

1. Calculating  $\nabla f$
2. Calculating  $\nabla^2 f, H_f$

# Brief recap of 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:

$$C = AB$$

$m \times p$      $m \times n$      $n \times p$

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

$$O(n^3)$$

$$O(n^2)$$

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

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

$$z = Ax$$

$$O(n^2)$$

is given by:

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

Finally, just to remind:

- $\underline{C = AB}$      $\underline{C^\top = B^\top A^\top}$
- $\underline{AB \neq BA}$      $\underline{P \times M} \quad \underline{P \times N} \quad \underline{N \times M}$
- $e^A = \sum_{k=0}^{\infty} \frac{1}{k!} A^k$
- $e^{A+B} \neq e^A e^B$  (but if  $A$  and  $B$  are commuting matrices, which means that  $AB = BA$ ,  $e^{A+B} = e^A e^B$ )
- $\langle x, Ay \rangle = \langle A^\top x, y \rangle$

$$\langle z^\top, p^\top B \rangle = \langle z^\top B^\top, p^\top \rangle$$

## Gradient

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

$$f(x) = C^\top X$$

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

$\in \mathbb{R}^n$

$$f(X) = \det X$$

$$\text{tr } X$$

Hessian

$$\frac{\partial f}{\partial X} = \left( \frac{\partial f}{\partial x_{ij}} \right)_{1 \times n}$$

MATRIX

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

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

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  $f(x) : \mathbb{R}^n \rightarrow \mathbb{R}^m$ :

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

## Summary

INPUT	OUTPUT	$f(x) : X \rightarrow Y;$ gradient	$\frac{\partial f(x)}{\partial x} \in G$
$\mathbb{R}$ scal	$\mathbb{R}$ sc	$\mathbb{R}$	$f'(x)$ (derivative)
$\mathbb{R}^n$ vec	$\mathbb{R}$ sc	$\mathbb{R}^n$	$\frac{\partial f}{\partial x_i}$ (gradient)
$\mathbb{R}^n$ vec	$\mathbb{R}^m$ vec	$\mathbb{R}^{m \times n}$	$\frac{\partial f_i}{\partial x_j}$ (jacobian)
$\mathbb{R}^{m \times n}$ mat	$\mathbb{R}$ sc	$\mathbb{R}^{m \times n}$	$\frac{\partial f}{\partial x_{ij}}$ $i, j$ matrix

$\left( \frac{\partial f}{\partial x_{ij}} \right)_{1 \times n}$  ← matrix

$\nabla f$

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

$$df = f(x+dx) - f(x) \rightarrow 0 \quad X \in \mathbb{R}^n \quad X+dx \quad \|dx\|_2 \rightarrow 0$$

$$\lim_{dx \rightarrow 0} \frac{df}{dx} = f'(x)$$

The idea 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.

$$f(x) = C^T X$$

**Differentials**  $df = C^T(X+dx) - C^T X = C^T dx = \langle C, dx \rangle$

$$\Rightarrow \boxed{\nabla f = C}$$

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

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

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).

- $dA = 0$
- $d(\alpha X) = \alpha(dX)$
- $d(AXB) = A(dX)B$
- $d(X + Y) = dX + dY$
- $d(X^\top) = (dX)^\top$
- $d(XY) = (dX)Y + X(dY)$
- $d\langle X, Y \rangle = \langle dX, Y \rangle + \langle X, dY \rangle$
- $d\left(\frac{X}{\phi}\right) = \frac{\phi dX - (d\phi)X}{\phi^2}$
- $d(\det X) = \det X \langle X^{-\top}, dX \rangle$
- $d(\text{tr } X) = \langle I, dX \rangle$
- $df(g(x)) = \frac{df}{dg} \cdot dg(x)$
- $H = (J(\nabla f))^T$
- $d(X^{-1}) = -X^{-1}(dX)X^{-1}$

$$\langle A, B \rangle = \text{tr}(A^\top B) = \text{tr}(B^\top A)$$

$$\langle A, A \rangle = \text{tr}(A^\top A) = \|A\|_F^2$$

$$\|A\|_F = \sqrt{\sum_{ij} a_{ij}^2}$$

← we'll derive it manually

## References

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

$$\nabla f = ? \quad \nabla^2 f = ?$$

Gradient

Hessians

## Examples

### Example 1

$$1. df = \dots \quad \nabla f$$

$$2. df = \langle \dots, dx \rangle$$

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

$$1) f = \frac{1}{2} \langle x, Ax \rangle + \langle b, x \rangle + c$$

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

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

$$= \frac{1}{2} \left( \underbrace{\langle Ax, dx \rangle}_{\text{---}} + \underbrace{\langle x, Adx \rangle}_{\text{---}} \right) + \langle b, dx \rangle = \frac{1}{2} (\langle Ax, dx \rangle + \langle A^T x, dx \rangle) + \langle b, dx \rangle =$$

$$= \langle \frac{1}{2}(A + A^T)x + b, dx \rangle \Rightarrow \boxed{\nabla f = \frac{1}{2}(A + A^T)x + b}$$

