

Семинар 2: Матрично-векторное дифференцирование

1 Матрично-векторное дифференцирование

1.1 Теория

Для вычисления большинства производных, которые возникают на практике, достаточно лишь небольшой таблицы стандартных производных и правил преобразования. Удобнее всего оказывается работать в терминах «дифференциала» — с ним можно не задумываться о промежуточных размерностях, а просто применять стандартные правила.

Замечание: В этом разделе описана сама техника матрично-векторного дифференцирования. Для более подробного описания математической теории, лежащей в основе этой техники, см. раздел А.

Правила преобразования	Таблица стандартных производных
$dA = 0$ $d(\alpha X) = \alpha(dX)$ $d(AXB) = A(dX)B$ $d(X + Y) = dX + dY$ $d(X^T) = (dX)^T$ $d(XY) = (dX)Y + X(dY)$ $d\langle X, Y \rangle = \langle dX, Y \rangle + \langle X, dY \rangle$ $d\left(\frac{X}{\phi}\right) = \frac{\phi dX - (d\phi)X}{\phi^2}$	$d\langle A, X \rangle = \langle A, dX \rangle$ $d\langle Ax, x \rangle = \langle (A + A^T)x, dx \rangle$ $d\langle Ax, x \rangle = 2\langle Ax, dx \rangle \quad (\text{если } A = A^T)$ $d(\text{Det}(X)) = \text{Det}(X)\langle X^{-T}, dX \rangle$ $d(X^{-1}) = -X^{-1}(dX)X^{-1}$

Здесь A, B — фиксированные матрицы; α — фиксированный скаляр; X, Y — произвольные дифференцируемые матричные функции (согласованные по размерностям, чтобы все операции имели смысл); ϕ — произвольная дифференцируемая скалярная функция.

Одним из самых важных является правило **производной композиции**. Пусть $g(Y)$ и $f(X)$ — две дифференцируемые функции, и мы знаем выражения для их дифференциалов: $dg(Y)$ и $df(X)$. Чтобы посчитать производную композиции $\phi(X) := g(f(X))$, как и в скалярном случае, нужно:

- взять выражение посчитанного дифференциала $dg(Y)$;
- подставить в него вместо Y значение $f(X)$, а вместо dY значение $df(X)$.

Пример
<p>Рассмотрим функцию $\phi(x) := \ln\langle Ax, x \rangle$, где $A \in \mathbb{S}_{++}^n$. В данном случае</p> $g(y) := \ln(y), \quad dg(y) = \frac{dy}{y}; \quad f(x) := \langle Ax, x \rangle, \quad df(x) = 2\langle Ax, dx \rangle.$ <p>Подставляем формально в $dg(y)$ вместо y выражение для $f(x) = \langle Ax, x \rangle$, а вместо dy выражение для $df(x) = 2\langle Ax, dx \rangle$:</p> $d\phi(x) = \frac{2\langle Ax, dx \rangle}{\langle Ax, x \rangle} \quad (\text{В нотации с «D»-большим: } D\phi(x)[h] = \frac{2\langle Ax, h \rangle}{\langle Ax, x \rangle}).$

Обычно, все возникающие на практике матрично-векторные функции составлены с помощью табличных функций и стандартных операций над ними. Благодаря универсальности приведённых правил,

дифференцировать сколь угодно сложные функции такого типа становится настолько же просто, как и дифференцировать одномерные функции.

Полученное в конце концов выражение нужно привести к одному из канонических видов:

Выход \ Вход	Скаляр	Вектор	Матрица
Скаляр	$df(x) = f'(x)dx$ ($f'(x)$: скаляр; dx : скаляр)	—	—
Вектор	$df(x) = \langle \nabla f(x), dx \rangle$ ($\nabla f(x)$: вектор; dx : вектор)	$df(x) = J_f(x)dx$ ($J_f(x)$: матрица; dx : вектор)	—
Матрица	$df(X) = \langle \nabla f(X), dX \rangle$ ($\nabla f(X)$: матрица; dX : матрица)	—	—

Случаи, отмеченные «—», нас интересовать не будут. Объект $\nabla f(x)$ (вектор для функции векторного аргумента и матрица для функции матричного аргумента) называется **градиентом**. Матрица $J_f(x)$ называется **матрицей Якоби**.

Найти *вторую производную* функции $f(X)$ можно по следующему «алгоритму»:

- посчитать первую производную функции; зафиксировать в выражении для $df(X)$ приращение dX — обозначить его как dX_1 ;
- посчитать производную для функции $g(X) = df(X)$, считая dX_1 фиксированным (константа). Новое приращение обозначать dX_2 .

Пример

Ввернёмся к функции $\phi(x) = \ln \langle Ax, x \rangle$, где $A \in \mathbb{S}_{++}^n$. Мы уже посчитали её первую производную: $d\phi(x) = \frac{2\langle Ax, dx \rangle}{\langle Ax, x \rangle}$. Обозначим dx за dx_1 и рассмотрим новую функцию:

$$g(x) = \frac{2\langle Ax, dx_1 \rangle}{\langle Ax, x \rangle}$$

Найдём производную $g(x)$, считая, что dx_1 — константный вектор:

$$\begin{aligned} d^2\phi(x) &= d\left(\frac{2\langle Ax, dx_1 \rangle}{\langle Ax, x \rangle}\right) = \frac{d(2\langle Ax, dx_1 \rangle)\langle Ax, x \rangle - 2\langle Ax, dx_1 \rangle d\langle Ax, x \rangle}{\langle Ax, x \rangle^2} \\ &= \frac{2\langle Adx_1, dx_2 \rangle \langle Ax, x \rangle - 2\langle Ax, dx_1 \rangle 2\langle Ax, dx_2 \rangle}{\langle Ax, x \rangle^2} = \left\langle \left(\frac{2A}{\langle Ax, x \rangle} - \frac{4Ax x^T A}{\langle Ax, x \rangle^2} \right) dx_1, dx_2 \right\rangle. \end{aligned}$$

(В нотации с D -большим: $D^2\phi(x)[h_1, h_2] = \left\langle \left(\frac{2A}{\langle Ax, x \rangle} - \frac{4Ax x^T A}{\langle Ax, x \rangle^2} \right) h_1, h_2 \right\rangle$.)

Для второй производной каноническая форма для скалярной функции векторного аргумента

$$d^2f(x) = \langle \nabla^2 f(x) dx_1, dx_2 \rangle.$$

Матрица $\nabla^2 f(x)$ называется **гессианом**. Для дважды непрерывно дифференцируемых функций гессиан является симметричной матрицей.

1.2 Задачи

Задача 1 (Квадратичная функция). Найти первую и вторую производные $df(x)$ и $d^2f(x)$, а также градиент $\nabla f(x)$ и гессиан $\nabla^2 f(x)$ функции

$$f(x) := \frac{1}{2} \langle Ax, x \rangle - \langle b, x \rangle + c, \quad x \in \mathbb{R}^n,$$

где $A \in \mathbb{S}^n$, $b \in \mathbb{R}^n$, $c \in \mathbb{R}$.

Решение. Найдем первую производную:

$$\boxed{df(x)} = d\left(\frac{1}{2}\langle Ax, x \rangle - \langle b, x \rangle + c\right) = \frac{1}{2}d\langle Ax, x \rangle - d\langle b, x \rangle = \frac{1}{2}2\langle Ax, dx \rangle - \langle b, dx \rangle = \langle Ax - b, dx \rangle.$$

Заметим, что $df(x)$ уже записан в канонической форме $df(x) = \langle \nabla f(x), dx \rangle$, поэтому

$$\boxed{\nabla f(x) = Ax - b}.$$

Теперь найдём вторую производную:

$$\boxed{d^2 f(x)} = d\langle Ax - b, dx_1 \rangle = \langle d(Ax - b), dx_1 \rangle = \langle d(Ax), dx_1 \rangle = \langle Adx_2, dx_1 \rangle.$$

Чтобы найти гессиан, приведем $d^2 f(x)$ к канонической форме $d^2 f(x) = \langle \nabla^2 f(x) dx_1, dx_2 \rangle$:

$$d^2 f(x) = \langle Adx_1, dx_2 \rangle \Rightarrow \boxed{\nabla^2 f(x) = A}.$$

Задача 2. Найти первую и вторую производные $df(x)$ и $d^2 f(x)$, а также градиент $\nabla f(x)$ и гессиан $\nabla^2 f(x)$ функции

$$f(x) := \frac{1}{2}\|Ax - b\|_2^2, \quad x \in \mathbb{R}^n,$$

где $A \in \mathbb{R}^{m \times n}$, $b \in \mathbb{R}^m$.

