

MLE, MAP & BAYESIAN

PART I: COIN FLIP

# Bayes Rule

$$P(A|B)P(B) = P(B|A)P(A)$$

A = Parameters ( $\theta$ )

B = Data (D)

$$P(\theta|D) = \frac{P(D|\theta)P(\theta)}{P(D)}$$

↑ Posterior

← Likelihood

← Prior

← Evidence

# Online Learning

No data

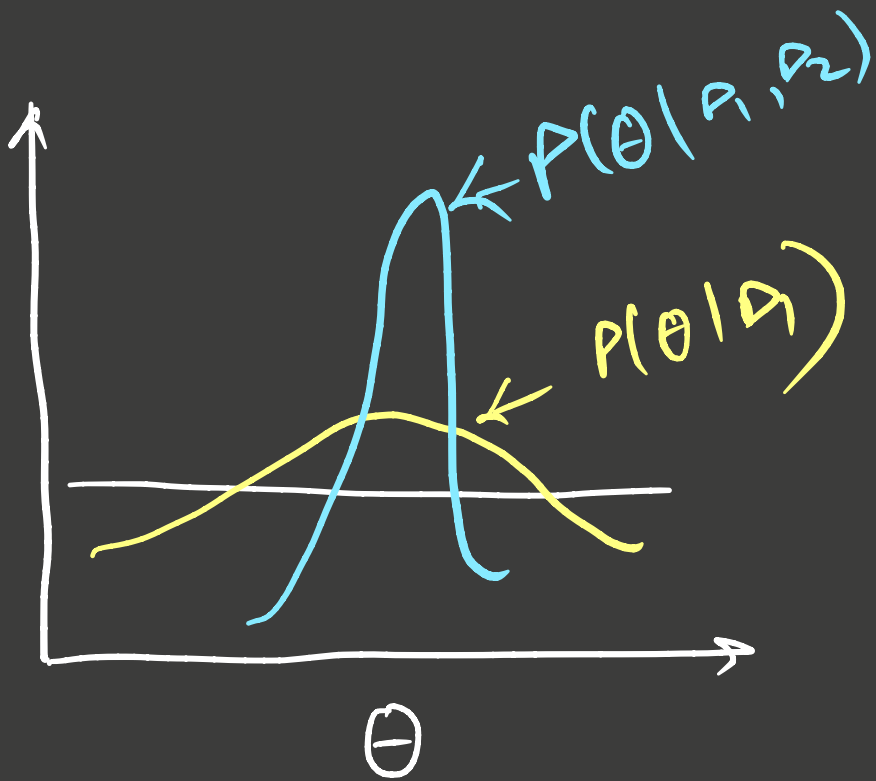
observed  $D_1$

observed  $D_2$

$P(\theta)$

$P(\theta | D_1)$

$$P(\theta | D_2) = \frac{P(D_2 | \theta) \text{Prior}}{P(D_2)}$$



Assume Coin Toss MULTIPLE TIMES

OBSERVATION  $= \{H, H\}$

θ) what is  $p(\text{Head})$ ?

# Maximum Likelihood Estimation (MLE)

$$6H, 4T$$

$$P(H) = .6; P(T) = .4$$

$$n_H, n_T$$

$$P(H) = \frac{n_H}{n_H + n_T}$$

← comes from an MLE estimate

$$\text{let } P(H) = \theta$$

$$\text{Likelihood} = P(D|\theta)$$

$$\hat{\theta}_{MLE} = \underset{\theta}{\operatorname{argmax}} P(D|\theta)$$

To prove

$$\hat{\theta}_{MLE} = \frac{n_H}{n_H + n_T}$$

$$D = n_H, n_T \quad \text{where } D = \{D_1, \dots, D_m\} \leftarrow \{H, T\}$$

$$P(H) = \theta \Rightarrow P(T) = 1 - \theta$$

$$P(D_1, D_2, \dots, D_m | \theta) = P(D_1 | \theta) P(D_2 | \theta) \dots P(D_m | \theta)$$

