

# Matrix calculus

## Useful definitions and notations

We will treat all vectors as column vectors by default.

### Matrix and vector multiplication

Let  $A$  be  $m \times n$ , and  $B$  be  $n \times p$ , and let the product  $AB$  be

then  $C$  is a  $m \times p$  matrix, with element  $(i, j)$  given by

$$\alpha^{(n)} - B^T$$

$$C = AB \quad m \times p \quad m \times n \times p$$

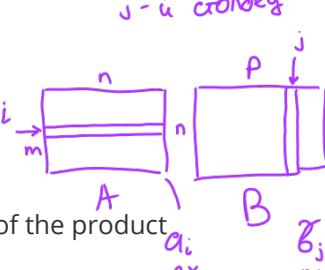
$$c_{ij} = \sum_{k=1}^n a_{ik} b_{kj}$$

$$c_{ij} = \alpha_i^T \cdot B_j$$



свойство

$i$ -ая строка  
 $j$ -ий столбец



$O(n^3)$

Множен

$O(n^{2.7...})$

Использовать

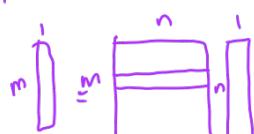
$O(n^{2.4})$

не торк

Python: Let  $A$  be  $m \times n$ , and  $x$  be  $n \times 1$ , then the typical element of the product

$$z = A @ (B @ (C @ z))$$

is given by



$$z = Ax$$

$$z_i = \sum_{k=1}^n a_{ik} x_k$$

свойство  
 $O(n^2)$

Finally, just to remind:

- $C = AB \quad C^T = B^T A^T$
- $AB \neq BA$
- $e^A = \sum_{k=0}^{\infty} \frac{1}{k!} A^k$
- $\langle y, x \rangle = y^T x$
- $\langle x, y \rangle = x^T y = \sum_{i=1}^n x_i y_i$
- $\langle x, Ay \rangle = \langle A^T x, y \rangle$

$$(AB)^T = B^T \cdot A^T$$

$$\langle y, x \rangle = y^T x$$

$$\langle x, y \rangle = x^T y = \sum_{i=1}^n x_i y_i$$

### Gradient

Gradient Let  $f(x) : \mathbb{R}^n \rightarrow \mathbb{R}$ , then vector, which contains all first order partial derivatives:

$$\nabla f(x) = \frac{df}{dx} = \begin{pmatrix} \frac{\partial f}{\partial x_1} \\ \frac{\partial f}{\partial x_2} \\ \vdots \\ \frac{\partial f}{\partial x_n} \end{pmatrix} \in \mathbb{R}^n$$

$$f(x_1, \dots, x_n)$$

### Hessian

Let  $f(x) : \mathbb{R}^n \rightarrow \mathbb{R}$ , then matrix, containing all the second order partial derivatives:

$$f''(x) = \frac{\partial^2 f}{\partial x_i \partial x_j} = \begin{pmatrix} \frac{\partial^2 f}{\partial x_1 \partial x_1} & \frac{\partial^2 f}{\partial x_1 \partial x_2} & \cdots & \frac{\partial^2 f}{\partial x_1 \partial x_n} \\ \frac{\partial^2 f}{\partial x_2 \partial x_1} & \frac{\partial^2 f}{\partial x_2 \partial x_2} & \cdots & \frac{\partial^2 f}{\partial x_2 \partial x_n} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial^2 f}{\partial x_n \partial x_1} & \frac{\partial^2 f}{\partial x_n \partial x_2} & \cdots & \frac{\partial^2 f}{\partial x_n \partial x_n} \end{pmatrix} \in \mathbb{R}^{n \times n}$$

But actually, Hessian could be a tensor in such a way:  $(f(x) : \mathbb{R}^n \rightarrow \mathbb{R}^m)$  is just 3d tensor, every slice is just hessian of corresponding scalar function ( $H(f_1(x)), H(f_2(x)), \dots, H(f_m(x))$ )

## Jacobian

The extension of the gradient of multidimensional  $f(x) : \mathbb{R}^n \rightarrow \mathbb{R}^m$ :

$$f'(x) = \frac{df}{dx^T} = \begin{pmatrix} \frac{\partial f_1}{\partial x_1} & \frac{\partial f_1}{\partial x_2} & \cdots & \frac{\partial f_1}{\partial x_n} \\ \frac{\partial f_2}{\partial x_1} & \frac{\partial f_2}{\partial x_2} & \cdots & \frac{\partial f_2}{\partial x_n} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial f_m}{\partial x_1} & \frac{\partial f_m}{\partial x_2} & \cdots & \frac{\partial f_m}{\partial x_n} \end{pmatrix} = J \in \mathbb{R}^{m \times n} \in \mathbb{R}^{n \times m}$$

## Summary

$$f(x) : X \rightarrow Y; \quad \frac{\partial f(x)}{\partial x} \in G$$

$$\det X$$

X	Y	G	Name
$\mathbb{R}$	$\mathbb{R}$	$\mathbb{R}$	$f'(x)$ (derivative)
$\mathbb{R}^n$	$\mathbb{R}$	$\mathbb{R}^n$	$\frac{\partial f}{\partial x_i}$ (gradient)
$\mathbb{R}^n$	$\mathbb{R}^m$	$\mathbb{R}^{m \times n}$	$\frac{\partial f_i}{\partial x_j}$ (jacobian)
$\mathbb{R}^{m \times n}$	$\mathbb{R}$	$\mathbb{R}^{m \times n}$	$\frac{\partial f}{\partial x_{ij}}$

named gradient of  $f(x)$ . This vector indicates the direction of steepest ascent. Thus, vector  $-\nabla f(x)$  means the direction of the steepest descent of the function in the point. Moreover, the gradient vector is always orthogonal to the contour line in the point.

