

MAXIMUM A-POSTERIORI ESTIMATION

ALEKS JOVCIC

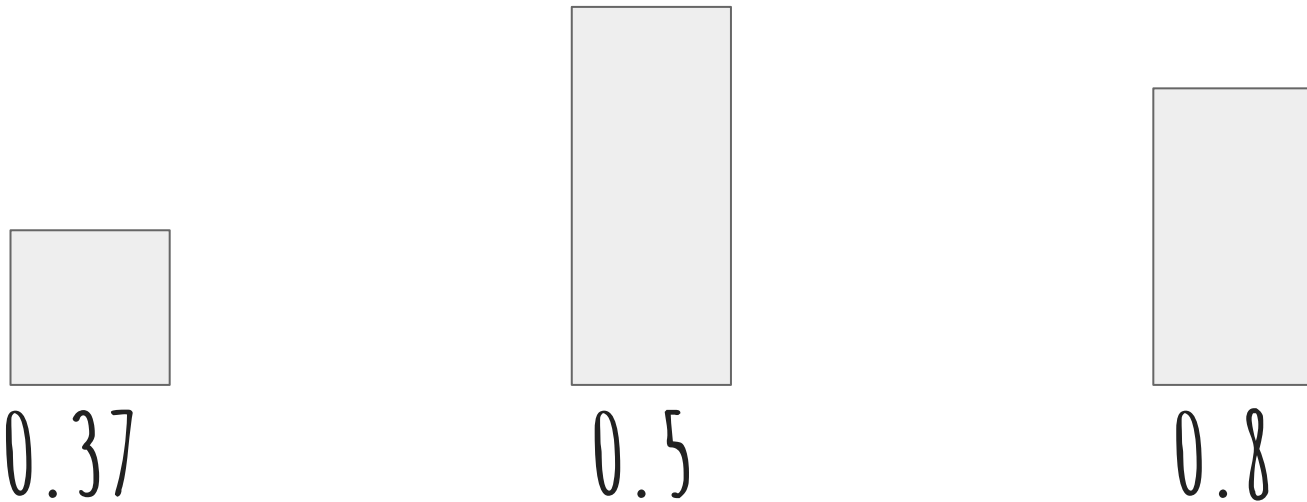
SLIDES BY JOSHUA FAN & ALEX TSUN

AGENDA

- THE BETA RANDOM VARIABLE
- MAP ESTIMATION
- MAP EXAMPLE

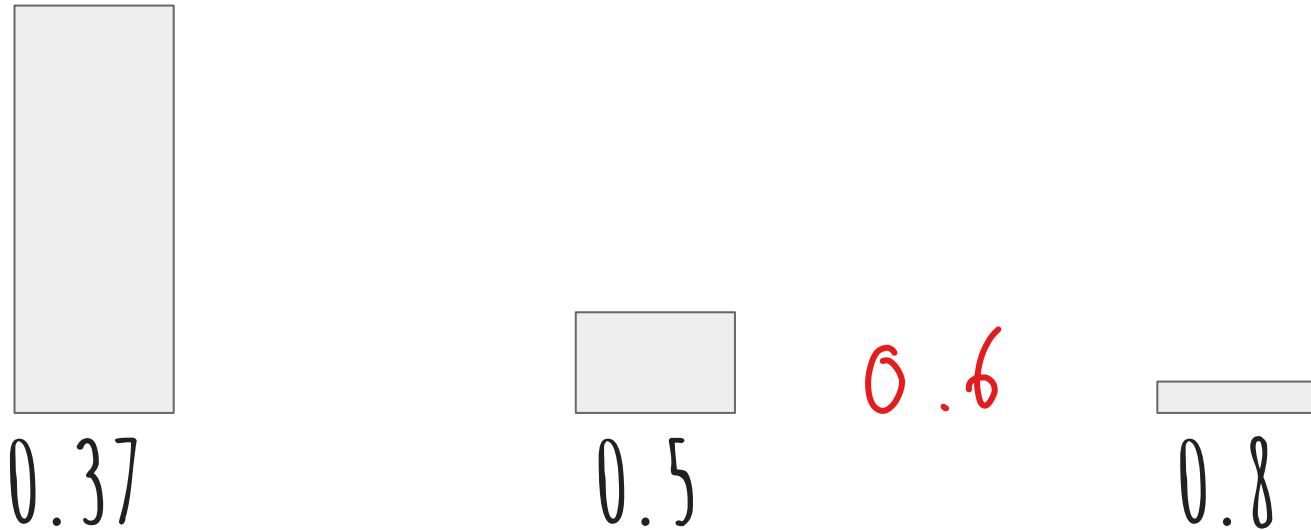
BETA RANDOM VARIABLE (INTUITION)

Suppose you want to model your belief on the unknown probability X of heads. You could assign a probability distribution as follows:



BETA RANDOM VARIABLE (INTUITION)