Решение. Найдем первую производную:

$$\boxed{df(x)} = d\left(\frac{1}{2}\|Ax - b\|_2^2\right) = \{d(\|x\|_2^2) = d\langle x, x \rangle = 2\langle x, dx \rangle\} = \frac{1}{2}2\langle Ax - b, d(Ax - b) \rangle = \langle Ax - b, Adx \rangle.$$

Чтобы найти градиент, приведем $df(x)$ к канонической форме $df(x) = \langle \nabla f(x), dx \rangle$:

$$df(x) = \langle A^T(Ax - b), dx \rangle \Rightarrow \boxed{\nabla f(x) = A^T(Ax - b)}.$$

Теперь найдём вторую производную:

$$\boxed{d^2 f(x)} = d\langle Ax - b, Adx_1 \rangle = \langle d(Ax - b), Adx_1 \rangle = \langle Adx_2, Adx_1 \rangle = \langle dx_2, A^T Adx_1 \rangle.$$

Чтобы найти гессиан, приведем $d^2 f(x)$ к канонической форме $d^2 f(x) = \langle \nabla^2 f(x) dx_1, dx_2 \rangle$:

$$d^2 f(x) = \langle A^T Adx_1, dx_2 \rangle \Rightarrow \boxed{\nabla^2 f(x) = A^T A}.$$

Задача 3 (Куб евклидовой нормы). Найти первую и вторую производные $df(x)$ и $d^2 f(x)$, а также градиент $\nabla f(x)$ и гессиан $\nabla^2 f(x)$ функции

$$f(x) := \frac{1}{3}\|x\|_2^3, \quad x \in \mathbb{R}^n.$$

Решение. Найдем первую производную:

$$\boxed{df(x)} = d\left(\frac{1}{3}\|x\|_2^3\right) = \frac{1}{3}d\langle x, x \rangle^{3/2} = \frac{1}{3} \frac{\mathcal{J}}{2} \langle x, x \rangle^{1/2} d\langle x, x \rangle = \frac{1}{2} \|x\|_2 (\mathcal{J} \langle x, dx \rangle) = \boxed{\|x\|_2 \langle x, dx \rangle}.$$

Чтобы найти градиент, приведем $df(x)$ к канонической форме $df(x) = \langle \nabla f(x), dx \rangle$:

$$df(x) = \langle \|x\|_2 x, dx \rangle \Rightarrow \boxed{\nabla f(x) = \|x\|_2 x}.$$

Теперь найдем вторую производную:

$$\begin{aligned} \boxed{d^2 f(x)} &= d(\|x\|_2 \langle x, dx_1 \rangle) = \underbrace{d(\|x\|_2)}_{=d(\langle x, x \rangle^{1/2})} \langle x, dx_1 \rangle + \|x\|_2 d\langle x, dx_1 \rangle \\ &= \left(\frac{1}{2} \langle x, x \rangle^{-1/2} (\mathcal{J} \langle x, dx_2 \rangle)\right) \langle x, dx_1 \rangle + \|x\|_2 \langle dx_2, dx_1 \rangle \\ &= \boxed{\|x\|_2^{-1} \langle x, dx_2 \rangle \langle x, dx_1 \rangle + \|x\|_2 \langle dx_2, dx_1 \rangle}. \end{aligned}$$

Чтобы найти гессиан, приведем $d^2 f(x)$ к канонической форме $d^2 f(x) = \langle \nabla^2 f(x) dx_1, dx_2 \rangle$:

$$\begin{aligned} d^2 f(x) &= \|x\|_2^{-1} \langle dx_1, x \rangle \langle x, dx_2 \rangle + \|x\|_2 \langle dx_1, dx_2 \rangle \\ &= \langle (\|x\|_2^{-1} x x^T + \|x\|_2 I_n) dx_1, dx_2 \rangle \Rightarrow \boxed{\nabla^2 f(x) = \|x\|_2^{-1} x x^T + \|x\|_2 I_n}. \end{aligned}$$

Отметим, что полученная формула для гессиана (и второй производной) верна только при $x \neq 0$, поскольку значение $\|x\|_2^{-1}$ не определено для $x = 0$. Такое ограничение возникло из-за того, что в самом начале мы воспользовались правилом произведения, и у нас возникла производная $d(\|x\|_2)$, которая не существует в точке $x = 0$. Тем не менее, можно показать, что рассматриваемая функция f является всюду дважды непрерывно дифференцируемой, и ее вторая производная в точке $x = 0$ равна нулю. Таким образом, можно сказать, что полученная формула, на самом деле, верна для любых значений x , с оговоркой, что в точке $x = 0$ значение $\|x\|_2^{-1} x x^T$ надо понимать как 0 (предел при $x \rightarrow 0$).

Задача 4 (Евклидова норма). Найти первую и вторую производные $df(x)$ и $d^2 f(x)$, а также градиент $\nabla f(x)$ и гессиан $\nabla^2 f(x)$ функции

$$f(x) := \|x\|_2, \quad x \in \mathbb{R}^n \setminus \{0\}.$$

Решение. Найдем первую производную:

$$\boxed{df(x)} = d(\|x\|_2) = d(\langle x, x \rangle^{1/2}) = \frac{1}{2} \langle x, x \rangle^{-1/2} d\langle x, x \rangle = \frac{1}{2} \|x\|_2^{-1} \mathcal{J} \langle x, dx \rangle = \boxed{\|x\|_2^{-1} \langle x, dx \rangle}.$$

Чтобы найти градиент, приведем $df(x)$ к канонической форме $df(x) = \langle \nabla f(x), dx \rangle$:

$$df(x) = \langle \|x\|_2^{-1} x, dx \rangle \Rightarrow \boxed{\nabla f(x) = \|x\|_2^{-1} x}.$$

Теперь найдем вторую производную:

$$\begin{aligned} \boxed{d^2 f(x)} &= d(\|x\|_2^{-1} \langle x, dx_1 \rangle) = d(\|x\|_2^{-1}) \langle x, dx_1 \rangle + \|x\|_2^{-1} d\langle x, dx_1 \rangle \\ &= -\|x\|_2^{-2} d(\|x\|_2) \langle x, dx_1 \rangle + \|x\|_2^{-1} \langle dx_2, dx_1 \rangle \\ &= -\|x\|_2^{-2} (\|x\|_2^{-1} \langle x, dx_2 \rangle) \langle x, dx_1 \rangle + \|x\|_2^{-1} \langle dx_2, dx_1 \rangle \\ &= \boxed{\|x\|_2^{-1} \langle dx_2, dx_1 \rangle - \|x\|_2^{-3} \langle x, dx_2 \rangle \langle x, dx_1 \rangle}. \end{aligned}$$

Чтобы найти гессиан, приведем $d^2 f(x)$ к канонической форме $d^2 f(x) = \langle \nabla^2 f(x) dx_1, dx_2 \rangle$:

$$\begin{aligned} d^2 f(x) &= \|x\|_2^{-1} (\langle dx_1, dx_2 \rangle - \|x\|_2^{-2} \langle dx_1, x \rangle \langle x, dx_2 \rangle) \\ &= \langle \|x\|_2^{-1} (I_n - \|x\|_2^{-2} x x^T) dx_1, dx_2 \rangle \Rightarrow \boxed{\nabla^2 f(x) = \|x\|_2^{-1} (I_n - \|x\|_2^{-2} x x^T)}. \end{aligned}$$

Задача 5 (Логистическая функция). Найти первую и вторую производные $df(x)$ и $d^2 f(x)$, а также градиент $\nabla f(x)$ и гессиан $\nabla^2 f(x)$ функции

$$f(x) := \ln(1 + \exp(\langle a, x \rangle)), \quad x \in \mathbb{R}^n,$$

где $a \in \mathbb{R}^n$.