## General concept

### Naive approach

The basic idea of naive approach is to reduce matrix\vector derivatives to the well-known scalar derivatives.

### Matrix notation of a function

$$f(x) = c^T x$$

### Scalar notation of a function

$$f(x) = \sum_{i=1}^n c_i x_i$$

### Matrix notation of a gradient

$$\nabla f(x) = c$$

$$\frac{\partial f(x)}{\partial x_k} = c_k$$

$$\frac{\partial f(x)}{\partial x_k} = \frac{\partial (\sum_{i=1}^n c_i x_i)}{\partial x_k}$$

$$\begin{aligned} \sum_{i=1}^n \frac{\partial (c_i x_i)}{\partial x_k} &= \\ &= \sum_{i=1}^n c_i \underbrace{\frac{\partial x_i}{\partial x_k}}_{1, i=k} = 0, i \neq k \end{aligned}$$

One of the most important practical trick here is to separate indices of sum ( $i$ ) and partial derivatives ( $k$ ). Ignoring this simple rule tends to produce mistakes.

## Guru approach

The guru approach implies formulating a set of simple rules, which allows you to calculate derivatives just like in a scalar case. It might be convenient to use the differential notation here.

## Differentials

After obtaining the differential notation of  $df$  we can retrieve the gradient using following formula:

*example*  
 $\downarrow$   
 $df(x) = \langle \nabla f(x), dx \rangle$

*Пусть*  
 $f(x) = c^T x = \langle c, x \rangle$   
 $df = d(\langle c, x \rangle) = \langle c, dx \rangle$

$$\nabla f = c$$

Then, if we have differential of the above form and we need to calculate the second derivative of the matrix\vector function, we treat "old"  $dx$  as the constant  $dx_1$ , then calculate  $d(df)$

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

## Properties

Let  $A$  and  $B$  be the **constant** matrices, while  $X$  and  $Y$  are the variables (or matrix functions).

- $dA = 0$
- ✓  $d(\alpha X) = \alpha(dX)$
- ✓  $d(AXB) = A(dX)B$
- ✓  $d(X + Y) = dX + dY$
- ✓  $d(X^\top) = (dX)^\top$
- $d(XY) = (dX)Y + X(dY)$
- $d\langle X, Y \rangle = \langle dX, Y \rangle + \langle X, dY \rangle$
- $d\left(\frac{X}{\phi}\right) = \frac{\phi dX - (d\phi)X}{\phi^2}$
- $d(\det X) = \det X \langle X^{-\top}, dX \rangle$
- $d\text{tr } X = \langle I, dX \rangle$

$f \in \mathbb{R}$     $df \in \mathbb{R}$   
 $x \in \mathbb{R}^n$     $dx \in \mathbb{R}^n$   
 $df = \lim_{\Delta x \rightarrow 0} (f(x+\Delta x) - f(x))$

$$\varphi : \mathbb{R}^n \rightarrow \mathbb{R}$$

- $df(g(x)) = \frac{df}{dg} \cdot dg(x)$
- $H = (J(\nabla f))^T$
- $d(X^{-1}) = -X^{-1}(dX)X^{-1}$

## References

- [Good introduction](#)
- [The Matrix Cookbook](#)
- [MSU seminars \(Rus.\)](#)
- [Online tool](#) for analytic expression of a derivative.
- [Determinant derivative](#)

Example 1  $f: \mathbb{R}^n \rightarrow \mathbb{R}$

Guru Approach

$$\text{Find } \nabla f(x), \text{ if } f(x) = \frac{1}{2}x^T(Ax) + b^T x + c.$$

$$\begin{aligned} \textcircled{1} \quad df &= d\left(\frac{1}{2}\langle x, Ax \rangle + \langle b, x \rangle + c\right) = \frac{1}{2} d(\langle x, Ax \rangle) + d(\langle b, x \rangle) + dc = \\ &= \frac{1}{2} (\langle dx, Ax \rangle + \langle x, d(Ax) \rangle) + \langle b, dx \rangle = \frac{1}{2} \langle Ax, dx \rangle + \frac{1}{2} \langle x, Adx \rangle + \langle b, dx \rangle = \\ &= \frac{1}{2} \langle Ax, dx \rangle + \frac{1}{2} \langle A^T x, dx \rangle + \langle b, dx \rangle = \langle \frac{1}{2}Ax + \frac{1}{2}A^T x + b, dx \rangle \\ \Rightarrow \boxed{\nabla f = \frac{1}{2}(A+A^T)x + b} \quad &\text{если } A \in \mathbb{S}_{++} \quad (\text{если } A = A^T) \quad \nabla f = Ax + b. \end{aligned}$$

$$f(x) = \langle c, x \rangle$$

$$df = \langle c, dx \rangle$$

$$d(\langle c, x \rangle) = \langle dc, x \rangle + \langle c, dx \rangle = \langle c, dx \rangle$$