$$\theta) \text{ what is } P(D_1 | \theta)? = \begin{cases} P(D_1 = H) \text{ given } \theta & \text{if } D_1 = H = \theta \\ P(D_1 = T) \text{ given } \theta & \text{if } D_1 = T = 1 - \theta \end{cases}$$

$$P(D|\theta) = (P(\text{head}))^{n_H} (P(\text{tail}))^{n_T} \leftarrow \text{Bernoulli}$$

$$P(D|\theta) = \theta^{n_H} (1-\theta)^{n_T}$$

(Not  
binomial)

$$\text{Log-likelihood (LL)} = \log P(D|\theta) = n_H \log \theta + n_T \log(1-\theta)$$

$$\frac{\partial LL}{\partial \theta} = 0 \Rightarrow \frac{n_H}{\theta} + \frac{n_T}{1-\theta} = 0$$

$$\Rightarrow \hat{\theta}_{MLE} = \frac{n_H}{n_H + n_T}$$

6)  $ZH, 0T$

$$P(H) = 1; P(T) = 0$$



# Maximum A posteriori (MAP)

\* MLE doesn't have a notion of prior knowledge

\* Overfitting

MAP overcomes these

$P(\theta)$  : Prior probability model.

$$P(\theta|D) = \frac{P(D|\theta)P(\theta)}{P(D)}$$

Posteriori

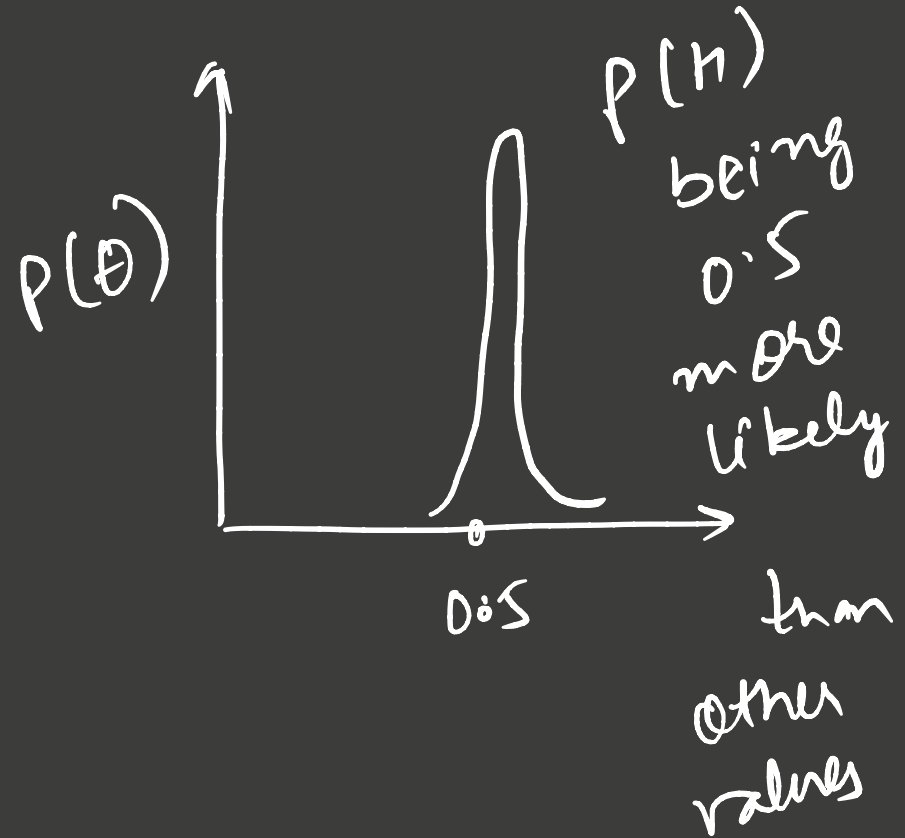
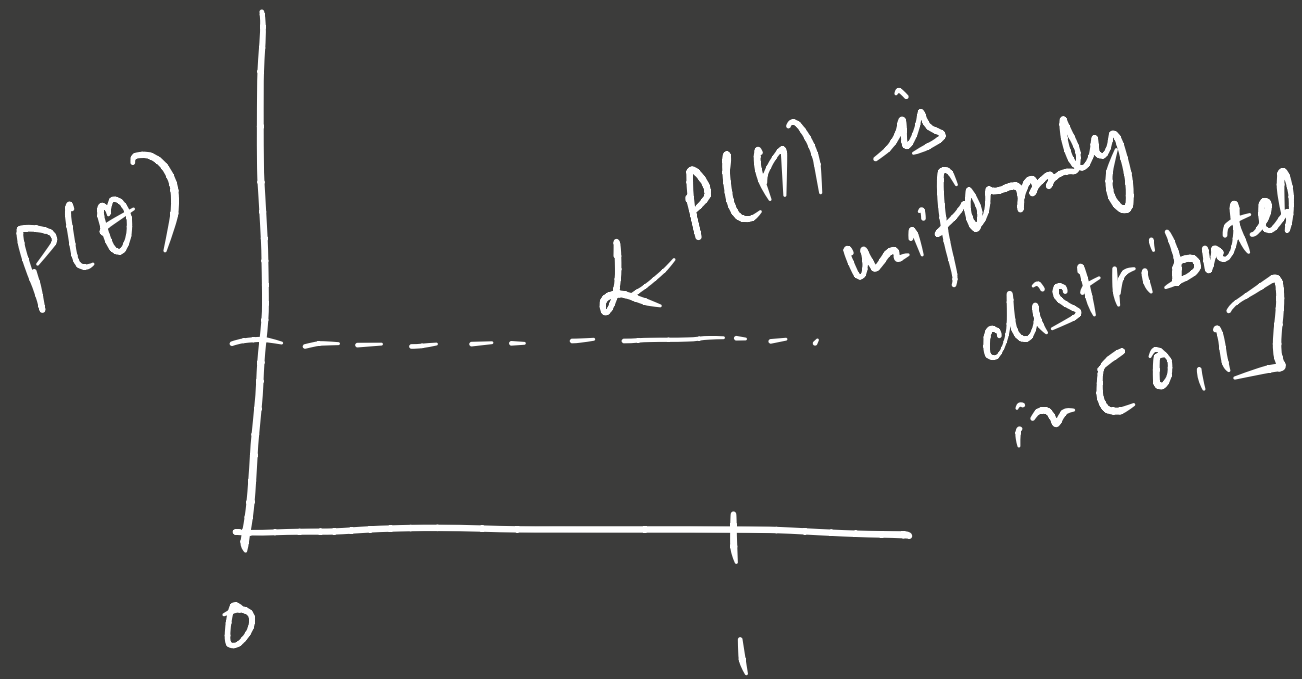
Goal for MAP

