# Linear Algebra for Computer Science

## Lecture 32
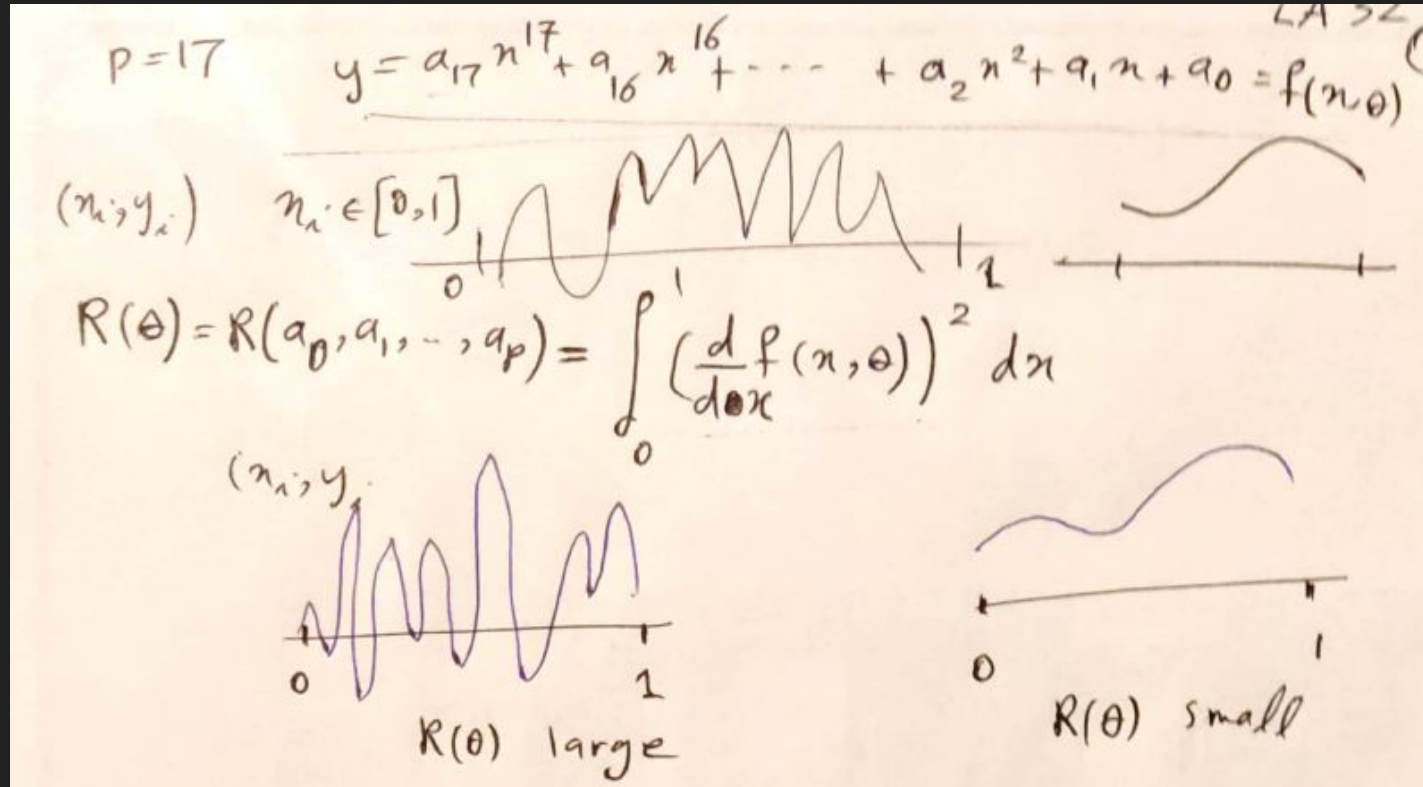
## Regularization, the Gradient Vector

# How to avoid overfitting

- Limit model complexity (e.g. degree of polynomial)
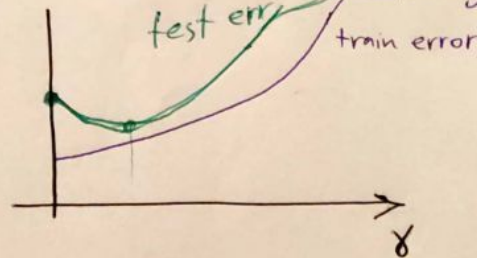- **Regularization**