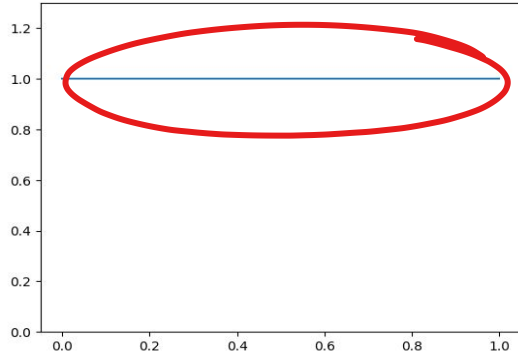
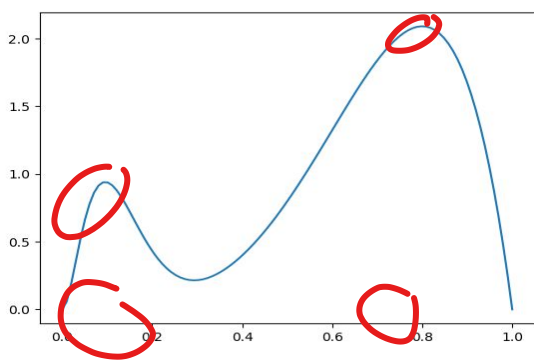
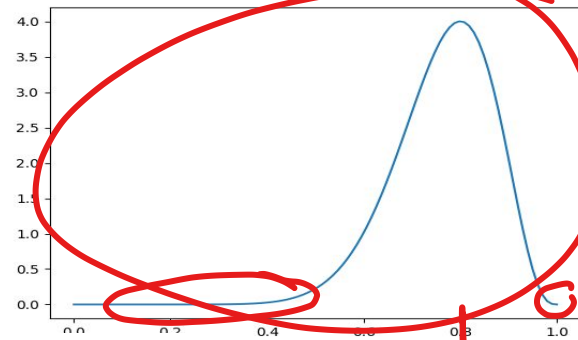
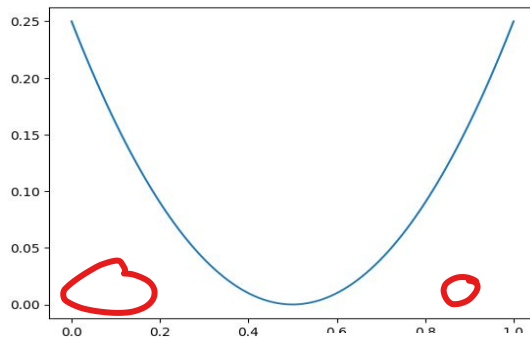
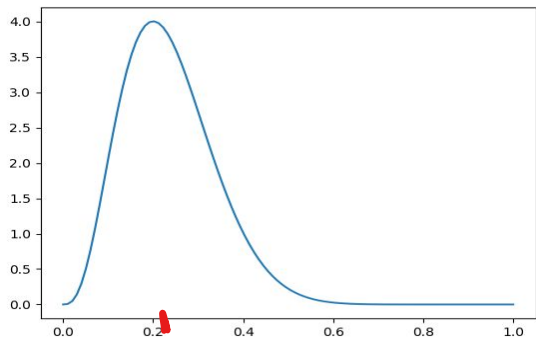
Suppose you want to model your belief on the unknown probability X of heads. You could assign a probability distribution as follows:



BETA RANDOM VARIABLE (INTUITION)

What if you wanted to be open to any value in $[0,1]$?

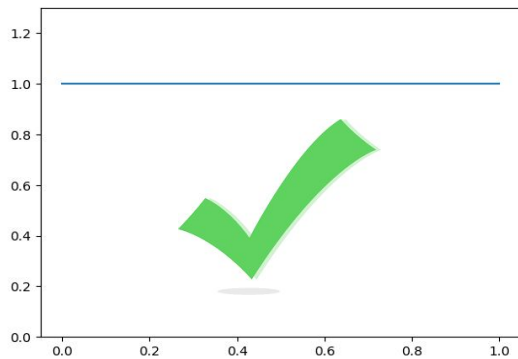
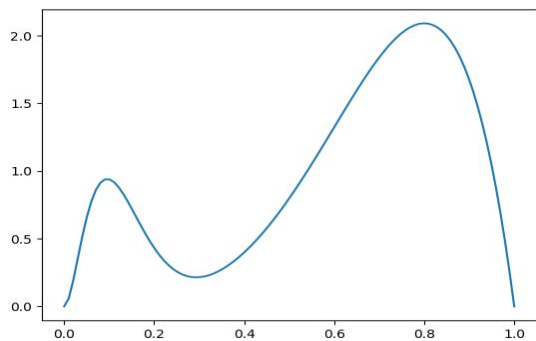
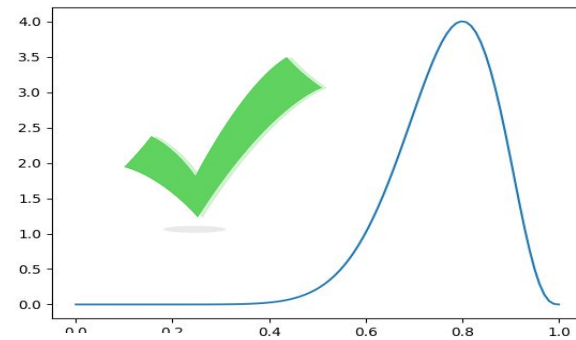
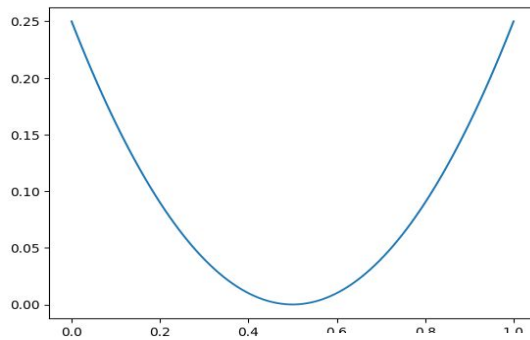
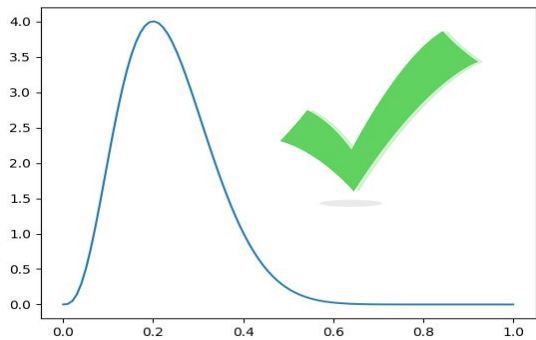
Need to use a *continuous* RV (with range $[0,1]$)!



BETA RANDOM VARIABLE (INTUITION)

What if you wanted to be open to any value in $[0,1]$?

