

Limitation of MLE

- Over-reliance on data sample D . If data is limited, estimates can be very wrong.
 - Example, Bernoulli p could be zero if no 1s in 10 trials.
- No indication on the uncertainty of the estimated parameters.
 - Example, for a Bernoulli parameters whether estimation is made from two with 50% heads or 1000 examples with 50% heads, the estimated parameter is the same.
 D_1 $|D_1| = 2$ $N_1(D_1) = 1$, $N - N = 1$ $\hat{p}_1 = 0.5$
 $\rightarrow D_2$ $|D_2| = 1000$ $N_1(D_2) = 500$ $\hat{p}_2 = 0.5$
- No mechanism to specify human's prior knowledge of the parameters.

Example of limitations of MLE

- Suppose a toss a coin 10 times and get

H, H, H, H, H, H, H, H, H, H

Estimate p? MLE: $\hat{p} = \# \text{ ones} / N = 1$

What is your guess on the probability p of head?

- Suppose you want to form a music band, and you are looking for bass guitarist. You ask 7 random batchmates: "Can you play the bass guitar" and you get answers

N, N, N, N, N, N, N D

What fraction of batchmates play bass guitar?

MLE: 0

Do you have a different guess?

$p \approx 0.01$

0.01

Bayesian estimation

- Treat the parameters as a random variable which has a distribution.
- Step 1: Humans specify their prior knowledge of the values of the parameters as a distribution $f_{pr}(\theta)$
 - Example: $f_{pr}(\theta) \sim U(0,1)$ where θ denotes the parameter p of a Bernoulli
 - Example for Gaussian:

Temperature of

CPU on your laptop $T \sim G(\theta, \sigma^2)$

$$\rightarrow f_{pr}(\theta) \sim N(\underline{30}, \underline{10})$$

Also called prior probability

Bayesian estimation

- Calculate the posterior distribution of parameters after observing data D following Bayes rule

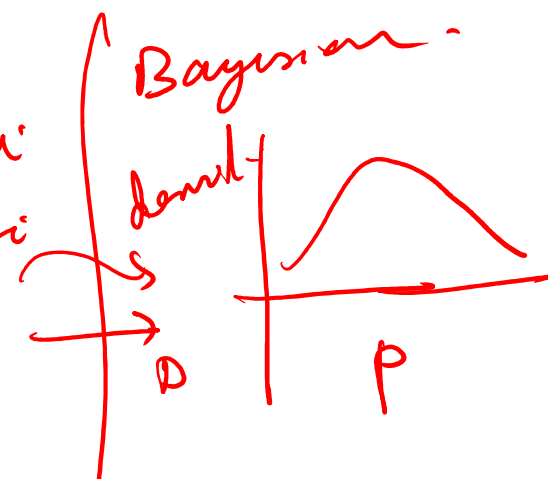
$$\cancel{f(D|\theta) = f(\theta)f(D|\theta) / \int_{\theta} f(\theta)f(D|\theta)}$$

$$f_{\text{post}}(\theta|D) = \frac{f(D|\theta)f_{\text{pr}}(\theta)}{\int_{\theta'} f(D|\theta')f_{\text{pr}}(\theta')} \quad \checkmark$$

Posterior probability

MLE p of Bernoulli

$$D \rightarrow \hat{p} = 0.6$$



Using Bayesian estimates

$$f(\theta|D) \equiv f_{\theta}(\theta|D)$$

- Exact Bayesian probability computation:

- Given a new x, calculate $f(X|D)$

$$f(x|D) = \int_{\theta} \underbrace{f(x|\theta)}_{\text{Binomial eg.}} f_{\theta}(\theta|D) d\theta$$