$$\hat{\theta}_{\text{MAP}} = \underset{\theta}{\operatorname{argmax}} P(\theta|D)$$

$$\hat{\theta}_{\text{MAP}} = \underset{\theta}{\operatorname{argmax}} P(D|\theta)P(\theta)$$

# Some examples of $P(\theta)$

$$P(H) = \theta$$



# Beta distribution

$$\text{Beta}(\theta | a, b) = \frac{\Gamma(a+b)}{\Gamma(a)\Gamma(b)} \theta^{a-1} (1-\theta)^{b-1}$$



Similar to

$$\theta^{nn} (1-\theta)^{nT}$$

$$\Gamma(n) = (n-1)!$$

(Natural nos)

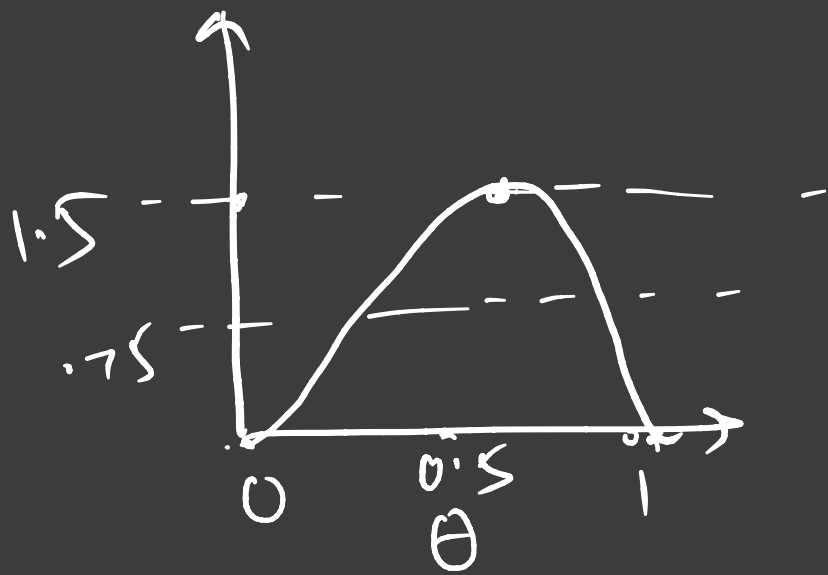
$$\Gamma(z) = \int_0^{\infty} x^{z-1} e^{-x} dx$$

$$\text{Beta}(\theta | 1, 1) = \frac{\Gamma(1+1)}{\Gamma(1)\Gamma(1)} \theta^{1-1} (1-\theta)^{1-1}$$

$$= 1$$



$$\text{Beta}(\theta | 2, 2) = \frac{\Gamma(4) \theta(1-\theta)}{\Gamma(2)\Gamma(2)} = 6\theta(1-\theta)$$



$$\text{Beta}(\theta | a=2, b=1)$$

$\Rightarrow$  Indicate higher probability  
of heads  
compared to  
tails.

$$D = n_H, n_T$$

$$P(\theta) = \text{Beta}(\theta | a, b) = \frac{\Gamma(a+b) \theta^{a-1} (1-\theta)^{b-1}}{\Gamma(a)\Gamma(b)}$$

$$\theta). \hat{\theta}_{MAP} = \underset{\theta}{\text{argmax}} P(D | \theta) P(\theta)$$

$$= \underset{\theta}{\text{argmax}} \theta^{n_H} (1-\theta)^{n_T} \theta^{a-1} (1-\theta)^{b-1} \times K$$

↓  
constant

$$= \underset{\theta}{\text{argmax}} \theta^{n_H+a-1} (1-\theta)^{n_T+b-1}$$

$$\hat{\theta}_{MAP} = \frac{n_H + a - 1}{n_H + n_T + a + b - 2}$$

## Conjugate Prior

$$P(\theta|D) = \frac{P(D|\theta)P(\theta)}{P(D)}$$

$P(\theta)$  is CONJUGATE TO  $P(D|\theta)$

if.

$P(\theta|D)$  and  $P(\theta)$  are from  
same distribution family.

