

# Maximum Likelihood Estimation (MLE)

$\mathcal{X}$  : training data set

$\theta$  : parameters

$$\theta_{MLE}^* = \arg \max_{\theta} \underline{p(\mathcal{X}|\theta)}$$

$$p(\mathcal{X}|\theta) = \prod_{i=1}^N p(x_i|\theta)$$

# Maximum a Posteriori Estimation (MAP)

$$\theta_{MAP}^* = \arg \max_{\theta} p(\theta|\mathcal{X})$$

$$= \arg \max_{\theta} \frac{\underline{p(\mathcal{X}|\theta)} \overbrace{p(\theta)}}{p(\mathcal{X})}$$

our prior belief about  $\theta$

# Parametric Regression:

$$\underbrace{y}_{\text{observations}} = \underbrace{f(x)}_{\text{underlying process}} + \underbrace{\epsilon}_{\text{noise}}$$

$$\textcircled{\text{I}} \quad p(\epsilon) \sim N(\epsilon; 0, \sigma^2)$$

$$\textcircled{\text{II}} \quad p(y|x) \sim N(y; g(x|\theta), \sigma^2)$$

$$y = f(x) + \epsilon$$

$$y = \underbrace{g(x|\theta)}_{\text{constant}} + \underbrace{\epsilon}_{\text{random variable}}$$

approximation

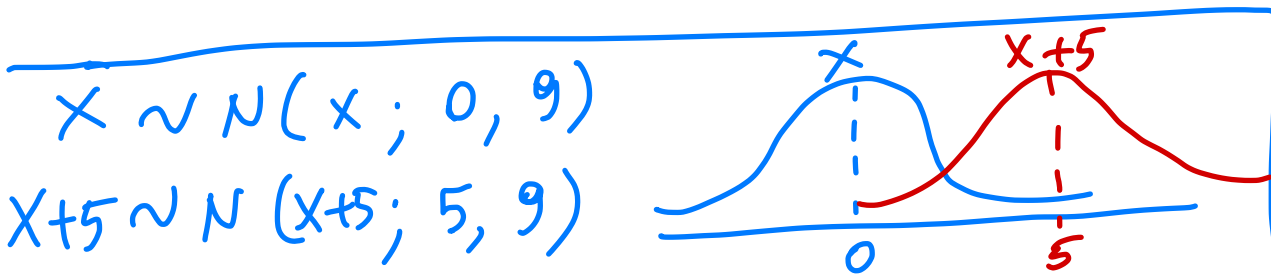
Learning problem:

approximate  $f(x)$  with  $g(x|\theta)$   
parameters

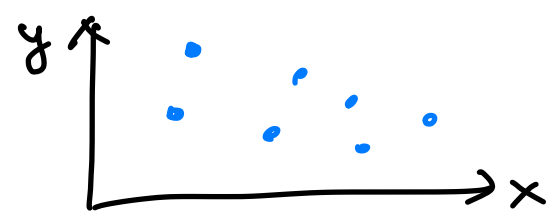
$$\begin{aligned} E[X] &= \delta \quad \text{constant} \\ E[X+c] &= \delta+c \\ \text{VAR}[X] &= k^2 \\ \text{VAR}[X+c] &= k^2 \end{aligned}$$

$$\begin{aligned} \underline{\underline{E[y|x]}} &= E[g(x|\theta) + \epsilon] \\ &= g(x|\theta) + \underbrace{E[\epsilon]}_0 \end{aligned}$$

$$\begin{aligned} \underline{\underline{\text{VAR}[y|x]}} &= \text{VAR}[g(x|\theta) + \epsilon] \\ &= 0 + \underbrace{\text{VAR}[\epsilon]}_{\sigma^2} \\ &= \sigma^2 \end{aligned}$$



$$\mathcal{X} = \{ (x_i, y_i) \}_{i=1}^N \quad x_i \in \mathbb{R}^1 \quad y_i \in \mathbb{R}^1$$



$$(x_i, y_i) \sim p(x_i, y_i)$$

$\xrightarrow{\text{i.i.d.}}$

$$\begin{aligned} p(x, y) &= p(x|y) \cdot p(y) \\ p(x, y) &= p(y|x) p(x) \end{aligned}$$

$$p(x_1, y_1, x_2, y_2, \dots, x_N, y_N) = \prod_{i=1}^N p(x_i, y_i)$$

Likelihood  $\Rightarrow L(\theta | \mathcal{X}) = \prod_{i=1}^N p(x_i, y_i)$

$$= \prod_{i=1}^N [p(y_i | x_i) p(x_i)]$$

log likelihood  $\Rightarrow \log \left[ \prod_{i=1}^N [p(y_i | x_i) p(x_i)] \right]$

$$\log(a \cdot b) = \log(a) + \log(b)$$

constant

$$= \sum_{i=1}^N \left[ \log[p(y_i | x_i)] + \log[p(x_i)] \right]$$

$$= + \sum_{i=1}^N \log[p(y_i | x_i)]$$

maximize  $\sum_{i=1}^N \log[p(y_i | x_i)] \Rightarrow \text{maximize } \sum_{i=1}^N \log \left[ \mathcal{N}(\underbrace{y_i}_{\text{R.V.}}; \underbrace{g(x_i | \theta)}_{\text{mean}}, \underbrace{\sigma^2}_{\text{variance}}) \right]$

