

DATASET:  $\{x_i\}_{i=1}^N$ ,  $x_i \in \{0, 1\}$

Example:  $\{0, 1, 0, 1, 0\} = D$

DATA - DRIVEN  
APPROACH

MAXIMUM LIKELIHOOD  
ESTIMATION  
(MLE)

$$P(x=1) = \frac{2}{5}$$

PRIOR - INDUCED  
APPROACH

MAXIMUM A POSTERIORI  
(MAP)

- 1) Define a set of plausible values for the underlying model parameters  
(in this example, types of coin)

$\{\text{fair}, \text{2-headed}, \text{coin flips}\}$   
 $\quad \quad \quad \text{Heads w/ 30\%}$

prior belief

- 2) TEST each hypothesis.

$$\begin{aligned} P(\text{fair} | \{0, 1, 0, 1, 0\}) &= \\ &= \frac{P(\{0, 1, 0, 1, 0\} | \text{fair}) \cdot P(\text{fair})}{P(\{0, 1, 0, 1, 0\})} \end{aligned}$$

Bayes' theorem

$$P(\text{fair} | D) = \frac{\overbrace{P(D|\text{fair}) \cdot P(\text{fair})}^{\text{P(D)}}}{P(D)}$$

$$= \frac{\overbrace{P(0|\text{fair}) \cdot P(1|\text{fair}) \cdot P(0|\text{fair}) \cdot P(1|\text{fair}) \cdot P(0|\text{fair}) \cdot P(\text{fair})}^{\text{P(D)}}}{P(D)}$$

↗  
Conditional  
independence

$$= \frac{(1/2)^3 (1/2)^2 (1/3)}{P(D)}$$

$$\begin{aligned} P(D) &= \overbrace{P(D|\text{fair}) \cdot P(\text{fair})}^{\text{P(D)}} + \overbrace{P(D|\text{2-headed}) \cdot P(\text{2-headed})}^{\text{P(D)}} \\ &\quad + \overbrace{P(D) \text{ flips Heads } 30\% \cdot P(\text{flips Heads } 30\%)}^{\text{P(D)}} \\ &= (1/2)^3 (1/2)^2 (1/3) + (0)^3 (1)^2 (1/3) + (7/10)^3 (3/10)^2 (1/3) \end{aligned}$$

$$P(\text{fair} | D) = A$$

$$P(\text{2-headed} | D) = B$$

$$P(\text{coin flip, H w/ 30\%} | D) = C$$

If  $A$  is the largest, then

pick prior choice "fair".

$$2) P(x=1 | \text{fair}) = \frac{1}{2}$$

Parameter Estimation of MLE  
and MAP w/ converge to

Same Solution if:  
is large enough.

- 1) Sample size
- 2) prior in MAP is uniformly-distributed.

MLE

$$\omega = \arg \max_{\omega} \underline{P(x|\omega)}$$

MAP

$$\begin{aligned}\omega &= \arg \max_{\omega} \underline{P(x|\omega)} \cdot \underline{P(\omega)} \\ &= \arg \max_{\omega} P(\omega|x) \cdot P(x) \\ &\propto \arg \max_{\omega} P(\omega|x)\end{aligned}$$

$\omega_{MLE} = \omega_{MAP}$  if  $P(\omega)$  is a constant.