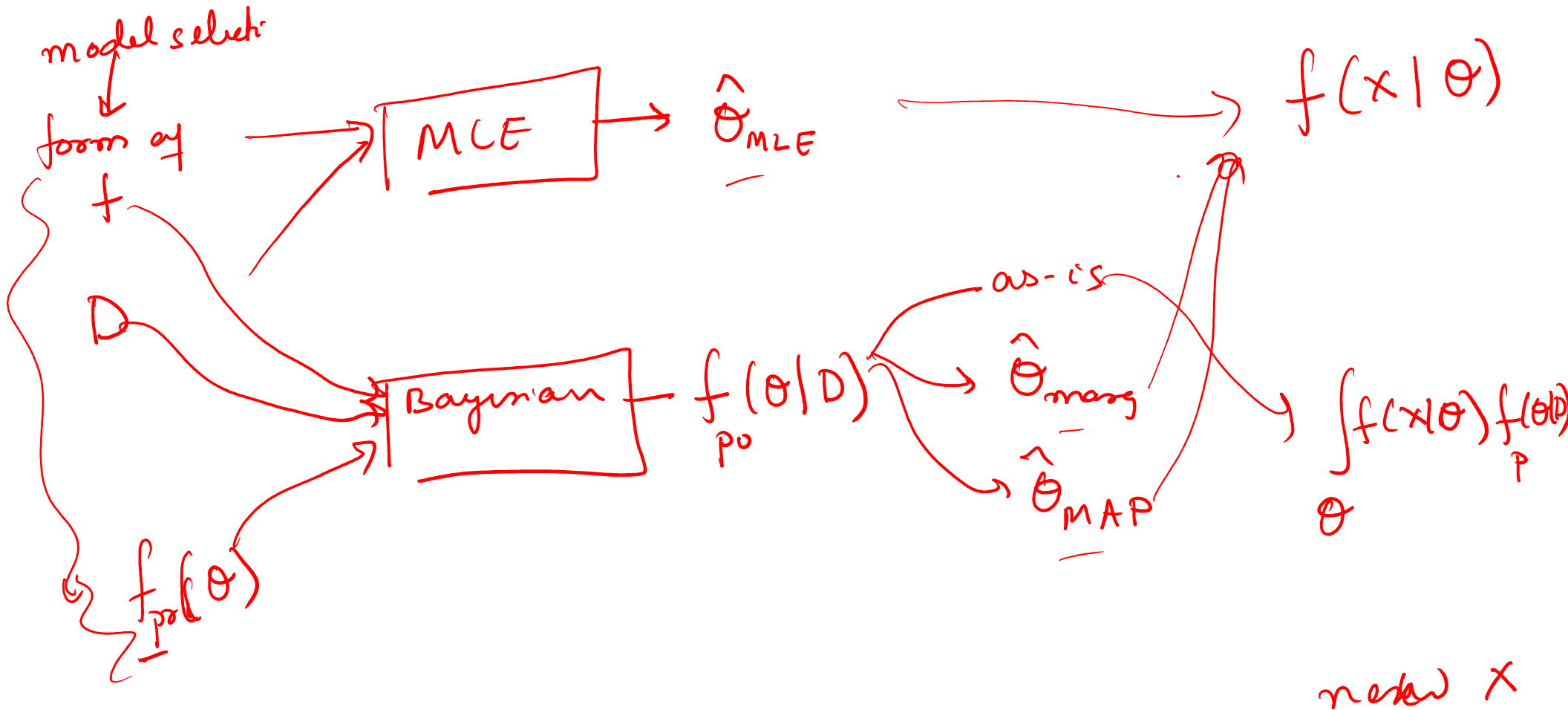
- Expected value of parameters: calculate expected value of $f(\theta|D)$

$$\hat{\theta}_{\text{marg}} = E[\theta]_{f_{\theta}(\theta|D)} = \int \theta \cdot f_{\theta}(\theta|D) d\theta$$

- MAP estimate: use $\max_{\theta} f(\theta|D)$

$$\hat{\theta}_{\text{MAP}} = \arg \max_{\theta} f(\theta|D)$$

Overall pipeline for MLE Vs Bayesian



Example:
Bayesian estimation of
Bernoulli/Binomial parameter p

Bayesian estimation of Bernoulli parameter

$$p \in [0, \dots, 1]$$

- Choose a prior distribution over parameter θ or p of Bernoulli

- $f_{\text{pr}}(\theta) \sim U(0,1)$

- Data D has n are ones and remaining $N-n = \underline{m}$ are 0s.
 $D \equiv \{x_1, x_2, \dots, x_N\}$ eg: $\{0, 1, 1, \dots, 0, \dots, 0\}$
 $D \equiv \{n, m\}$
 $\begin{matrix} \uparrow & \uparrow \\ \text{\#1s} & \text{\#0s} \end{matrix}$
 $f(D|\theta) \equiv \prod_{i=1}^n f(x_i|\theta) = \theta^n (1-\theta)^{\underbrace{N-n}_m}$

