

Методы оптимизации

КОНСПЕКТ ТЕОРИИ

Содержание

1	Matrix calculus and linear algebra	3
1.1	Матричное дифференцирование	3
1.2	Псевдообратная матрица	3
1.3	Сингулярное разложение	4
2	Convex sets and functions, projections	5
2.1	Выпуклые множества	5
2.2	Проекции	5
2.3	Выпуклые функции	5
3	Conjugate sets	7
4	Conjugate functions	8
4.1	Сопряженные функции	8
4.2	Сопряженная норма	8
5	Subgradient and subdifferential	9
6	Karush-Kuhn-Tucker conditions	10
7	Duality	11
7.1	Построение двойственной задачи	11
7.2	Связь двойственной задачи и условий ККТ	12
7.3	Теорема Фенхеля-Рокафеллара	12
7.4	Задачи линейного программирования	13
8	Maximum likelihood estimation	14
8.1	Постановка задачи	14
8.2	Линейная регрессия	14
8.3	Логистическая регрессия	15

Источники:

1. семинары Даниила Меркулова в 776 группе (source: fmin.xyz);
2. лекции Осипенко К.Ю. на ФУПМе;
3. “Convex Optimization”, Stephen Boyd.

Используемые обозначения:

- $\nabla f(x)$ — градиент функции $f : \mathbb{R}^p \rightarrow \mathbb{R}$ по вектору x . По умолчанию является столбцом $(p \times 1)$.
- $H(x)$ — гессиан функции $f : \mathbb{R}^p \rightarrow \mathbb{R}$. Является матрицей $(p \times p)$.
- $J(x)$ — якобиан функции $f : \mathbb{R}^p \rightarrow \mathbb{R}^m$. Является матрицей $(m \times p)$.
- $f'(X)$ — производная функции $f : \mathbb{R}^{m \times n} \rightarrow \mathbb{R}$ по матрице X :

$$f'(X) = \left(\frac{\partial f}{\partial x_{ij}} \right)_{\substack{i=1, \dots, m \\ j=1, \dots, n}} = (m \times n)$$

- $\langle x, y \rangle = x^T y$ — скалярное произведение векторов x и y .
- $\langle X, Y \rangle = \text{tr}(A^T B) = \sum_{i=1}^n \sum_{j=1}^m a_{ij} b_{ij}$ — скалярное произведение матриц X и Y (одинакового размера).
- $X^{-T} = (X^{-1})^T = (X^T)^{-1}$ для квадратной невырожденной матрицы X .
- $I = E$ — единичная матрица.
- \mathbb{S}_+^n — симметричные положительно полуопределенные матрицы порядка n .
- \mathbb{S}_{++}^n — симметричные положительно определенные матрицы порядка n .
- $A \succ 0$ — матрица A положительно определена.
- $A \succeq 0$ — матрица A положительно полуопределена.
- $\mathbb{R}_+^n = \{x \in \mathbb{R}^n \mid x_i \geq 0\}$ — векторы с неотрицательными компонентами.
- $\mathbb{R}_{++}^n = \{x \in \mathbb{R}^n \mid x_i > 0\}$ — векторы с положительными компонентами.
- $\|\cdot\|_2$ — евклидова норма вектора.
- $\text{int} A, \text{relint} A$ — внутренность A , относительная внутренность A .
- $B_r(a)$ — открытый шар радиуса r с центром в точке a .
- $\text{conv} A$ — выпуклая оболочка множества A .
- $\text{aff} A$ — аффинная оболочка множества A .
- $\text{cone} A$ — коническая оболочка множества A .
- $\text{lin} A$ — линейная оболочка множества A .
- $A + B = \{a + b \mid a \in A, b \in B\}$ — сумма Минковского множеств A и B .
- $\mathbb{E}\xi$ — матожидание случайной величины ξ (expected value).
- $\mathbb{V}\xi$ — дисперсия случайной величины ξ (variance).

1 Matrix calculus and linear algebra

1.1 Матричное дифференцирование

Для нахождения первой и второй производной (градиента, якобиана, гессиана и т.п.) многомерных функций используется запись дифференциала функции, следующая из формулы Тейлора:

- $f : \mathbb{R}^n \rightarrow \mathbb{R}$

$$df = \langle \nabla f, dx \rangle$$
$$d^2 f = \langle H dx_1, dx_2 \rangle,$$

где dx_1 — дифференциал x при первом дифференцировании, а dx_2 — при втором.

- $f : \mathbb{R}^n \rightarrow \mathbb{R}^m$

$$df = J dx,$$

где $J = (m \times n)$ — якобиан функции f .

- $f : \mathbb{R}^{m \times n} \rightarrow \mathbb{R}$

$$df = \langle f'(X), dX \rangle$$

Производные старших порядков также можно получить таким образом, но они уже будут тензорами ранга ≥ 3 , поэтому они не будут иметь матричного представления.

Некоторые свойства дифференцирования матриц:

1. $d(X^T) = (dX)^T$
2. $d(XY) = (dX)Y + X(dY)$
3. $d(\det X) = \det X \langle X^{-T}, dX \rangle$
4. $d(\text{tr} X) = \langle I, dX \rangle$

1.2 Псевдообратная матрица

Опр. Пусть $A \in \mathbb{R}^{m \times n}$. Псевдообратной матрицей (Moore-Penrose inverse) к матрице A называется

$$A^\dagger = \lim_{\alpha \rightarrow 0} (A^T A + \alpha I)^{-1} A^T = \lim_{\alpha \rightarrow 0} A^T (A A^T + \alpha I)^{-1}$$

Оба предела всегда существуют и равны.