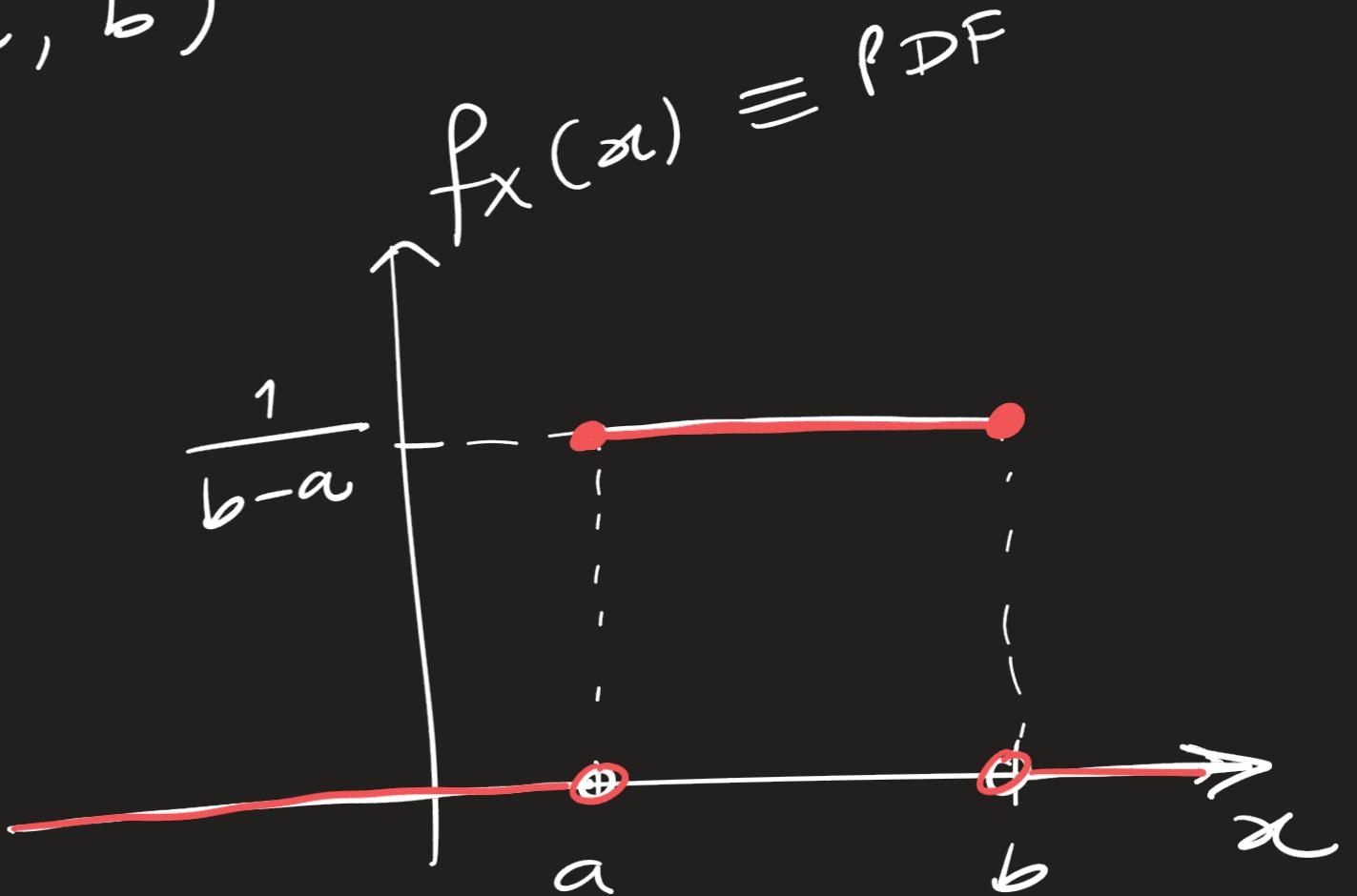
example Uniform  $(-\infty, +\infty)$ .

}  $\rightarrow$  if  $P(\omega)$  is selected  
 incorrectly, MLE will  
 always outperform MAP.  
 } Under  
 But if  $P(\omega)$  is selected  
 correctly then MAP outperforms  
 MLE.  
 small-ish  
 sample  
 sizes.

} it needs a lot more  
 data to catch up  
 with MLE.

$$X \sim \text{Uniform}(a, b)$$

$$f_X(x) = \frac{1}{b-a}$$



$\mu$  = learnable parameter (represents prob.  
of flipping Heads)

$x_i$  = sample =  $\{0, 1\}$ ,  $\{x_i\}_{i=1}^N$  i.i.d.

$$P(x=0|\mu) = 1-\mu$$

$$P(x=1|\mu) = \mu$$

In general, we can write the  
data likelihood for sample

$x$  as:

$$P(x|\mu) = \mu^x \cdot (1-\mu)^{1-x}$$

BERNOULLI  
Distribution

$$X \sim \text{BERNOULLI}(\mu)$$

## MAXIMUM LIKELIHOOD ESTIMATION (MLE)

- ① SET UP THE OBSERVED DATA LIKELIHOOD  
FOR DATASET  $\{x_i\}_{i=1}^N$  i.i.d.

$$\boxed{\begin{aligned} \mathcal{L}^o &= P(\{x_1, x_2, \dots, x_N\} | \mu) \\ &= \prod_{i=1}^N P(x_i | \mu) \end{aligned}}$$

$$\begin{aligned} \text{i.i.d. samples} &= \prod_{i=1}^N \mu^{x_i} \cdot (1-\mu)^{1-x_i} \\ &= \mu^{\sum_{i=1}^N x_i} \cdot (1-\mu)^{\sum_{i=1}^N (1-x_i)} \\ &= \mu^{\sum_{i=1}^N x_i} \cdot (1-\mu)^{N - \sum_{i=1}^N x_i} \end{aligned}$$

Illustration:

$$\begin{aligned} &\mu \cdot (1-\mu) \cdot \mu \\ &= \mu^2 \cdot (1-\mu) \end{aligned}$$

② Log-likelihood:

$$\begin{aligned} \mathcal{L} &= \ln \mathcal{L}^* \\ &= \left( \sum_{i=1}^N x_i \right) \cdot \ln(\mu) + \left( N - \sum_{i=1}^N x_i \right) \cdot \ln(1-\mu) \end{aligned}$$

③ DERIVE w.r.t.  $\mu$ :

$$\frac{\partial \mathcal{L}}{\partial \mu} = 0 \Leftrightarrow \left( \sum_{i=1}^N x_i \right) \cdot \frac{1}{\mu} + \left( N - \sum_{i=1}^N x_i \right) \cdot \left( \frac{-1}{1-\mu} \right) = 0$$

$$\Leftrightarrow (1-\mu) \cdot \left( \sum_{i=1}^N x_i \right) - \mu \cdot \left( N - \sum_{i=1}^N x_i \right) = 0$$

$$\Leftrightarrow \sum_{i=1}^N x_i - \cancel{\mu \cdot \sum_{i=1}^N x_i} - \mu \cdot N + \cancel{\mu \cdot \sum_{i=1}^N x_i} = 0$$

$$\Leftrightarrow \boxed{\mu_{MLE} = \frac{\sum_{i=1}^N x_i}{N}}$$

## Maximum A Posteriori (MAP)

It starts by encoding prior beliefs  
for  $\mu$  in the form of a  
probabilistic model.