Решение. Найдем первую производную:

$$\begin{aligned} \boxed{df(x)} &= d(\ln(1 + \exp(\langle a, x \rangle))) = \left\{ d(\ln(x)) = \frac{dx}{x} \right\} = \frac{d(1 + \exp(\langle a, x \rangle))}{1 + \exp(\langle a, x \rangle)} = \frac{d(\exp(\langle a, x \rangle))}{1 + \exp(\langle a, x \rangle)} \\ &= \{d(\exp(x)) = \exp(x) dx\} = \frac{\exp(\langle a, x \rangle) d\langle a, x \rangle}{1 + \exp(\langle a, x \rangle)} = \frac{\exp(\langle a, x \rangle) \langle a, dx \rangle}{1 + \exp(\langle a, x \rangle)} = \frac{\langle a, dx \rangle}{1 + \exp(\langle a, x \rangle)} \\ &= \boxed{\sigma(\langle a, x \rangle) \langle a, dx \rangle}. \end{aligned}$$

Здесь введено обозначение $\sigma : \mathbb{R} \rightarrow \mathbb{R}$ для *сигмоидной функции*:

$$\boxed{\sigma(x) := \frac{1}{1 + \exp(-x)}}.$$

Чтобы найти градиент, приведем $df(x)$ к канонической форме $df(x) = \langle \nabla f(x), dx \rangle$:

$$df(x) = \langle \sigma(\langle a, x \rangle) a, dx \rangle \Rightarrow \boxed{\nabla f(x) = \sigma(\langle a, x \rangle) a}.$$

Таким образом, градиент $\nabla f(x)$ — это вектор, коллинеарный вектору a с коэффициентом $\sigma(\langle a, x \rangle) \in (0, 1)$. В зависимости от точки x меняется лишь длина вектора $\nabla f(x)$, но не его направление.

Теперь найдем вторую производную:

$$\begin{aligned} \boxed{d^2 f(x)} &= d(\sigma(\langle a, x \rangle) \langle a, dx_1 \rangle) = d(\sigma(\langle a, x \rangle)) \langle a, dx_1 \rangle = \{d(\sigma(x)) = \sigma'(x) dx\} = (\sigma'(\langle a, x \rangle) d\langle a, x \rangle) \langle a, dx_1 \rangle \\ &= \sigma'(\langle a, x \rangle) \langle a, dx_2 \rangle \langle a, dx_1 \rangle = \{\sigma'(x) = \sigma(x)(1 - \sigma(x))\} \\ &= \boxed{\sigma(\langle a, x \rangle) (1 - \sigma(\langle a, x \rangle)) \langle a, dx_2 \rangle \langle a, dx_1 \rangle}. \end{aligned}$$

Чтобы найти гессиан, приведем $d^2 f(x)$ к канонической форме $d^2 f(x) = \langle \nabla^2 f(x) dx_1, dx_2 \rangle$:

$$\begin{aligned} d^2 f(x) &= \sigma(\langle a, x \rangle) (1 - \sigma(\langle a, x \rangle)) \langle dx_1, a \rangle \langle a, dx_2 \rangle \\ &= \langle (\sigma(\langle a, x \rangle) (1 - \sigma(\langle a, x \rangle)) a a^T) dx_1, dx_2 \rangle \Rightarrow \boxed{\nabla^2 f(x) = \sigma(\langle a, x \rangle) (1 - \sigma(\langle a, x \rangle)) a a^T}. \end{aligned}$$

Заметим, что гессиан $\nabla^2 f(x)$ — это одноранговая матрица, пропорциональная матрице aa^T с коэффициентом $\sigma(\langle a, x \rangle) (1 - \sigma(\langle a, x \rangle)) \in (0, 0.25)$. Точка x влияет лишь на коэффициент пропорциональности.

Задача 6 (Логарифм определителя). Найти первую и вторую производные $df(X)$ и $d^2 f(X)$, а также градиент $\nabla f(X)$ функции

$$f(X) := \ln(\text{Det}(X)),$$

заданной на множестве \mathbb{S}_{++}^n в пространстве \mathbb{S}^n .

Решение. Найдём первую производную:

$$\boxed{df(X)} = d(\ln \text{Det}(X)) = \left\{ d(\ln(x)) = \frac{dx}{x} \right\} = \frac{d(\text{Det}(X))}{\text{Det}(X)} = \frac{\text{Det}(X) \langle X^{-1}, dX \rangle}{\text{Det}(X)} = \langle X^{-1}, dX \rangle.$$

Заметим, что $df(X)$ уже и так записан в канонической форме $df(X) = \langle \nabla f(X), dX \rangle$.¹ Поэтому,

$$\boxed{\nabla f(X) = X^{-1}}.$$

Теперь найдем вторую производную:

$$\boxed{d^2 f(X)} = d\langle X^{-1}, dX_1 \rangle = \langle d(X^{-1}), dX_1 \rangle = \langle -X^{-1}(dX_2)X^{-1}, dX_1 \rangle = -\langle X^{-1}(dX_2)X^{-1}, dX_1 \rangle.$$

В итоге получилась билинейная форма от приращений dX_1 и dX_2 в пространстве матриц.

Рассмотрим

$$D^2 f(X)[H, H] = -\langle X^{-1}HX^{-1}, H \rangle.$$

Покажем, что $D^2 f(X)[H, H]$ имеет отрицательный знак для всех $X \in \mathbb{S}_{++}^n$ и $H \in \mathbb{S}^n$, т. е. что функция f является вогнутой функцией. Действительно, раскладывая $X^{-1} = X^{-1/2}X^{-1/2}$, перепишем $D^2 f(X)[H, H]$ в следующем виде:

$$D^2 f(X)[H, H] = -\langle X^{-1/2}HX^{-1/2}, X^{-1/2}HX^{-1/2} \rangle = -\|X^{-1/2}HX^{-1/2}\|_F^2.$$

Отсюда видно, что $D^2 f(X)[H, H]$, действительно, имеет отрицательный знак.

Задача 7. Найти производную $df(X)$ и градиент $\nabla f(X)$ функции

$$f(X) := \|AX - B\|_F, \quad X \in \mathbb{R}^{k \times n},$$

где $A \in \mathbb{R}^{m \times k}$, $B \in \mathbb{R}^{m \times n}$.

Решение. Вычислим отдельно $d(\|X\|_F)$:

$$\begin{aligned} d(\|X\|_F) &= d(\langle X, X \rangle^{1/2}) = \left\{ d(x^{1/2}) = \frac{1}{2}x^{-1/2}dx \right\} = \frac{1}{2}\langle X, X \rangle^{-1/2}d\langle X, X \rangle \\ &= \frac{1}{2}\|X\|_F^{-1}2\langle X, dX \rangle = \|X\|_F^{-1}\langle X, dX \rangle. \end{aligned}$$

Теперь используем полученную формулу, чтобы найти $df(X)$:

$$\begin{aligned} \boxed{df(X)} &= d(\|AX - B\|_F) = \|AX - B\|_F^{-1}\langle AX - B, d(AX - B) \rangle \\ &= \|AX - B\|_F^{-1}\langle AX - B, AdX \rangle. \end{aligned}$$

Чтобы найти градиент, приведем $df(X)$ к канонической форме $df(X) = \langle \nabla f(X), dX \rangle$:

$$df(X) = \langle \|AX - B\|_F^{-1}A^T(AX - B), dX \rangle \quad \Rightarrow \quad \boxed{\nabla f(X) = \|AX - B\|_F^{-1}A^T(AX - B)}.$$

¹В этом примере мы работаем в пространстве симметричных матриц \mathbb{S}^n , поэтому знак транспонирования можно опустить.

Задача 8. Найти производную $df(X)$ и градиент $\nabla f(X)$ функции

$$f(X) := \text{Tr}(AXBX^{-1}), \quad X \in \mathbb{R}^{n \times n}, \quad \text{Det}(X) \neq 0,$$

где $A, B \in \mathbb{R}^{n \times n}$.

Решение. Для удобства перепишем след через скалярное произведение:

$$f(X) = \langle I_n, AXBX^{-1} \rangle.$$

Найдем первую производную:

$$\begin{aligned} \boxed{df(X)} &= d\langle I_n, AXBX^{-1} \rangle = \langle I_n, d(AXBX^{-1}) \rangle = \langle I_n, (d(AXB))X^{-1} + (AXB)d(X^{-1}) \rangle \\ &= \langle I_n, (A(dX)B)X^{-1} + (AXB)(-X^{-1}(dX)X^{-1}) \rangle = \langle I_n, A(dX)BX^{-1} - AXBX^{-1}(dX)X^{-1} \rangle. \end{aligned}$$

Чтобы найти градиент, приведем $df(X)$ к канонической форме $df(X) = \langle \nabla f(X), dX \rangle$:

$$\begin{aligned} df(X) &= \langle I_n, A(dX)BX^{-1} \rangle - \langle I_n, AXBX^{-1}(dX)X^{-1} \rangle = \\ &= \langle A^T X^{-T} B^T, dX \rangle - \langle X^{-T} B^T X^T A^T X^{-T}, dX \rangle \\ &= \langle A^T X^{-T} B^T - X^{-T} B^T X^T A^T X^{-T}, dX \rangle \\ &\Rightarrow \boxed{\nabla f(X) = A^T X^{-T} B^T - X^{-T} B^T X^T A^T X^{-T}}. \end{aligned}$$

Задача 9. Рассмотрим функцию скалярного аргумента

$$\phi(\alpha) := f(x + \alpha p), \quad \alpha \in \mathbb{R},$$

где $x, p \in \mathbb{R}^n$, $f : \mathbb{R}^n \rightarrow \mathbb{R}$ — дважды непрерывно дифференцируемая функция. Найдите первую и вторую производные $\phi'(\alpha)$ и $\phi''(\alpha)$ и выразите их через градиент $\nabla f(\cdot)$ и гессиан $\nabla^2 f(\cdot)$.

