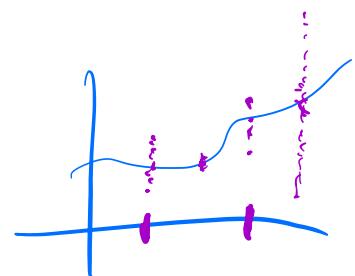


# Maximum Likelihood Estimation (MLE)

How about using a different learner from ERM?

$$f_{\text{Bayes}} = \underset{\{f \in \mathcal{F} \mid f: X \rightarrow Y\}}{\operatorname{argmin}} L(f) \quad E[e(f(\vec{x}), y)]$$

$$\begin{aligned} f_{\text{Bayes}}(\vec{x}) &= E[Y \mid \vec{X} = \vec{x}] \\ &= \int_y y \cdot P_{Y|\vec{X}}(y \mid \vec{x}) dy \end{aligned}$$



Let's use the dataset  $D$  to estimate  $P_{Y|\vec{X}}$

## MLE Basics

$$D = (Z_1, \dots, Z_n) \in \mathbb{Z}^n, \quad P_D, P_D$$

$Z_i$  are i.i.d. with  $P_Z$  and pmf or pdf  $p_Z$

If we have a fixed dataset  $D = (Z_1, \dots, Z_n)$   
how can learn what  $p$  is

$\Rightarrow$  pick the  $p$  that makes the data most likely

$$P_D(D) = P_D(z_1, \dots, z_n) \stackrel{\text{independent}}{=} P_{z_1}(z_1) P_{z_2}(z_2) \dots P_n(z_n)$$

identically distributed

$$= p_z(z_1) \dots p(z_n) \rightarrow \text{i.i.d.}$$

$$P_{\text{MLE}} = \underset{p \in \mathcal{H}}{\operatorname{argmax}} \prod_{i=1}^n p(z_i) \approx p_z$$

$\downarrow P_{\vec{x}, y}$

We pick the pmf/pdf that makes the data under this prob. dist. most likely.

Ex:  $Z_i \sim \text{Bernoulli}(\alpha^*)$  is the  $i$ -th flip of an unfair coin

$$p_z(z) = (\alpha^*)^z (1-\alpha^*)^{(1-z)}$$

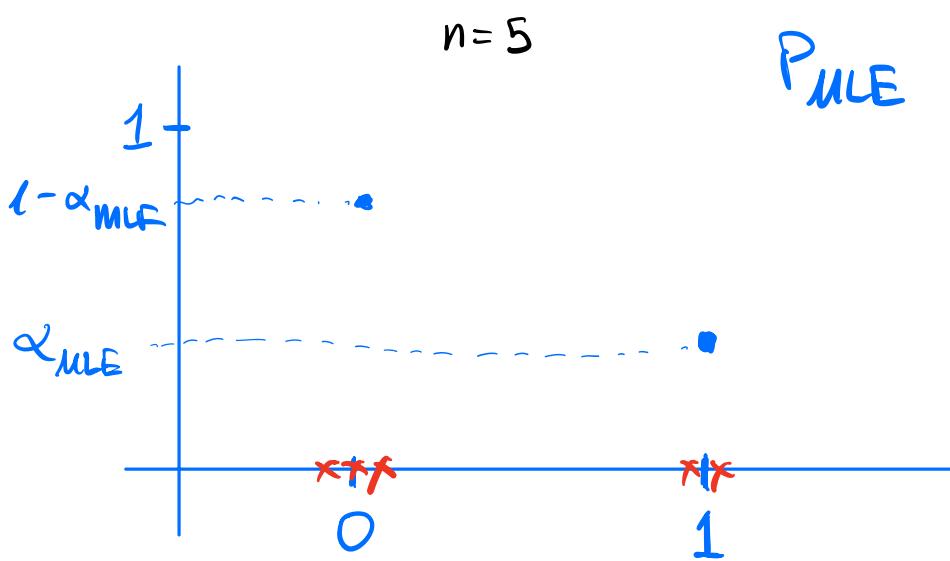
$$\mathcal{H} = \{p \mid p: \mathbb{Z} \rightarrow [0,1] \text{ and } p(z) = \alpha^z (1-\alpha)^{1-z}, \alpha \in [0,1]\}$$

Auf  $p_\alpha \in \mathcal{H}$  has the form

$$p_\alpha(z) = \alpha^z (1-\alpha)^{1-z} \text{ for some } \alpha \in [0,1]$$

$$P_{\text{MLE}} = P_{\alpha_{\text{MLE}}} \approx p_z \quad \text{likelihood} = P_D(D | \alpha)$$

$$\text{where } \alpha_{\text{MLE}} = \underset{\alpha \in [0,1]}{\operatorname{argmax}} \underbrace{\prod_{i=1}^n p_\alpha(z_i)}_{= p(z_i | \alpha)}$$



Ex:  $Z_i \sim \mathcal{N}(\mu^*, 1)$  is the  $i$ -th persons height

$$P_Z(z) = \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{(z-\mu^*)^2}{2}\right)$$

$$\mathcal{H} = \{p \mid p: \mathbb{R} \rightarrow [0, \infty) \text{ and } p(z) = \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{(z-\mu)^2}{2}\right), \mu \in \mathbb{R}\}$$

