# Mathematics for AI
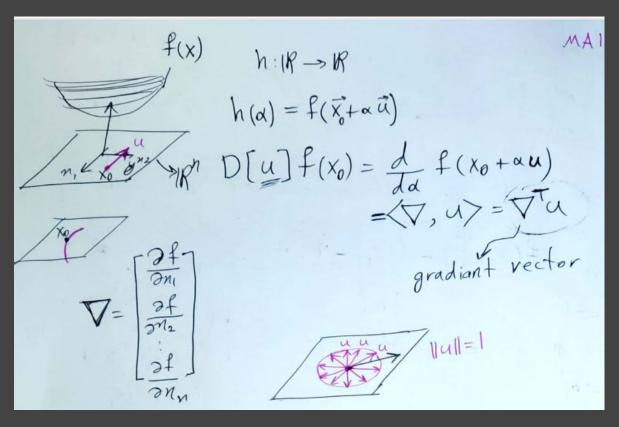
## Lecture 16

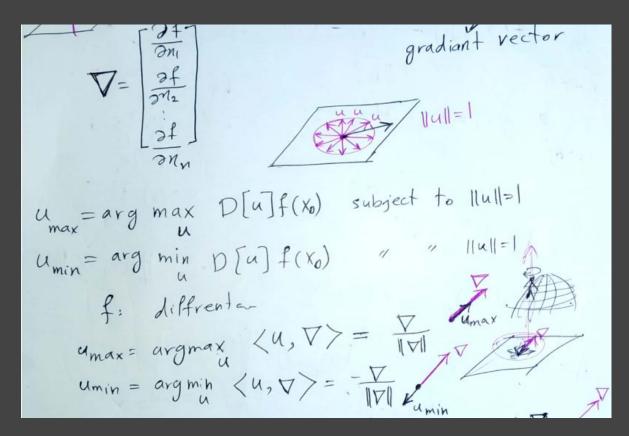# The gradient vector

# The gradient and steepest directions

# Moving perpendicular to gradient



$$D[u]f(x_0) = 0 \Rightarrow \langle u, \nabla \rangle = 0 \Rightarrow u \perp \nabla$$

$u_{min}$

$h = 1400m$

# Example

$$f(x) = f\left(\begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix}\right) = x_1 x_2 + x_3 \sin x_2 + x_1 x_2 \; e^{x_3}$$

$$f: \mathbb{R}^3 \longrightarrow \mathbb{R}$$

$$\nabla f = \begin{bmatrix} \dfrac{\partial f}{\partial x_1} \\[2mm] \dfrac{\partial f}{\partial x_2} \\[2mm] \dfrac{\partial f}{\partial x_3} \end{bmatrix} = \begin{bmatrix} x_2 + x_2 \, e^{x_3} \\[2mm] x_1 + x_3 \cos x_2 + x_1 e^{x_3} \\[2mm] \sin x_2 + x_1 x_2 \, e^{x_3} \end{bmatrix} \in \mathbb{R}^3$$

# Definition of differentiability



$$f(x_0) + \vec{m}^T(x - x_0)$$

$$f: \mathbb{R}^n \longrightarrow \mathbb{R} \quad \text{is differentiable at } x_0 \Rightarrow \exists \, \vec{m} \in \mathbb{R}^n$$

$$\lim_{\vec{h} \to \vec{0}} \frac{f(\vec{x}_0 + \vec{h}) - f(\vec{x}_0) - \vec{m}^T h}{\|\vec{h}\|} = 0$$

LA 16 Ⅱ

# Limits in higher dimensions

$$g: \mathbb{R}^m \longrightarrow \mathbb{R}^n$$

$$\lim_{\vec{x} \to \vec{x_0}} g(\vec{x}) = \vec{y_0} \qquad \forall \varepsilon \ \exists \delta \quad \|x - x_0\| < \delta \implies \|g(x) - \vec{y_0}\| < \varepsilon$$