Решение. В этой задаче нужно постоянно помнить, что дифференцирование выполняется по α , а x — постоянный вектор.

Найдем первую производную:

$$\begin{aligned} d\phi(\alpha) &= d_\alpha(f(x + \alpha p)) = \{df(x) = \langle \nabla f(x), dx \rangle\} = \langle \nabla f(x + \alpha p), d_\alpha(x + \alpha p) \rangle \\ &= \langle \nabla f(x + \alpha p), (d\alpha)p \rangle = \langle \nabla f(x + \alpha p), p \rangle d\alpha. \end{aligned}$$

Здесь последнее равенство следует из того, что $d\alpha$ — это скаляр. Заметим, что мы представили $d\phi(\alpha)$ в канонической форме $d\phi(\alpha) = \phi'(\alpha)d\alpha$. Значит,

$$\boxed{\phi'(\alpha) = \langle \nabla f(x + \alpha p), p \rangle}.$$

Теперь найдем вторую производную:

$$\begin{aligned} d^2\phi(\alpha) &= d_\alpha(\langle \nabla f(x + \alpha p), p \rangle d\alpha_1) = \langle d_\alpha \nabla f(x + \alpha p), p \rangle d\alpha_1 = \{d\nabla f(x) = \nabla^2 f(x)dx\} \\ &= \langle \nabla^2 f(x + \alpha p) d_\alpha(x + \alpha p), p \rangle d\alpha_1 = \langle \nabla^2 f(x + \alpha p)(d\alpha_2)p, p \rangle d\alpha_1 \\ &= \langle \nabla^2 f(x + \alpha p)p, p \rangle d\alpha_1 d\alpha_2. \end{aligned}$$

Таким образом, из канонической формы $d^2\phi(\alpha) = \phi''(\alpha)\alpha_1\alpha_2$, получаем

$$\boxed{\phi''(\alpha) = \langle \nabla^2 f(x + \alpha p)p, p \rangle}.$$

Задача 10. Рассмотрим функцию скалярного аргумента

$$\phi(\alpha) := \|r(x + \alpha p)\|_2, \quad \alpha \in \mathbb{R}_+, \quad r(x + \alpha p) \neq 0,$$

где $x, p \in \mathbb{R}^n$, $r : \mathbb{R}^n \rightarrow \mathbb{R}^m$ — дифференцируемое отображение. Найдите производную $\phi'(\alpha)$ и выразите ее через матрицу Якоби $J_r(\cdot)$.

Решение. В этой задаче, как и в предыдущей, нужно постоянно помнить, что дифференцирование выполняется по α , а x — постоянный вектор.

Найдем первую производную:

$$\begin{aligned} d\phi(\alpha) &= d_\alpha(\|r(x + \alpha p)\|_2) = \left\{ d\|x\|_2 = \frac{\langle x, dx \rangle}{\|x\|_2} \right\} = \frac{\langle r(x + \alpha p), d_\alpha r(x + \alpha p) \rangle}{\|r(x + \alpha p)\|_2} = \{dr(x) = J_r(x)dx\} \\ &= \frac{\langle r(x + \alpha p), J_r(x + \alpha p)d_\alpha(x + \alpha p) \rangle}{\|r(x + \alpha p)\|_2} = \frac{\langle r(x + \alpha p), J_r(x + \alpha p)(d\alpha)p \rangle}{\|r(x + \alpha p)\|_2} \\ &= \frac{\langle r(x + \alpha p), J_r(x + \alpha p)p \rangle}{\|r(x + \alpha p)\|_2} d\alpha. \end{aligned}$$

Отсюда

$$\boxed{\phi'(\alpha) = \frac{\langle r(x + \alpha p), J_r(x + \alpha p)p \rangle}{\|r(x + \alpha p)\|_2}}.$$

2 Условия оптимальности

2.1 Теория

Рассмотрим функцию $f : X \rightarrow \mathbb{R}$, где $X \subseteq U$ — подмножество нормированного линейного пространства, например, $U = \mathbb{R}^n$ — векторы или $U = \mathbb{R}^{n \times m}$ — матрицы.

Определение 2.1 (Локальные экстремумы). Точка $x \in X$ называется *точкой локального минимума* функции f , если существует шар некоторого радиуса $r > 0$, с центром в этой точке: $W = \{z \in U : \|z - x\| < r\}$, и выполнено:

$$f(x) \leq f(z) \quad \text{для любого } z \in W \cap X.$$

Если для всех $z \in W$ кроме $z = x$ в определении выполняется строгое неравенство, то локальный минимум называется *строгим локальным минимумом*.

Если неравенство выполнено в другую сторону, то точка x называется *локальным максимумом*. Локальные минимумы и максимумы вместе называются *локальными экстремумами*.

Если функция является дифференцируемой, то отыскать её локальные экстремумы иногда удаётся с помощью следующих утверждений.

Утверждение 2.1 (условие оптимальности первого порядка). Пусть для функции $f : X \rightarrow \mathbb{R}$ точка x является *точкой локального экстремума*.

Тогда если функция непрерывно-дифференцируема в окрестности этой точки, то её производная в этой точке равна нулю:

$$df(x) = 0.$$

Замечание 2.1. Равенство $df(x) = 0$ означает, что $Df(x)[h] = 0$ для любого $h \in U$.

Замечание 2.2. Напомним, что для функции векторного аргумента $f : \mathbb{R}^n \rightarrow \mathbb{R}$ ее производную всегда можно представить с помощью градиента: $Df(x)[h] = \langle \nabla f(x), h \rangle$. В этом случае равенство производной нулю эквивалентно равенству нулю градиента:

$$df(x) = 0 \quad \Leftrightarrow \quad \langle \nabla f(x), h \rangle = 0 \quad \text{для любого } h \in \mathbb{R}^n \quad \Leftrightarrow \quad \nabla f(x) = 0.$$

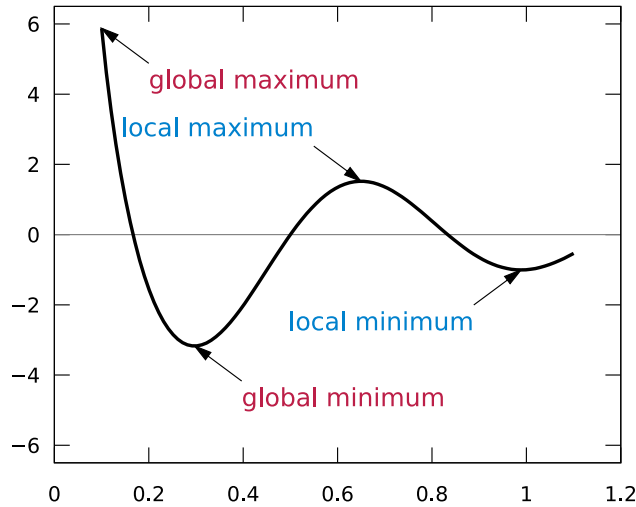


Рис. 1: Пример экстремумов из Википедии.

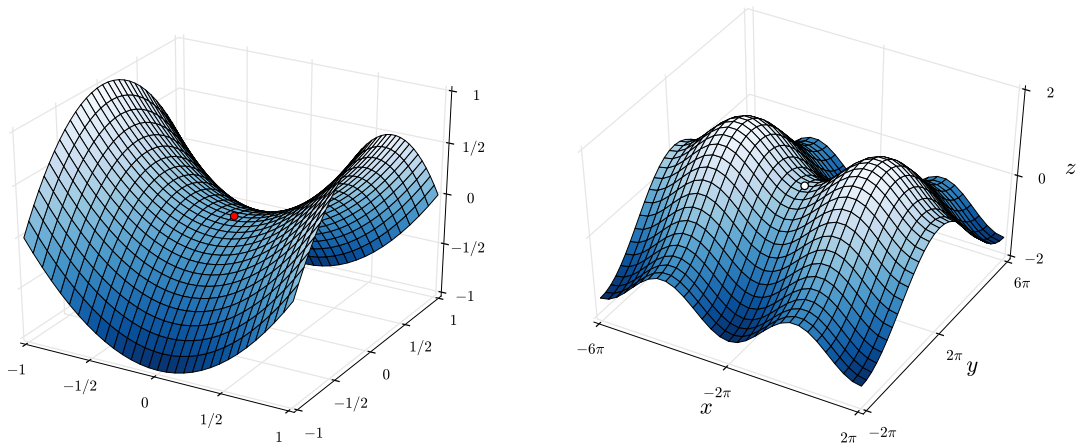


Рис. 2: Примеры седловых точек (слева — красная точка, справа — белая) из Википедии.