Bernoulli likelihood, gamma is conjugate



# Relationship b/w MAP & MLE

$\hat{\theta}_{MAP}$

$\hat{\theta}_{MLE}$

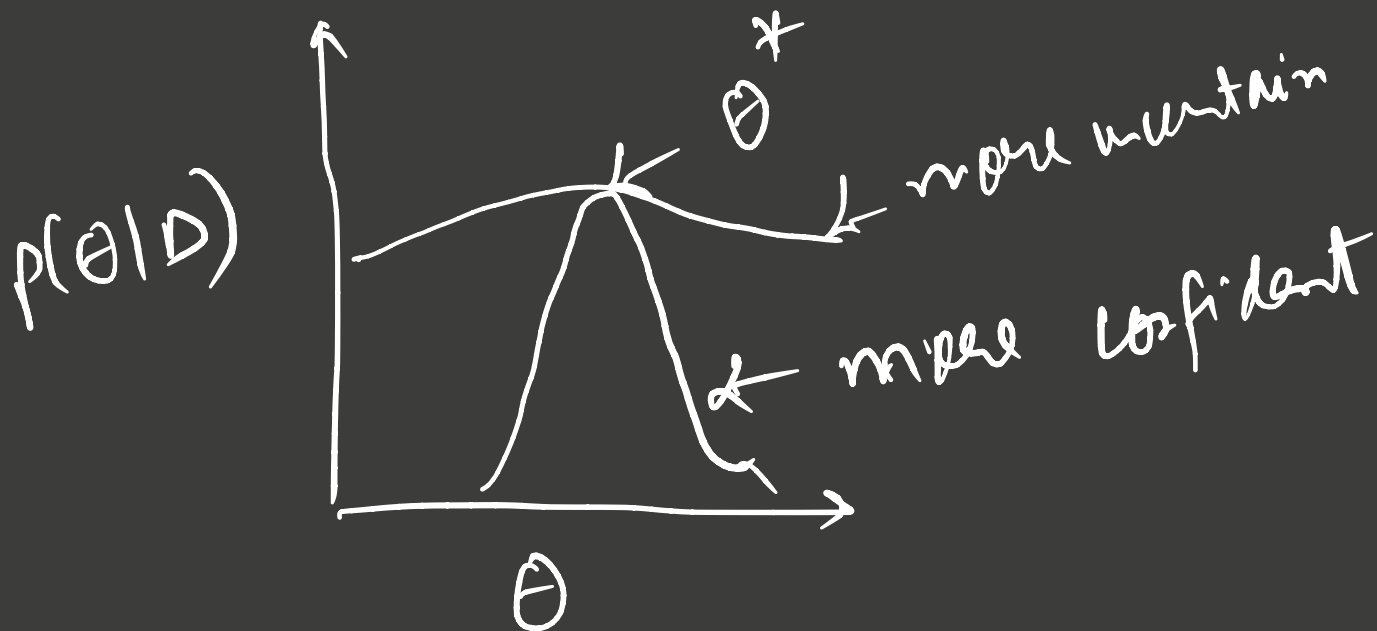
$$\hat{\theta}_{MAP} = \underset{\theta}{\operatorname{argmax}} P(D|\theta) P(\theta) \rightarrow \text{if this is uniform}$$

$$= \hat{\theta}_{MLE} \text{ (for uniform prior)}$$

# Fully Bayesian

\* MLE & MAP

don't give you uncertainty  
(or full distribution)



Predictive distribut<sup>2</sup>

$P(\text{Next coin} = H \mid \text{Data})?$

(1) What  $\theta$  to use?

use all possible values of  $\theta$

$$P(\text{Next} = H \mid D) = \int P(\text{Next} = H, \theta \mid D) d\theta$$

$$H \rightarrow 1$$

$$T \rightarrow 0$$

$$P(\text{Next} = c \mid \theta) = \theta^c (1-\theta)^{1-c}$$

or  
 $P(c \mid \theta)$

↑ why this

$$L = 0$$

$$P(\text{tails} \mid \theta) = \theta^0 (1-\theta)^1$$

$$P(\text{Next} = c \mid D, a, b) = \int P(\text{Next} = c, \theta \mid D, a, b) d\theta$$

~  
Beta

distribut<sup>2</sup>  
for prior

↑ why?

SUM RULE

$$\therefore P(a) = \int P(a, y) dy$$

$$= \int P(N_{\text{ent}} = c \mid \theta) P(\theta \mid D, a, b) d\theta$$

↑ why? ∵ once  $\theta$  is known,  $D, a, b$  don't influence  $P(N_{\text{ent}} = c)$

$$= \int \theta^c (1-\theta)^{1-c} \frac{\Gamma(n_H + n_T + a + b)}{\Gamma(n_H + a) \Gamma(n_T + b)} \theta^{n_H + a - 1} (1-\theta)^{n_T + b - 1} d\theta$$

$$= \frac{\Gamma(n_H + n_T + a + b)}{\Gamma(n_H + a) \Gamma(n_T + b)} \int \theta^{c + n_H + a - 1} (1-\theta)^{n_T + b - c} d\theta$$

$$= \frac{\Gamma(n_H + n_T + a + b) \Gamma(c + n_H + a) \Gamma(n_T + b - c + 1)}{\Gamma(n_H + a) \Gamma(n_T + b) \Gamma(1 + n_H + a + n_T + b)}$$