# Regularization



$$p = 17 \qquad y = a_{17} x^{17} + a_{16} x^{16} + \cdots + a_2 x^2 + a_1 x + a_0 = f(x, \theta)$$

$(x_i, y_i) \qquad x_i \in [0, 1]$

$$R(\theta) = R(a_0, a_1, \ldots, a_p) = \int_0^1 \left( \frac{d}{dx} f(x, \theta) \right)^2 dx$$

$(x_i, y_i)$

$R(\theta)$ large

$R(\theta)$ small

# Regularization



$$R(\theta) = \int_0^1 \left( \frac{d}{dx} \sum_{i=0}^{p} a_i x^i \right)^2 dx = \int_0^1 \left( \sum_{i=0}^{p} i\, a_i x^{i-1} \right)^2 dx$$

$$\Rightarrow R(\theta) = \theta^T M \theta$$

a simpler choice $\Longleftarrow$ $R(\theta) = \theta^T \theta = \|\theta\|^2$

$$\min_{\theta} \ C(\theta)$$

$$\min_{\theta} \ C(\theta) + \gamma R(\theta)$$

hyper parameter

regularizer

test error    train error

# Regularization

$$\text{cost function:} \left( \sum_{i=1}^{N} f\left(x_i, \begin{bmatrix} a_1 \\ \vdots \\ a_n \end{bmatrix}\right) - y_i \right)^2 + \gamma\, R(a_1, \ldots, a_n)$$

$$C(\theta) = \left( \sum_{i=1}^{W} f(x_i, \theta) - y_i \right)^2 + \gamma\, \underbrace{R(\theta)}_{\text{Regularizer}}$$

Regularization

$\gamma$  Regularization parameter     $\gamma \uparrow$ emphasis on smoothness

(hyperparameter)      $\gamma \downarrow$      "      "  data

Simple Example:  $R(\theta) = R(a_0, a_1, \ldots, a_n) = a_0^2 + a_1^2 + \cdots + a_n^2$

$$= \|\theta\|^2$$

# Solving Regularized Polynomial Regression

$$C(\theta) = \| M\theta - y \|^2 \qquad \text{polynomial regression}$$

$$R(\theta) = \theta^T \theta = \| \theta \|^2$$

$$\theta^* = \underset{\theta}{\arg\min} \; C(\theta) + \gamma R(\theta)$$

$$= \underset{\theta}{\arg\min} \; \| M\theta - y \|^2 + \gamma \| \theta \|^2$$
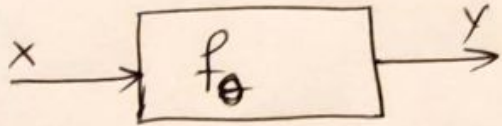
$$\approx \| M\theta - y \|^2 + \| \sqrt{\gamma} I \theta \|^2$$

$$= \underset{\theta}{\arg\min} \; \left\| \begin{bmatrix} M \\ \sqrt{\gamma} I \end{bmatrix} \theta - \begin{bmatrix} y \\ 0 \end{bmatrix} \right\|^2$$

$$\boxed{\| u \|^2 + \| v \|^2 = \left\| \begin{bmatrix} u \\ v \end{bmatrix} \right\|^2}$$

$$\| B\theta - z \|^2$$

$$\theta^* = (B^T B)^{-1} B^T z$$

# Nonlinear in input / Linear in parameters



so far $f(x, \theta) = \phi(x)^T \theta$     linear in $\theta$

$$f(x, \theta) = \tanh(x)\, a + (\sin x)b + (\log x)c + \exp(x)d$$

$$= \left[ \tanh(x),\ \sin x,\ \log x,\ \exp(x) \right] \begin{bmatrix} a \\ b \\ c \\ d \end{bmatrix}$$

$$C(\theta) = \left\| \phi \begin{bmatrix} \phi(x_1)^T \\ \phi(x_2)^T \\ \vdots \\ \phi(x_n)^T \end{bmatrix} \theta - \begin{bmatrix} y_1 \\ y_2 \\ | \\ y_n \end{bmatrix} \right\|^2$$