- Posterior distribution is:

$$\underbrace{f_{\text{po}}(\theta|D)} = \frac{\underbrace{f(\theta)}_{\text{pr}} f(D|\theta)}{\int_{\theta'} \underbrace{f(\theta') f(D|\theta')}_{\text{pr}}} = \frac{1 \cdot \theta^n (1-\theta)^m}{\int_{\theta'} \theta'^n (1-\theta')^m}$$

$$\underline{f_{po}(\theta|D)} = \frac{(1-\theta)^m (\theta)^n}{Z} \quad 0 \leq \theta \leq 1$$

$Z \leftarrow \text{normalizer.}$

mode of $f_{po}(\theta|D)$

$$\max_{\theta} (1-\theta)^m \theta^n$$

$$\equiv \frac{\partial}{\partial \theta} [(1-\theta)^m \theta^n] = -m(1-\theta)^{m-1} \theta^n + n(1-\theta)^m \theta^{n-1} = 0$$

$$\Rightarrow -m\theta + n(1-\theta) = 0$$

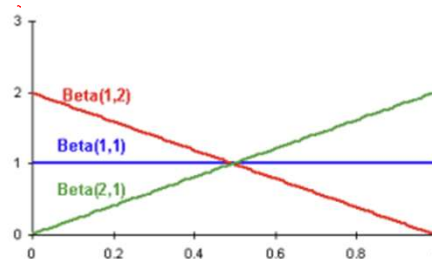
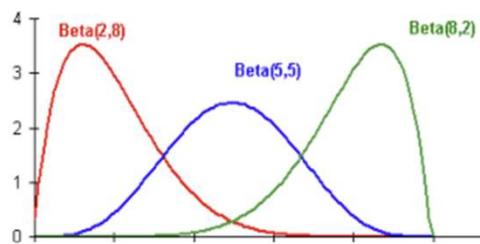
$$\theta(n+m) = n \Rightarrow \theta = \frac{n}{m+n}$$

Beta Random Variable *~ (Generic defn. of Beta)*

X is a **Beta Random Variable**: $X \sim \text{Beta}(a, b)$

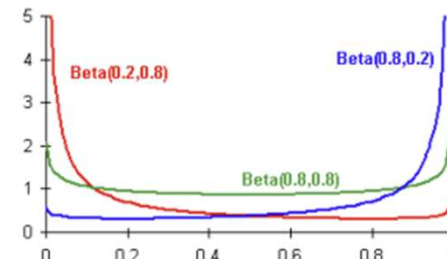
- Probability Density Function (PDF): (where $a, b > 0$)

$$f(x) = \begin{cases} \frac{1}{B(a,b)} x^{a-1} (1-x)^{b-1} & 0 < x < 1 \\ 0 & \text{otherwise} \end{cases}$$