Example 2

$$dc = \begin{pmatrix} 0 \\ \vdots \\ 0 \end{pmatrix}$$

Find  $\nabla f(x), f''(x)$ , if  $f(x) = -e^{-x^T x}$ .

$$\begin{aligned} 1) \quad f(x) &= -e^{-\langle x, x \rangle}. \quad df = d(-e^{-\langle x, x \rangle}) = -d(e^{-\langle x, x \rangle}) = \\ &= -e^{-\langle x, x \rangle} \cdot d(-\langle x, x \rangle) = -e^{-\langle x, x \rangle} \cdot (-1) d(\langle x, x \rangle) = e^{-\langle x, x \rangle} \cdot (\langle dx, x \rangle + \langle x, dx \rangle) = \\ &= \underbrace{e^{-\langle x, x \rangle}}_{\text{const}} \cdot 2\langle x, dx \rangle = \langle 2e^{-\langle x, x \rangle} \cdot x, dx \rangle \Rightarrow \boxed{\nabla f = 2e^{-\langle x, x \rangle} \cdot x} \end{aligned}$$

2) Как получаем  $\nabla f(x)$ ? (если  $dx = dx_1 = \text{const}$ )

$$2.1) g(x) = df(x)$$

$$2.2) \text{ Где } dg = d(\langle 2e^{-\langle x, x \rangle} \cdot x, dx_1 \rangle) = \langle d(2e^{-\langle x, x \rangle} \cdot x), dx_1 \rangle =$$

$$\begin{aligned} &= \langle 2d(e^{-\langle x, x \rangle} \cdot x), dx_1 \rangle = \langle 2\left(d(e^{-\langle x, x \rangle}) \cdot x + e^{-\langle x, x \rangle} \cdot dx\right), dx_1 \rangle = \\ &= \langle 2\left(e^{-\langle x, x \rangle} \cdot d(-\langle x, x \rangle) \cdot x + e^{-\langle x, x \rangle} \cdot dx\right), dx_1 \rangle \quad \text{⇒} \end{aligned}$$

$$X^T y = \langle X, y \rangle$$

~~$X^T y$~~   $X^T y$

ベクトル ベクトル  
nx1 nx1  
nxn nx1

$$\begin{pmatrix} x_1 \\ x_2 \end{pmatrix} \cdot \begin{pmatrix} y_1 \\ y_2 \end{pmatrix} = x_1 y_1 + x_2 y_2$$

$$\begin{aligned}
 & \stackrel{=} \langle 2 \cdot \bar{\ell}^{-\langle x, x \rangle} \left( -2 \underbrace{\langle x, dx \rangle}_{\text{nx1}} \cdot x + dx \right), dx \rangle = \\
 &= \langle 2 \cdot \bar{\ell}^{-\langle x, x \rangle} \left( -2 \underbrace{x^T dx}_{\text{nxn}} \cdot x + dx \right), dx_1 \rangle = \\
 &= \langle 2 \cdot \bar{\ell}^{-\langle x, x \rangle} \left( -2 \underbrace{x \cdot x^T dx}_{\text{nx1} \text{ nxn}} + dx \right), dx_1 \rangle = \\
 &= \langle 2 \cdot \bar{\ell}^{-\langle x, x \rangle} \left( -2 \underbrace{\mathbf{x} \mathbf{x}^T}_{\text{nxn}} \cdot dx + \mathbf{I} \cdot dx \right), dx_1 \rangle = \\
 &= \langle 2 \cdot \bar{\ell}^{-\langle x, x \rangle} \left( -2 \mathbf{x} \mathbf{x}^T + \mathbf{I} \right) dx, dx_1 \rangle = \\
 &= \langle dx, \underbrace{\left( 2 \cdot \bar{\ell}^{-\langle x, x \rangle} \left( \mathbf{I} - 2 \mathbf{x} \mathbf{x}^T \right) \right)^T}_{\text{dx}} dx_1 \rangle = \\
 &= \langle 2 \bar{\ell}^{-\langle x, x \rangle} \cdot \underbrace{\left( \mathbf{I} - 2 \mathbf{x} \mathbf{x}^T \right)}_{\nabla^2 f}, dx_1, dx \rangle
 \end{aligned}$$

$$dx = I \cdot dx$$

$$df = \langle H_f \cdot dx_1, dx \rangle$$

$$\begin{aligned}
 & \left( \mathbf{I} - 2 \mathbf{x} \mathbf{x}^T \right)^T = \\
 &= \mathbf{I}^T - 2 \left( \mathbf{x} \mathbf{x}^T \right)^T = \\
 &= \mathbf{I} - 2 \mathbf{x}^T \cdot \mathbf{x}^T = \\
 &= \mathbf{I} - 2 \mathbf{x} \mathbf{x}^T
 \end{aligned}$$

$$\nabla^2 f = \nabla^2 f = f''(x) = H_f = 2 \bar{\ell}^{-\langle x, x \rangle} (\mathbf{I} - 2 \mathbf{x} \mathbf{x}^T)$$

$$df = \langle \nabla f, dx \rangle = \boxed{g(x) = \langle \nabla f(x), dx_1 \rangle}$$

