nx) dx, \end{cases} \quad (3.7)$$

### 3.8 Частные производные

В математическом анализе **частная производная** — одно из обобщений понятия производной на случай функции нескольких переменных. **Частная производная** — это предел отношения приращения функции **по выбранной переменной** к приращению этой переменной, при стремлении этого приращения к нулю:

$$\frac{\partial f(a_1, \dots, a_n)}{\partial x_k} = \lim_{\Delta x_k \rightarrow 0} \frac{f(a_1, \dots, a_k + \Delta x_k, \dots, a_n) - f(a_1, \dots, a_k, \dots, a_n)}{\Delta x_k} \quad (3.8)$$

Например, если  $f(x, y, z) = (x^5, y^3, z)$ , то

$$\left\{ \begin{array}{l} \frac{\partial f(x, y, z)}{\partial x} = (5x^4, y^3, z) \\ \frac{\partial f(x, y, z)}{\partial y} = (x^5, 3y^2, z) \\ \frac{\partial f(x, y, z)}{\partial z} = (x^5, y^3, 1) \\ \frac{\partial^2 f(x, y, z)}{\partial x^2} = (20x^3, y^3, z) \\ \frac{\partial f(x, y, z)}{\partial y \partial x} = (5x^4, 3y^2, z) \end{array} \right. \quad (3.9)$$

### 3.9 Градиент

**Градиент функции**  $f(\mathbf{x}) = f(x_1, \dots, x_n)$ ,  $x \in \mathbb{R}^n$  – вектор *первых* частных производных:

$$\text{grad } f(\mathbf{x}) = \nabla f(\mathbf{x}) = \begin{pmatrix} \frac{\partial f(\mathbf{x})}{\partial x_1} \\ \frac{\partial f(\mathbf{x})}{\partial x_2} \\ \vdots \\ \frac{\partial f(\mathbf{x})}{\partial x_{n-1}} \\ \frac{\partial f(\mathbf{x})}{\partial x_n} \end{pmatrix}_{n \times 1}. \quad (3.10)$$

Показывает направление наискорейшего роста функции.

Оператор  $\nabla$  называется **оператором набла**:

$$\nabla = \frac{\partial}{\partial x_1} \vec{e}_1 + \frac{\partial}{\partial x_2} \vec{e}_2 + \dots + \frac{\partial}{\partial x_n} \vec{e}_n = \sum_{i=1}^n \frac{\partial}{\partial x_i} \vec{e}_i, \quad (3.11)$$

где  $\vec{e}_1, \vec{e}_2, \dots, \vec{e}_n$  – единичные векторы по осям  $x_1, x_2, \dots, x_n$  соответственно.

### 3.10 Лапласиан

$$\Delta = \frac{\partial^2}{\partial^2 x_1} \vec{e}_1 + \frac{\partial^2}{\partial^2 x_2} \vec{e}_2 + \dots + \frac{\partial^2}{\partial^2 x_n} \vec{e}_n = \sum_{i=1}^n \frac{\partial^2}{\partial^2 x_i} \vec{e}_i \quad (3.12)$$

– **оператор Лапласа** (лапласиан), с его помощью можно составить вектор *вторых*



частных производных:

$$\Delta f(\mathbf{x}) = \begin{pmatrix} \frac{\partial^2 f(\mathbf{x})}{\partial x_1^2} \\ \frac{\partial^2 f(\mathbf{x})}{\partial x_1^2} \\ \vdots \\ \frac{\partial^2 f(\mathbf{x})}{\partial x_{n-1}^2} \\ \frac{\partial^2 f(\mathbf{x})}{\partial x_n^2} \end{pmatrix}_{n \times 1}. \quad (3.13)$$

### 3.11 Матрица Якоби

Матрица Якоби – матрица *первых* частных производных для вектор-функции

$$\mathbf{f}(\mathbf{x}) = \begin{pmatrix} f_1(\mathbf{x}) \\ f_2(\mathbf{x}) \\ \vdots \\ f_{n-1}(\mathbf{x}) \\ f_n(\mathbf{x}) \end{pmatrix}_{n \times 1}, \quad \mathbf{x} \in \mathbb{R}^n:$$

$$\mathbf{J}_f = \begin{pmatrix} \frac{\partial f_1}{\partial x_1} & \cdots & \frac{\partial f_1}{\partial x_n} \\ \vdots & \ddots & \vdots \\ \frac{\partial f_m}{\partial x_1} & \cdots & \frac{\partial f_m}{\partial x_n} \end{pmatrix}_{m \times n} \quad (3.14)$$

Определитель матрицы Якоби – **якобиан** – может возникать, например, при переходе к новым координатам (от декартовых – к сферическим и т. д.).