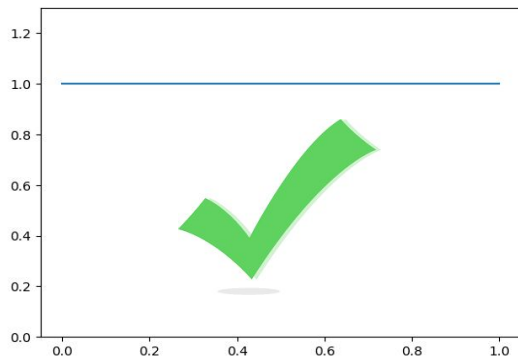
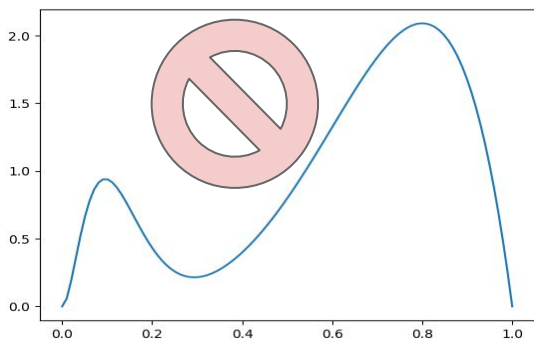
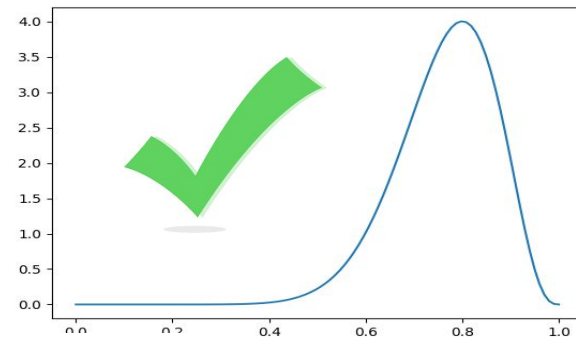
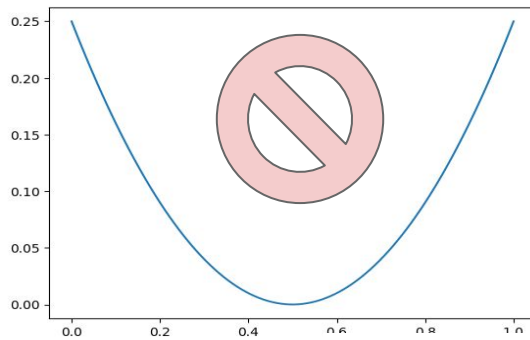
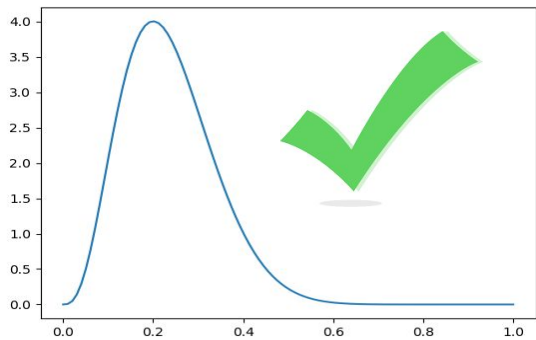
Need to use a *continuous* RV (with range $[0,1]$)!



BETA RANDOM VARIABLE (INTUITION)

What if you wanted to be open to any value in $[0,1]$?

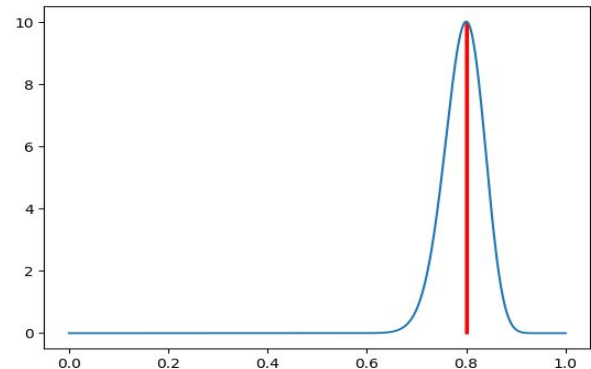
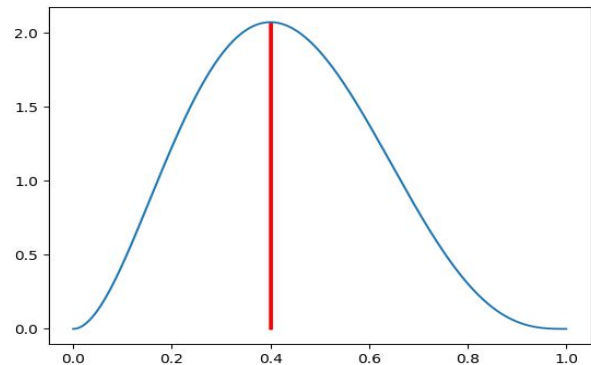
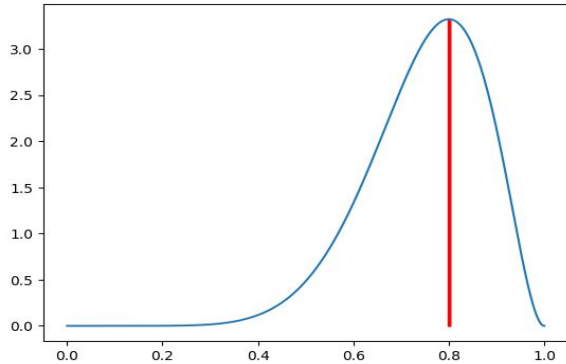
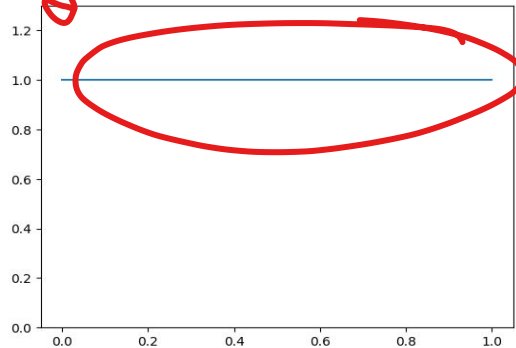
Need to use a *continuous* RV (with range $[0,1]$)!



BETA RANDOM VARIABLE (INTUITION)

If you flip a coin with unknown probability of heads X , what does your belief distribution look like if:

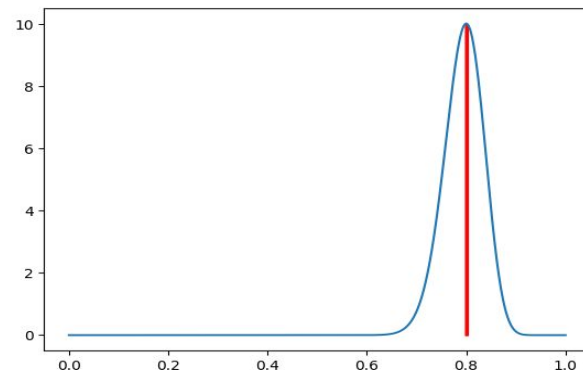
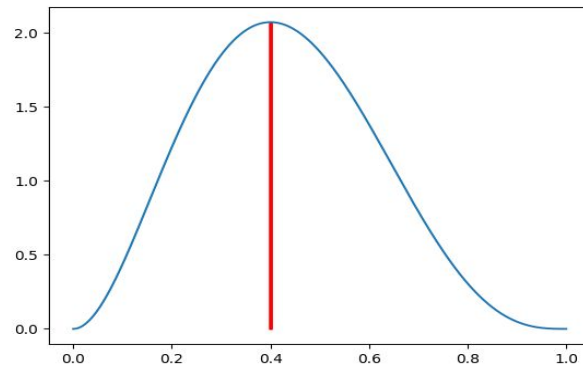
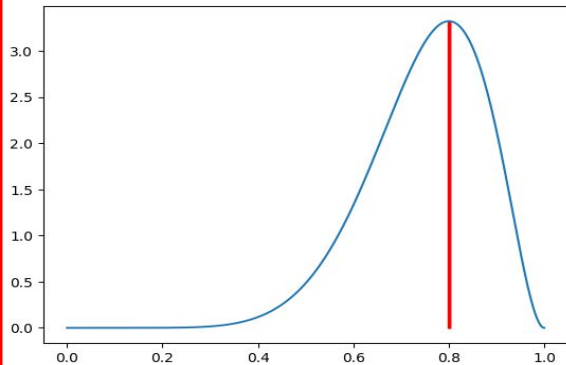
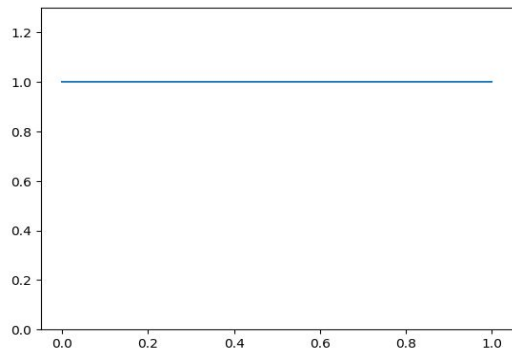
- You didn't observe anything?
- You observed 8 heads and 2 tails?
- You observed 80 heads and 20 tails?
- You observed 2 heads and 3 tails?



BETA RANDOM VARIABLE (INTUITION)

If you flip a coin with unknown probability of heads X , what does your belief distribution look like if:

- You didn't observe anything?
- You observed 8 heads and 2 tails?
- You observed 80 heads and 20 tails?
- You observed 2 heads and 3 tails?

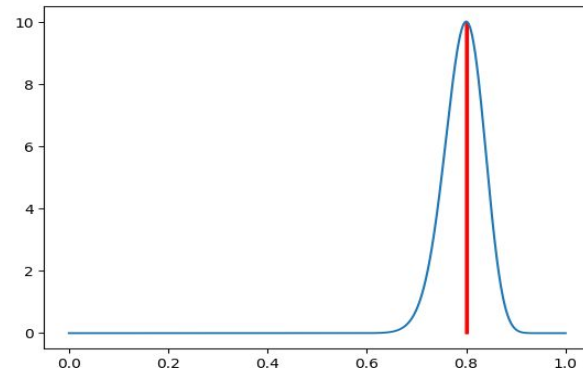
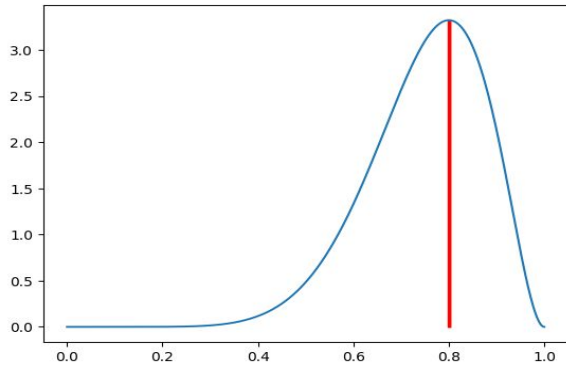
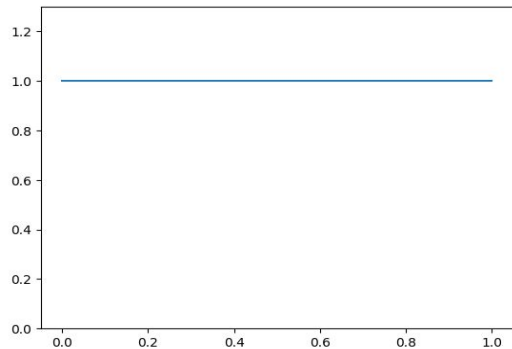
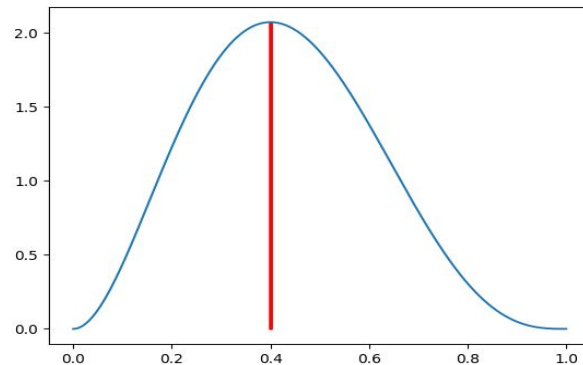


BETA RANDOM VARIABLE (INTUITION)

If you flip a coin with unknown probability of heads X , what does your belief distribution look like if:

- You didn't observe anything?
- You observed 8 heads and 2 tails?
- You observed 80 heads and 20 tails?
- You observed 2 heads and 3 tails?