$dx_1 = \text{const}$

$$dg \quad \langle d(\nabla f(x)), dx_1 \rangle$$

$$\langle \underbrace{(\cdot \cdot \cdot)}_{\text{meapura (Геометрия)}} \cdot dx_1, dx \rangle$$

$$\begin{pmatrix} x_1 \\ \vdots \\ x_n \end{pmatrix} f(x) = \frac{1}{2} x^T A x + b^T x + c$$

Naive approach

$$\textcircled{1} \quad f(x) = \frac{1}{2} \sum_{i=1}^n x_i (Ax)_i + \sum_{i=1}^n b_i x_i + c =$$

?

$$= \frac{1}{2} \sum_{i=1}^n \left( x_i \left( \sum_{j=1}^n a_{ij} x_j \right) \right) + \sum_{i=1}^n b_i x_i + c$$

$$\textcircled{2} \quad \frac{\partial f}{\partial x_m} = \frac{\partial}{\partial x_m} \left( \frac{1}{2} \sum_{i=1}^n \left( x_i \left( \sum_{j=1}^n a_{ij} x_j \right) \right) + \sum_{i=1}^n b_i x_i + c \right) =$$

$$= \frac{1}{2} \frac{\partial}{\partial x_m} \left( \sum_{i=1}^n \left( x_i \left( \sum_{j=1}^n a_{ij} x_j \right) \right) \right) + \frac{\partial}{\partial x_m} \left( \sum_{i=1}^n b_i x_i \right) + \frac{\partial}{\partial x_m} (c) =$$

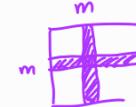
$$= \frac{1}{2} \sum_{i,j=1}^n \frac{\partial}{\partial x_m} \left( x_i a_{ij} x_j \right) + \underbrace{\sum_{i=1}^n \frac{\partial (b_i x_i)}{\partial x_m}}_{=}$$

$\frac{\partial x_i}{\partial x_m} = \begin{cases} 1 & , i=m \\ 0 & , i \neq m \end{cases}$

$x_1, x_2, x_3$   
 $\frac{\partial x_2}{\partial x_1} = 0$   
 $\frac{\partial x_2}{\partial x_2} = 1$

$$\frac{\partial}{\partial x_m} (a_{ij} x_i x_j) =$$

$$b_m \cdot 1 + \cancel{b_1 \cancel{b_2 \cdots \cancel{b_m}}}$$



$$= \begin{cases} 0 & i \neq m, j \neq m \\ a_{mj} x_i & i=m, j \neq m \\ a_{im} x_i & i \neq m, j=m \\ 2a_{mm} x_m & i=m, j=m \end{cases}$$

$$\frac{\partial (a_{mj} x_m x_i)}{\partial x_m}$$

$$\frac{\partial f}{\partial x_m} = \frac{1}{2} \sum_{i=1}^n (a_{mi} + a_{im}) x_i + b_m$$

$$\boxed{\nabla f = \frac{1}{2} (A + A^T) x + b}$$

$\langle Ax-b, d(Ax-b) \rangle =$   
 $= \langle d(Ax-b), Ax-b \rangle$

Example 3

$\|X\|_2 = \sqrt{\sum_{i=1}^n x_i^2}$      $\left\| \begin{pmatrix} x_1 \\ x_2 \end{pmatrix} \right\|_2 = \sqrt{x_1^2 + x_2^2}$   
 ↗  
 ↘  
 ↙

Find the gradient  $\nabla f(x)$  and hessian  $f''(x)$ , if  $f(x) = \frac{1}{2} \|Ax-b\|_2^2$

Бюджет (ЕБНУГОБА)

1)  $f(x) = \frac{1}{2} (Ax-b)(Ax-b)^T = \frac{1}{2} \langle Ax-b, Ax-b \rangle$   
 2)  $df = d\left(\frac{1}{2} \langle Ax-b, Ax-b \rangle\right) = \frac{1}{2} (\langle d(Ax-b), Ax-b \rangle + \langle Ax-b, d(Ax-b) \rangle)$   
 $= \frac{1}{2} 2 \langle d(Ax-b), Ax-b \rangle = \langle d(Ax) - d(b), Ax-b \rangle = \langle Adx, Ax-b \rangle =$   
 $= \langle Ax-b, Adx \rangle = \langle A^T(Ax-b), dx \rangle \quad \boxed{\nabla f = A^T(Ax-b)}$   
 ↗  
 ↙

3)  $g(x) = df(x) = \langle A^T(Ax-b), dx \rangle$   
 $dg = \langle d(A^TAx - A^Tb), dx \rangle = \langle A^T A dx, dx \rangle = \langle dx, (A^T A) dx \rangle =$   
 $K(x+y) = Kx + Ky$   
 $A^T(Ax-b) = A^T Ax - A^T b$

огнко и то же

Example 4

$\nabla f = \begin{matrix} n \times n \\ n \times m \\ m \times n \end{matrix} \quad \boxed{\nabla f = A^T A} \quad x \in \mathbb{R}^n$   
 ноль химичного оп.