### 3.12 Гессиан

Гессиан, или матрица Гессе – матрица *вторых* частных производных для функции  $f(\mathbf{x})$  (не векторной, как для случая с Якобианом!):

$$\mathbf{H}_f = \begin{pmatrix} \frac{\partial^2 f}{\partial x_1^2} & \frac{\partial^2 f}{\partial x_1 \partial x_2} & \cdots & \frac{\partial^2 f}{\partial x_1 \partial x_n} \\ \frac{\partial^2 f}{\partial x_2 \partial x_1} & \frac{\partial^2 f}{\partial x_2^2} & \cdots & \frac{\partial^2 f}{\partial x_2 \partial x_n} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial^2 f}{\partial x_n \partial x_1} & \frac{\partial^2 f}{\partial x_n \partial x_2} & \cdots & \frac{\partial^2 f}{\partial x_n^2} \end{pmatrix}. \quad (3.15)$$

Используются в методах второго порядка при численных методах оптимизации функций. Считать и хранить такую матрицу очень затратно при высоких размерностях, поэтому существуют разные приближенные методы (квази-ньютоновские).

### 3.13 Метрика (расстояние)

Функция  $\rho(x, y)$  – метрика и определяет расстояние между точками  $x, y$  в  $\mathbb{R}^n$ , если

- $\rho(x, y) = 0 \Leftrightarrow x = y$  – тождество
- $\rho(x, y) \geq 0$  – неотрицательность
- $\rho(x, y) = \rho(y, x)$  – симметричность
- $\rho(x, z) \leq \rho(x, y) + \rho(y, z)$  – неравенство треугольника

Пусть  $x, y \in \mathbb{R}^n$  – точки в  $n$ -мерном пространстве (например, евклидовом). Часто используются:

- **Евклидова метрика:**

$$\rho(x, y) = \sqrt{\sum_{i=1}^n (x_i - y_i)^2} \quad (3.16)$$

- **Манхеттенская метрика (метрика городских кварталов):**

$$\rho(x, y) = \sum_{i=1}^n |x_i - y_i| \quad (3.17)$$

- **Расстояние Левенштейна**<sup>4</sup>: метрика сходства между двумя строковыми последовательностями. Чем больше расстояние, тем более различны строки. Для двух одинаковых последовательностей расстояние равно нулю. По сути, это минимальное число односимвольных преобразований (удаления, вставки или замены), необходимых, чтобы превратить одну последовательность в другую.

Например, если сравнить слова «лупа», «пу~~п~~а», то они отличаются в 1 позиции, поэтому расстояние Левенштейна = 1 (нужно сделать 1 замену, чтобы слова стали равными).

- **Расстояние Вассерштейна**: определяет расстояние между вероятностными распределениями. Чем-то похоже на KL-дивергенцию и дивергенцию Йенсена-Шеннона, но имеет некоторые теоретические преимущества. Почитать можно в [вики](#) и [тут](#), [тут](#).

### 3.14 Норма

**Норма** задает «длину» вектора или функции. Обозначается через  $\|\cdot\|$ . Сначала общий случай: пусть  $\mathbf{x} \in \mathbb{R}^n, x = (x_1, \dots, x_n)$ . Тогда  $L_p$ -норма этого вектора:

$$\|\mathbf{x}\|_p = \sqrt[p]{\sum_{i=1}^n |x_i|^p} \quad (3.18)$$

---

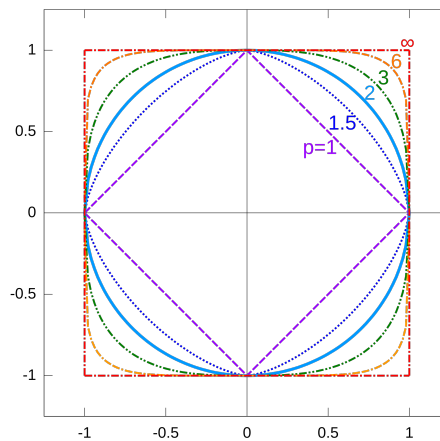
<sup>4</sup>Можно почитать на [хабре](#)