Замечание 2.3. Важно понимать, что утверждение 2.1 является лишь *необходимым* условием локального минимума. Например, для функции $f(x) = x^2$ производная равна нулю в единственной точке $x = 0$, и эта точка действительно является локальным минимумом функции (в данном случае и глобальным). Но для $f(x) = x^3$, в той же точке $x = 0$ — производная равна нулю, но сама точка не является ни локальным минимумом, ни локальным максимумом.

Определение 2.2. Точка x называется *стационарной точкой*, если производная в ней обращается в ноль: $df(x) = 0$. Стационарная точка, которая не является ни локальным минимумом, ни локальным максимумом, называется *седловой точкой*.

В случае функции двух переменных «седловая точка» действительно напоминает седло. Но мы используем это понятие для произвольных функций, в значении, указанном выше.

Для классификации стационарных точек удобным оказывается следующее утверждение.

Утверждение 2.2 (условия оптимальности второго порядка). Пусть функции $f : X \rightarrow \mathbb{R}$ дважды непрерывно дифференцируема в окрестности стационарной точки $x \in X$. Тогда:

- если $D^2f(x)[h, h] > 0$ для всех $h \in U, h \neq 0$, то x — строгий локальный минимум;
- если $D^2f(x)[h, h] < 0$ для всех $h \in U, h \neq 0$, то x — строгий локальный максимум;
- если существует $h, h' \in U$ такие, что: $D^2f(x)[h, h] > 0$ и $D^2f(x)[h', h'] < 0$, то x — седловая точка.

Также, справедливы и необходимые условия оптимальности второго порядка:

- если x — локальный минимум, то $D^2f(x)[h, h] \geq 0$ для всех $h \in U$;
- если x — локальный максимум, то $D^2f(x)[h, h] \leq 0$ для всех $h \in U$.

Замечание 2.4. Напомним, что для функции векторного аргумента $f : \mathbb{R}^n \rightarrow \mathbb{R}$ её вторую производную всегда можно представить с помощью матрицы — гессиана:

$$D^2f(x)[h, h] = \langle \nabla^2 f(x)h, h \rangle.$$

В этом случае утверждение 2.2 можно переписать в терминах знакоопределённости гессиана:

- если $\nabla^2 f(x) \succ 0$, то x — строгий локальный минимум;
- если $\nabla^2 f(x) \prec 0$, то x — строгий локальный максимум;
- если $\nabla^2 f(x)$ — неопределённая матрица (т.е. не выполнено ни $\nabla^2 f(x) \succeq 0$ ни $\nabla^2 f(x) \preceq 0$), то x — седловая точка.
- если x — локальный минимум, то $\nabla^2 f(x) \succeq 0$;
- если x — локальный максимум, то $\nabla^2 f(x) \preceq 0$.

Итак, в общем случае условие оптимальности первого порядка ($Df(x) = 0$) является только *необходимым* условием глобального минимума — точка x может быть не глобальным, а лишь локальным минимумом, или вообще локальным максимумом или седловой точкой. Тем не менее, для определённого класса функций — *класса выпуклых функций* — условие оптимальности первого порядка является не просто необходимым, но также и *достаточным* условием глобального минимума.

Утверждение 2.3 (Условие оптимальности первого порядка для выпуклой функции). Пусть X является открытым множеством, и пусть $f : X \rightarrow \mathbb{R}$ — выпуклая дифференцируемая функция. Точка $x^* \in X$ является глобальным минимумом функции f тогда и только тогда, когда $\nabla f(x^*) = 0$. Другими словами, любая стационарная точка выпуклой функции автоматически является глобальным минимумом.

2.2 Задачи

Задача 11. Рассмотрим функцию $f(x_1, x_2) = x_1^2 - x_2^2 + 2x_1$. Найти все её стационарные точки.

Решение. Найдем градиент и гессиан:

$$\nabla f(x_1, x_2) = \begin{bmatrix} 2x_1 + 2 \\ -2x_2 \end{bmatrix}$$

$$\nabla^2 f(x_1, x_2) = \begin{bmatrix} 2 & 0 \\ 0 & -2 \end{bmatrix} \quad - \quad \text{всюду постоянная матрица.}$$

Градиент обращается в ноль в единственной точке: $(-1, 0)$.

Гессиан — неопределённая матрица, значит стационарная точка $(-1, 0)$ — седловая.

Локальных экстремумов нет, функция неограниченная:

$$\inf_{x \in \mathbb{R}^2} f(x) = -\infty, \quad \sup_{x \in \mathbb{R}^2} f(x) = +\infty.$$

Задача 12 (Регрессия наименьших квадратов). Рассмотрим следующую задачу оптимизации:

$$\min_{x \in \mathbb{R}^n} \|Ax - b\|_2,$$

где $A \in \mathbb{R}^{m \times n}$, $b \in \mathbb{R}^m$, $\text{Rank}(A) = n$. Найдите множество ее решений и оптимальное значение целевой функции.

Решение. Прежде всего, перейдем от исходной негладкой задачи к эквивалентной ей гладкой:

$$\min_{x \in \mathbb{R}^n} \|Ax - b\|_2^2. \quad (2.1)$$

Заметим, что при таком переходе меняется лишь оптимальное значение целевой функции, а множество оптимальных решений остается неизменным.

Чтобы найти решения (2.1), воспользуемся условием оптимальности первого порядка:

$$\nabla(\|Ax - b\|_2^2) = 2A^T(Ax - b) = 0.$$

Заметим, что целевая функция в задаче (2.1) является выпуклой (как композиция выпуклой функции $x \mapsto \|x\|_2$ и аффинного преобразования $x \mapsto Ax - b$). Поэтому условие оптимальности первого порядка является не только необходимым, но также и достаточным условием того, что точка x является глобальным решением задачи (2.1). Итак, все решения задачи (2.1) — это решения следующей системы линейных уравнений, и только они:

$$A^T Ax = A^T b. \quad (2.2)$$

Уравнения (2.2) называются *нормальными уравнениями*. Заметим, что матрица $A^T A$ является обратной, поскольку она имеет полный ранг: согласно условию, $\text{Rank}(A) = n$, и из линейной алгебры известно, что $\text{Rank}(A^T A) = \text{Rank}(A)$. Отсюда заключаем, что нормальные уравнения (2.2) имеют единственное решение

$$x^* = (A^T A)^{-1} A^T b.$$

Таким образом, найденная точка x^* является единственным решением исходной задачи. Оптимальное значение целевой функции при этом равно

$$\boxed{\text{Opt}} := \|Ax^* - b\|_2 = \|A(A^T A)^{-1} A^T b - b\|_2.$$

Матрица $(A^T A)^{-1} A^T$ называется *псевдообратной* (по Муру-Пенроузу) и обозначается A^+ . Это прямое обобщение понятия обратной матрицы для неквадратных матриц. Если A квадратная и обратимая, то псевдообратная матрица A^+ совпадает с обратной матрицей A^{-1} .

Матрица $\Pi := A(A^T A)^{-1} A^T$ называется *проекционной матрицей* и проектирует заданный вектор b на линейную оболочку, натянутую на столбцы матрицы A .

Отметим, что задача (2.1) имеет решения всегда, даже если $\text{Rank}(A) < n$. Действительно, как было показано выше, задача (2.1) разрешима тогда и только тогда, когда разрешимы нормальные уравнения (2.2). Нормальные уравнения (2.2) разрешимы тогда и только тогда, когда $A^T b \in \text{Im}(A^T A)$, где $\text{Im}(A^T A)$ — образ матрицы $A^T A$. Из линейной алгебры известно, что $\text{Im}(A^T A) = \text{Im}(A^T)$. Значит, условие $A^T b \in \text{Im}(A^T A)$ всегда выполняется. Итак, задача (2.1) всегда разрешима, и в случае $\text{Rank}(A) < n$ имеет бесконечное множество решений. (Одно из возможных решений записывается через псевдообратную матрицу.)