$$\mu \sim \text{Beta}(\alpha, \beta) = \frac{\Gamma(\alpha + \beta)}{\Gamma(\alpha) \cdot \Gamma(\beta)} \cdot \mu^{\alpha-1} \cdot (1-\mu)^{\beta-1}$$

$\swarrow$

"constant for  
some  $\alpha$  and  $\beta$ .

continuous  
R.V.  $\Gamma(x) = (x-1)!$

$$\alpha, \beta > 0$$

① OBSERVED DATA LIKELIHOOD  $\times$  PRIOR

$$\mathcal{L}^o = P(\{x_1, \dots, x_N\} | \mu) \cdot P(\mu)$$

$$= \left[ \prod_{i=1}^N P(x_i | \mu) \right] \cdot [P(\mu)]$$

*independent*

$$= \left[ \mu^{\sum_{i=1}^N x_i} \cdot (1-\mu)^{N - \sum_{i=1}^N x_i} \right] \cdot \left[ \frac{\Gamma(\alpha + \beta)}{\Gamma(\alpha) \cdot \Gamma(\beta)} \mu^{\alpha-1} \cdot (1-\mu)^{\beta-1} \right]$$

*identically distributed*

$$\mathcal{L} = \frac{\mu^{\sum_{i=1}^N x_i} \cdot (1-\mu)^{N - \sum_{i=1}^N x_i} \cdot \mu^{\alpha-1} \cdot (1-\mu)^{\beta-1}}{\sum_{i=1}^N x_i + \alpha - 1 \cdot (1-\mu)^{N - \sum_{i=1}^N x_i + \beta - 1}}$$

② DATA Log-likelihood:

$$\begin{aligned} \mathcal{L} &= \ln \mathcal{L}^o \\ &= \left( \sum_{i=1}^N x_i + \alpha - 1 \right) \cdot \ln(\mu) + \left( N - \sum_{i=1}^N x_i + \beta - 2 \right) \cdot \ln(1-\mu) \end{aligned}$$

③ Deriv w.r.t.  $\mu$ :

$$\frac{\partial \mathcal{L}}{\partial \mu} = 0 \quad (\Rightarrow) \quad \mu_{MAP} = \frac{\sum_{i=1}^N x_i + \alpha - 1}{N + \alpha + \beta - 2}$$

In cases where the resulting  
posterior has the same  
probabilistic  
shape as the prior

prob. up + a constant  
factor).

⇒ they are said to  
have a conjugate prior  
relationship.

DATA  
Likelihood  $\times$  PRIOR  $\propto$  POSTERIOR

$$\overline{\mathcal{L}}^{\text{MAP}} = \left[ \mu^{\frac{\sum_{i=1}^n x_i + \alpha - 1}{N - \sum_{i=1}^n x_i + \beta - 1}} \quad (1-\mu)^{\frac{N - \sum_{i=1}^n x_i + \beta - 1}{N - \sum_{i=1}^n x_i + \beta - 1}} \right] \sim \text{Beta}\left(\alpha + \sum_{i=1}^n x_i, \beta + N - \sum_{i=1}^n x_i\right)$$

and prior

$$P(\mu) = C \cdot \left[ \mu^{\alpha-1} \quad (1-\mu)^{\beta-1} \right] \sim \text{Beta}(\alpha, \beta)$$

Constant

Bernoulli  $\times$  Beta  $\propto$  Beta

This forms a conjugate prior relationship

Online update of prior

$$\begin{aligned} \alpha^{(t+1)} &\leftarrow \alpha^{(t)} + \sum_{i=1}^n x_i && (\text{recursive update}) \\ \beta^{(t+1)} &\leftarrow \beta^{(t)} + N - \sum_{i=1}^n x_i \end{aligned}$$

## Pseudocode for Online Prior Update:

$t = 0$  (iteration)

- 1) Start w/ initial guess for prior parameters.

$$\alpha^{(t)}, \beta^{(t)}$$

- 2) Calculate the solution for parameters w/ current prior values:

$$\mu_{MAP} = \frac{\sum_{i=1}^n x_i + \alpha^{(t)} - 1}{N + \alpha^{(t)} + \beta^{(t)} - 2}$$

- 3) Update the prior parameters

$$\alpha^{(t+1)} \leftarrow \alpha^{(t)} + \sum_{i=1}^n x_i$$

$$\beta^{(t+1)} \leftarrow \beta^{(t)} + N - \sum_{i=1}^n x_i$$

- 4)  $t \leftarrow t + 1$ , go back to step ②.

Before   Next class ..

1) Review

Univariate

Bivariate

Multivariate  
(in general)

Gaussian distributions.