# What if model is not linear in parameters?



$$C(\theta) = C\left(\begin{bmatrix} \theta_1 \\ \theta_2 \\ \vdots \\ \theta_p \end{bmatrix}\right)$$

LA 32 (III)

$$\frac{\partial C}{\partial \theta_1}(\theta) = 0$$

$$\frac{\partial C(\theta)}{\partial \theta_2} = 0$$

$$\vdots$$

$$\frac{\partial C(\theta)}{\partial \theta_p} = 0$$

$n$ equations

$n$ unknowns $(\theta_1, \theta_2, \dots, \theta_p)$

- solve for $\theta$

- use gradient based optimization

# What if model is not linear in parameters?

1. Compute the partial derivatives and set them equal to zero
   a. n (nonlinear) equations in n unknowns
   b. Cannot be solved in most cases
2. Use Gradient-based optimization
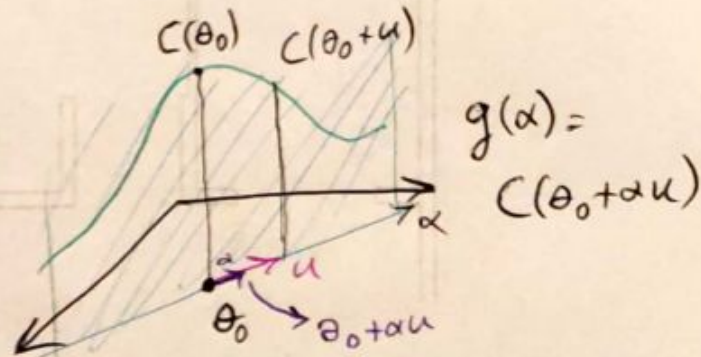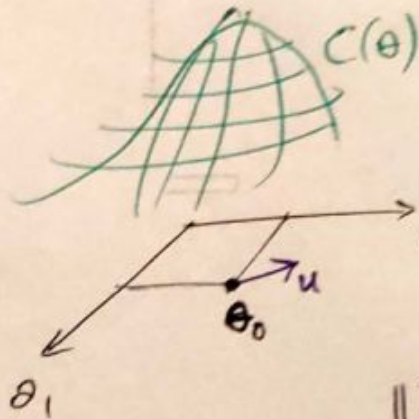
# 1D derivative



$$f'(x_0) = \frac{df}{dx}\bigg|_{x=x_0} = \;?$$

$$f'(x_0) = \lim_{\delta \to 0} \frac{f(x_0 + \delta) - f(x_0)}{\delta}$$

# N-dimentional functions: Directional Derivatives



$$C(\theta) = C(\theta_1, \theta_2, \underline{\phantom{xx}}, \theta_p)$$

$$C(\theta)$$

$$\theta_2$$

$$u$$

$$\theta_0$$

$$\theta_1$$

$$C(\theta_0) \quad C(\theta_0 + u)$$

$$g(\alpha) = C(\theta_0 + \alpha u)$$

$$\alpha$$

$$u$$

$$\theta_0 \quad \theta_0 + \alpha u$$

$$\|\vec{u}\| = 1$$
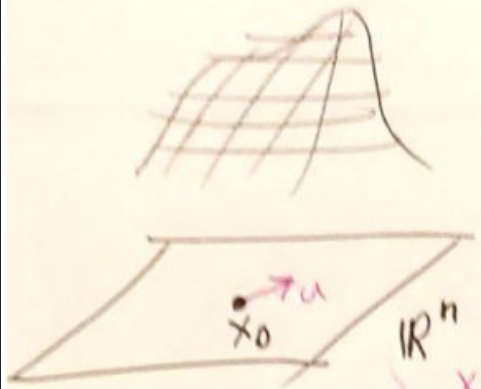
$$D[u]C(\theta_0) = \frac{C(\theta_0 + \alpha u) - C(\theta_0)}{\alpha} = \frac{d}{d\alpha} C(\theta_0 + \alpha u)\Big|_{\alpha = 0}$$

