

# *MLE & MAP*

**INSTRUCTOR: HONGJIE CHEN**  
**MAY 26TH 2022**

# *Deal with Data Sparsity*

- Estimate probabilities from sparse data
  - Maximum Likelihood Estimation (MLE)
  - Maximum A Posteriori estimation (MAP)
- Represent joint probability distribution other than using tables
  - Graphical models

# *A Toy Example to Explain MLE & MAP*

- Given a coin, estimate the probability that it turns up heads ( $X = 1$ ) or tails ( $X = 0$ )
- Test A: 3 flips, 2 heads ( $X = 1$ ), 1 tail ( $X = 0$ )
- Test B: 100 flips, 51 heads ( $X = 1$ ), 49 tails ( $X = 0$ )
- Test C: Keep flipping and develop an online learning algorithm that gives reasonable estimate for each flip

# Estimating Probabilities

- Maximum Likelihood Estimation (**MLE**)

- Choose parameter  $\theta$  that maximizes  $P(\text{data} | \theta)$

- $$\hat{\theta} = \frac{\alpha_1}{\alpha_1 + \alpha_0}$$

- Maximum A Posteriori estimation (**MAP**)

- Choose parameter  $\theta$  that maximizes  $P(\theta | \text{data})$

- $$\hat{\theta} = \frac{\alpha_1 + \beta_1 - 1}{(\alpha_1 + \beta_1 - 1) + (\alpha_0 + \beta_0 - 1)}$$

**Confusing?**  
**Prior:  $P(\theta)$**

# MLE - $P(\text{data} \mid \theta)$

- $\alpha_1$ : # head up,  $\alpha_0$ : # tail up
- $P(X = 1) = \theta$ ,  $P(X = 0) = 1 - \theta$
- Data  $D$ : 1, 0, 0, 1, 1
- $P(D \mid \theta)$ :  $\theta, 1 - \theta, 1 - \theta, \theta, \theta = \theta^{\alpha_1}(1 - \theta)^{\alpha_0}$  **i.i.d.~Bernoulli**
  
- Our goal:  $\hat{\theta} = \operatorname{argmax}_{\theta} P(D \mid \theta)$
- How to obtain the extreme value for a function?

# Derivative

- Our goal:  $\hat{\theta} = \operatorname{argmax}_{\theta} P(D | \theta)$
- $P(D | \theta) = \theta^{\alpha_1} (1 - \theta)^{\alpha_0}$
- Calculate  $\theta$  from  $\frac{d}{d\theta} P(D | \theta) = 0$  (Show me the maths!)
  - An easier route with  $\ln P(D | \theta)$
- $\hat{\theta}_{\text{MLE}} = \frac{\alpha_1}{\alpha_1 + \alpha_0}$  **Questions: How to get a more accurate  $\hat{\theta}_{\text{MLE}}$ ?**  
**Think about # flips**

# Prepare for MAP

- Recall Bayes theorem

- $$P(A | B) = \frac{P(B | A)P(A)}{P(B)}$$

- $$P(\theta | \text{data}) = \frac{P(\text{data} | \theta)P(\theta)}{P(\text{data})}$$

# MAP

- Bayes theorem

$$\underbrace{P(\theta | \text{data})}_{\text{Posteriori}} = \frac{\overbrace{P(\text{data} | \theta)}^{\text{Likelihood}} \underbrace{P(\theta)}_{\text{Prior}}}{P(\text{data})} \quad \text{No variable in data}$$

- Equivalently,  $P(\theta | \text{data}) \propto P(\text{data} | \theta)P(\theta)$



# Prior Probability: $P(\theta)$

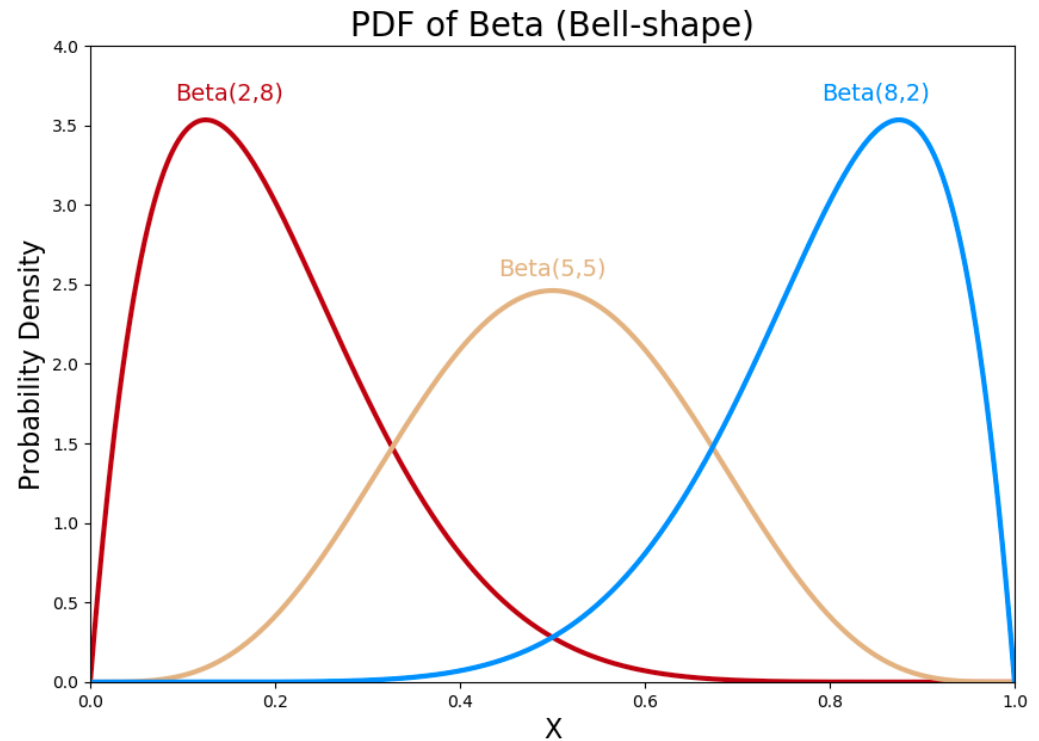
- Let's review what  $\theta$  and  $P(\theta)$  mean
  - $\theta$ : the probability a coin turns head up in one flipping
  - $P(\theta)$ : the distribution of this probability parameter