Если матрица A имеет полный ранг ($\text{rg } A = \min\{m, n\}$), то для A^\dagger есть алгебраическое выражение:

Случай	Алгебраическое выражение	Задача, решением которой является $A^\dagger b$
$A \in \mathbb{R}^{n \times n}$	$A^\dagger = A^{-1}$	$Ax = b$
$A \in \mathbb{R}^{m \times n}$ $m \geq n$	$A^\dagger = (A^T A)^{-1} A^T$	$\ Ax - b\ _2^2 \rightarrow \min_{x \in \mathbb{R}^n}$
$A \in \mathbb{R}^{m \times n}$ $m \leq n$	$A^\dagger = A^T (A A^T)^{-1}$	$\begin{cases} \ x\ _2^2 \rightarrow \min_{x \in \mathbb{R}^n} \\ Ax = b \end{cases}$

Свойства:

- $AA^\dagger A = A$
- $A^\dagger AA^\dagger = A^\dagger$
- $(A^\dagger)^\dagger = A$
- $(A^T)^\dagger = (A^\dagger)^T$
- $(\alpha A)^\dagger = \alpha^{-1} A^\dagger$
- $(AB)^\dagger = B^\dagger A^\dagger$, если A или B полного ранга

Если линейная система $Ax = b$ имеет решения, то все они задаются формулой

$$x = A^\dagger b + [I - A^\dagger A]w, \quad \forall w \in \mathbb{R}^n$$

1.3 Сингулярное разложение

Теорема о сингулярном разложении (SVD)

Пусть $A \in \mathbb{R}^{m \times n}$ — произвольная вещественная матрица ранга r . Тогда при $m \geq n$:

$$A = U \Sigma V^T,$$

- $U = (m \times m)$, $V = (n \times n)$ — ортогональные матрицы,
- $\Sigma = (m \times n)$ — матрица с r ненулевыми элементами на диагонали:

$$\sigma_j = \sqrt{\lambda_j} \quad \text{— сингулярные числа матрицы } A^T A \text{ в порядке убывания}$$

- Столбцы U — собственные векторы AA^T , столбцы V — собственные векторы $A^T A$.

Аналогично теорема формулируется для случая, когда $m \leq n$.

Если $A = U \Sigma V^T$, то псевдообратная матрица находится по формуле:

$$A^\dagger = V \Sigma^\dagger U^T, \quad \Sigma^\dagger = \text{diag} \left(\frac{1}{\sigma_1}, \dots, \frac{1}{\sigma_r}, 0, \dots, 0 \right)$$

2 Convex sets and functions, projections

2.1 Выпуклые множества

Способы проверить выпуклость множества:

- По определению:

Опр. Множество $S \subset \mathbb{R}^n$ называется *выпуклым*, если

$$\forall x, y \in S \quad \forall \lambda \in [0, 1] \rightarrow \lambda x + (1 - \lambda)y \in S;$$

- Пересечение любого (даже несчетного) числа выпуклых множеств — выпуклое множество;
- Образ и прообраз выпуклого множества при аффинном отображении ($f(x) = Ax + b$) — выпуклые множества.

2.2 Проекции

Опр. Проекцией точки $y \in \mathbb{R}^n$ на множество $S \subset \mathbb{R}^n$ называется точка $\pi = \pi_S(y) \in S$, если

$$\forall x \in S \rightarrow \|\pi - y\| \leq \|\pi - x\|$$

Свойства:

- Если проекция y на S существует, то

$$\pi_S(y) = \arg \min_{x \in S} \|x - y\|$$

- Проекция может вообще не существовать, быть единственна, или их может быть много.
- (*теорема Рисса о проекции*) Если множество $S \subset \mathbb{R}^n$ выпукло и замкнуто, то проекция любой точки на S существует и единственна, и выполнено:

$$\forall x \in S: \quad \langle \pi - y, x - \pi \rangle \geq 0 \quad \Longleftrightarrow \quad \pi_S(y) = \pi$$

- Если множество S открыто, а $y \notin S$, то проекции не существует.
- Если S — аффинное подпространство, то

$$\forall x \in S: \quad \langle \pi - y, x - \pi \rangle = 0 \quad \Longleftrightarrow \quad \pi_S(y) = \pi$$

2.3 Выпуклые функции

Способы проверить (нестрогую) выпуклость функции:

- По определению:

Опр. Функция $f: S \rightarrow \mathbb{R}$, определенная на *выпуклом* множестве $S \subseteq \mathbb{R}^n$, называется *выпуклой* на S , если

$$\forall x, y \in S \quad \forall \lambda \in [0, 1] \rightarrow f(\lambda x + (1 - \lambda)y) \leq \lambda f(x) + (1 - \lambda)f(y)$$

- Дифференциальный критерий 1-го порядка:

Пусть f дифференцируема на S . Тогда f выпукла на S тогда и только тогда, когда

$$\forall x, y \in S \rightarrow f(y) \geq f(x) + (\nabla f(x))^T (y - x),$$

то есть в каждой точке можно провести касательную гиперплоскость, являющуюся глобальной нижней оценкой.

- Дифференциальный критерий 2-го порядка:

Пусть f дважды дифференцируема на S . Тогда f выпукла тогда и только тогда, когда

$$\forall x \in \text{relint} S \rightarrow H(x) = \nabla^2 f(x) \succeq 0,$$

то есть гессиан f является положительно полуопределенной матрицей.