# General Directional Derivative

$f(x)$

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

$x_0$    $u$    $\mathbb{R}^n$

$x_0 \in \mathbb{R}^n,\ u \in \mathbb{R}^n$

$D[u]f(x_0) = \dfrac{d}{d\alpha} f(x_0 + \alpha u)\Big|_{\alpha=0}$

$\|u\| = 1$

Let $u$ be any vector (not just a unit vector $\|u\|=1$)

$D[u]f(x_0) = \dfrac{d}{d\alpha} f(x_0 + \alpha u)\Big|_{\alpha=0}$    directional derivative

# Scaling the direction

$$D[\beta u] f(x_0) = \frac{d}{d\alpha} f(x_0 + \alpha \beta u)\Big|_{\alpha = 0}$$

$$= \lim_{\alpha \to 0} \frac{\left(f(x_0 + \overbrace{\alpha \beta}^{\gamma} u) - f(x_0)\right) \beta}{\alpha \underbrace{\beta}_{\gamma}}$$

$$= \lim_{\gamma \to 0} \frac{f(x_0 + \gamma u) - f(x_0)}{\gamma} \beta$$

$$\Rightarrow D[\beta u] f(x_0) = \beta \, D[u] f(x_0)$$

# Directional Derivative is linear in the direction variable

$$\Rightarrow D[\beta u] f(x_0) = \beta\, D[u] f(x_0)$$

for differentiable functions

$$D[u+v] f(x_0) = D[u] f(x_0) + D[v] f(x_0)$$

$$\Rightarrow \text{For } \cancel{\text{differ}} \text{ a differentiable function } f: \mathbb{R}^n \to \mathbb{R}$$

$$D[u] f(x_0) = D[u] f \Big|_{x=x_0} \quad \text{is linear in } \vec{u}.$$

# The Gradient Vector

$\Rightarrow$ For ~~differ~~ a differentiable function $f: \mathbb{R}^n \to \mathbb{R}$

$$D[u]f(x_0) = D[u]f\Big|_{x=x_0} \quad \text{is linear in } \vec{u}.$$

$$D[u]f(x_0) = m^T u = \nabla^T u = \nabla(x_0)^T u$$

$D[\cdot]f\Big|_{x_0} : \mathbb{R}^n \to \mathbb{R}$  $\in \mathbb{R}^n$  $m \in \mathbb{R}^n$  $\hookrightarrow$ gradient vector

# The Gradient Vector and Partial Derivatives

# The Gradient Vector and Partial Derivatives



$$\nabla = \begin{bmatrix} \partial f / \partial x_1 \\ \partial f / \partial x_2 \\ \vdots \\ \partial f / \partial x_n \end{bmatrix}$$

# How to derive the gradient?



How to ~~comput~~ calculate $\nabla$

حل حد حد : 1- find $\partial f / \partial x_1 \quad \dfrac{\partial f}{\partial x_2} \quad - \quad \dfrac{\partial f}{\partial x_n}$

2- arrange in a vector

$$\nabla = \begin{bmatrix} \partial f / \partial x_1 \\ \partial f / \partial x_2 \\ | \\ \partial f / \partial x_n \end{bmatrix}$$

# Example: Least Squares



Least squares

$$x^* = \text{argmin} \; \|Ax - b\|^2 \qquad x^* = (A^T A)^{-1} A^T b$$

$$A = \begin{bmatrix} c_1 & c_2 & \cdots & c_n \end{bmatrix} = \begin{bmatrix} r_1^T \\ r_2^T \\ \vdots \\ r_m^T \end{bmatrix} = \begin{bmatrix} a_{11} & a_{12} & \cdots & a_{1n} \\ & & & \\ a_{m1} & a_{m2} & \cdots & a_{mn} \end{bmatrix}$$

$$x = \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix}$$

$$\|Ax - b\|^2 = \left\| \begin{bmatrix} r_1^T \\ r_2^T \\ \vdots \\ r_m^T \end{bmatrix} x - \begin{bmatrix} b_1 \\ b_2 \\ \vdots \\ b_m \end{bmatrix} \right\|^2$$

$$= \left\| \begin{bmatrix} r_1^T x - b_1 \\ r_2^T x - b_2 \\ \vdots \\ r_m^T x - b_m \end{bmatrix} \right\|^2 = \sum_{i=1}^{m} (r_i^T x - b_i)^2 = f(x)$$