Задача 13. Для следующей задачи оптимизации найдите ее множество решений и оптимальное значение целевой функции:

$$\min_{X \in \mathbb{S}_{++}^n} \{\langle S, X \rangle - \ln \text{Det}(X)\},$$

где $S \in \mathbb{S}_{++}^n$.

Решение. Введем обозначение $f(X) := \langle S, X \rangle - \ln \text{Det}(X)$. Запишем условие оптимальности первого порядка:

$$\nabla f(X) = S - X^{-1} = 0.$$

Заметим, что в нашем случае целевая функция f является выпуклой на рассматриваемом множестве \mathbb{S}_{++}^n как сумма двух выпуклых функций: $X \mapsto \langle S, X \rangle$ и $X \mapsto -\ln \text{Det}(X)$. Поэтому условие оптимальности первого порядка является не только необходимым, но также и достаточным условием глобального минимума.

Из уравнения оптимальности находим $X = S^{-1}$. Итак, оптимальное решение единственное — это

$$\boxed{X^* = S^{-1}}.$$

Соответствующее оптимальное значение целевой функции равно

$$\boxed{\text{Opt}} := f(X^*) = \langle S, S^{-1} \rangle - \ln \text{Det}(S^{-1}) = \boxed{n + \ln \text{Det}(S)}.$$

А Производные: теория

А.1 Определение

Начнём с напоминания понятия производной.

Для функции одной переменной $f : \mathbb{R} \rightarrow \mathbb{R}$ её производная в точке x обозначается $f'(x)$ и определяется из равенства:

$$f(x+h) = f(x) + f'(x)h + o(h) \quad \text{для всех достаточно малых } h.$$

Другими словами, зафиксировав некоторую точку x , мы хотим приблизить изменение функции $f(x+h) - f(x)$ в окрестности этой точки с помощью линейной функции по h , и $f'(x)h$ — наилучший способ это сделать.

Рассмотрим теперь более общую ситуацию.

Пусть U и V суть конечномерные линейные пространства с нормами. Основными примерами таких пространств для нас будут служить числа: \mathbb{R} , векторы: \mathbb{R}^n и матрицы: $\mathbb{R}^{n \times m}$, а также их комбинации (декартовы произведения).

Рассмотрим функцию $f : X \rightarrow V$, где $X \subseteq U$.

Определение А.1 (Дифференцируемость). Пусть $x \in X$ — внутренняя точка множества X , и пусть $L : U \rightarrow V$ — линейный оператор. Будем говорить, что функция f дифференцируема в точке x с производной L , если для всех достаточно малых $h \in U$ справедливо следующее разложение:

$$f(x+h) = f(x) + L[h] + o(\|h\|). \quad (\text{А.1})$$

Если для любого линейного оператора $L : U \rightarrow V$ функция f не является дифференцируемой в точке x с производной L , то будем говорить, что f не является дифференцируемой в точке x . Если точка x не является внутренней точкой множества X , то оставим понятие дифференцируемости функции f в точке x неопределённым.

Замечание А.1. Выражение $o(\|h\|)$ имеет стандартное значение:

$$f(x+h) - f(x) - L[h] = o(\|h\|) \quad \stackrel{\text{def}}{\iff} \quad \lim_{h \rightarrow 0} \frac{\|f(x+h) - f(x) - L[h]\|}{\|h\|} = 0.$$

Замечание А.2. Поскольку рассматриваемые пространства U и V являются конечномерными (а в конечномерном пространстве все нормы топологически эквивалентны), то не имеет значения, какие конкретно нормы используются в данном выше определении: если функция f является дифференцируемой в точке x с производной L для одного выбора норм, то f также будет дифференцируемой в точке x с производной L для любого другого выбора норм.

Утверждение А.1. Предположим, что функция f дифференцируема в точке x с производной L_1 и также дифференцируема в точке x с производной L_2 . Тогда $L_1 = L_2$.

Таким образом, если функция f является дифференцируемой в точке x , то её производная L определяется единственным образом. Будем обозначать её символом $df(x)$.

Замечание А.3. Объект df зависит от двух параметров: точка $x \in X$, в которой мы аппроксимируем функцию, и приращение $h \in U$, которое откладывается от зафиксированной точки:

$$df : X \times U \rightarrow V, \quad \text{линейный по второму аргументу — по «}h\text{»}.$$

Замечание А.4. Встречаются разные обозначения производной функции f в точке x :

$$Df(x)[h] \equiv df(x)[h] \equiv Df(x)[\Delta x] \equiv df(x)[\Delta x].$$

Все они обозначают одно и то же. При работе с определением производной, удобно явным образом указывать приращение (h или Δx) в квадратных скобках. При вычислении производных на практике, пользуясь уже известными посчитанными производными и свойствами пересчёта, приращение в квадратных скобках обычно не пишут: $df(x)$ или даже просто df , когда понятно, о чём идёт речь.

Итак, производная функции в точке x — это линейный оператор $df(x)$ который лучше всего аппроксимирует приращение функции:

$$f(x+h) - f(x) \approx Df(x)[h].$$

Ещё одним известным и важным понятием является производная функции по направлению. Оказывается, зная производную функции f мы можем легко посчитать её производную вдоль любого направления h .

Утверждение А.2. Пусть f дифференцируема в точке x . Выберем произвольное направление h . Тогда:

$$Df(x)[h] = \frac{\partial f(x)}{\partial h} := \lim_{t \rightarrow +0} \frac{f(x+th) - f(x)}{t}$$

То есть, чтобы посчитать $\frac{\partial f(x)}{\partial h}$ — производную функции f вдоль направления h , достаточно применить $df(x)[\cdot]$ к этому направлению.

Набор векторов

$$e_i = (0, \dots, 0, \underset{i}{1}, 0, \dots, 0) \in \mathbb{R}^n, \quad i = 1, \dots, n$$

называется *стандартным базисом* в \mathbb{R}^n .

Если для некоторого i у функции существует (двусторонняя) производная вдоль направления e_i , то её называют *частной производной по i -ой координате*:

$$\frac{\partial f(x)}{\partial x_i} := \lim_{t \rightarrow 0} \frac{f(x+te_i) - f(x)}{t} = Df(x)[e_i].$$

Обратите внимание, что функция может быть недифференцируемой, даже если у неё существуют производные по всем направлениям.

Пример А.1. Рассмотрим функцию $f(x) = \|x\|_2$. Найдём её производную по направлению h в точке $x = 0$:

$$\left. \frac{\partial \|x\|_2}{\partial h} \right|_{x=0} = \lim_{t \rightarrow +0} \frac{\|0+th\|_2 - \|0\|_2}{t} = \lim_{t \rightarrow +0} \frac{t\|h\|_2}{t} = \|h\|_2.$$

Если бы функция $f(x) = \|x\|_2$ была дифференцируема в нуле, то по утверждению 2:

$$Df(0)[h] = \frac{\partial f(0)}{\partial h} = \|h\|_2,$$

но функция $\|h\|_2$ не является линейной, что противоречит тому, что производная — линейный оператор. Значит $\|x\|_2$ не дифференцируема в нуле, хотя и имеет производные вдоль всех направлений.

Градиент функции, матрица Якоби.

- В случае $U = \mathbb{R}^n$, $V = \mathbb{R}$ линейную функцию $Df(x)[h]$ всегда можно представить с помощью скалярного произведения с некоторым вектором:

$$Df(x)[h] = \langle a_x, h \rangle \quad \text{где } a_x \in \mathbb{R}^n \text{ — разный для каждого } x.$$

Вектор a_x называется *градиентом* функции f в точке x и обозначается $\nabla f(x)$.

В стандартном базисе градиент функции представляется в виде вектора из частных производных:

$$\nabla f(x) = \left(\frac{\partial f}{\partial x_1}(x), \dots, \frac{\partial f}{\partial x_n}(x) \right) \in \mathbb{R}^n.$$

Как и все векторы у нас, этот вектор — вектор-столбец.