# Beta Distribution

- We choose Beta distribution as the prior

- $$P(\theta) = \frac{\theta^{\beta_1-1}(1-\theta)^{\beta_0-1}}{B(\beta_1, \beta_0)} \sim \text{Beta}(\beta_1, \beta_0)$$

**A normalization constant**



# Why Beta Distribution as Prior

- Reasonable, and easy to compute

$$P(\theta) = \frac{\theta^{\beta_1-1}(1-\theta)^{\beta_0-1}}{B(\beta_1, \beta_0)}, \quad P(D|\theta) = \theta^{\alpha_1}(1-\theta)^{\alpha_0}$$

$$\begin{aligned} P(\theta|D) &\propto P(D|\theta)P(\theta) \\ &= \frac{\theta^{\alpha_1+\beta_1-1}(1-\theta)^{\alpha_0+\beta_0-1}}{B(\alpha_1+\beta_1, \alpha_0+\beta_0)} \Rightarrow \hat{\theta}_{\text{MAP}} = \operatorname{argmax}_{\theta} P(\theta|D) = \frac{\alpha_1 + \beta_1 - 1}{(\alpha_1 + \beta_1 - 1) + (\alpha_0 + \beta_0 - 1)} \end{aligned}$$

As  $N \rightarrow \infty$ , prior is dominated when  $N = (\alpha_0 + \alpha_1) \gg \beta_0 + \beta_1$

# *Questions about Prior*

- Are other prior probability distributions possible?
- How to set parameters for the prior, e.g., for Beta distribution?

# Thoughts

- Are other prior probability distributions possible?

**Yes. substitute  $P(\theta)$  and calculate derivative to get  $\hat{\theta}_{\text{MAP}}$**

- How to set parameters for the prior, e.g., for Beta distribution?

**Reparametrize  $\beta_1, \beta_0$  to other parameters such as mean and variance  $\mu, \sigma^2$ , conduct pilot experiments to obtain  $\mu$  and  $\sigma^2$**

$$\nu := \beta_1 + \beta_0 = \frac{\mu(1 - \mu)}{\sigma^2} - 1, \quad \begin{cases} \beta_1 = \mu\nu \\ \beta_0 = (1 - \mu)\nu \end{cases}$$

**New parameters to estimate :(**

# Conjugate Prior

- If the posterior distribution  $P(\theta | D)$  and the prior distribution  $P(\theta)$  are in the same probability distribution family.

Likelihood:  $P(D | \theta) = \theta^{\alpha_1}(1 - \theta)^{\alpha_0} \sim \text{Binomial}(\alpha_1, \alpha_0)$

Prior:  $P(\theta) = \frac{\theta^{\beta_1-1}(1 - \theta)^{\beta_0-1}}{B(\beta_1, \beta_0)} \sim \text{Beta}(\beta_1, \beta_0)$

Posterior:  $P(\theta | D) = \frac{\theta^{\alpha_1+\beta_1-1}(1 - \theta)^{\alpha_0+\beta_0-1}}{B(\alpha_1 + \beta_1, \alpha_0 + \beta_0)} \sim \text{Beta}(\alpha_1 + \beta_1, \alpha_0 + \beta_0)$

**For Binomial likelihood, conjugate prior is Beta distribution**

# Conjugate Prior for Multinomial

- Let's use a dice instead of a coin ( $k$  outcomes instead of 2)

$$\text{Likelihood: } P(D | \theta) = \theta_1^{\alpha_1} \theta_2^{\alpha_2} \dots \theta_k^{\alpha_k} \sim \text{Multinomial}(\theta_1, \theta_2, \dots, \theta_k)$$

$$\text{If the prior is a Dirichlet distribution: } P(\theta) = \frac{\prod_{i=1}^k \theta_i^{\beta_i - 1}}{B(\beta_1, \dots, \beta_k)} \sim \text{Dirichlet}(\beta_1, \dots, \beta_k)$$

Then the posterior is a Dirichlet distribution:

$$P(\theta | D) \sim \text{Dirichlet}(\alpha_1 + \beta_1, \dots, \alpha_k + \beta_k)$$

**For Multinomial likelihood, conjugate prior is Dirichlet distribution**