- Ограничение на прямую:

Пусть $f : S \rightarrow \mathbb{R}$, $S \subseteq \mathbb{R}^n$ — выпукло. Пусть $x \in S$, $v \in \mathbb{R}^n$. Определим на выпуклом множестве $T = \{t \mid x + tv \in S\} \subseteq \mathbb{R}$ функцию числового аргумента g :

$$g : T \rightarrow \mathbb{R}, \quad g(t) = f(x + tv)$$

Тогда функция f выпукла на S тогда и только тогда, когда

$$\forall x \in S, v \in \mathbb{R}^n \text{ функция } g \text{ выпукла на } T.$$

Чтобы проверить *строгую выпуклость* нужно во всех критериях поменять знаки неравенств на строгие.

Способы проверить μ -сильную выпуклость:

- По определению:

Опр. Функция $f : S \rightarrow \mathbb{R}$, определенная на *выпуклом* множестве $S \subseteq \mathbb{R}^n$, называется μ -сильно выпуклой на S или просто *сильно выпуклой*, если

$$\exists \mu > 0 \quad \forall x, y \in S \quad \forall \lambda \in [0, 1] \rightarrow f(\lambda x + (1 - \lambda)y) \leq \lambda f(x) + (1 - \lambda)f(y) - \mu \lambda(1 - \lambda) \|x - y\|^2$$

- Дифференциальный критерий 1-го порядка:

Пусть f дифференцируема на S . Тогда f сильно выпукла на S тогда и только тогда, когда

$$\exists \mu > 0 \quad \forall x, y \in S \rightarrow f(y) \geq f(x) + (\nabla f(x))^T (y - x) + \frac{\mu}{2} \|y - x\|^2,$$

то есть в каждой точке можно провести касательную параболу, являющуюся глобальной нижней оценкой.

- Дифференциальный критерий 2-го порядка:

Пусть f дважды дифференцируема на S . Тогда f сильно выпукла тогда и только тогда, когда

$$\exists \mu > 0 \quad \forall x \in \text{relint} S \rightarrow H(x) = \nabla^2 f(x) \succeq \mu I,$$

то есть матрица $(\nabla^2 f(x) - \mu I)$ является положительно полуопределенной матрицей.

Опр. Функция f называется *вогнутой*, если функция $(-f)$ выпукла.

3 Conjugate sets

Опр. Пусть $S \subset \mathbb{R}^n$. *Сопряженным или двойственным множеством* ко множеству S называется

$$S^* = \{y \in \mathbb{R}^n \mid \langle x, y \rangle \geq -1 \quad \forall x \in S\}$$

Опр. Пусть $S \subset \mathbb{R}^n$ — конус. *Сопряженным конусом* называется

$$K^* = \{y \in \mathbb{R}^n \mid \langle x, y \rangle \geq 0 \quad \forall x \in K\}$$

Свойства:

- S^* всегда выпукло, замкнуто и содержит 0.
- $S^* = \bigcap_{x \in S} \{y \mid \langle x, y \rangle \geq -1\}$ — пересечение полупространств.
- $S^* = (\overline{S})^*$, $S^* = (\text{conv} S)^*$, $S^* = (S \cup \{0\})^*$
- $S^{**} = \overline{\text{conv}(S \cup \{0\})}$
- $\left(\bigcup_{\alpha} S_{\alpha}\right)^* = \bigcap_{\alpha} S_{\alpha}^*$
- Для конуса K и произвольного множества S : $(S + K)^* = S^* \cap K^*$
- Для конусов K_1, K_2 , имеющих внутреннюю точку: $(K_1 \cap K_2)^* = K_1^* + K_2^*$

Теорема. Сопряженным ко множеству

$$S = \text{conv}(x_1, \dots, x_k) + \text{cone}(x_{k+1}, \dots, x_m)$$

является полиэдр (многогранник)

$$S^* = \{p \mid \langle p, x_i \rangle \geq -1, \quad i = \overline{1, k}, \quad \langle p, x_j \rangle \geq 0, \quad j = \overline{k+1, m}\}$$

4 Conjugate functions

4.1 Сопряженные функции

Опр. Пусть $f : \mathbb{R}^n \rightarrow \mathbb{R}$. Сопряженной к функции f функцией называется

$$f^*(y) = \sup_{x \in \mathbb{R}^n} (\langle y, x \rangle - f(x))$$

Областью определения f^* является множество таких y , что супремум в определении выше конечен.

Свойства:

- f^* — выпуклая функция;
- $f^{**} = f \iff f$ — выпуклая замкнутая функция (Теорема Фенхеля-Моро);
- Неравенство Фенхеля-Юнга: $f(x) + f^*(y) \geq \langle y, x \rangle$
- Если $f(x) \leq g(x)$, то $f^*(y) \geq g^*(y)$
- Если $f(x, y) = f_1(x) + f_2(y)$ и функции f_1, f_2 выпуклы, то $f^*(p, q) = f^*(p) + f^*(q)$

Как искать сопряженную функцию к дифференцируемой функции $f(x)$:

1. Положить $g(x, y) = \langle y, x \rangle - f(x)$.
2. Определить, при каких y $\sup_{x \in \mathbb{R}} g(x, y)$ конечен — это область определения f^* .
3. Найти максимум $g(x, y)$ по x : $\nabla_x g(x, y) = y - \nabla f(x) = 0$.
Часто получается (но не всегда), что все значения y , при которых уравнение $y = \nabla f(x)$ разрешимо относительно x , есть область определения f^* .
4. Выразить x через y и подставить в $g(x, y)$ — это выражение для $f^*(y)$.