$\mathcal{N}(y_i; g(x_i | \theta), \sigma^2)$

maximize  $\sum_{i=1}^N \log \left[ \frac{1}{\sqrt{2\pi\sigma^2}} \cdot \exp \left[ -\frac{(y_i - g(x_i|\theta))^2}{2\sigma^2} \right] \right]$

*(Note: In the original image, the terms  $\frac{1}{\sqrt{2\pi\sigma^2}}$  and  $\exp$  are marked with 'x' and a red arrow points from 'constant' to the exponent's denominator.)*

maximize  $\sum_{i=1}^N \left[ -\frac{[y_i - g(x_i|\theta)]^2}{2\sigma^2} \right]$

*(Note: In the original image,  $2\sigma^2$  is crossed out with a red arrow pointing to 'constant'.)*

maximize  $\sum_{i=1}^N -[y_i - g(x_i|\theta)]^2$

minimize  $\sum_{i=1}^N [y_i - \hat{y}_i]^2$

*(Note: In the original image,  $\hat{y}_i$  is indicated by a purple arrow pointing to  $g(x_i|\theta)$  in the expression above.)*

minimize  $\sum_{i=1}^N (y_i - \hat{y}_i)^2 = \sum_{i=1}^N e_i^2$

*(Note: In the original image,  $e_i$  is indicated by a purple arrow pointing to  $y_i - \hat{y}_i$ .)*

$$g(x_i|\theta) = w_0 + w_1 \cdot x_i$$

$$\theta = \{w_0, w_1\}$$

$$g(x_i|\theta) = w_0 + w_1 \cdot x_i + w_2 x_i^2$$

$$\theta = \{w_0, w_1, w_2\}$$

minimize  $\sum_{i=1}^N [y_i - g(x_i|\theta)]^2$

$g(x_i|\theta) = w_0 + w_1 x_i$   
 $\theta^* = \{w_0^*, w_1^*\} = ?$

$\text{Error}[\theta|x] = \sum_{i=1}^N [y_i - (w_0 + w_1 x_i)]^2$

$\frac{\partial \text{Error}[\theta|x]}{\partial w_0} = \sum_{i=1}^N [y_i - (w_0 + w_1 x_i)] \cdot (1) = 0$

$\frac{\partial \text{Error}[\theta|x]}{\partial w_1} = \sum_{i=1}^N [y_i - (w_0 + w_1 x_i)] \cdot (x_i) = 0$

$\sum_{i=1}^N y_i = \sum_{i=1}^N w_0 + \sum_{i=1}^N w_1 x_i$   
 $\sum_{i=1}^N (y_i \cdot x_i) = \sum_{i=1}^N w_0 x_i + \sum_{i=1}^N w_1 x_i^2$

$\bar{A}^T A \theta = \bar{A}^T b$   
 $\theta = A^{-1} \cdot b$

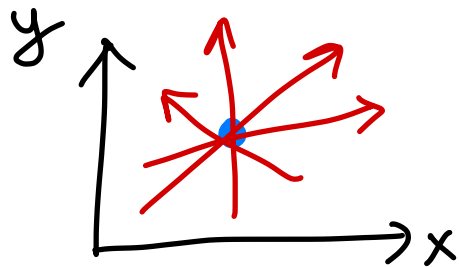
Exercise: Show this  
 2x2 matrix is invertible  
 if  $N \geq 2$ .

$\begin{bmatrix} N & \sum_{i=1}^N x_i \\ \sum_{i=1}^N x_i & \sum_{i=1}^N x_i^2 \end{bmatrix} \begin{bmatrix} w_0 \\ w_1 \end{bmatrix} = \begin{bmatrix} \sum_{i=1}^N y_i \\ \sum_{i=1}^N (y_i \cdot x_i) \end{bmatrix}$

$$G = \begin{bmatrix} w_0 \\ w_1 \end{bmatrix} = \begin{bmatrix} N & \sum_{i=1}^N x_i \\ \sum_{i=1}^N x_i & \sum_{i=1}^N x_i^2 \end{bmatrix}^{-1} \cdot \begin{bmatrix} \sum_{i=1}^N y_i \\ \sum_{i=1}^N (y_i \cdot x_i) \end{bmatrix}$$

this matrix is invertible if  $N \geq 2$ .