# Example

$$x \to x_0$$

~~$A \in \mathbb{R}^{m \times n}$~~ $\quad x \in \mathbb{R}^n \quad \mathrm{Diag}(x) = \begin{bmatrix} x_1 & & & \\ & x_2 & \cdots & \\ & & & x_n \end{bmatrix}$

$$\mathrm{Diag}(x) = \mathrm{Diag}\left( \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix} \right) = \begin{bmatrix} x_1 & 0 & 0 & 0 \\ 0 & x_2 & 0 & 0 \\ 0 & 0 & \ddots & 0 \\ 0 & 0 & 0 & x_n \end{bmatrix}$$

~~$f(x)$~~ $\quad f: \mathbb{R}^n \to \mathbb{R} \quad f(x) = u^T \mathrm{Diag}(x) V \quad$ for

constant vectors $\quad u, v \in \mathbb{R}^n$.

$$u^T \mathrm{Diag}(x) V = \begin{bmatrix} u_1 & u_2 & \cdots & u_n \end{bmatrix} \begin{bmatrix} x_1 & & \emptyset \\ & x_2 & \\ \emptyset & & x_n \end{bmatrix} \begin{bmatrix} v_1 \\ v_2 \\ \vdots \\ v_n \end{bmatrix} = \sum_{i=1}^{n} x_i \cdot u_i \cdot v_i$$

$$\frac{\partial f}{\partial x_k} = \frac{\partial}{\partial x_k} \sum_{i=1}^{n} x_i u_i v_i = u_k v_k$$

$$\nabla = \begin{bmatrix} \frac{\partial f}{\partial x_1} \\ \frac{\partial f}{\partial x_2} \\ \vdots \\ \frac{\partial f}{\partial x_n} \end{bmatrix} = \begin{bmatrix} u_1 v_1 \\ u_2 v_2 \\ \vdots \\ u_n v_n \end{bmatrix} = \mathrm{Diag}(u) v$$

$$= \mathrm{Diag}(v) u$$

$$= u \odot v$$

# Hadamard Product

Hadamard product (element-wise product)

$$u \odot v = \begin{bmatrix} u_1 \\ u_2 \\ \vdots \\ u_n \end{bmatrix} \odot \begin{bmatrix} v_1 \\ v_2 \\ \vdots \\ v_n \end{bmatrix} = \begin{bmatrix} u_1 v_1 \\ u_2 v_2 \\ \vdots \\ u_n v_n \end{bmatrix} = \underbrace{Diag(u)}_{n \times n} v = Diag(v) u$$

# Example

$f: \mathbb{R}^n \to \mathbb{R} \qquad f(x) = x^T A x$ for a constant matrix $A \in \mathbb{R}^{n \times n}$.

$$f(x) = x^T A x = \underbrace{[x_1 \; x_2 \cdots x_n]}_{1 \times n} \underbrace{\begin{bmatrix} a_{11} & a_{12} \cdots & a_{12} \\ a_{21} & a_{22} - & a_{2n} \\ \vdots & & \vdots \\ a_{n1} & a_{n2} & a_{nn} \end{bmatrix}}_{n \times n} \overbrace{\begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix}}^{A} \in \mathbb{R}$$

$$= \sum_{i=1}^{n} \sum_{j=1}^{n} x_i x_j a_{ij}$$

$$\frac{\partial f}{\partial x_k} \left( \sum_{i=1}^{n} \sum_{j=1}^{n} a_{ij} x_i x_j \right)$$

$$\frac{\partial}{\partial x_k} \left( a_{kk} \underbrace{x_k x_k}_{x_k^2} + \sum_{\substack{j=1 \\ j \neq k}}^{n} a_{kj} x_k x_j + \sum_{\substack{i=1 \\ i \neq k}}^{n} a_{ik} x_i x_k + \underbrace{\sum_{\substack{i=1 \\ i \neq k}}^{n} \sum_{\substack{j=1 \\ j \neq k}}^{n} a_{ij} x_i x_j}_{0} \right)$$