4.2 Сопряженная норма

Опр. Пусть $(X, \|\cdot\|)$ — линейное нормированное пространство. Сопряженным пространством X^* называется множество всех линейных непрерывных функционалов на X .

Действие функционала $y \in X^*$ на элементе $x \in X$ обозначается $\langle y, x \rangle$.

Опр. Сопряженной нормой (dual norm) на X^* называется функция $\|\cdot\|_* : X^* \rightarrow \mathbb{R}$:

$$\|y\|_* = \sup_{x \neq 0} \frac{|\langle y, x \rangle|}{\|x\|} = \sup_{\|x\| \leq 1} |\langle y, x \rangle| = \inf \{L > 0 \mid |\langle y, x \rangle| \leq L\|x\| \quad \forall x \in X\}$$

Свойства:

- $(X^*, \|\cdot\|_*)$ — линейное нормированное пространство.
- Неравенство Коши-Буняковского-Шварца: $\langle y, x \rangle \leq \|y\|_* \cdot \|x\|$.
- Сопряженным пространством ко множеству столбцов \mathbb{R}^n является множество всех строк \mathbb{R}^n , а $\langle y, x \rangle$ является обычным скалярным произведением.
- Сопряженной нормой к $\|\cdot\|_p$ на \mathbb{R}^n является $\|\cdot\|_q$, где $\frac{1}{p} + \frac{1}{q} = 1$, $p > 1$.
- Сопряженной нормой к $\|\cdot\|_1$ является $\|\cdot\|_\infty$, сопряженной нормой к $\|\cdot\|_\infty$ является $\|\cdot\|_1$.
- Норма $\|\cdot\|_2$ самосопряжена.

Сопряженная норма не является сопряженной функцией для $f(x) = \|x\|$. Сопряженной функцией будет

$$f^*(y) = \begin{cases} 0 & , \|y\|_* \leq 1; \\ +\infty & , \text{ иначе.} \end{cases}$$

5 Subgradient and subdifferential

Опр. Пусть $f : S \rightarrow \mathbb{R}$, $S \subset \mathbb{R}^n$. Вектор g называется *субградиентом* функции f в точке x_0 , если

$$\forall x \in S \rightarrow f(x) - f(x_0) \geq \langle g, x - x_0 \rangle$$

Опр. Множество всех субградиентов f в точке x_0 называется *субдифференциалом* функции f в точке x_0 и обозначается $\partial f(x_0) \equiv \partial f_S(x_0)$.

Свойства:

- Если $x_0 \in \text{relint} S$, то $\partial f_S(x_0)$ — выпуклое компактное множество;
- Выпуклая функция f дифференцируема в $x_0 \iff \partial f(x_0) = \{\nabla f(x_0)\}$;
- Если $\forall x \in S \partial f_S(x_0) \neq \emptyset$, то f выпукла на S ;
- Если $\alpha \geq 0$, то $\partial(\alpha f)(x) = \alpha \partial f(x)$;
- Если f выпукла, то $\partial(f(Ax + b))(x) = A^T \partial f(Ax + b)$.

Теорема Моро-Рокафеллара.

Пусть $f : E \rightarrow \mathbb{R}$, $g : G \rightarrow \mathbb{R}$ — выпуклые функции, $x_0 \in E \cap G$, $E \cap \text{int} G \neq \emptyset$. Тогда

$$\partial(f + g)(x_0) = \partial f(x_0) + \partial g(x_0)$$

Теорема Дубовицкого-Милютин.

Пусть $f_i : E_i \rightarrow \mathbb{R}$, $i = \overline{1, m}$ — выпуклые функции, $x_0 \in \text{int} \left(\bigcap_{i=1}^m E_i \right)$, а функция $f : \bigcap_{i=1}^m E_i \rightarrow \mathbb{R}$:

$$f(x) = \max \{f_1(x), \dots, f_m(x)\}$$

Тогда

$$\partial f(x_0) = \text{conv} \left(\bigcup_{j \in I} \partial f_j(x_0) \right), \quad I = \{j \mid f_j(x_0) = f(x_0)\}$$

Теорема о субдифференциале сложной функции.

Пусть $g_i : S \rightarrow \mathbb{R}$, $i = \overline{1, m}$ — выпуклые функции, $\varphi : U \rightarrow \mathbb{R}$ — неубывающая (по всем переменным) выпуклая функция, причем $U \supset (g_1(S), \dots, g_m(S))$, $U \subset \mathbb{R}^m$ — открытое множество. Тогда при $f(x) = \varphi(g_1(x), \dots, g_m(x))$:

$$\partial f(x) = \bigcup_{p \in \partial \varphi(u)} \left(\sum_{i=1}^m p_i \partial g_i(x) \right), \quad u = (g_1(x), \dots, g_m(x))$$

В частности, если φ дифференцируема в точке u , то

$$\partial f(x) = \sum_{i=1}^m \frac{\partial \varphi}{\partial u_i}(u) \partial g_i(x)$$

6 Karush-Kuhn-Tucker conditions

Рассматривается задача оптимизации (*задача математического программирования*):

$$\begin{cases} f(x) \longrightarrow \min_{x \in \mathbb{R}^n} \\ g_i(x) \leq 0, & i = \overline{1, m} \\ h_j(x) = 0, & j = \overline{1, p} \end{cases} \iff f(x) \longrightarrow \min_{x \in S}, \quad (*)$$

где множество S задается ограничениями.

