

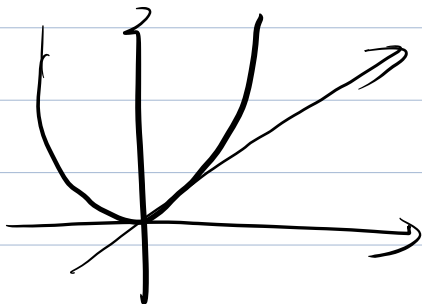
Youtube Link:

<https://www.youtube.com/watch?v=oGZK3yGF-6k>

Gradients

Suppose  $y = x^T x$

$\Rightarrow dy = (2x)^T dx$  [dot product, small change of  $(dx)$  (vector)]



$$\underline{y - y_0 = 2x_0^T (x - x_0)}$$

Linearisation  $\Rightarrow$  plane

Matrix and vector product rule

$$d(AB) = (dA) \cdot B + A \cdot (dB) \text{ (usually do not commute)}$$

Specially,  $d(x^T x) = dx^T \cdot x + x^T \cdot dx$ , since dot product can commute,  $\Rightarrow d(x^T \cdot x) = 2x^T \cdot dx$

slight generally  $d(u^T \cdot v) = du^T \cdot v + u^T \cdot dv$   
(note  $du^T \cdot v = u^T \cdot dv$ )

Gradients Derivation

$$f(x) = (Ax - b)^T (Ax - b)$$

$$df(x) = (Adx)^T (Ax - b) + (Ax - b)^T (Adx)$$

$$(\text{dot product}) = 2(Ax - b)^T A dx$$

$$= \underline{(2A^T(Ax - b))^T \cdot dx}$$

Why this form?

geometric of gradient  $\Rightarrow$  [hyperplane] a dot product of  $\boxed{0^T \cdot dx}$

$$\Rightarrow \text{setting } \frac{df}{dx} = 0 \Rightarrow A^T A x = A^T b \quad (\text{least square})$$

$$\nabla x^T x = 2x$$

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

Trace

① a linear transformation from matrix to real

② Cyclic property  $\text{trace}(ABC) = \text{trace}(BCA) = \text{trace}(CAB)$

③  $\text{tr}(A^T) = \text{tr}(A)$

$$\begin{aligned} \textcircled{4} \quad x^T A y &= \text{tr}(x^T A y) = \text{tr}(y x^T A) \\ &= \text{tr}((x y^T)^T \cdot A) \\ &= \text{sum}((x y^T) \cdot^* A) \end{aligned}$$

If  $f$  is scalar and linear on  $m \times n$  matrices, then there exists an  $m \times n$  matrix  $U$ , such that  $f(A) = \sum_{ij} U_{ij} \cdot A_{ij} = \text{sum}(U \cdot^* A) = \text{tr}(U^T A)$

Gradients of functions from matrices to scalars

$$\text{eg: } f(A) = \text{trace}(A^2)$$

$$\begin{aligned} df &= \text{trace}(A \cdot dA + dA \cdot A) \\ &= \text{tr}((2A^T)^T \cdot dA) \quad (\text{cyclic rule}) \end{aligned}$$

$$\text{gradient} = 2A^T$$

$$\text{in general } df = \text{tr}(U^T dA) \Rightarrow \text{gradient} = U, \text{ no indices needed}$$

Side: Jacobian and Gradients  
vector function  $\mathbb{R}^n$  to  $\mathbb{R}^m$

$$f = \begin{pmatrix} f_1(x_1, \dots, x_n) \\ f_2(x_1, \dots, x_n) \\ \vdots \\ f_m(x_1, \dots, x_n) \end{pmatrix} \quad \begin{array}{l} n \text{ input scalars} \\ m \text{ output scalars} \end{array}$$

$$\text{Jacobian} = \begin{pmatrix} \frac{\partial f_1}{\partial x_1} & \dots & \frac{\partial f_1}{\partial x_n} \\ \vdots & & \vdots \\ \frac{\partial f_m}{\partial x_1} & \dots & \frac{\partial f_m}{\partial x_n} \end{pmatrix}$$

eg:  $f(A) = (AX - b)^T (AX - b)$

$$\begin{aligned} df &= (dAX)^T (AX - b) + (AX - b)^T (dAX) \\ &= \text{tr} [2(AX - b)^T (dA)x] \\ &= \text{tr} [2x(AX - b)^T \cdot dA] \\ &= \text{tr} [(2AX - b)x^T]^T \cdot dA \end{aligned}$$

So the gradient is  $\underline{2(AX - b) \cdot x^T}$   
(elementary sum)

## Jacobian Matrix

If  $x \in \mathbb{R}^n$  and  $f(x)$  is a differentiable function whose values are in  $\mathbb{R}^m$ ,

then the linearization is given by Jacobian matrix

$$df = J \cdot dx \quad \text{where } J_{ij} = \partial f_i / \partial x_j$$

eg:  $f(x) = h(Wx - b)$  ( $h$  is a scalar function)

$$df = h'(Wx - b) \underset{\text{element-wise}}{*} (W dx) \quad (\text{chain rule})$$

$$J = \text{diagonal}(h'(Wx - b)) * W$$

Relationship to columns (any dimension)

$$df = J dx$$

We have a small area around  $x$  is transformed  
to a small area scaled by  $|J|$ .

Since  $J$  depends on  $x$  for nonlinear transformation,  
the scaling is not uniform.

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

$$x \cdot A = I$$

$$d(x \cdot A) = d(I) = 0$$

$$dx \cdot A + x \cdot dA = 0$$

$$\begin{aligned} dx &= -x \cdot dA \cdot A^{-1} \\ &= -A^{-1} \cdot dA \cdot A^{-1}. \end{aligned}$$