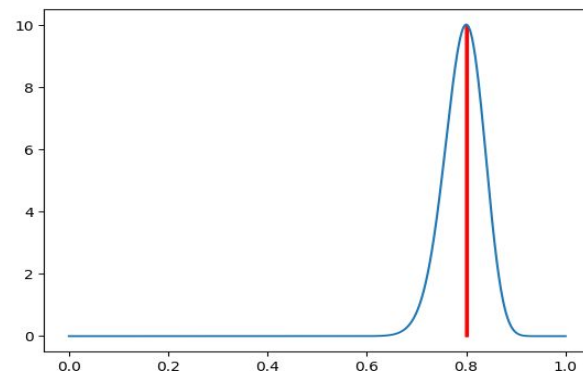
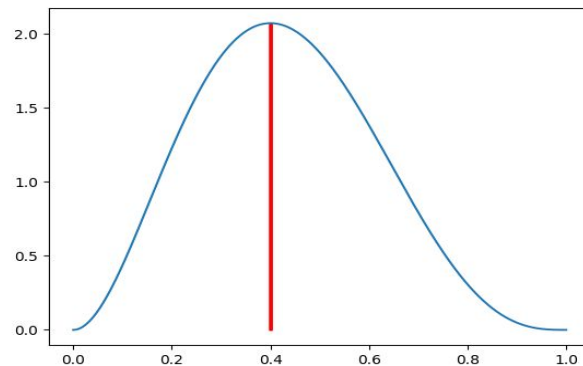
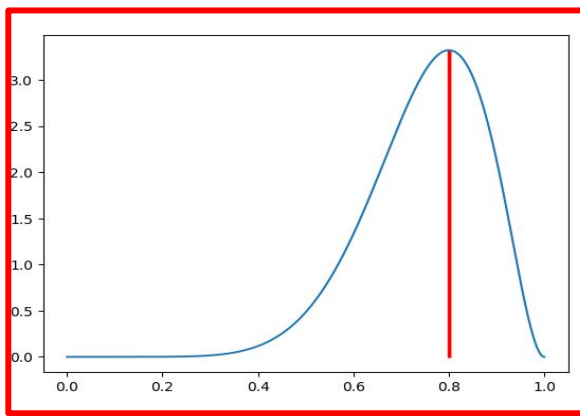
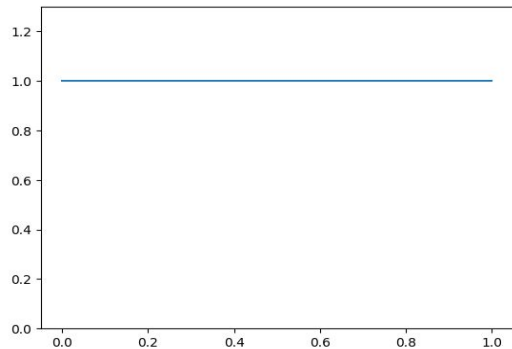
8/10



BETA RANDOM VARIABLE (INTUITION)

If you flip a coin with unknown probability of heads X , what does your belief distribution look like if:

- You didn't observe anything?
- You observed 8 heads and 2 tails?
- You observed 80 heads and 20 tails?
- You observed 2 heads and 3 tails?

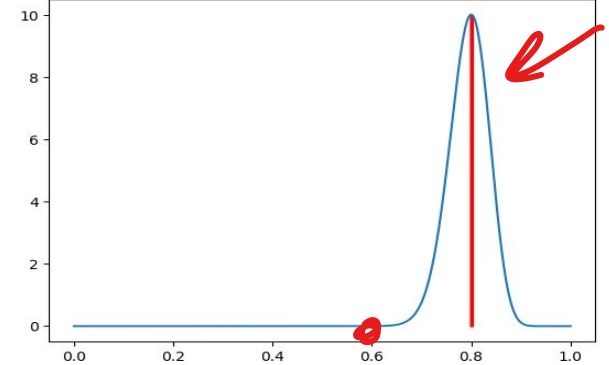
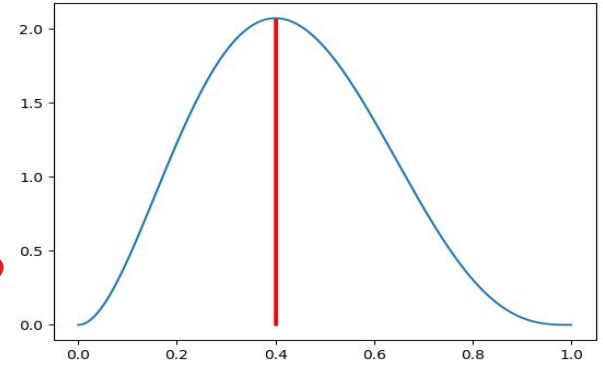
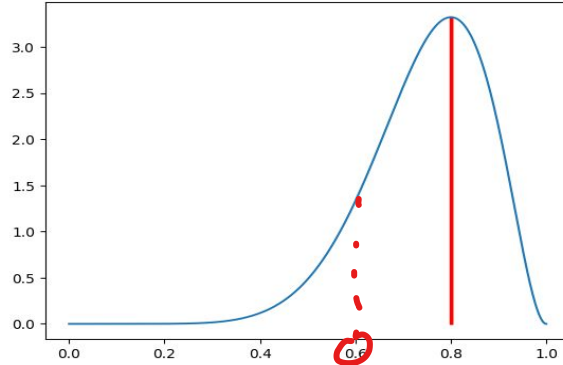
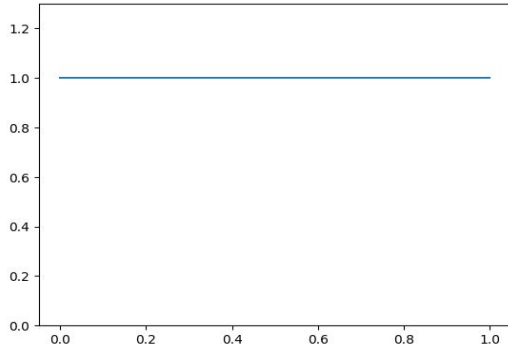


BETA RANDOM VARIABLE (INTUITION)

If you flip a coin with unknown probability of heads X , what does your belief distribution look like if:

- You didn't observe anything?
- You observed 8 heads and 2 tails?
- You observed 80 heads and 20 tails?
- You observed 2 heads and 3 tails?