$$B(a,b) = \int_0^1 x^{a-1} (1-x)^{b-1} dx$$

$$B(z_1, z_2) = \frac{\Gamma(z_1) \Gamma(z_2)}{\Gamma(z_1 + z_2)}$$



$z_1, z_2 \equiv \text{integers}$
$$\frac{(z_1-1)! (z_2-1)!}{(z_1+z_2-1)!}$$

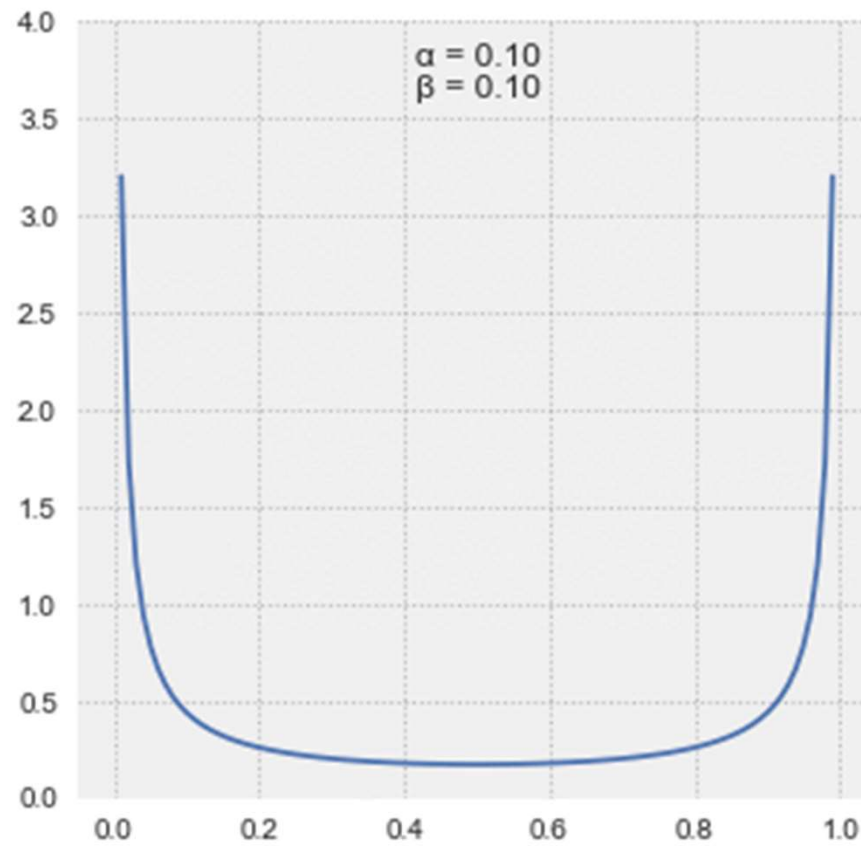
- Symmetric when $a = b$

$$E[X] = \frac{a}{a+b}$$

$$\text{Var}(X) = \frac{ab}{(a+b)^2 (a+b+1)}$$

Mode: $\frac{a-1}{a+b-1}$

The shapes of the Beta distribution



Beta distribution is the
distribution of probabilities

.

More properties of Beta distributions

- Uniform distribution $U(0,1) = B(1,1)$

$$f(\theta | a, b) = \frac{\theta^{a-1} (1-\theta)^{b-1}}{B(a, b)} = \frac{1 \cdot 1}{B(1, 1)} = 1$$

- Relationship between Beta and Gamma distribution

- Let $Y = G(a, 1)$ and $W = G(b, 1)$

$$f(y | a, 1) = \frac{e^{-y} y^{a-1}}{\Gamma(a)}$$

$$f(x | \alpha, \lambda) = \frac{\lambda e^{-\lambda x} (\lambda x)^{\alpha-1}}{\Gamma(\alpha)}$$
$$f(w | b, 1) = \frac{e^{-w} w^{b-1}}{\Gamma(b)}$$

- The $X = Y/(Y+W)$ follows a Beta distribution $B(a, b)$

$$X = \frac{Y}{Y+W} \quad \text{then} \quad X \sim B(a, b)$$

Expected value of the posterior of Binomial

$$f_{p_0}(\theta | D) \equiv \frac{\theta^n (1-\theta)^m}{Z} \equiv B(a=n+1, b=m+1)$$

$$D \equiv \begin{matrix} \{n, m\} \\ \uparrow \quad \uparrow \\ \#1s \quad \#0s \end{matrix}$$

$$\hat{\theta}_{\text{marg}} = \frac{n+1}{n+m+2}$$

Laplace smoothing.

Bass guitar example:

$$\hat{\theta}_{\text{marg}} = \frac{1}{9}$$

Contrast with MLE

$$\hat{\theta}_{\text{MLE}} \equiv \frac{n}{m+n}$$

$$\hat{\theta}_{\text{MLE}} = \frac{0}{7}$$