Замечания:

- можно убрать ограничения типа $h_j(x) = 0$, заменив их на $h_j(x) \leq 0$ и $-h_j(x) \leq 0$;
- под записью $f(x) \longrightarrow \min$ понимается нахождение нижней грани (инфимума);
- далее все функции f, g_i и h_j считаем достаточно гладкими.

Опр. Функцией Лагранжа для задачи (*) называется

$$L(x, \lambda, \mu) = \lambda_0 f(x) + \sum_{j=1}^p \lambda_j h_j(x) + \sum_{i=1}^m \mu_i g_i(x)$$

Теорема Каруша-Куна-Таккера.

Пусть x^* — решение задачи (*).

Тогда $\exists (\lambda^*, \mu^*) = (\lambda_0^*, \dots, \lambda_p^*, \mu_1^*, \dots, \mu_m^*) \neq \mathbf{0}$ такой, что выполнены условия Каруша-Куна-Таккера:

1. $x^* \in S$ (x^* — допустимая точка, т.е. выполнены ограничения);
2. $\lambda_0^* \geq 0, \mu_i^* \geq 0, i = \overline{1, m}$ (неотрицательность);
3. $\nabla_x L(x^*, \lambda^*, \mu^*) = 0$ (минимальность);
4. $\mu_i^* \cdot g_i(x^*) = 0, i = \overline{1, m}$ (дополняющая нежесткость).

Опр. Задача (*) называется *выпуклой*, если функции f, g_i выпуклы, а ограничения-равенства либо отсутствуют, либо являются аффинными (имеют вид $Ax = b$).

Опр. Задача (*) называется *регулярной*, если для нее условия ККТ выполнены при $\lambda_0^* > 0$.

Если задача (*) регулярна, то в лагранжиане можно положить $\lambda_0 = 1$:

$$L(x, \lambda, \mu) = f(x) + \sum_{j=1}^p \lambda_j h_j(x) + \sum_{i=1}^m \mu_i g_i(x),$$

то есть регулярность в некотором смысле означает невырожденность задачи.

Условие регулярности Слейтера. Пусть в задаче (*)

- функции g_i выпуклы, а h_j либо отсутствуют, либо аффинны;
- $\exists \tilde{x} \in S: g_i(\tilde{x}) < 0 \iff \text{relint}(S) \neq \emptyset$.

Тогда задача (*) является регулярной.

Условие регулярности через двойственность. Если в задаче (*) присутствует сильная двойственность, то задача (*) регулярна.

Если задача (*) выпукла и регулярна, то условия Каруша-Куна-Таккера являются *необходимыми и достаточными условиями глобального минимума*.

7 Duality

7.1 Построение двойственной задачи

Опр. Функция $f : \mathbb{R}^n \rightarrow \bar{\mathbb{R}}$ называется *собственной*, если f не принимает значения $-\infty$ и $f \not\equiv +\infty$.

Рассматривается задача оптимизации (*primal problem*) собственной функции $f : \mathbb{R}^n \rightarrow \mathbb{R} \cup \{+\infty\}$:

$$\begin{cases} f(x) \longrightarrow \min_{x \in \mathbb{R}^n} \\ g_i(x) \leq 0, & i = \overline{1, m} \\ h_j(x) = 0, & j = \overline{1, p} \end{cases} \iff f(x) \longrightarrow \min_{x \in S}, \quad (P)$$

где множество S задается ограничениями. Оптимальное значение этой задачи будем обозначать p^* .

Ей соответствует лагранжиан

$$L(x, \lambda, \mu) = f(x) + \sum_{j=1}^p \lambda_j h_j(x) + \sum_{i=1}^m \mu_i g_i(x)$$

Обозначим

$$F(x) = \max_{\substack{\lambda \in \mathbb{R}^p \\ \mu \in \mathbb{R}_+^m}} L(x, \lambda, \mu) = \begin{cases} f(x), & x \in S; \\ +\infty, & x \notin S. \end{cases}$$

Второе равенство легко доказывается.

Часто в качестве нотации вместо \sup и \inf используется \max и \min , но подразумеваются верхняя и нижняя грани соответственно.

Тогда исходную задачу можно записать в виде:

$$f(x) \longrightarrow \min_{x \in S} \iff F(x) \longrightarrow \min_{x \in \mathbb{R}^n}$$

Опр. Двойственной функцией (*dual function*) к задаче (P) называется функция $g : \mathbb{R}^p \times \mathbb{R}_+^m \rightarrow \mathbb{R} \cup \{-\infty\}$

$$g(\lambda, \mu) = \min_{x \in \mathbb{R}^n} L(x, \lambda, \mu)$$

Не стоит путать понятия сопряженной функции и двойственной функции.

Несложно видеть, что

$$g(\lambda, \mu) \leq L(x, \lambda, \mu) \leq F(x) \quad \forall x \in \mathbb{R}^n, \lambda \in \mathbb{R}^p, \mu \in \mathbb{R}_+^m,$$

поэтому

$$\max_{\substack{\lambda \in \mathbb{R}^p \\ \mu \in \mathbb{R}_+^m}} g(\lambda, \mu) \leq \min_{x \in \mathbb{R}^n} F(x) = \min_{x \in S} f(x) \quad (1)$$

Опр. Двойственной задачей (*dual problem*) к задаче (P) называется задача

$$\begin{cases} g(\lambda, \mu) \longrightarrow \max \\ \mu_i \geq 0, & i = \overline{1, m} \end{cases} \quad (D)$$

Оптимальное значение задачи (D) будем обозначать d^* .