- В случае $U = \mathbb{R}^{n \times m}$, $V = \mathbb{R}$ линейную функцию $df(x)[H]$ всегда можно представить с помощью скалярного произведения с некоторой матрицей:

$$df(x)[H] = \langle A_x, H \rangle, \quad A_x, H \in \mathbb{R}^{n \times m}.$$

Эта матрица также называется *градиентом функции* в точке x : $\nabla f(x) = A_x$ и в стандартном базисе (из матриц, у которых все нули, кроме одной единички) записывается в виде матрицы частных производных:

$$\nabla f(x) = \left(\frac{\partial f}{\partial x_{ij}}(x) \right)_{i=1, j=1}^{n, m}$$

- В случае $U = \mathbb{R}^m$, $V = \mathbb{R}^n$, линейный оператор $df(x)[\cdot]$, зафиксировав базисы, всегда можно представить матрицей:

$$Df(x)[h] = J_x h, \quad J_x \in \mathbb{R}^{n \times m}.$$

Матрица J_x называется *матрицей Якоби* функции f . В стандартном базисе она состоит из частных производных:

$$J_x = \left(\frac{\partial f_i}{\partial x_j}(x) \right)_{i=1, j=1}^{n, m}.$$

Утверждение А.3 (Дифференциальное исчисление). Пусть U и V — векторные пространства, X — подмножество U , $x \in X$ — внутренняя точка X . Справедливы следующие свойства:

- (Производная константы) Пусть $f : X \rightarrow V$ — постоянная функция, т. е. найдется $v \in V$, что $f(x') = v$ для всех $x' \in X$. Тогда f дифференцируема в точке x , и $df(x) = 0$.
- (Производная тождественной функции) Пусть $f : X \rightarrow V$ — тождественная функция, т. е. $f(x') = x'$ для всех $x' \in X$. Тогда f дифференцируема в точке x , и ее производная — также тождественная функция: $Df(x)[h] = h$ для всех $h \in U$.
- (Линейность) Пусть $f : X \rightarrow V$ и $g : X \rightarrow V$ — функции. Если f и g дифференцируемы в точке x , а $c_1, c_2 \in \mathbb{R}$ — числа, то функция $(c_1 f + c_2 g)$ также дифференцируема в точке x , и

$$d(c_1 f + c_2 g)(x) = c_1 df(x) + c_2 dg(x).$$

- (Правило произведения) Пусть $\alpha : X \rightarrow \mathbb{R}$ и $f : X \rightarrow V$ — функции. Если α и f дифференцируемы в точке x , то функция αf также дифференцируема в точке x , и

$$D(\alpha f)(x)[h] = (D\alpha(x)[h])f(x) + \alpha(x)(Df(x)[h])$$

для всех $h \in U$.

- (Правило композиции) Пусть Y — подмножество V , $f : X \rightarrow Y$ — функция. Также пусть W — векторное пространство, $g : Y \rightarrow W$ — функция. Если f дифференцируема в точке x , и g дифференцируема в точке $f(x)$, то их композиция $(g \circ f) : X \rightarrow W$ (определенная как $(g \circ f)(x) = g(f(x))$) также будет дифференцируема в точке x , и

$$D(g \circ f)(x) = Dg(f(x)) \left[df(x) \right] \quad \text{или, более подробно,} \quad D(g \circ f)(x)[h] = Dg(f(x)) \left[Df(x)[h] \right].$$

- (Правило частного) Пусть $\alpha : X \rightarrow \mathbb{R}$ и $f : X \rightarrow V$ — функции. Если α и f дифференцируемы в точке x , и если α не обращается в ноль на X , то функция $(1/\alpha)f$ также дифференцируема в точке x , и

$$D\left(\frac{1}{\alpha}f\right)(x)[h] = \frac{\alpha(x)(Df(x)[h]) - (D\alpha(x)[h])f(x)}{\alpha(x)^2}.$$

для всех $h \in U$.

Доказательство. Первые четыре свойства доказываются по определению, а последнее выводится из правил произведения и композиции. \square

Заметим, что правило произведения в утверждении А.3 установлено только в случае, когда одна из функций является скалярной. Это понятно, поскольку в векторных пространствах определено лишь умножение на скаляр, а не на произвольный элемент векторного пространства. Тем не менее, в некоторых частных случаях правило произведения остается верным даже если обе функции являются не скалярными. Например, справедливо следующее утверждение.

Утверждение А.4. Пусть U — векторное пространство, X — подмножество U , $x \in X$ — внутренняя точка X . Пусть $f : X \rightarrow \mathbb{R}^{m \times n}$ и $g : X \rightarrow \mathbb{R}^{n \times k}$ — матрично-значные функции. Предположим, что f и g дифференцируемы в точке x . Тогда функция fg также дифференцируема в точке x , и

$$D(fg)(x)[h] = (Df(x)[h])g(x) + f(x)(Dg(x)[h]).$$

для всех $h \in U$. (Здесь под операцией умножения подразумевается матричное умножение, поэтому порядок множителей имеет значение.)

А.2 Вторая производная

Пусть функция $f : X \rightarrow V$ дифференцируема в каждой точке $x \in X \subseteq U$.

Рассмотрим производную функции f при фиксированном приращении $h_1 \in U$ как функцию от x :

$$g(x) = Df(x)[h_1].$$

Определение А.2. Если для функции g в некоторой точке x существует производная, то она называется *второй производной* функции f в точке x :

$$D^2f(x)[h_1, h_2] := Dg(x)[h_2].$$

Можно показать, что $D^2f(x)[h_1, h_2]$ является билинейной функцией по h_1 и h_2 .

По аналогии определяются третья: $D^3f(x)[h_1, h_2, h_3]$, четвёртая и производные более высоких порядков.

Если производная $df(x)$ является непрерывной функцией по x , то говорят, что f — *непрерывно дифференцируема*. Если вторая производная $D^2f(x)$ непрерывна по x , то тогда f — *дважды непрерывно дифференцируема*.

Для функций $f : \mathbb{R}^n \rightarrow \mathbb{R}$ вторую производную, как и любую билинейную форму, можно представить с помощью матрицы:

$$D^2f(x)[h_1, h_2] = \langle H_x h_1, h_2 \rangle, \quad H_x \in \mathbb{R}^{n \times n}.$$

Матрица H_x называется *гессианом* функции f в точке x и обозначается обычно $\nabla^2 f(x)$. В стандартном базисе эта матрица состоит из вторых частных производных:

$$\nabla^2 f(x) = \left(\frac{\partial^2 f}{\partial x_i \partial x_j}(x) \right)_{i=1, j=1}^{n, n}$$

Для дважды непрерывно дифференцируемой функции её гессиан — симметричная матрица:

$$\nabla^2 f(x) \in \mathbb{S}^n.$$

А.3 Формула Тейлора

Для дважды непрерывно-дифференцируемой функции справедлива формула Тейлора:

$$f(x+h) = f(x) + Df(x)[h] + \frac{1}{2}D^2f(x)[h, h] + o(\|h\|^2).$$

Для функции $f : \mathbb{R}^n \rightarrow \mathbb{R}$ её можно записать, используя градиент и гессиан:

$$f(x+h) = f(x) + \langle \nabla f(x), h \rangle + \frac{1}{2} \langle \nabla^2 f(x)h, h \rangle + o(\|h\|^2).$$

Если функция имеет непрерывные производные до порядка k включительно, то формулу Тейлора можно записать до k -ой производной:

$$f(x+h) = f(x) + Df(x)[h] + \frac{1}{2!}D^2f(x)[h, h] + \frac{1}{3!}D^3f(x)[h, h, h] + \dots + \frac{1}{k!}D^k f(x)[h, \dots, h] + o(\|h\|^k).$$

А.4 Подсчет табличных производных

Замечание: всюду в дальнейшем $\|\cdot\|$ обозначает (для краткости) евклидову норму для векторов и спектральную (операторную) норму для матриц.

Пример А.2 (Линейная функция). Пусть $c \in \mathbb{R}^n$, и пусть $f : \mathbb{R}^n \rightarrow \mathbb{R}$ — функция $f(x) := \langle c, x \rangle$. Покажем, что f является дифференцируемой в произвольной точке $x \in \mathbb{R}^n$ и найдем ее производную $df(x) : \mathbb{R}^n \rightarrow \mathbb{R}$. Для этого зафиксируем произвольное приращение аргумента $h \in \mathbb{R}^n$ и вычислим соответствующее приращение функции:

$$f(x+h) - f(x) = \langle c, x+h \rangle - \langle c, x \rangle = \langle c, h \rangle.$$

Заметим, что отображение $h \mapsto \langle c, h \rangle$ является линейным. Значит, для функции f справедливо разложение (А.1) с $Df(x)[h] := \langle c, h \rangle$. Таким образом, функция f является дифференцируемой в произвольной точке $x \in \mathbb{R}^n$ с производной $Df(x)[h] = \langle c, h \rangle$.

