

Lec 04 : Linear Regression (Contd.)

$$\min_w E(w, D) = \|Xw - y\|^2$$

1-d fn: $(ax - y)^2$

d-dim $\|Ax - y\|^2 \rightarrow A$ is PSD $\frac{d^2 f}{dx^2} \geq 0$

$$W^* = (X^T X)^{-1} X^T y$$

(pseudoinverse)

$$\nabla_w E = 0 \Rightarrow$$

$$X^T (Xw) - X^T y = 0$$

$$X^T y$$

$$a_{11}x_1 + \dots + a_{1d}x_d = b_1$$

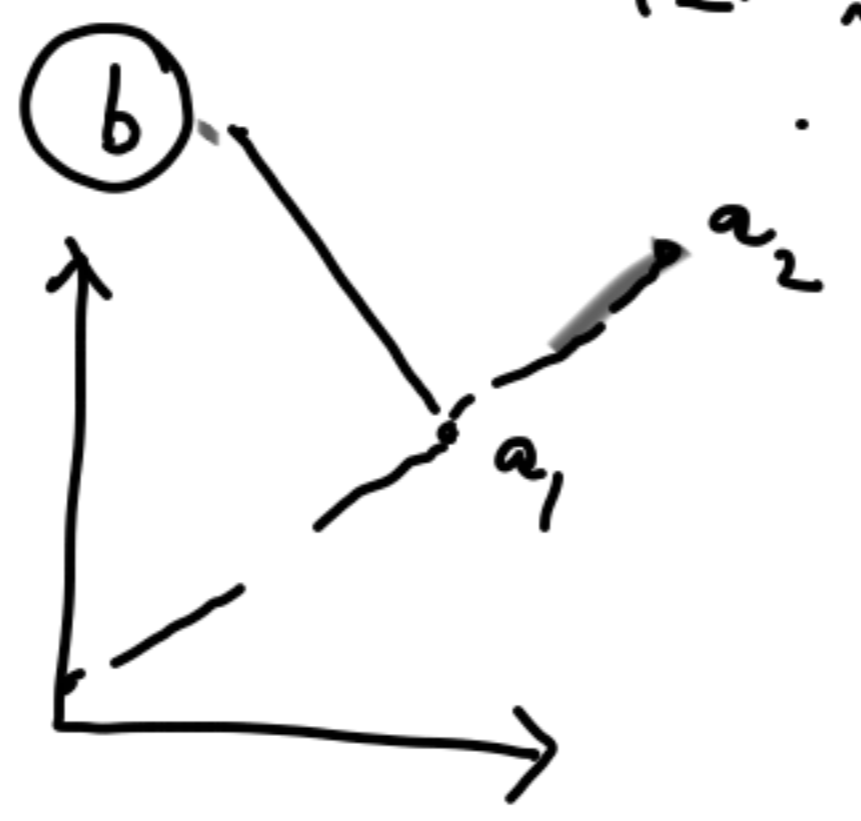
⋮

$$a_{m1}x_1 + \dots + a_{md}x_d = b_m$$

$$Ax = b \rightarrow$$

$$\begin{bmatrix} | & & | \\ a_1 & \dots & a_d \\ | & & | \end{bmatrix} \begin{bmatrix} x_1 \\ \vdots \\ x_d \end{bmatrix} = x_1 \underline{a_1} + x_2 \underline{a_2} + \dots + x_d \underline{a_d} \sim b$$

$$\begin{bmatrix} 1 \\ 2 \\ 4 \end{bmatrix} x = \begin{bmatrix} 3 \\ 4 \end{bmatrix} \begin{bmatrix} 3 \\ 6 \end{bmatrix} a_1$$