Опр. Говорят, что в задаче (P) присутствует *слабая двойственность* (*weak duality*), если $d^* \leq p^*$.

Из (1) следует, что слабая двойственность есть всегда.

Опр. Говорят, что в задаче (P) присутствует *сильная двойственность* (*strong duality*), если $d^* = p^*$.

Сильная двойственность есть не всегда, однако есть некоторые достаточные условия, гарантирующие ее наличие.

Опр. Разница $p^* - d^*$ называется *разрывом двойственности* (*duality gap*).

Условие Слейтера.

В задаче (P) есть сильная двойственность, если (P) — выпуклая задача и $\text{relint}(S) \neq \emptyset$, т.е.

$$\exists \tilde{x} \in \mathbb{R}^n : g_i(\tilde{x}) < 0$$

7.2 Связь двойственной задачи и условий ККТ

Рассмотрим выпуклую задачу, для которой выполнено условие Слейтера

$$\begin{cases} f(x) \longrightarrow \min_{x \in \mathbb{R}^n} \\ g_i(x) \leq 0, & i = \overline{1, m} \\ Ax = b. \end{cases}$$

Для нее условия Каруша-Куна-Таккера являются необходимыми и достаточными условиями глобального минимума, и наблюдается сильная двойственность.

При таких условиях точки x^*, λ^*, μ^* — решение системы из условий ККТ тогда и только тогда, когда

- x^* — точка оптимума прямой задачи;
- (λ^*, μ^*) — точка оптимума двойственной задачи.

7.3 Теорема Фенхеля-Рокафеллара

Рассмотрим задачу оптимизации

$$f(x) + g(Ax) \longrightarrow \min_{x \in E \cap A^{-1}(G)},$$

где $x \in \mathbb{R}^n$, а $A \in \mathbb{R}^{m \times n}$ — матрица линейного отображения из \mathbb{R}^n в \mathbb{R}^m .

Эквивалентная задача:

$$\begin{cases} f(x) + g(y) \longrightarrow \min \\ Ax = y \end{cases} \quad (*)$$

Можно считать, что f и g равны $+\infty$ вне множеств E и G соответственно, то есть $E = \text{dom } f$, $G = \text{dom } g$.

Лагранжиан:

$$L(x, y, \lambda) = f(x) + g(y) + \lambda^T (Ax - y)$$

Несложно видеть, что двойственная функция выражается через сопряженные:

$$g_d(\lambda) = -f^*(-A^T \lambda) - g^*(\lambda)$$

Тогда задача, двойственная к $(*)$ имеет вид

$$-f^*(-A^T \lambda) - g^*(\lambda) \longrightarrow \max_{\lambda \in \mathbb{R}^m}$$

Теорема Фенхеля-Рокафеллара.

1. Пусть $f: \mathbb{R}^n \rightarrow \overline{\mathbb{R}}$, $g: \mathbb{R}^m \rightarrow \overline{\mathbb{R}}$ — собственные функции, $A \in \mathbb{R}^{m \times n}$, а p^* и d^* — значения оптимумов прямой и двойственной задач:

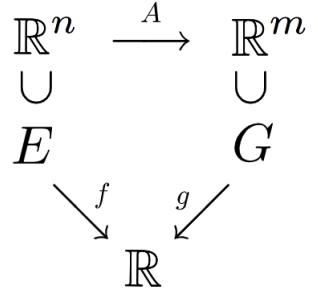
$$p^* = \min_{x \in \mathbb{R}^n} [f(x) + g(Ax)], \quad d^* = \max_{\lambda \in \mathbb{R}^m} [-f^*(-A^T \lambda) - g^*(\lambda)]$$

Тогда $p^* \geq d^*$. (Это мы доказали, построив двойственную задачу.)

2. Кроме того, пусть функции f и g выпуклы, и $A(\text{relint } E) \cap \text{relint } G \neq \emptyset$. Тогда $p^* = d^*$.

При этом, если $p^* = d^* < +\infty$, то точки x^* и λ^* являются точками оптимума тогда и только тогда, когда

$$-A^T \lambda^* \in \partial f(x^*), \quad \lambda^* \in \partial g(Ax^*)$$



7.4 Задачи линейного программирования

Форма задачи линейного программирования	Прямая задача (P)	Двойственная задача (D)
нормальная	$\begin{cases} c^T x \longrightarrow \min \\ Ax \geq b \\ x \geq 0 \end{cases}$	$\begin{cases} y^T b \rightarrow \max \\ y^T A \leq c^T \\ y^T \geq 0 \end{cases}$
общая	$\begin{cases} c^T x \longrightarrow \min \\ Ax \geq b \end{cases}$	$\begin{cases} y^T b \rightarrow \max \\ y^T A = c^T \\ y^T \geq 0 \end{cases}$
каноническая	$\begin{cases} c^T x \longrightarrow \min \\ Ax = b \\ x \geq 0 \end{cases}$	$\begin{cases} y^T b \rightarrow \max \\ y^T A \leq c^T \end{cases}$

Задачу линейного программирования (ЛП) в одной форме можно свести к другой, то есть все три формы эквивалентны.

Везде подразумевается, что заданы столбцы $c \in \mathbb{R}^n$, $b \in \mathbb{R}^m$ и матрица $A \in \mathbb{R}^{m \times n}$, а оптимум ищется по векторам $x \in \mathbb{R}^n$, $y \in \mathbb{R}^m$.

Теорема о двойственности.