$$= 2 a_{kk} x_k + \sum_{\substack{j=1 \\ j \neq k}}^{n} a_{kj} x_j + \sum_{\substack{i=1 \\ i \neq k}}^{n} a_{ik} x_i$$

$$= 2 a_{kk} x_k + \sum_{\substack{j=1 \\ j \neq k}}^{n} a_{kj} x_j + \sum_{\substack{i=1 \\ i \neq k}}^{n} a_{ik} x_i$$

$$= \sum_{j=1}^{n} a_{kj} x_j + \sum_{i=1}^{n} a_{ik} x_i = [a_{k1} \; a_{k2} \cdots a_{kn}] \underbrace{\begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix}}_{x} + [a_{1k} \; a_{2k} \cdots a_{nk}] \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix}$$

$$\frac{\partial f}{\partial x_k} = a[k,:]^T x + a[:,k]^T x$$

$a[k,:]^T$ → k-th row of A
$a[:,k]^T$ → k-th column of A

$$\nabla = \begin{bmatrix} \frac{\partial f}{\partial x_1} \\ \frac{\partial f}{\partial x_2} \\ \vdots \\ \frac{\partial f}{\partial x_n} \end{bmatrix} = \begin{bmatrix} a_{1:}^T x + a_{:1}^T x \\ a_{2:}^T x + a_{:2}^T x \\ \vdots \\ a_{n:}^T x + a_{:n}^T x \end{bmatrix} = (A + A^T) x$$

# Example

$$f : \mathbb{R}^n \longrightarrow \mathbb{R} \qquad x \longmapsto \underline{x^T A x}$$

$$\nabla f = \underbrace{(A + A^T)}_{n \times n} \underbrace{x}_{n \times 1} \in \mathbb{R}^n$$

$$A \text{ symmetric} \Rightarrow \nabla f = 2Ax$$

# Example: Least Squares

least squares problem $\quad Ax = b, \; A \in \mathbb{R}^{m \times n}$

$\quad m \geqslant n$

$$\left[\begin{array}{c} A \end{array}\right]\left[x\right] = \left[b\right]$$

$A$ has full column rank

$\text{rank}(A) = n$

least squares problem $\qquad$ *geometric method*

$x^* = \underset{x}{\text{argmin}} \; \|Ax - b\|^2 \implies x = (A^TA)^{-1}A^Tb$

$f(x) = \|Ax - b\|^2 \qquad \nabla f = 0 \implies x = \checkmark$

$f(x) = \left\| \begin{bmatrix} a_{1:}^T x - b_1 \\ a_{2:}^T x - b_2 \\ a_{n:}^T x - b_n \end{bmatrix} \right\|^2 = \sum_{i=1}^{n} \left( a_{i:}^T x - b_i \right)^2$

$= \sum_{i=1}^{n} \left( a_{i1}x_1 + a_{i2}x_2 + \cdots + a_{in}x_n - b_i \right)^2$

$\qquad\qquad\qquad\qquad\qquad\qquad |^2$

# Example: Least Squares

$$\frac{\partial f}{\partial x_k} = \frac{\partial}{\partial x_k} \sum_{i=1}^{n} \left( a_{i1} x_1 + a_{i2} x_2 + \cdots + a_{in} x_n - b_i \right)^2$$

$$= \sum_{i=1}^{n} 2 a_{ik} \left( a_{i1} x_1 + a_{i2} x_2 + \cdots + a_{in} x_n - b_i \right)$$

$$= 2 \sum_{i=1}^{n} a_{ik} \left( a_{i:}^{T} x - b_i \right)$$

$$\frac{\partial f}{\partial x_k} = 2 \begin{bmatrix} a_{1k} & a_{2k} & \cdots & a_{nk} \end{bmatrix} \begin{bmatrix} a_1^T x - b_1 \\ a_2^T x - b_2 \\ \vdots \\ a_n^T x - b_n \end{bmatrix} = 2 a_{:k}^T (Ax - b)$$

$$\nabla = \begin{bmatrix} \frac{\partial f}{\partial x_1} \\ \frac{\partial f}{\partial x_2} \\ \vdots \\ \partial f / \partial x_n \end{bmatrix} = 2 \begin{bmatrix} a_{:1}^T \\ a_{:2}^T \\ \vdots \\ a_{:n}^T \end{bmatrix} (Ax - b) = 2 A^T (Ax - b)$$