Calculate:  $\frac{\partial}{\partial X} \sum \text{eig}(X)$ ,  $\frac{\partial}{\partial X} \prod \text{eig}(X)$ ,  $\frac{\partial}{\partial X} \text{tr}(X)$ ,  $\frac{\partial}{\partial X} \det(X)$   
 $\forall x \in \mathbb{R}^n$

$f(x) = \boxed{\|x\|^2} = \langle x, x \rangle$   
 $\sum_{i=1}^n x_i^2 = \sum_{i=1}^n x_i \cdot x_i$   
 $\frac{1}{2} \left\| \begin{matrix} t \\ t \end{matrix} \right\|_2^2 = \frac{1}{2} t^T t =$   
 $= \frac{1}{2} \langle t, t \rangle$   
 $\|Ax\|_2^2 \geq 0$

$x^T (A^T A) x = c \geq 0$   
 $x^T A^T A x$   
 $(Ax)^T A x$

$\text{tr}(X) = \sum_{i=1}^n x_{ii}$   
 $\frac{\partial \text{tr} X}{\partial X_{kp}} = \frac{\partial \left( \sum_{i=1}^n x_{ii} \right)}{\partial X_{kp}} = \sum_{i=1}^n \frac{\partial x_{ii}}{\partial X_{kp}}$   
 $\begin{cases} 1, & k=p=i \\ 0, & \text{where} \end{cases}$

Example 5

$\Rightarrow \frac{\partial \text{tr} X}{\partial X} = I \rightarrow \begin{pmatrix} 1 & & \\ & \ddots & \\ & & 1 \end{pmatrix}$

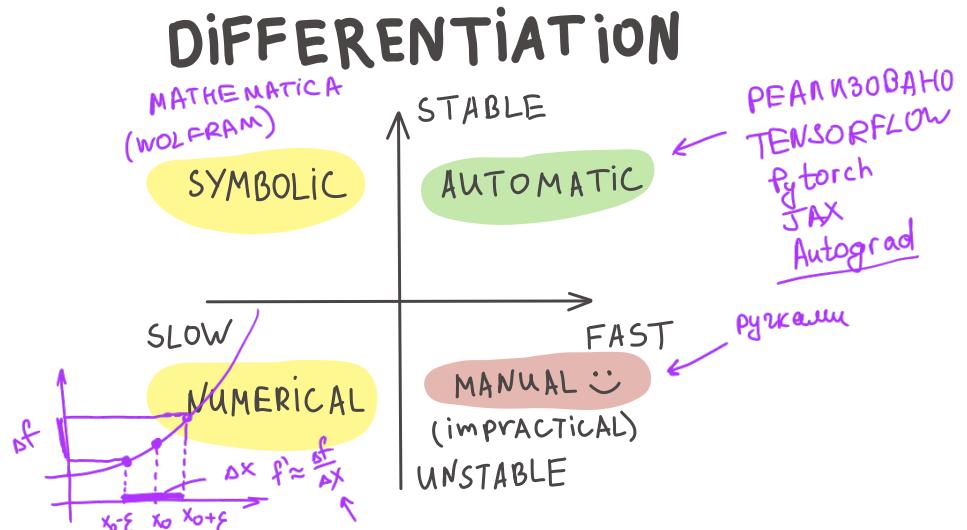
$\boxed{\text{tr}(X) = \sum_{i=1}^n \lambda_i(X)}$   
 $\frac{\partial \left( \sum_{i=1}^n \lambda_i(X) \right)}{\partial X} = I$

$\det X = \prod_{i=1}^n \lambda_i(X)$   
 $\frac{\partial \left( \prod_{i=1}^n \lambda_i(X) \right)}{\partial X} = \frac{\partial (\det X)}{\partial X} = \det \cdot X^{-1}$

$\mathbb{R}^{n \times n}$  Rational  
 $\text{CGR}$   
 $\begin{matrix} \text{c} \in \mathbb{R}^1 \\ \text{c} \in \mathbb{R}^{n \times 1} \end{matrix}$   
 $\text{df} = d(\langle S, X \rangle - \ln \det X) = d(\langle S, X \rangle) - d(\ln \det X) =$   
 $= \langle S, dX \rangle - \frac{d(\det X)}{\det X} = \rightarrow \text{uz tablizki}$   
 $\langle x, y \rangle = \langle y, x \rangle$   
 $\langle x, x \rangle = 0 \Leftrightarrow x=0$   
 $\langle dx + dy, dz \rangle =$   
 $= 2\langle x, dz \rangle + \langle y, dz \rangle$   
 $\langle x, x \rangle = \|x\|_F^2$   
 $\boxed{\text{Find } \nabla f(X), \text{ if } f(X) = \langle S, X \rangle - \log \det X}$   
 $S \in \mathbb{R}^{n \times n} = \text{const}$   
 $\begin{matrix} \text{metric} & \langle X, Y \rangle = \text{tr}(X^T Y) = \sum_{i=1}^n (XY)_{ii} \\ \text{mn} & \text{mn} \\ \text{n} \times \text{m} & \text{m} \times \text{n} \end{matrix}$   
 $\det X \neq 0$