$\Rightarrow$  If there is a single data point ( $N=1$ )

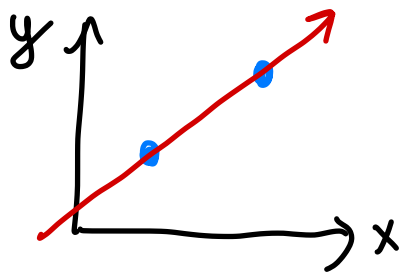


$$A = \begin{bmatrix} 1 & x_1 \\ x_1 & x_1^2 \end{bmatrix}$$

$$\Rightarrow \det(A) = 1 \cdot x_1^2 - x_1 \cdot x_1 = 0$$

$\Downarrow$   
not invertible

$\Rightarrow$  If there are two data points ( $N=2$ )



$$A = \begin{bmatrix} 2 & (x_1 + x_2) \\ (x_1 + x_2) & x_1^2 + x_2^2 \end{bmatrix}$$

$$\Rightarrow \det(A) = 2x_1^2 + 2x_2^2 - x_1^2 - x_2^2 - 2x_1x_2 = (x_1 - x_2)^2$$

$> 0$   
 $\Downarrow$   
invertible

# Polynomial Regression!

$$w_0 = w_0 \cdot x_i^0$$

$k^{\text{th}}$  order polynomial

$k=1 \Rightarrow$  linear regression

$$g(x_i) | w_0, w_1, w_2, \dots, w_k = w_0 + w_1 x_i^1 + w_2 x_i^2 + \dots + w_k x_i^k$$

$$\begin{bmatrix} \boxed{N} & \sum_{i=1}^N x_i^1 & \sum_{i=1}^N x_i^2 & \dots & \sum_{i=1}^N x_i^k \\ \sum_{i=1}^N x_i & \sum_{i=1}^N x_i^2 & \sum_{i=1}^N x_i^3 & \dots & \sum_{i=1}^N x_i^{k+1} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ \sum_{i=1}^N x_i^k & \sum_{i=1}^N x_i^{k+1} & \dots & \sum_{i=1}^N x_i^{2k} \end{bmatrix}$$

$$A = D^T \cdot D$$

A is invertible if  $N \geq k+1$ .

$$\begin{bmatrix} \boxed{w_0} \\ \boxed{w_1} \\ \vdots \\ w_k \end{bmatrix}$$

$\theta$

$$\begin{bmatrix} \sum_{i=1}^N (y_i \cdot x_i^0) \\ \sum_{i=1}^N (y_i \cdot x_i^1) \\ \vdots \\ \sum_{i=1}^N (y_i \cdot x_i^k) \end{bmatrix}$$

$b$

$$\theta = \bar{A}^{-1} \cdot b$$

$K=0 \Rightarrow$  constant fit

$$N.wo = \sum_{i=1}^N y_i \Rightarrow \boxed{w_0^* = \frac{\sum_{i=1}^N y_i}{N}} \quad \bar{y}$$

↪ sample mean

$$D = \begin{bmatrix} 1 & x_1 & - & - & x_1^K \\ 1 & x_2 & - & - & x_2^K \\ \vdots & \vdots & \ddots & & \vdots \\ 1 & x_N & - & - & x_N^K \end{bmatrix}$$

$$D^T = \begin{bmatrix} 1 & 1 & - & - & 1 \\ x_1 & x_2 & - & - & x_N \\ \vdots & \vdots & \ddots & & \vdots \\ x_1^K & x_2^K & - & - & x_N^K \end{bmatrix}$$

$D^T$

$$D = \begin{bmatrix} 1 & x_1 & - & - & x_1^K \\ 1 & x_2 & - & - & x_2^K \\ \vdots & \vdots & \ddots & & \vdots \\ 1 & x_N & - & - & x_N^K \end{bmatrix}$$

$D$

$$= \begin{bmatrix} N & \sum_{i=1}^N x_i & - & - & \sum_{i=1}^N x_i^K \\ \sum_{i=1}^N x_i & \sum_{i=1}^N x_i^2 & - & - & \sum_{i=1}^N x_i^3 \\ \vdots & \vdots & \ddots & & \vdots \\ \sum_{i=1}^N x_i^K & \sum_{i=1}^N x_i^{K+1} & - & - & \sum_{i=1}^N x_i^{2K} \end{bmatrix}$$

$$D^T \cdot D = A$$

if  $N < K+1$ ,  $D^T D$  is not invertible.  
if  $N \geq K+1$ ,  $D^T D$  is invertible.