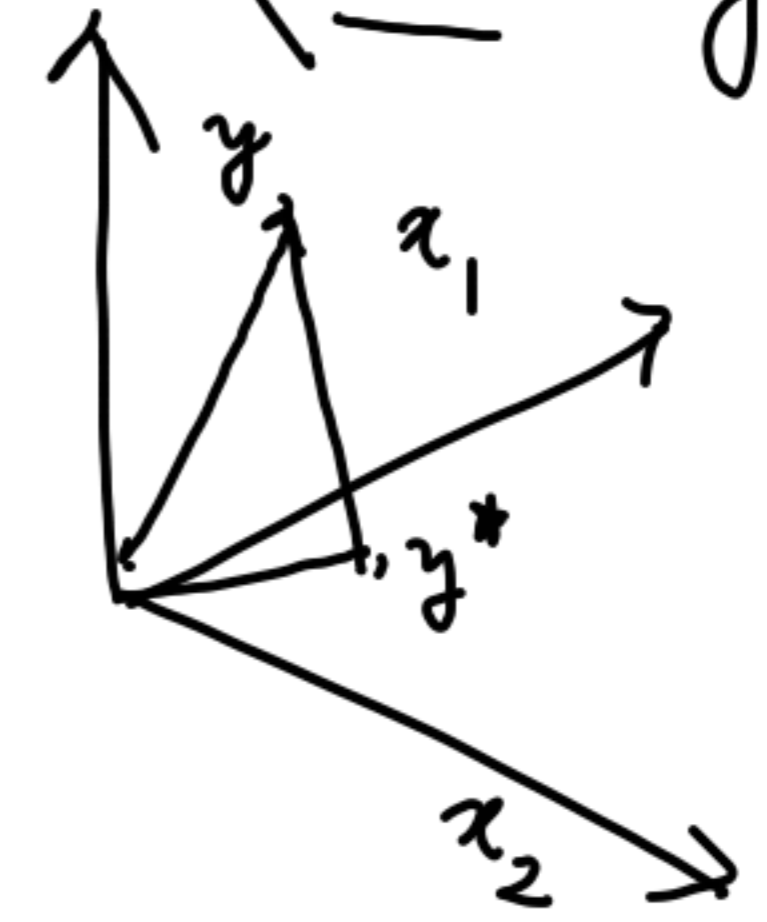
$$\text{Case 1: } \min \|Xw - y\|^2$$

$$Xw^* = y$$

Case 2: can't find w^*

$$X^T Xw - X^T y = 0$$

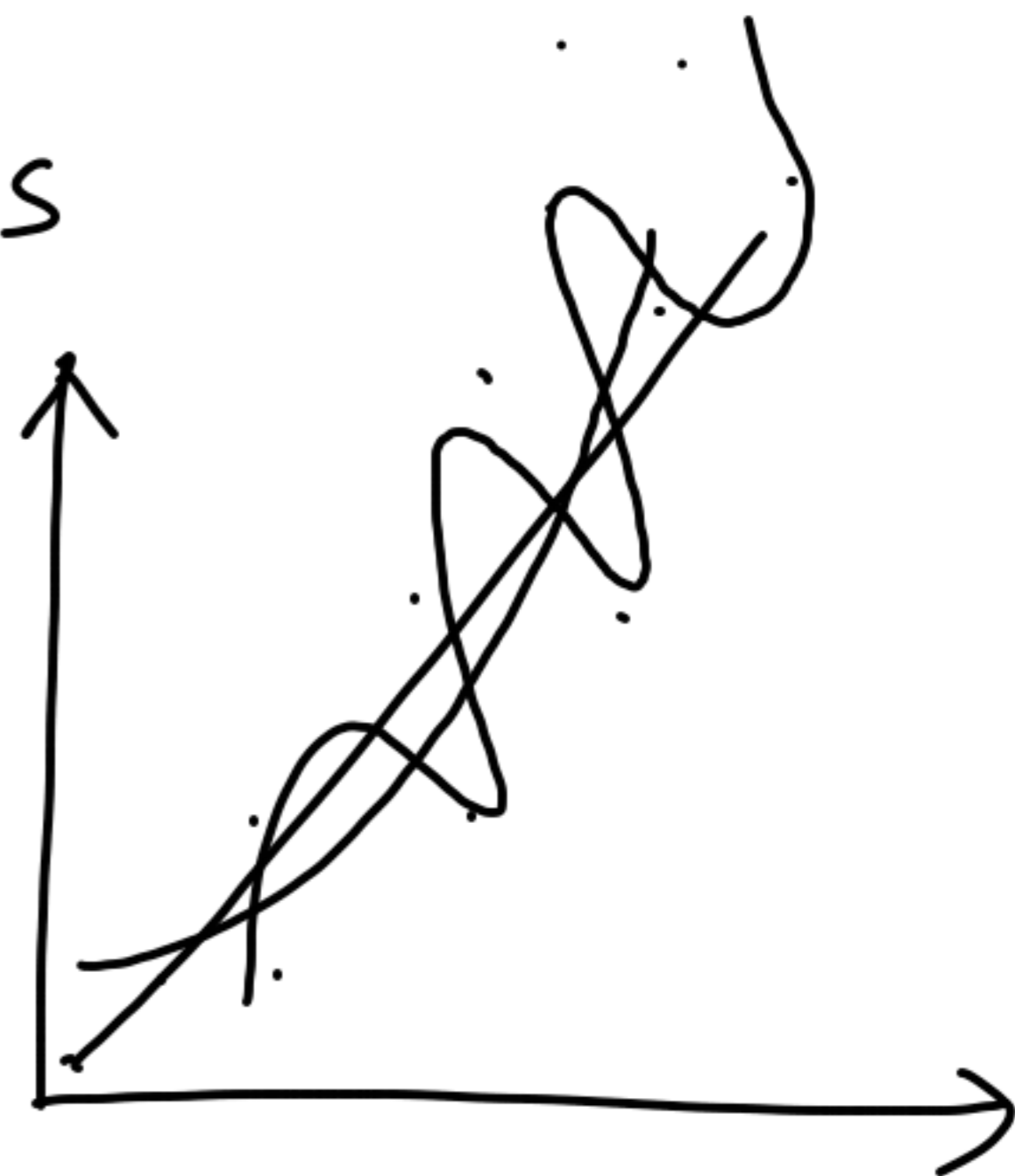
$$\Rightarrow X^T (Xw^* - y) = 0$$



Regression model with basis functions

$$\hat{y}_i = W_0 + W_1 x_i + \left(W_2 x_i^2 + W_3 x_i^3 + \dots + W_m x_i^m \right)$$