## Automatic differentiation

### Idea



Automatic differentiation is a scheme that allows you to compute a value of gradient of function with a cost of computing function itself only twice.  $\leftarrow \text{BACKWARD} = 2 \text{ FORWARD}$

### Chain rule

We will illustrate some important matrix calculus facts for specific cases

### Univariate chain rule

Suppose, we have the following functions  $R : \mathbb{R} \rightarrow \mathbb{R}$ ,  $L : \mathbb{R} \rightarrow \mathbb{R}$  and  $W \in \mathbb{R}$ . Then

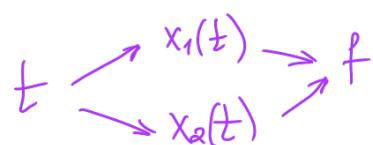
$$\begin{aligned} R(L(W)) & \quad \frac{\partial R}{\partial W} = \frac{\partial R}{\partial L} \frac{\partial L}{\partial W} \\ \frac{\partial R}{\partial W} & \end{aligned}$$

Производная сложной функции

## Multivariate chain rule

$$x_1(t) = \sin t$$

$$x_2(t) = e^{2t}$$



The simplest example:

$$\frac{\partial}{\partial t} f(x_1(t), x_2(t)) = \frac{\partial f}{\partial x_1} \frac{\partial x_1}{\partial t} + \frac{\partial f}{\partial x_2} \frac{\partial x_2}{\partial t}$$

$$f = 3x_1(t) + 5x_2(t) =$$

$$= \sin t + e^{2t}$$

$$\frac{\partial f}{\partial t} = 3 \cdot \cos t + 5 \cdot 2e^{2t}$$

$$J_{ij} = \frac{\partial f_i}{\partial x_j}$$

Now, we'll consider  $f: \mathbb{R}^n \rightarrow \mathbb{R}$ :

$$\frac{\partial}{\partial t} f_i = \sum_{j=1}^n \frac{\partial f_i}{\partial x_j} \frac{\partial x_j}{\partial t} \quad \frac{\partial}{\partial t} f(x_1(t), \dots, x_n(t)) = \frac{\partial f}{\partial x_1} \frac{\partial x_1}{\partial t} + \dots + \frac{\partial f}{\partial x_n} \frac{\partial x_n}{\partial t}$$

But what if we will add another dimension  $f: \mathbb{R}^n \rightarrow \mathbb{R}^m$ , than the  $j$ -th output of  $f$  will be:

$$J \in \mathbb{R}^{m \times n}$$

$$\frac{\partial f_i}{\partial x_j}$$

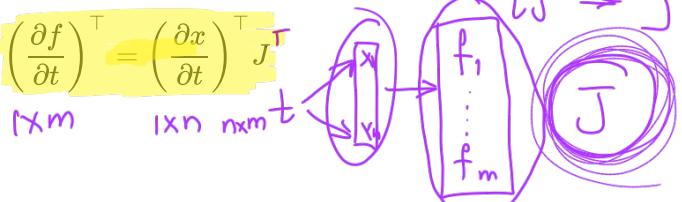
$$\frac{\partial}{\partial t} f_j(x_1(t), \dots, x_n(t)) = \sum_{i=1}^n \frac{\partial f_j}{\partial x_i} \frac{\partial x_i}{\partial t} = \sum_{i=1}^n J_{ji} \frac{\partial x_i}{\partial t},$$

$$j_i \rightarrow J$$

where matrix  $J \in \mathbb{R}^{m \times n}$  is the jacobian of the  $f$ . Hence, we could write it in a vector way:

$$\frac{\partial f}{\partial t} = J \frac{\partial x}{\partial t}$$

$$\left( \frac{\partial f}{\partial t} \right)^T = \left( \frac{\partial x}{\partial t} \right)^T J^T$$



## Backpropagation

The whole idea came from the applying chain rule to the computation graph of primitive operations