### Example 2

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

$$1. f = -e^{-\langle x, x \rangle} \quad \nabla f$$

$$2. df = d(-e^{-\langle x, x \rangle}) = -d(e^{-\langle x, x \rangle}) = -e^{-\langle x, x \rangle} \cdot d(-\langle x, x \rangle) =$$

$$= f \cdot -(\langle dx, x \rangle + \langle x, dx \rangle) = -f \cdot \langle 2x, dx \rangle = \langle -2f x, dx \rangle$$

$$\Rightarrow \boxed{\nabla f = 2 \bar{e}^{-\langle x, x \rangle} \cdot x} \in \mathbb{R}^n$$

$$df = -f \cdot \langle 2x, dx_1 \rangle$$

$$d^2f = d(df) = d(-f \cdot \langle 2x, dx_1 \rangle) = -d(f \cdot \langle 2x, dx_1 \rangle) =$$

$$= -\underline{df} \cdot \langle 2x, dx_1 \rangle - f \cdot d(\langle 2x, dx_1 \rangle) =$$

$$= +f \langle 2x, dx \rangle \cdot \langle 2x, dx_1 \rangle - f \langle 2dx, dx_1 \rangle =$$

$$= 2f [2x^T dx \cdot x^T dx_1 - dx^T \cdot dx_1] = \stackrel{df = \langle H dx_1, dx_2 \rangle}{\langle (x^T dx_1 \cdot x^T), dx \rangle}$$

$$= 2f [2 \underbrace{x^T dx_1 \cdot x^T dx}_{} - \underbrace{dx_1^T \cdot dx}_{}] = \langle (x^T dx_1 \cdot x^T)^T, dx \rangle$$

$$= 2f \cdot \langle 2 \underbrace{x^T dx_1}_{} x - \underbrace{dx_1}_{}, dx \rangle = \langle dx_1, dx \rangle$$

$$= 2f \langle 2x x^T dx_1 - dx_1, dx \rangle =$$

$$= 2f \langle (2x x^T - I_n) dx_1, dx \rangle$$

$n \times 1 \times n$

$$I_n \cdot dx_1 = dx_1$$

$$d^2f = \langle 2f \cdot (2x x^T - I) dx_1, dx_2 \rangle$$

$$\Rightarrow \boxed{\nabla^2 f = 2f \cdot (2x x^T - I)}$$

$$df = \langle \frac{1}{2}(A+A^T)x + b, dx_1 \rangle$$

$$d^2f = d(df) = d\left(\langle \frac{1}{2}(A+A^T)x + b, dx_1 \rangle\right) = \langle \frac{1}{2}d[(A+A^T)x + b], dx_1 \rangle =$$

$$= \langle \frac{1}{2}(A+A^T)dx_3, dx_1 \rangle = \langle dx_2, \frac{1}{2}(A+A^T)dx_1 \rangle =$$

$$= \langle \frac{1}{2}(A+A^T)dx_1, dx_2 \rangle \Rightarrow \boxed{\nabla^2 f = \frac{1}{2}(A+A^T)}$$

$$d(\ln \varphi) = \frac{d\varphi}{\varphi}$$

### Example 3

$$\ln x = \log_e x$$

Find  $\nabla f(X)$ , if  $f(X) = \langle S, X \rangle - \log \det X$ .

$$1. df = d(\langle S, X \rangle - \ln \det X) = d(\langle S, X \rangle) - d(\ln \det X) =$$

$$= \langle S, dX \rangle - \frac{d(\det X)}{\det X} = \langle S, dX \rangle - \frac{\cancel{\det X} \cdot \langle X^{-T}, dx \rangle}{\cancel{\det X}} =$$

$$= \langle S - X^{-T}, dx \rangle \Rightarrow \boxed{\nabla f = S - X^{-T}}$$

### Example 4

Find  $\nabla f(X)$ , if  $f(X) = \ln \langle Ax, x \rangle$ ,  $A \in \mathbb{S}_{++}^n$

$$f = \ln \langle Ax, x \rangle \quad A \in S_{++}^n \quad A = A^T$$

$$df = d(\ln \langle Ax, x \rangle) = \frac{d(\langle Ax, x \rangle)}{\langle Ax, x \rangle} = \frac{\langle Adx, x \rangle + \langle Ax, dx \rangle}{\langle Ax, x \rangle}$$

$$= \frac{\cancel{\langle A^T x, dx \rangle} + \langle Ax, dx \rangle}{\langle Ax, x \rangle} = \left( \frac{2Ax}{\langle Ax, x \rangle}, dx \right)$$

$$\Rightarrow \nabla f = \frac{2Ax}{\langle Ax, x \rangle}$$

$f(X) = \det X$ . What is  $df = ?$

$\det X \neq 0$

$\nabla f = ?$

$$df = f(X+dx) - f(X)$$

$$df = \langle \nabla f, dx \rangle$$

$$\|dx\| \rightarrow 0 \quad I = \begin{pmatrix} 1 & 0 \\ 0 & \dots \end{pmatrix}$$

$$df = \det(X+dx) - \det X = \det X(I + X^{-1}dx) - \det X =$$

$$= \det X \cdot \det(I + X^{-1}dx) - \det X = \det X [\det(I + X^{-1}dx) - 1] \quad \textcircled{1}$$

$$\det(I + X^{-1}dx) = \prod_{i=1}^n \lambda_i(I + X^{-1}dx) = \prod_{i=1}^n (1 + \lambda_i(X^{-1}dx)) =$$

$$= 1 + \sum_{i=1}^n \lambda_i(X^{-1}dx) + O(\|dx\|) \approx 1 + \left( \sum_{i=1}^n \lambda_i(X^{-1}dx) \right) \quad \begin{matrix} \text{When } n \\ \|dx\| \rightarrow 0 \end{matrix}$$

$$\text{tr}(X^{-1}dx)$$

$$\textcircled{1} \quad \det X [1 + \text{tr}(X^{-1}dx) - 1] = \det X \cdot \text{tr}(X^{-1}dx) \quad \textcircled{2}$$

$$\textcircled{3} \quad \langle \det X \cdot X^{-T}, dx \rangle$$

$$\boxed{\nabla f = \det X \cdot X^{-T}}$$

$\text{tr}$

$$\begin{aligned} \langle X^T, dx \rangle &= \\ &= \text{tr}(X^{-1}dx) \end{aligned}$$