При  $p = 1$  получаем манхеттенскую норму ( $L_1$ -норму):

$$\|x\|_1 = \sum_{i=1}^n |x_i|, \quad (3.19)$$

при  $p = 2$  получаем манхеттенскую норму ( $L_2$ -норму):

$$\|x\|_2 = \sqrt{\sum_{i=1}^n x_i^2}, \quad (3.20)$$



Разные нормы, иллюстрация единичных окружностей

**Дополнительно** Функциональные пространства, Банаховы и Гильбертовы пространства, метрические пространства, нормированные пространства, метрика, норма, функциональные последовательности, непрерывность и гладкость, класс функций  $C^n$  —  $n$  раз непрерывно-дифференцируемых отображений, понятие отображения: биективность, функции многих переменных, кратные интегралы,  $\int \dots \int$ , криволинейные интегралы  $\oint$ , поверхностные интегралы, векторные поля (в.п.), ротор в.п., дивергенция в.п.

На этот случай есть крутая лит-ра<sup>5 6</sup>, но это уже больше функкан.

<sup>5</sup>Кутателадзе С. С. «Основы функционального анализа»

<sup>6</sup>Колмогоров, Фомин. «Элементы теории функций и функционального анализа»



## Глава 4

# Теория вероятностей и математическая статистика

Литература<sup>1 2</sup>

### 4.1 Основные понятия

**Случайная величина** — это результат случайного эксперимента. Обозначают часто как  $\xi$  (кси). Например, случайная величина — подбрасывание монетки. **Вероятность** изменяется на промежутке  $p \in [0, 1]$ , где  $p = 0$  — событие точно не произойдет,  $p = 1$  — сто процентов произойдет. Вероятность «50/50» —  $p = 0.5$ .

### 4.2 Типы случайных величин: дискретные и непрерывные

- **Дискретные:** подбрасывание монеты, количество звонков в день, ...
- **Непрерывные:** те, которые задаются функцией распределения и соответствующей функцией плотности распределения

### 4.3 Плотность распределения случайной величины

**Плотность распределения** — любая функция, которая удовлетворяет специальным условиям: это нормированная функция, то есть если проинтегрировать ее от  $-\infty$  до  $+\infty$  (то есть по всему вещественному промежутку в одномерном  $\mathbb{R}^1$  случае), то получим единицу:

$$\int_{-\infty}^{+\infty} f(x)dx = 1 \quad (4.1)$$

---

<sup>1</sup> «Теория вероятностей и математическая статистика» Буре, Парилина

<sup>2</sup> «Методы прикладной статистики в R и Excel» Буре, Парилина

## 4.4 Некоторые распределения

### 4.4.1 Нормальное (гауссовское)

$$f(x) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(x-\mu)^2}{2\sigma^2}} \quad (4.2)$$

— плотность распределения, причем параметр  $\mu$  — **математическое ожидание** (среднее значение), медиана и мода распределения, а параметр  $\sigma$  — **среднеквадратическое отклонение**,  $\sigma^2$  — **дисперсия** распределения.

Если  $\mu = 0, \sigma = 1$ , то это **стандартное нормальное распределение**. Рассмотрим чуть подробнее его плотность:

$$f(x) = \frac{1}{\sqrt{2\pi}} e^{-\frac{x^2}{2}} \quad (4.3)$$

Поскольку

$$\int_{-\infty}^{+\infty} e^{-x^2/2} dx = \sqrt{2\pi}, \quad (4.4)$$

то

$$\int_{-\infty}^{+\infty} f(x) dx = \frac{1}{\sqrt{2\pi}} \int_{-\infty}^{+\infty} e^{-\frac{x^2}{2}} dx = \frac{1}{\sqrt{2\pi}} \sqrt{2\pi} = 1 \quad (4.5)$$

Вот почему плотность нормального распределения выглядит именно так, откуда там взялся коэффициент с  $\pi$ . И это еще раз проливает свет на определение плотности вероятностного распределения.