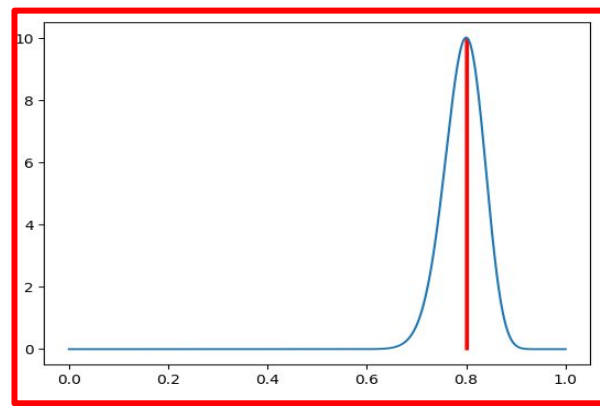
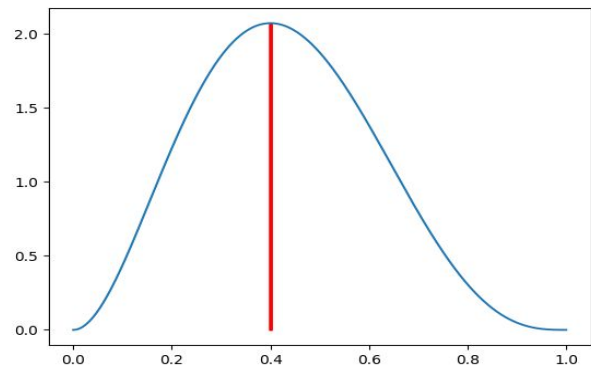
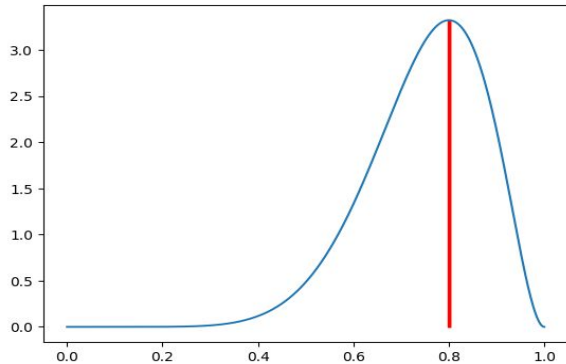
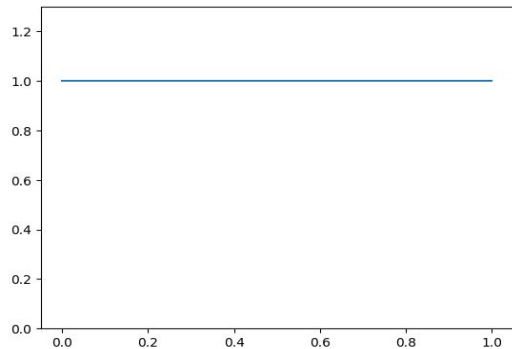
$\frac{80}{100}$



BETA RANDOM VARIABLE (INTUITION)

If you flip a coin with unknown probability of heads X , what does your belief distribution look like if:

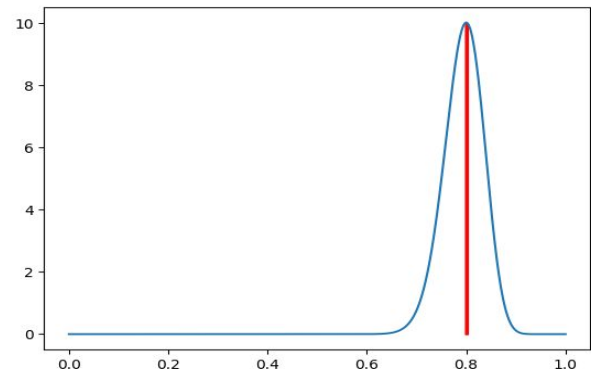
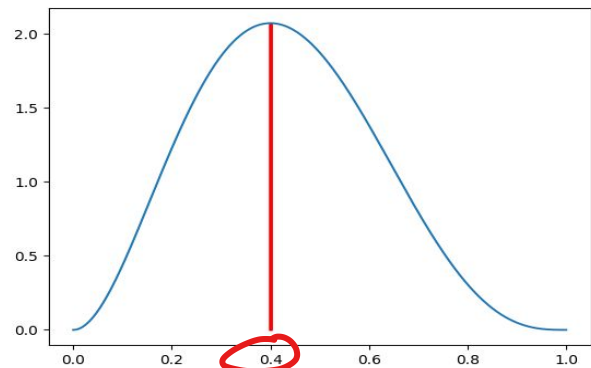
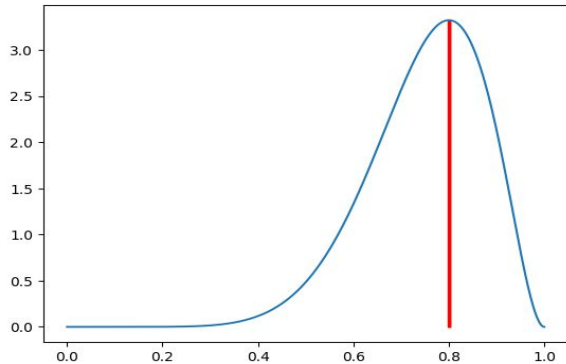
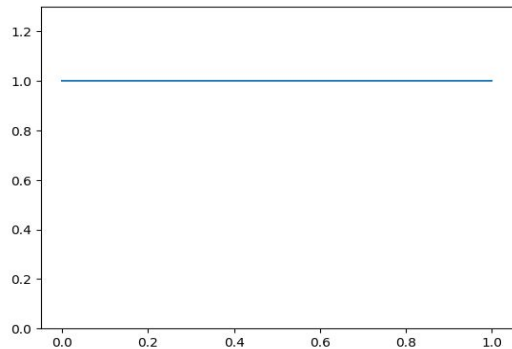
- You didn't observe anything?
- You observed 8 heads and 2 tails?
- You observed 80 heads and 20 tails?
- You observed 2 heads and 3 tails?



BETA RANDOM VARIABLE (INTUITION)

If you flip a coin with unknown probability of heads X , what does your belief distribution look like if:

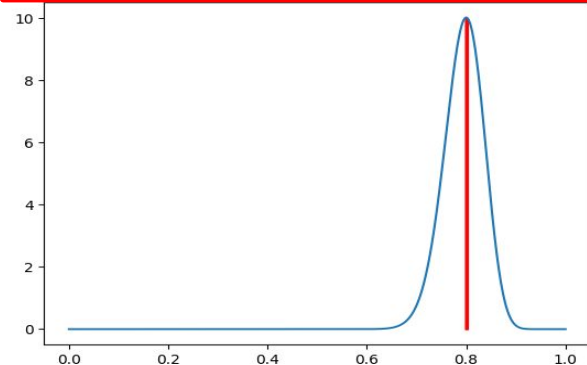
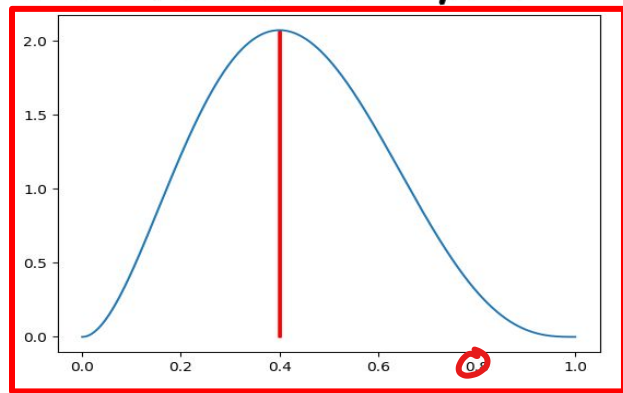
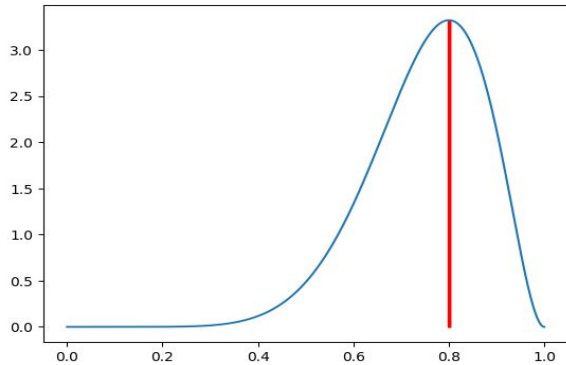
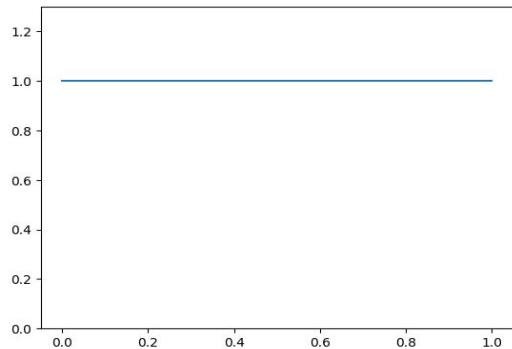
- You didn't observe anything?
- You observed 8 heads and 2 tails?
- You observed 80 heads and 20 tails?
- You observed 2 heads and 3 tails?



BETA RANDOM VARIABLE (INTUITION)

If you flip a coin with unknown probability of heads X , what does your belief distribution look like if:

- You didn't observe anything?
- You observed 8 heads and 2 tails?
- You observed 80 heads and 20 tails?
- You observed 2 heads and 3 tails?



THE BETA RANDOM VARIABLE

Beta RV: $X \sim \text{Beta}(\alpha, \beta)$, if and only if X has the following pdf:

$$f_X(x) = \begin{cases} \frac{1}{B(\alpha, \beta)} x^{\alpha-1} (1-x)^{\beta-1}, & 0 \leq x \leq 1 \\ 0, & \text{otherwise} \end{cases}$$

THE BETA RANDOM VARIABLE

Beta RV: $X \sim \text{Beta}(\alpha, \beta)$, if and only if X has the following pdf:

$$f_X(x) = \begin{cases} \frac{1}{B(\alpha, \beta)} x^{\alpha-1} (1-x)^{\beta-1}, & 0 \leq x \leq 1 \\ 0, & \text{otherwise} \end{cases}$$

X is typically the belief distribution about some unknown probability of success, where we pretend we saw $\alpha - 1$ successes and $\beta - 1$ failures ahead of time. Hence, the mode, $\arg \max_{x \in [0,1]} f_X(x)$, is

$$\text{mode}[X] = \frac{\alpha - 1}{(\alpha - 1) + (\beta - 1)}$$

BETA RANDOM VARIABLE EXAMPLES

If you flip a coin with unknown probability of heads X , what does your belief distribution look like if:

- You didn't observe anything?

$$\alpha - 1 = 0 \quad , \quad \alpha = 1 \quad X \sim \text{Beta}(1, 1)$$
$$\beta - 1 = 0 \quad \beta = 1$$

$$x^0 (1-x)^0 = 1$$

BETA RANDOM VARIABLE EXAMPLES

If you flip a coin with unknown probability of heads X , what does your belief distribution look like if:

- You didn't observe anything?

$$\underline{Beta(0 + 1, 0 + 1)} \equiv \underline{Beta(1, 1)} \equiv \underline{Unif(0, 1)} \rightarrow \underline{\text{mode} = ???}$$

BETA RANDOM VARIABLE EXAMPLES

Beta (α, β)

If you flip a coin with unknown probability of heads X , what does your belief distribution look like if:

- You didn't observe anything?

$Beta(0 + 1, 0 + 1) \equiv Beta(1, 1) \equiv Unif(0, 1) \rightarrow \text{mode} = ???$

- You observed 8 heads and 2 tails?

$X \sim Beta(9, 3)$

BETA RANDOM VARIABLE EXAMPLES

If you flip a coin with unknown probability of heads X , what does your belief distribution look like if:

- You didn't observe anything?

$$Beta(0 + 1, 0 + 1) \equiv Beta(1, 1) \equiv Unif(0, 1) \rightarrow \text{mode} = ???$$

- You observed 8 heads and 2 tails?

$$Beta(8 + 1, 2 + 1) \equiv Beta(9, 3) \rightarrow \text{mode} = \frac{(9 - 1)}{(9 - 1) + (3 - 1)} = \frac{8}{10}$$

BETA RANDOM VARIABLE EXAMPLES

If you flip a coin with unknown probability of heads X , what does your belief distribution look like if:

- You didn't observe anything?

$$\text{Beta}(0 + 1, 0 + 1) \equiv \text{Beta}(1, 1) \equiv \text{Unif}(0, 1) \rightarrow \text{mode} = ???$$

- You observed 8 heads and 2 tails?

$$\text{Beta}(8 + 1, 2 + 1) \equiv \text{Beta}(9, 3) \rightarrow \text{mode} = \frac{(9 - 1)}{(9 - 1) + (3 - 1)} = \frac{8}{10}$$

- You observed 80 heads and 20 tails?

$$\text{Beta}(80 + 1, 20 + 1) \equiv \text{Beta}(81, 21)$$

BETA RANDOM VARIABLE EXAMPLES

If you flip a coin with unknown probability of heads X , what does your belief distribution look like if:

- You didn't observe anything?

$$\text{Beta}(0 + 1, 0 + 1) \equiv \text{Beta}(1, 1) \equiv \text{Unif}(0, 1) \rightarrow \text{mode} = ???$$

- You observed 8 heads and 2 tails?

$$\text{Beta}(8 + 1, 2 + 1) \equiv \text{Beta}(9, 3) \rightarrow \text{mode} = \frac{(9 - 1)}{(9 - 1) + (3 - 1)} = \frac{8}{10}$$

- You observed 80 heads and 20 tails?

$$\text{Beta}(80 + 1, 20 + 1) \equiv \text{Beta}(81, 21) \rightarrow \text{mode} = \frac{(81 - 1)}{(81 - 1) + (21 - 1)} = \frac{80}{100}$$

BETA RANDOM VARIABLE EXAMPLES

If you flip a coin with unknown probability of heads X , what does your belief distribution look like if:

- You didn't observe anything?

$$Beta(0 + 1, 0 + 1) \equiv Beta(1, 1) \equiv Unif(0, 1) \rightarrow \text{mode} = ???$$

- You observed 8 heads and 2 tails?

$$Beta(8 + 1, 2 + 1) \equiv Beta(9, 3) \rightarrow \text{mode} = \frac{(9 - 1)}{(9 - 1) + (3 - 1)} = \frac{8}{10}$$

- You observed 80 heads and 20 tails?

$$Beta(80 + 1, 20 + 1) \equiv Beta(81, 21) \rightarrow \text{mode} = \frac{(81 - 1)}{(81 - 1) + (21 - 1)} = \frac{80}{100}$$

- You observed 2 heads and 3 tails?