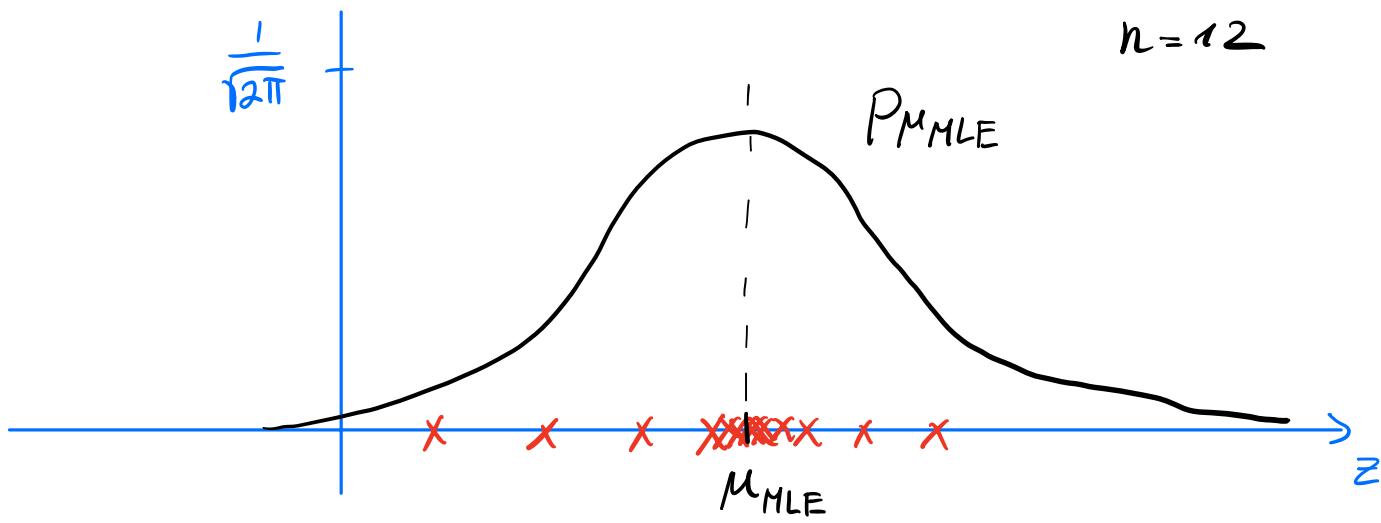
$$P_{MLE} = P_{\mu_{MLE}} \approx P_Z$$

$= P(z; \mu)$

where  $\mu_{MLE} = \underset{\mu \in \mathbb{R}}{\operatorname{argmax}} \prod_{i=1}^n P_\mu(z_i)$

and  $p(\cdot | \mu) \in \mathcal{H}$

all possible inputs



## Calculating $\mu_{MLE}$ :

$$\mu_{MLE} = \underset{\mu \in \mathbb{R}}{\operatorname{argmax}} \prod_{i=1}^n P(z_i | \mu)$$

have the same  
 maximizer  
 → log is monotonically  
 increasing  
 $\log = \log_e = \ln$

$$= \underset{\mu \in \mathbb{R}}{\operatorname{argmax}} \log \left( \prod_{i=1}^n P(z_i | \mu) \right)$$

$$= \underset{\mu \in \mathbb{R}}{\operatorname{argmax}} \sum_{i=1}^n \log(P(z_i | \mu))$$

$$= \underset{\mu \in \mathbb{R}}{\operatorname{argmin}} - \sum_{i=1}^n \log(P(z_i | \mu))$$

negative log-likelihood =  $-\log p(D | \mu)$

$$= \underset{\mu \in \mathbb{R}}{\operatorname{argmin}} - \sum_{i=1}^n \log \left( \frac{1}{\sqrt{2\pi}} \exp \left( -\frac{(z_i - \mu)^2}{2} \right) \right)$$

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

$$= \underset{\mu \in \mathbb{R}}{\operatorname{argmin}} - \sum_{i=1}^n \left[ \log \left( \frac{1}{\sqrt{2\pi}} \right) + \log \left( \exp \left( -\frac{(z_i - \mu)^2}{2} \right) \right) \right]$$

$\log(\exp(x)) = x$

$$= \underset{\mu \in \mathbb{R}}{\operatorname{argmin}} - \sum_{i=1}^n \left[ \log \left( \frac{1}{\sqrt{2\pi}} \right) - \frac{(z_i - \mu)^2}{2} \right]$$

$$= \underset{\mu \in \mathbb{R}}{\operatorname{argmin}} \sum_{i=1}^n \frac{(z_i - \mu)^2}{2}$$

$\underbrace{g(\mu)}$

$$\frac{d}{d\mu} g(\mu) = -\frac{2}{2} \sum_{i=1}^n (z_i - \mu) = -\sum_{i=1}^n (z_i - \mu) = 0$$

$$\Rightarrow \sum_{i=1}^n \mu = \sum_{i=1}^n z_i$$

$$n\mu = \sum_{i=1}^n z_i$$

$$\Rightarrow \mu_{MLE} = \mu = \frac{1}{n} \sum_{i=1}^n z_i \rightarrow^* \text{sample mean (but with } \frac{1}{n} \text{ instead of } \frac{1}{n-1})$$

If you want to find the pdf of a normal dist that makes the data most likely, calculate the sample mean\* over the data.

## Estimating $P_{Y|\vec{X}}$

$$D = ((\vec{X}_1, Y_1), \dots, (\vec{X}_n, Y_n)) \in (\mathcal{X} \times \mathcal{Y})^n, P_D, p_D$$

$(\vec{X}_i, Y_i)$  are i.i.d with  $P_{\vec{X}, Y}$  and  $p_{\vec{X}, Y}$   
product rule

$$P_{\vec{X}, Y}(\vec{x}, y) \stackrel{\text{product rule}}{=} P_{Y|\vec{X}}(y | \vec{x}) P_{\vec{X}}(\vec{x})$$

$$\text{Assume } Y_i | \vec{X}_i = \vec{x}_i \sim N(\vec{x}_i^T \vec{w}^*, 1)$$