### 4.4.2 Экспоненциальное

Случайная величина  $X$  имеет экспоненциальное распределение с параметром  $\lambda > 0$ , если её плотность вероятности имеет вид:

$$f_X(x) = \begin{cases} \lambda e^{-\lambda x}, & x \geq 0, \\ 0, & x < 0. \end{cases} \quad (4.6)$$

### 4.4.3 Бернулли

Случайная величина  $X$  имеет распределение Бернулли, если она принимает всего два значения: 1 и 0 с вероятностями  $p$  и  $q \equiv 1 - p$  соответственно. Таким образом:  $\mathbb{P}(X = 1) = p$ ,  $\mathbb{P}(X = 0) = q$ . Принято говорить, что событие  $\{X = 1\}$  соответствует «успеху», а событие  $\{X = 0\}$  — «неудаче». Эти названия условные, и в зависимости от конкретной задачи могут быть заменены на противоположные.

### 4.4.4 Пуассона

Выберем фиксированное число  $\lambda > 0$  и определим дискретное распределение, задаваемое следующей функцией вероятности:

$$p(k) \equiv \mathbb{P}(Y = k) = \frac{\lambda^k}{k!} e^{-\lambda} \quad (4.7)$$

## 4.5 Математическое ожидание

Идейно, **мат ожидание случайной величины** — это что ты с точки зрения математики ожидаешь получить, случайно сэмплируя величины из распределения при стремлении количества сэмплов к бесконечности. Обозначают символом  $E, \mathbb{E}$ .

Например, бросая монетку бесконечное число раз, в пределе получим, что орел/решка выпадают в среднем одинаковое кол-во раз. То есть мат ожидание такой случайной величины (она, кстати, Бернулевская, дискретная) равно 0.5.

Для непрерывных распределений мат ожидание — это интеграл

$$\mathbb{E}\xi = \int_{-\infty}^{+\infty} x f(x) dx, \quad (4.8)$$

для дискретных — сумма

$$\mathbb{E}\xi = \sum_{i=1}^n a_i p_i, \quad (4.9)$$

$a_i, p_i$  — значение случайной величины и ее вероятность соответственно,  $n$  — количество экспериментов.

## 4.6 Дисперсия

**Дисперсия** — это мера разброса данных. Чем больше дисперсия, тем больше разброс значений, и наоборот. Для расчета дисперсии в математике используется формула, которая учитывает разницу каждого значения от среднего значения и возводит эту разницу в квадрат:

$$\mathbb{D}\xi = \mathbb{E}(\xi - \mathbb{E}\xi)^2 = \mathbb{E}\xi^2 - (\mathbb{E}\xi)^2 = \int_{-\infty}^{+\infty} x^2 f(x) dx - \left( \int_{-\infty}^{+\infty} x f(x) dx \right)^2 \quad (4.10)$$

В дискретном случае:

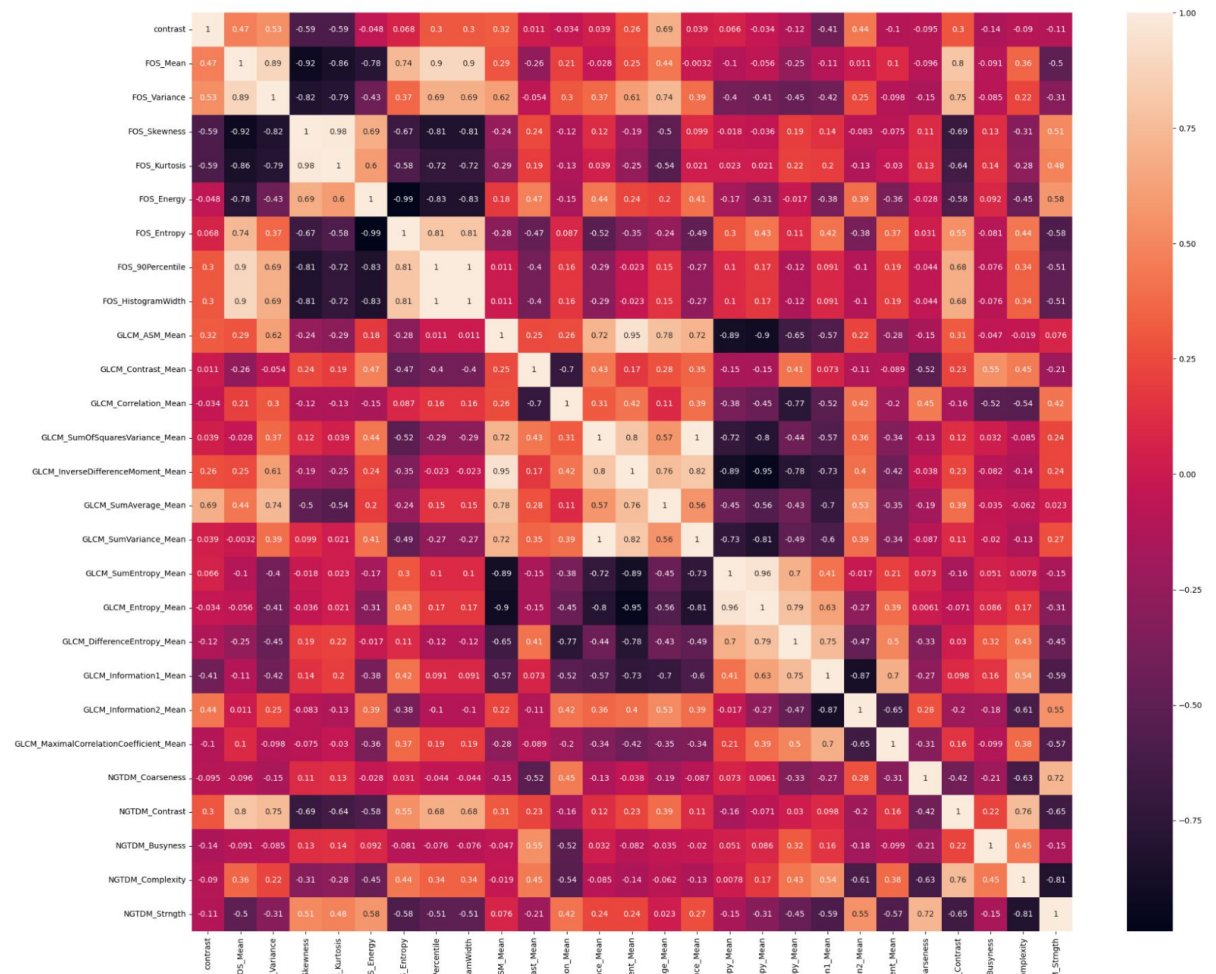
$$\mathbb{D}\xi = \mathbb{E}(\xi - \mathbb{E}\xi)^2 = \mathbb{E}\xi^2 - (\mathbb{E}\xi)^2 = \sum_{i=1}^n a_i^2 p_i - \left( \sum_{i=1}^n a_i p_i \right)^2 \quad (4.11)$$

## 4.7 Корреляция

Бывает разной: линейной, ранговой, нелинейной. Корреляции Пирсона (линейная), Спирмена (ранговая), Кенделла (ранговая) и другие ( $\varphi_k$  — нелинейная). По-разному показывают, как переменные зависят друг от друга.

## 4.8 Выборка, генеральная совокупность

**Генеральная совокупность** — это вообще все возможные существующие на свете объекты, которые нам интересно рассматривать. Например, можно рассматривать множество всех людей на планете как генеральную совокупность. Тогда выборкой будет называться (определение по самому слову) некоторое подмножество,



Пример матрицы корреляций Пирсона между разными признаками

в которое входят отобранные (например, случайно) люди. Можно выбрать конкретных людей. Выборка светловолосых и темноволосых людей, выборка очень низких и очень высоких, выборка из всех петербуржцев и так далее.

Обычно выборку обозначают за  $X$ .