- Для оптимальных значений прямой и двойственной задач линейного программирования возможны следующие 4 случая:

	1	2	3	4
значение (P)	c	\emptyset	$-\infty$	\emptyset
значение (D)	c	$+\infty$	\emptyset	\emptyset

где $c \in \mathbb{R}$, а \emptyset означает, что допустимое множество задачи пусто.

- Пусть \hat{x} и \hat{y}^T — допустимые точки задач (P) и (D). Тогда

$$\left. \begin{array}{l} \hat{x} - \text{решение } (P) \\ \hat{y}^T - \text{решение } (D) \end{array} \right\} \iff c^T \hat{x} = \hat{y}^T b$$

То есть если значение хотя бы одной из задач (P) или (D) конечно, то значения обеих задач конечны и совпадают, то есть присутствует сильная двойственность.

Теорема о двойственности верна для задач линейного программирования в любой форме.

8 Maximum likelihood estimation

8.1 Постановка задачи

Дано: выборка x_1, \dots, x_m — независимые измерения случайного вектора $X \in \mathbb{R}^n$.

Задача: найти распределение случайного вектора X .

Сначала делается общая гипотеза о том, распределение какого класса имеет случайный вектор X . То есть мы предполагаем, что X имеет плотность распределения $p(x | \theta)$, где $\theta \in \mathbb{R}^k$ — набор параметров.

Например, мы можем предположить, что X имеет нормальное распределение. Тогда $\theta = (\mu, \sigma)$ и

$$p(x | \mu, \sigma) = \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{(x - \mu)^2}{2\sigma^2}\right)$$

Найдем такие параметры θ^* , что вероятность исходной выборки при $\theta = \theta^*$ максимальна. В этом и заключается суть метода максимального правдоподобия.

Опр. *Функцией правдоподобия (likelihood function)* называется вероятность исходной выборки:

$$L(\theta) = \prod_{i=1}^m p(x_i | \theta)$$

Почти всегда удобно перейти к логарифму этой функции, и иногда именно его называют функцией правдоподобия.

Опр. *Функцией правдоподобия (log-likelihood function)* называется

$$L(\theta) = \log \prod_{i=1}^m p(x_i | \theta) = \sum_{i=1}^m \log p(x_i | \theta)$$

Тогда оптимальные параметры:

$$\theta^* = \arg \max_{\theta} L(\theta) = \arg \max_{\theta} \sum_{i=1}^m \log p(x_i | \theta)$$

Находить их можно, например, приравняв градиент функции правдоподобия к нулю: $\nabla_{\theta} L(\theta^*) = 0$.

8.2 Линейная регрессия

Дано: точки $a_1, \dots, a_m \in \mathbb{R}^n$ ($m > n$) и измерения $b_1, \dots, b_m \in \mathbb{R}$ в этих точках.

Задача: найти наилучшее линейное приближение $b \approx \theta^T A$, где матрица $A \in \mathbb{R}^{n \times m}$ имеет столбцы a_i , а $\theta \in \mathbb{R}^n$ — вектор параметров.

Если в модель хочется добавить смещение: $b \approx \theta^T A + \theta_0$, то можно добавить еще одну величину $a_{m+1} = 1$, всегда равную единице. Тогда задача сведется к описанному выше случаю.

Рассмотрим два способа решения этой задачи.

1. Метод наименьших квадратов

Предположим, что слово “наилучшее” означает, что сумма квадратов отклонений наименьшая. Тогда можно сформулировать задачу оптимизации:

$$\sum_{i=1}^m (\theta^T a_i - b_i)^2 = \|\theta^T A - b\|_2^2 \rightarrow \min_{\theta}$$

В случае полноранговой матрицы A , приравняв градиент по θ к нулю, получаем, что решение задается псевдообратной матрицей:

$$\theta^* = (A^T)^{\dagger} b^T = A(A^T A)^{-1} b^T$$

2. Метод максимального правдоподобия

Сделаем гипотезу, что измерения b_i не просто зависят линейно от a_i , но и имеют некоторый шум ξ_i :

$$b_i = \theta^T a_i + \xi_i$$

Предположим, что ξ_i — независимые значения одной и той же случайной величины $\xi \sim \mathcal{N}(0, \sigma^2)$, то есть шум нормальный:

$$p(x | \sigma) = \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{x^2}{2\sigma^2}\right)$$

Таким образом, функция правдоподобия:

$$L(\theta, \sigma) = \sum_{i=1}^m p(\xi_i | \sigma) = \sum_{i=1}^m p(b_i - \theta^T a_i | \sigma) = -\frac{m}{2} \ln(2\pi\sigma^2) - \frac{1}{2\sigma^2} \sum_{i=1}^m (b_i - \theta^T a_i)^2$$

Заметим, что в этом примере параметры θ являются параметрами распределения, но все равно входят в функцию правдоподобия.

Максимизация этой функции по θ равносильна минимизации $\|b - \theta^T A\|_2^2$, что как раз и есть метод наименьших квадратов.

Итак, мы показали, что следующие два подхода эквивалентны:

- искать такие параметры θ^* , что сумма квадратов отклонений минимальна;
- искать такие параметры θ^* , что невязки ξ_i как можно лучше описываются нормальным распределением.

Аналогичным образом, можно показать, что эквивалентны следующие подходы:

- искать такие параметры θ^* , что сумма модулей отклонений минимальна;
- искать такие параметры θ^* , что невязки ξ_i как можно лучше описываются распределением Лапласа.

8.3 Логистическая регрессия

Решается задача бинарной классификации.