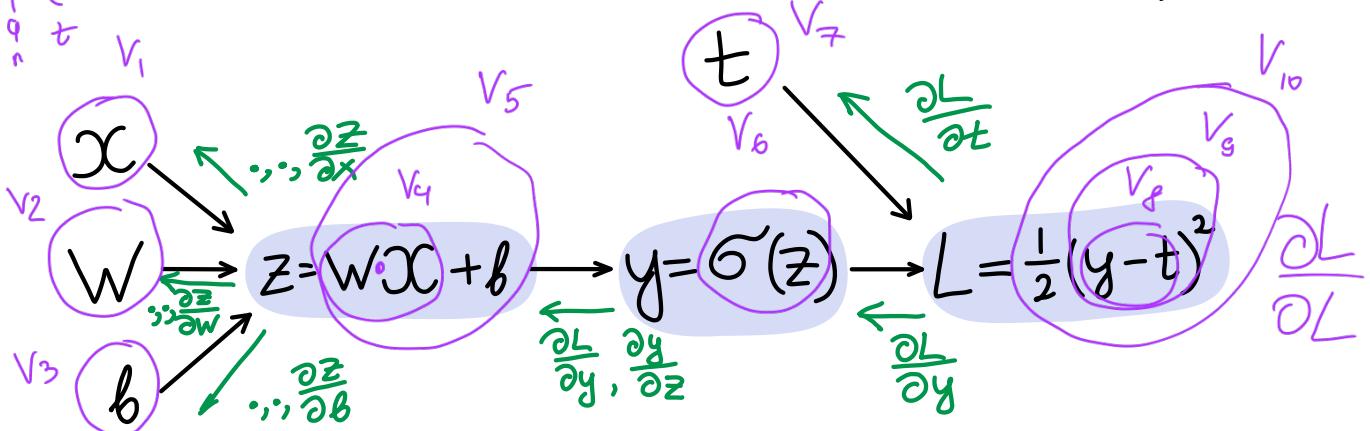
$$\left( \frac{\partial x}{\partial t} \right)^T = \left( \frac{\partial f}{\partial t} \right)^T \cdot J$$

$$\frac{\partial x}{\partial t} = J^T \frac{\partial f}{\partial t}$$

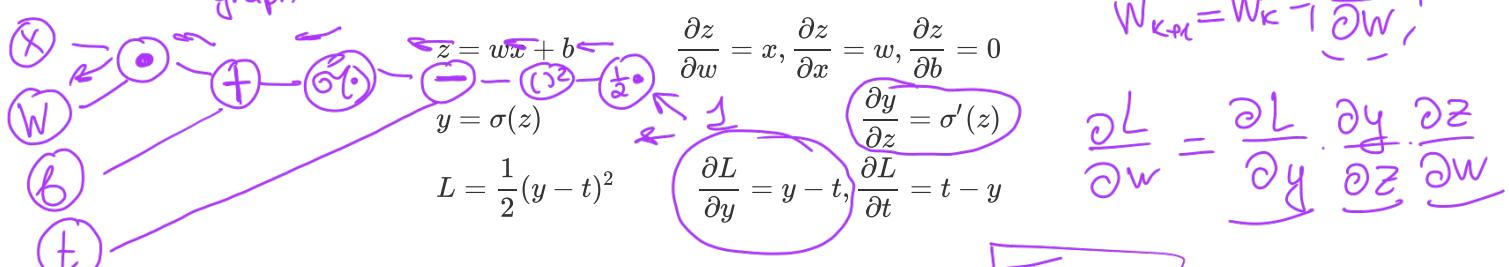
$$J_{m \times n}$$

## FORWARD PASS (COMPUTE LOSS)

→



## BACKWARD PASS (compute derivatives)



$$\frac{\partial z}{\partial w} = x, \frac{\partial z}{\partial x} = w, \frac{\partial z}{\partial b} = 0$$

$$\frac{\partial y}{\partial z} = \sigma'(z)$$

$$\frac{\partial L}{\partial y} = y - t, \frac{\partial L}{\partial t} = t - y$$

$$W_{k \times m} = W_k - \frac{\partial L}{\partial W}$$

$$\frac{\partial L}{\partial w} = \frac{\partial L}{\partial y} \cdot \frac{\partial y}{\partial z} \cdot \frac{\partial z}{\partial w}$$

$$J \cdot x$$

сразу  
запускать коечку

All frameworks for automatic differentiation construct (implicitly or explicitly) computation graph. In deep learning we typically want to compute the derivatives of the loss function  $L$  w.r.t. each intermediate parameters in order to tune them via gradient descent. For this purpose it is convenient to use the following notation:

$$\overline{v_i} = \frac{\partial L}{\partial v_i}$$

Let  $v_1, \dots, v_N$  be a topological ordering of the computation graph (i.e. parents come before children).  $v_N$  denotes the variable we're trying to compute derivatives of (e.g. loss).

### Forward pass:

- For  $i = 1, \dots, N$ :
  - Compute  $\overline{v_i}$  as a function of its parents.

### Backward pass:

- $\overline{v_N} = 1$
- $\frac{\partial L}{\partial L} = 1$
- For  $i = N - 1, \dots, 1$ :

◦ Compute derivatives  $\overline{v_i} = \sum_{j \in \text{Children}(v_i)} \overline{v_j} \frac{\partial v_j}{\partial v_i}$

$$\frac{\partial L}{\partial v_i} = J \cdot \frac{\partial L}{\partial v_{\text{children}}}$$

Note, that  $\overline{v_j}$  term is coming from the children of  $\overline{v_i}$ , while  $\frac{\partial v_j}{\partial v_i}$  is already precomputed effectively.

## Jacobian vector product

The reason why it works so fast in practice is that the Jacobian of the operations are already developed in effective manner in automatic differentiation frameworks. Typically, we even do not construct or store the full Jacobian, doing matvec directly instead.