**Дополнительно** Ковариации, ковариационная матрица, многомерные случайные величины, моменты, центральная предельная теорема,  $\gamma$ - и  $\beta$ -распределения, статистические тесты, алгебра множеств, счетная аддитивность,  $\sigma$ -алгебры, аксиоматика Колмогорова, расстояния между распределениями



## Глава 5

# Численные методы. Методы оптимизации

Литература<sup>1 2</sup>

### 5.1 Определение функционала

Короче, **функционал** — обобщение функций: это функция, принимающая функцию (не путать с программированием).

### 5.2 Что такое минимизация функционала

Пусть целевая функция имеет вид:

$$F(\vec{x}) : \mathbb{X} \rightarrow \mathbb{R} \quad (5.1)$$

И задача оптимизации задана следующим образом

$$F(\vec{x}) \rightarrow \min_{\vec{x} \in \mathbb{X}}, \quad (5.2)$$

то есть надо минимизировать  $F(x)$  на множестве всех  $x \in \mathbb{X}$ .

$$\min_{\vec{x} \in \mathbb{X}} F(\vec{x}) \quad (5.3)$$

— функционал качества / задача оптимизации / **минимизация функционала** / минимизация целевой функции и проч.

С максимизацией все так же с точностью до знака.

### 5.3 Градиент. Градиентный спуск

**Градиент** — направление наискорейшего роста функции. **Антиградиент** — градиент с минусом, то есть направление наискорейшего уменьшения функции

---

<sup>1</sup> «Численные методы», Вержбицкий

<sup>2</sup> Интернет, репозитории Github

Основная идея метода заключается в том, чтобы идти в направлении наискорейшего спуска, а это направление задаётся антиградиентом  $-\nabla F$ :

$$\vec{x}_{n+1} = \vec{x}_n - \lambda \nabla F(\vec{x}_n) \quad (5.4)$$

За счет этого из заданной точки мы движемся в направлении наискорейшего убывания функции, таким образом попадем в конечном итоге в желаемый минимум.

Проблема в том, что нужно правильно подобрать  $\lambda$  — **градиентный шаг**. Если сделать его слишком большим, то шаги будут чересчур большими и можем перепрыгнуть через точку минимума; если сделать слишком маленьким, то будет очень долго сходиться алгоритм. Обычно берут что-то вроде  $\lambda = 10^{-3}$ , но это тоже гиперпараметр, поэтому только практикой можно понять, какой  $\lambda$  подходит лучше всего.

## Глава 6

# Теоретическая информатика

### 6.1 Энтропия Шеннона

**Информационная энтропия** — мера неопределённости некоторой системы (в статистической физике или теории информации), в частности, непредсказуемость появления какого-либо символа первичного алфавита.

Например, в последовательности букв, составляющих какое-либо предложение на русском языке, разные буквы появляются с разной частотностью, поэтому неопределённость появления для некоторых букв меньше, чем для других.

Информационная двоичная энтропия, при отсутствии информационных потерь, рассчитывается по **формуле Хартли**:

$$I = \log_2 N, \quad (6.1)$$

где  $N$  — мощность алфавита,  $I$  — количество информации в каждом символе сообщения. Для случайной величины  $x$ , принимающей  $n$  независимых случайных значений  $x_i$  с вероятностями  $p_i$  ( $i = 1, \dots, n$ ), формула Хартли переходит в **формулу Шеннона**:

$$H(x) = - \sum_{i=1}^n p_i \log_2 p_i. \quad (6.2)$$

Здесь  $-\log_2 p_i$  означает измеряемое **в битах** (отсюда  $-\log_2 p_i$ ) количество информации, содержащейся в том событии, что случайная величина приняла значение  $x_i$ ,  $H(x)$  — количество информации, которое в среднем приходится на одно событие (для предложений на русском языке — количество информации в среднем на одну букву).

Если измерять информацию не **в битах**, то основание логарифма поменяется на любое другое.

**Дополнительно** Условная энтропия, перекрестная (совместная) энтропия, пропускная способность канала, основы кодирования, дерево Хаффмана, расстояние Левенштейна, дифференциальная энтропия