Пример А.3 (Квадратичная форма). Пусть $A \in \mathbb{R}^{n \times n}$, и пусть $f : \mathbb{R}^n \rightarrow \mathbb{R}$ — функция $f(x) := \langle Ax, x \rangle$. Зафиксируем произвольную точку $x \in \mathbb{R}^n$ и произвольное приращение аргумента $h \in \mathbb{R}^n$ и вычислим соответствующее приращение функции:

$$f(x+h) - f(x) = \langle A(x+h), x+h \rangle - \langle Ax, x \rangle = \langle (A+A^T)x, h \rangle + \langle Ah, h \rangle.$$

Заметим, что отображение $h \mapsto \langle (A+A^T)x, h \rangle$ является линейным, а $\langle Ah, h \rangle = o(\|h\|)$, поскольку для всех $h \in \mathbb{R}^n$ справедлива следующая цепочка неравенств:

$$|\langle Ah, h \rangle| \leq \|h\| \|Ah\| \leq \|A\| \|h\|^2.$$

Здесь первое неравенство следует из неравенства Коши-Буняковского; второе неравенство следует из согласованности матричной и векторной норм. Таким образом, функция f дифференцируема в произвольной точке $x \in \mathbb{R}^n$ с производной $Df(x)[h] = \langle (A+A^T)x, h \rangle$.

Пример А.4 (Обратная матрица). Пусть $S := \{X \in \mathbb{R}^{n \times n} : \text{Det}(X) \neq 0\}$ — множество всех квадратных невырожденных матриц размера n . Рассмотрим функцию $f : S \rightarrow S$, которая для каждой матрицы $X \in S$ возвращает ее обратную: $f(X) := X^{-1}$. Покажем, что f является дифференцируемой в любой точке $X \in S$. Для этого зафиксируем произвольное достаточно малое приращение аргумента $H \in \mathbb{R}^{n \times n}$ (удовлетворяющее $X+H \in S$ и $\|H\| < 1/\|X^{-1}\|$) и рассмотрим соответствующее приращение функции:

$$f(X+H) - f(X) = (X+H)^{-1} - X^{-1} = (X(I_n + X^{-1}H))^{-1} - X^{-1} = ((I_n + X^{-1}H)^{-1} - I_n)X^{-1}.$$

Оценим отдельно $(I_n + X^{-1}H)^{-1}$. Для этого разложим эту матрицу в ряд Неймана:²

$$(I_n + X^{-1}H)^{-1} = I_n - X^{-1}H + \sum_{k=2}^{\infty} (-X^{-1}H)^k.$$

Заметим, что ряд, стоящий в правой части последнего равенства, является абсолютно сходящимся, поскольку $\|X^{-1}H\| < 1$ в силу достаточной малости H . Покажем, что сумма этого ряда есть $o(\|H\|)$:

$$\left\| \sum_{k=2}^{\infty} (-X^{-1}H)^k \right\| \leq \sum_{k=2}^{\infty} \|(-X^{-1}H)^k\| \leq \sum_{k=2}^{\infty} \|X^{-1}\|^k \|H\|^k = \frac{\|X^{-1}\|^2 \|H\|^2}{1 - \|X^{-1}\| \cdot \|H\|}.$$

Здесь первое неравенство следует из неравенства треугольника для нормы; второе неравенство следует из субмультипликативности нормы; далее вычисляется сумма геометрического ряда. Таким образом,

$$(I_n + X^{-1}H)^{-1} = I_n - X^{-1}H + o(\|H\|).$$

Подставляя это выражение в полученную выше формулу для приращения функции, получаем

$$f(X + H) - f(X) = -X^{-1}HX^{-1} + o(\|H\|).$$

Таким образом, функция f дифференцируема в произвольной точке $X \in S$ с производной $df(X)[H] = -X^{-1}HX^{-1}$.

Замечание А.5. Выведенную формулу для производной функции X^{-1} можно очень просто получить с помощью следующего трюка. Рассмотрим дифференциал единичной матрицы $d(I_n)$. С одной стороны, поскольку матрица постоянная, $d(I_n) = 0$. С другой стороны, по правилу произведения, $dI_n = d(XX^{-1}) = (dX)X^{-1} + Xd(X^{-1})$. Приравняв выражения, получим $d(X^{-1}) = -X^{-1}(dX)X^{-1}$, или, в другой форме, $d(X^{-1})[H] = -X^{-1}HX^{-1}$. Заметим, однако, что приведённое рассуждение не является полным доказательством тождества, так как предполагает, но не доказывает существование дифференциала $d(X^{-1})$.

Пример А.5 (Определитель матрицы). Пусть $f : \mathbb{R}^{n \times n} \rightarrow \mathbb{R}$ — функция $f(X) := \text{Det}(X)$. Рассмотрим произвольную точку $X \in \mathbb{R}^{n \times n}$ и произвольное приращение аргумента $H \in \mathbb{R}^{n \times n}$. Будем предполагать, что матрица X обратима. Выпишем соответствующее приращение функции:

$$f(X+H) - f(X) = \text{Det}(X+H) - \text{Det}(X) = \text{Det}(X(I_n + X^{-1}H)) - \text{Det}(X) = \text{Det}(X)(\text{Det}(I_n + X^{-1}H) - 1).$$

Оценим отдельно $\text{Det}(I_n + X^{-1}H)$. Для этого воспользуемся тем, что определитель матрицы равен произведению ее собственных значений. Пусть $\lambda_1(X^{-1}H), \dots, \lambda_n(X^{-1}H)$ — собственные значения матрицы $X^{-1}H$ (пронумерованные в произвольном порядке и, возможно, комплексные). Заметим, что собственными значениями матрицы $I_n + X^{-1}H$ будут $1 + \lambda_1(X^{-1}H), \dots, 1 + \lambda_n(X^{-1}H)$. Поэтому

$$\text{Det}(I_n + X^{-1}H) = \prod_{i=1}^n [1 + \lambda_i(X^{-1}H)] = 1 + \sum_{i=1}^n \lambda_i(X^{-1}H) + \left(\sum_{1 \leq i < j \leq n} \lambda_i(X^{-1}H)\lambda_j(X^{-1}H) + \dots \right),$$

где многоточие скрывает сумму всевозможных троек $\lambda_i(X^{-1}H)\lambda_j(X^{-1}H)\lambda_k(X^{-1}H)$, всевозможных четверок и т. д. Заметим, что выражение, стоящее в скобках, представляет из себя величину $o(\|H\|)$. Это следует из неравенства треугольника и того факта, что для произвольной матрицы $A \in \mathbb{R}^{n \times n}$ все

²Имеется в виду разложение $(I_n - A)^{-1} = \sum_{k=0}^{\infty} A^k$, справедливое для любой матрицы $A \in \mathbb{R}^{n \times n}$, такой, что $\|A\| < 1$.

Эта формула является обобщением известной формулы для суммы геометрического ряда: $(1 - q)^{-1} = \sum_{k=0}^{\infty} q^k$ для любого $|q| < 1$.

ее собственные значения не превосходят по модулю ее нормы $\|A\|$. (Действительно, пусть $\lambda \in \mathbb{C}$ — собственное значение матрицы A , и пусть $x \in \mathbb{C}^n \setminus \{0\}$ — соответствующий собственный вектор: $Ax = \lambda x$. Тогда $|\lambda|\|x\| = \|Ax\| \leq \|A\|\|x\|$.) Таким образом,

$$\text{Det}(I_n - X^{-1}H) = 1 + \sum_{i=1}^n \lambda_i(X^{-1}H) + o(\|H\|) = 1 + \text{Tr}(X^{-1}H) + o(\|H\|).$$

Подставляя полученное выражение в полученную выше формулу для приращения функции, получаем

$$f(X + H) - f(X) = \text{Det}(X) \text{Tr}(X^{-1}H) + o(\|H\|).$$

Таким образом, для любой обратимой матрицы $X \in \mathbb{R}^{n \times n}$ функция f дифференцируема в точке X с производной $df(x)[H] = \text{Det}(X) \text{Tr}(X^{-1}H) = \text{Det}(X) \langle X^{-T}, H \rangle$.

Замечание А.6. Можно показать, что рассматриваемая функция $f(X) = \text{Det}(X)$ будет дифференцируемой всюду на $\mathbb{R}^{n \times n}$, а не только на подмножестве обратимых матриц. Общая формула для производной в этом случае называется *формулой Якоби* и выглядит следующим образом: $df(X)[H] = \text{Tr}(\text{Adj}(X)H)$, где $\text{Adj}(X)$ — присоединенная матрица к X . Заметим, что если X — невырожденная матрица, тогда $\text{Adj}(X) = \text{Det}(X)X^{-1}$ и формула Якоби переходит в доказанную формулу $df(X)[H] = \text{Det}(X) \text{Tr}(X^{-1}H)$.