# Example: Least Squares

$$\nabla f = \vec{0} \implies 2A^{\top}(Ax - b) = \vec{0}$$

$$\implies A^{\top}Ax - A^{\top}b = 0 \implies A^{\top}Ax = A^{\top}b$$

$$\boxed{x = (A^{\top}A)^{-1}A^{\top}b}$$

least squares solution

# Easier method of calculating gradient

$$f(x) \qquad D[u]f = \langle \nabla, u \rangle = \nabla^T u$$

1- derive the directional derivative of $\underline{f}$ for an arbitrary direction $u \in \mathbb{R}^n$.

2- write the solution in form of $\langle z, u \rangle$

3 - $z$ is the gradient vector.

# Inner product

$$\langle Ax, y \rangle \qquad \begin{array}{l} x \in R^m \\ y \in \mathbb{R}^n \\ A \in \mathbb{R}^{n \times m} \end{array} \qquad \langle Ax, y \rangle = \langle x, A^\top y \rangle$$

product

# Inner product for matrices

$$A \in \mathbb{R}$$

$$\langle A, B \rangle \quad A, B \in \mathbb{R}^{m \times n}. \quad \text{Matrix } \cancel{\text{dot}} \text{ inner product}$$

$$\langle A, B \rangle = \sum_{i=1}^{m} \sum_{j=1}^{n} A_{ij} B_{ij}$$

$$\langle A, B \rangle = \text{trace}(A^T B) = \text{trace}(A B^T)$$

$$\langle A, BC \rangle = \langle B^T A, C \rangle = \langle A C^T, B \rangle$$

# Compute gradient (easy way)

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

$$D[u]f = \frac{d}{d\alpha} f(x+\alpha u)\Big| = \frac{d}{d\alpha}(\vec{x} + \alpha \vec{u})^T A (x+\alpha u)$$
$$\Big|_{\alpha=0}$$

$$\left[\frac{d}{d\alpha}(x+\alpha u)\right]^T A (x+\alpha u) + (x+\alpha u)^T A \left[\frac{d}{d\alpha}(x+\alpha u)\right]$$

$$u^T A (x+\alpha u) + (x+\alpha u)^T A u\Big|_{\alpha=0} \Rightarrow \underbrace{u^T A x}_{1\times1} + \underbrace{x^T A u}_{1\times1}$$

$$= x^T A^T u + x^T A u = (x^T A^T + x^T A) u = (Ax + A^T x)^T u \equiv$$

$$= (Ax + A^T x)^T u \Rightarrow \langle \underbrace{Ax + A^T x}_{\nabla}, u \rangle$$

$$\Rightarrow \nabla f = Ax + A^T x = (A + A^T)x$$

# Compute Gradient (easy way) least squares

$$f(x) = \|Ax - b\|^2 = (Ax-b)^\top (Ax-b)$$

$$\frac{d}{d\alpha} f(x+\alpha u) = \frac{d}{d\alpha} \left. (A(x+\alpha u) - b)^\top (A(x+\alpha u) - b) \right\}\Big|_{\alpha=0}$$

$$= \left. (Au)^\top (A(x+\alpha u)) + (A(x+\alpha u) - b)^\top (Au) \right|_{\alpha=0}$$

$$= \alpha(Au)^\top (Ax-b) + (Ax-b)^\top Au$$

$$= 2(Ax-b)^\top Au = (2A^\top(Ax-b))^\top u$$

$$= \langle 2A^\top(Ax-b), u \rangle$$

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