Дано: точки (векторы признаков) $x_1, \dots, x_m \in \mathbb{R}^n$ и значения $y_1, \dots, y_m \in \{0, 1\}$ бинарной функции в этих точках.

Задача: построить функцию $\varphi : \mathbb{R}^n \rightarrow [0, 1]$, которая по набору признаков x будет давать вероятность $\varphi(x)$ того, что $y = 1$.

Предположим, что вероятность того, что $y = 1$ подчиняется *логистической функции* или *сигмоиде*:

$$\mathbb{P}\{y = 1 | x\} = \sigma(u) = \frac{1}{1 + e^{-u}},$$

где $u = u(x) \in \mathbb{R}$ — некоторая величина, характеризующая выборку. В логистической регрессии, как в линейной регрессии, используется линейная комбинация:

$$u = \theta_0 + \theta_1 x^{(1)} + \dots + \theta_n x^{(n)} = \theta_0 + \theta^T x$$

Таким образом:

$$\mathbb{P}\{y = 1 | x\} = \sigma(\theta^T x), \quad \mathbb{P}\{y = 0 | x\} = 1 - \sigma(\theta^T x)$$

Здесь и далее будем без ограничения общности опускать коэффициент θ_0 .

Можно записать компактно:

$$\mathbb{P}\{y | x\} = [\sigma(\theta^T x)]^y \cdot [1 - \sigma(\theta^T x)]^{1-y}, \quad y \in \{0, 1\}$$

Задача сводится к тому, что найти наилучшие коэффициенты θ . Найдем их *методом максимального правдоподобия*:

$$\prod_{i=1}^m \mathbb{P}\{y = y_i \mid x_i\} \longrightarrow \max_{\theta}$$

Функция правдоподобия:

$$L(\theta) = \sum_{i=1}^m \ln \mathbb{P}\{y = y_i \mid x_i\} = \sum_{i=1}^m \left[y_i \ln \sigma(\theta^T x_i) + (1 - y_i) \ln (1 - \sigma(\theta^T x_i)) \right] \longrightarrow \max_{\theta}$$

Можно упростить это выражение. Распишем $\sigma(u) = \frac{e^u}{1 + e^u}$, $1 - \sigma(u) = \frac{1}{1 + e^u}$, и тогда

$$L(\theta) = \sum_{i=1}^m \left[y_i (u_i - \ln(1 + e^{u_i})) - (1 - y_i) \ln(1 + e^{u_i}) \right] = \sum_{i=1}^m \left[y_i \cdot \theta^T x_i - \ln(1 + \exp(\theta^T x_i)) \right] \longrightarrow \max_{\theta}$$

Максимизация функции правдоподобия эквивалентна минимизации *логистической функции ошибки (log-loss function)*, это частный случай функции *кросс-энтропии* при числе классов $M = 2$:

$$\text{Log_loss}(y, p) = - \sum_{i=1}^m \left[y_i \ln p_i + (1 - y_i) \ln(1 - p_i) \right] \longrightarrow \min_{\theta},$$

где $p_i = \sigma(\theta^T x_i)$ — предсказываемые моделью вероятности, y_i — реальные значения.

Ясно, что бинарная энтропия достигает минимума, когда все $p_i = y_i$, но в логистической регрессии мы ищем их в особом виде, зависящем от коэффициентов θ , поэтому оптимальные значения будут другими:

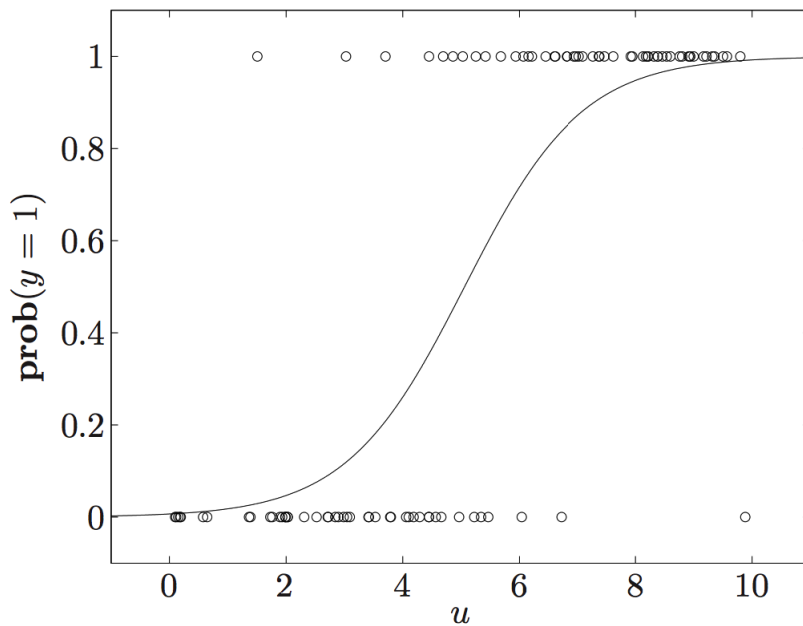


Иллюстрация из книги “Convex optimization”, Boyd, §7.1.1.

Кружочками отмечены реальные пары (u_i, y_i) из выборки, где $u = \theta^T x$, при этом параметры θ выбраны оптимальными. На графике также изображена сигмоида. Ее значения в каждой точке $\sigma(u_i) = p_i$ — предсказываемые моделью вероятности.

Итак, конечной моделью будет

$$\varphi(x) = \sigma(\theta^T x) = \frac{1}{1 + \exp(-\theta^T x)}.$$