### Example: element-wise exponent

$$y = \exp(z) \quad J = \text{diag}(\exp(z)) \quad \bar{z} = \bar{y}J$$

See the examples of Vector-Jacobian Products from autodidact library:

```
defvjp(anp.add,
       lambda g, ans, x, y : unbroadcast(x, g),
       lambda g, ans, x, y : unbroadcast(y, g))
defvjp(anp.multiply,
       lambda g, ans, x, y : unbroadcast(x, y * g),
       lambda g, ans, x, y : unbroadcast(y, x * g))
defvjp(anp.subtract,
       lambda g, ans, x, y : unbroadcast(x, g),
       lambda g, ans, x, y : unbroadcast(y, -g))
defvjp(anp.divide,
       lambda g, ans, x, y : unbroadcast(x, g / y),
       lambda g, ans, x, y : unbroadcast(y, -g * x / y**2))
defvjp(anp.true_divide,
       lambda g, ans, x, y : unbroadcast(x, g / y),
       lambda g, ans, x, y : unbroadcast(y, -g * x / y**2))
```

$$\begin{aligned} z_1 &\rightarrow \exp(z_1) O(n) \\ &\vdots \\ z_n &\rightarrow \exp(z_n) O(n) \end{aligned}$$

$$J = \begin{pmatrix} \exp(z_1) & & \\ & \exp(z_2) & \\ & & \ddots & \exp(z_n) \end{pmatrix}$$

$J_x = \begin{pmatrix} x_1 & \dots & x_n \end{pmatrix} \cdot y$

$J_y = \begin{pmatrix} y_1 & \dots & y_n \end{pmatrix} \cdot \begin{pmatrix} x_1 & \dots & x_n \end{pmatrix}$

$$\begin{aligned} y &= Ax \\ &\text{m} \times n \quad n \times 1 \\ J &= A \\ &\text{m} \times n \end{aligned}$$

$$J \cdot g = O(n^2)$$

## Hessian vector product

Interesting, that the similar idea could be used to compute Hessian-vector products, which is essential for second order optimization or conjugate gradient methods. For a scalar-valued function  $f : \mathbb{R}^n \rightarrow \mathbb{R}$  with continuous second derivatives (so that the Hessian matrix is symmetric), the Hessian at a point  $x \in \mathbb{R}^n$  is written as  $\partial^2 f(x)$ . A Hessian-vector product function is then able to evaluate

$$v \mapsto \partial^2 f(x) \cdot v$$

for any vector  $v \in \mathbb{R}^n$ .

The trick is not to instantiate the full Hessian matrix: if  $n$  is large, perhaps in the millions or billions in the context of neural networks, then that might be impossible to store. Luckily, `grad` (in the jax/autograd/pytorch/tensorflow) already gives us a way to write an efficient Hessian-vector product function. We just have to use the identity

$$\partial^2 f(x)v = \partial[x \mapsto \partial f(x) \cdot v] = \partial g(x),$$

where  $g(x) = \partial f(x) \cdot v$  is a new scalar-valued function that dots the gradient of  $f$  at  $x$  with the vector  $v$ . Notice that we're only ever differentiating scalar-valued functions of vector-valued arguments, which is exactly where we know `grad` is efficient.

```
import jax.numpy as jnp

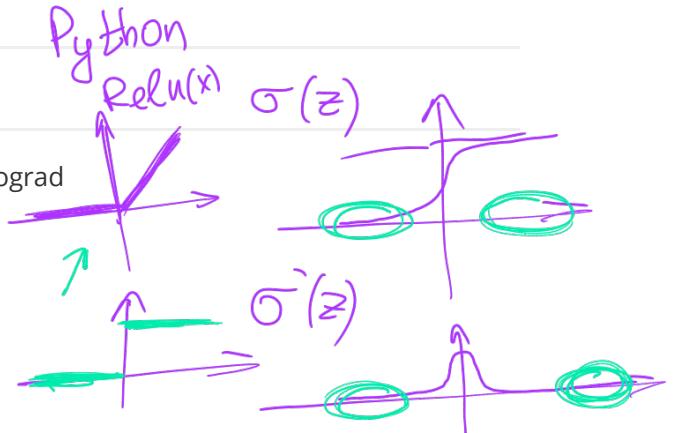
def hvp(f, x, v):
    return grad(lambda x: jnp.vdot(grad(f)(x), v))(x)
```

## Code

### Materials

- [Autodidact](#) - a pedagogical implementation of Autograd
- [CSC321](#) Lecture 6
- [CSC321](#) Lecture 10
- [Why](#) you should understand backpropagation :)
- [JAX autodiff cookbook](#)

200 CPOK koga



## Convex set

In this chapter variety of convex sets and related definitions are described.

НЕ УСПЕЛЫ

## Convex function

## Convex function

The function  $f(x)$ , which is defined on the convex set  $S \subseteq \mathbb{R}^n$ , is called **convex**  $S$ , if:

$$f(\lambda x_1 + (1 - \lambda)x_2) \leq \lambda f(x_1) + (1 - \lambda)f(x_2)$$

for any  $x_1, x_2 \in S$  and  $0 \leq \lambda \leq 1$ .

If above inequality holds as strict inequality  $x_1 \neq x_2$  and  $0 < \lambda < 1$ , then function is called **strictly convex**  $S$ .