1-D data: x_i 's scalars



$$y^* = \phi w^*$$

$$\phi_0(x_i) = 1$$

$$\phi_1(x_i) = x_i$$

$$\vdots$$

$$\phi_m(x_i) = x_i^m$$

$$\hat{y}_i = \sum_{j=0}^m W_j \phi_j(x_i)$$

$$\hat{y}_i = \sum_{j=0}^m \phi_j(x_i) w_j$$

$m \gg d$

$$x_i = \begin{bmatrix} x_{i1} \\ \vdots \\ x_{id} \end{bmatrix}$$

$$\begin{aligned}
 \phi &= \begin{bmatrix} \phi_0(x_1) & \dots & \phi_m(x_1) \\ \vdots & & \vdots \\ \phi_0(x_n) & \dots & \phi_m(x_n) \end{bmatrix}_{n \times (m+1)} \\
 y &= \begin{bmatrix} y_1 \\ \vdots \\ y_n \end{bmatrix} \\
 \text{argmin } &\| \phi w - y \|^2 \\
 w^* &= (\phi^T \phi)^{-1} \phi^T y \\
 w &= \begin{bmatrix} w_0 \\ w_1 \\ \vdots \\ w_m \end{bmatrix}_{(m+1) \times 1}
 \end{aligned}$$

$$y_i = w^T x_i + \epsilon_i$$

noisy linear model

parameters: w

noise: $\epsilon_i \sim N(0, \sigma^2)$ i.i.d.

$$\text{Cov}(\epsilon_i, \epsilon_j) = 0 \quad \forall i \neq j$$

Probabilistic model of linear regression

$$D = \left\{ (x_i, y_i) \right\}_{i=1}^n$$

• Estimate w from this probabilistic model

• Maximum likelihood estimation (MLE)

MLE

A set of independent and identically distributed (i.i.d.) observations $\{y_1, \dots, y_n\}$ are generated by a probabilistic model parametrized by θ

$$y_j \sim \underbrace{P(y|\theta)}_{\text{Likelihood}}$$

$$P(y; \theta)$$

Log likelihood: $\log(P(y|\theta))$

- Mathematically nicer
- Numerical advantage

$$\theta_{MLE} = \operatorname{argmax}_{\theta} L(\theta) = \operatorname{argmax}_{\theta} \sum_{j=1}^n \log P(y_j|\theta)$$

$$P(\underline{y}|\underline{x}, \theta) \\ = \prod_{i=1}^n P(y_i|x_i, \theta)$$

Coin toss: Toss a coin n times, each is a binary RV
with Bernoulli dist

$$P(y_j | \theta) = \theta^{y_j} (1-\theta)^{1-y_j}$$

$$\begin{aligned} L(\theta) &= \log P(y | \theta) = \log \left(\prod_{j=1}^n P(y_j | \theta) \right) \\ &= \sum_{j=1}^n \log(P_j | \theta) = \sum_{j=1}^n \left(y_j \log \theta + (1-y_j) \log(1-\theta) \right) \\ &\Rightarrow \hat{\theta}_{MLE} = \frac{1}{n} \sum_{j=1}^n y_j \end{aligned}$$

MLE for regression

$$y_j = w^T x_j + \epsilon_j \sim N(0, \sigma^2)$$

$$y_j \sim N(w^T x_j, \sigma^2)$$

$$P(y_j | x_j, w) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left\{-\frac{(y_j - w^T x_j)^2}{2\sigma^2}\right\}$$

$$L(w) = \text{const.} \cdot \prod_{j=1}^n \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left\{-\frac{(y_j - w^T x_j)^2}{2\sigma^2}\right\}$$

$$\arg\max_w L(w)$$

$$w_{MLE}^* = \arg\min_w \sum_{j=1}^n (y_j - w^T x_j)^2$$

argmin $\|Xw - y\|^2$

$$W^* = (X^T X)^{-1} X^T y$$

Algorithmic way to optimize

Gradient Descent

- $w \leftarrow w_0$
- Repeat until convergence

$$w \leftarrow w - \eta \nabla_w E$$

$\|\nabla E(w)\| < \epsilon$
learning rate

$$\|Xw - y\|^2 = E$$