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

# Example: Least Squares

$$= \left\| \begin{bmatrix} r_1^T x - b_1 \\ r_2^T x - b_2 \\ \vdots \\ r_m^T x - b_m \end{bmatrix} \right\|^2 = \sum_{i=1}^{m} \left( r_i^T x - b_i \right)^2 = f(x)$$

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

$$\frac{\partial f}{\partial x_k} = \sum_{i=1}^{m} 2 \; a_{ik} \; \left( r_i^T x - b \right)$$

$$= 2 \sum_{i=1}^{m} a_{ik} \left( r_i^T x - b \right) \qquad \begin{bmatrix} r_1^T x - b \\ r_2^T x - b \\ \vdots \\ r_m^T x - b \end{bmatrix}$$

$$= 2 \begin{bmatrix} a_{1k} & a_{2k} & \cdots & a_{mk} \end{bmatrix}$$

$$= 2 \, c_k^T \left( \begin{bmatrix} r_1^T \\ r_2^T \\ \vdots \\ r_m^T \end{bmatrix} x - b \right)$$

$$\frac{\partial f}{\partial x_k} = 2 \, c_k^T \left( A x - b \right)$$

# Example: Least Squares

$$\frac{\partial f}{\partial x_k} = 2 c_k^T (Ax - b)$$

$$\nabla = \begin{bmatrix} \frac{\partial f}{\partial x_1} \\ \frac{\partial f}{\partial x_2} \\ \vdots \\ \frac{\partial f}{\partial x_n} \end{bmatrix} = 2 \begin{bmatrix} c_1^T (Ax-b) \\ c_2^T (Ax-b) \\ \vdots \\ c_n^T (Ax-b) \end{bmatrix} = 2 \underbrace{\begin{bmatrix} c_1^T \\ c_2^T \\ \vdots \\ c_n^T \end{bmatrix}}_{A^T} (Ax-b) = 2 A^T (Ax-b)$$

$$= 2 A^T (Ax-b)$$

# Example: Least Squares



$$f(x) = \|Ax - b\|^2 \qquad A \in \mathbb{R}^{m \times n}$$

$$x^* = \arg\min_x f(x)$$

$$\nabla f(x) = 2A^\top(Ax - b) = 2\underbrace{A^\top}_{n \times m}\left(\underbrace{A}_{m \times n}\underbrace{x}_{n \times 1} - \underbrace{b}_{m \times 1}\right)$$

$$\underbrace{\qquad}_{m \times 1}$$

$$\underbrace{\qquad}_{n \times 1} \in \mathbb{R}^n$$

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

$$\Rightarrow A^\top A x - A^\top b = 0$$

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

# Derive Gradient: Second Method

1- Calculate the diretional derivative

$$D[u]\, f(x_0) = \frac{d}{d\alpha} f(x + \alpha u)\Big|_{\alpha = 0} = g(x, u)$$

2- Write $D[u]f$ in the form of $\nabla^T u$

$$f(x) = \|Ax - b\|^2 = (Ax - b)^T (Ax - b) \qquad \langle \nabla, u \rangle \; \overset{\longleftarrow}{\text{ضرب داخلی}}$$

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

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

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

$$\overset{\alpha = 0}{=} (Au)^T (Ax - b) + (Ax - b)^T A u$$

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

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

# Derive Gradient: Second Method

$$D[u]\, f(x) = 2(Ax - b)^T A\, u = \left(2 A^T (Ax - b)\right)^T u$$

$$= \langle \nabla, u \rangle = \nabla^T u$$

$$\Rightarrow \boxed{\nabla = 2 A^T (Ax - b)}$$

# Final Project