BETA RANDOM VARIABLE EXAMPLES

If you flip a coin with unknown probability of heads X , what does your belief distribution look like if:

- You didn't observe anything?

$$\text{Beta}(0 + 1, 0 + 1) \equiv \text{Beta}(1, 1) \equiv \text{Unif}(0, 1) \rightarrow \text{mode} = ???$$

- You observed 8 heads and 2 tails?

$$\text{Beta}(8 + 1, 2 + 1) \equiv \text{Beta}(9, 3) \rightarrow \text{mode} = \frac{(9 - 1)}{(9 - 1) + (3 - 1)} = \frac{8}{10}$$

- You observed 80 heads and 20 tails?

$$\text{Beta}(80 + 1, 20 + 1) \equiv \text{Beta}(81, 21) \rightarrow \text{mode} = \frac{(81 - 1)}{(81 - 1) + (21 - 1)} = \frac{80}{100}$$

- You observed 2 heads and 3 tails?

$$\text{Beta}(2 + 1, 3 + 1) \equiv \text{Beta}(3, 4) \rightarrow \text{mode} = \frac{(3 - 1)}{(3 - 1) + (4 - 1)} = \frac{2}{5}$$

AGENDA

- THE BETA RANDOM VARIABLE
- MAP ESTIMATION
- MAP EXAMPLE

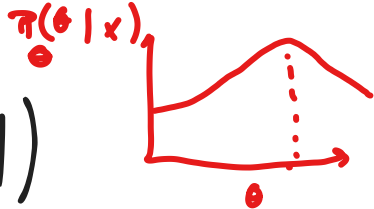
MAXIMUM A POSTERIORI (INTUITION)



In maximum likelihood estimation, we use iid samples $\mathbf{x} = (x_1, \dots, x_n)$ from some distribution with unknown parameter(s) θ , in order to estimate θ .

$$\hat{\theta}_{MLE} = \arg \max_{\theta} L(\mathbf{x} | \theta) = \arg \max_{\theta} \prod_{i=1}^n f_X(x_i; \theta)$$

MAXIMUM A POSTERIORI (INTUITION)



In maximum likelihood estimation, we use iid samples $x = (x_1, \dots, x_n)$ from some distribution with unknown parameter(s) θ , in order to estimate θ .

$$\hat{\theta}_{MLE} = \arg \max_{\theta} L(x | \theta) = \arg \max_{\theta} \prod_{i=1}^n f_X(x_i; \theta)$$

Shouldn't we maximize " $P(\theta | x)$ " instead? Doesn't make sense unless Θ is a rv.

MAXIMUM A POSTERIORI (INTUITION)

$f_X(x)$



In maximum likelihood estimation, we use iid samples $x = (x_1, \dots, x_n)$ from some distribution with unknown parameter(s) θ , in order to estimate θ .

$$\hat{\theta}_{MLE} = \arg \max_{\theta} L(x | \theta) = \arg \max_{\theta} \prod_{i=1}^n f_X(x_i; \theta)$$

Shouldn't we maximize " $P(\theta | x)$ " instead? Doesn't make sense unless Θ is a rv.

Maximum a Posteriori (MAP) Estimation Idea: Actually, unknown parameter(s) is a random variable θ . We have a prior distribution $\pi_{\theta}(\theta)$ and posterior distribution (given data) $\pi_{\theta}(\theta | x)$.

$\pi_{\Theta}(\theta)$



MAXIMUM A POSTERIORI (INTUITION)

In maximum likelihood estimation, we use iid samples $x = (x_1, \dots, x_n)$ from some distribution with unknown parameter(s) θ , in order to estimate θ .

$$\hat{\theta}_{MLE} = \arg \max_{\theta} L(x | \theta) = \arg \max_{\theta} \prod_{i=1}^n f_X(x_i; \theta)$$

Shouldn't we maximize " $P(\theta | x)$ " instead? Doesn't make sense unless Θ is a rv.

Maximum a Posteriori (MAP) Estimation Idea: Actually, unknown parameter(s) is a random variable θ . We have a prior distribution $\pi_{\theta}(\theta)$ and posterior distribution (given data) $\pi_{\theta}(\theta | x)$. By Bayes' Theorem,

$$\pi_{\theta}(\theta | x) = \frac{L(x | \theta)\pi_{\theta}(\theta)}{P(x)}$$



MAXIMUM A POSTERIORI (INTUITION)

In maximum likelihood estimation, we use iid samples $x = (x_1, \dots, x_n)$ from some distribution with unknown parameter(s) θ , in order to estimate θ .

$$\hat{\theta}_{MLE} = \arg \max_{\theta} L(x | \theta) = \arg \max_{\theta} \prod_{i=1}^n f_X(x_i; \theta)$$

Shouldn't we maximize " $P(\theta | x)$ " instead? Doesn't make sense unless Θ is a rv.

Maximum a Posteriori (MAP) Estimation Idea: Actually, unknown parameter(s) is a random variable θ . We have a prior distribution $\pi_{\theta}(\theta)$ and posterior distribution (given data) $\pi_{\theta}(\theta | x)$. By Bayes' Theorem,

$$\pi_{\theta}(\theta | x) = \frac{L(x | \theta)\pi_{\theta}(\theta)}{P(x)} \propto \underbrace{L(x | \theta)\pi_{\theta}(\theta)}$$



MAXIMUM A POSTERIORI (INTUITION)

In maximum likelihood estimation, we use iid samples $\mathbf{x} = (x_1, \dots, x_n)$ from some distribution with unknown parameter(s) θ , in order to estimate θ .

$$\hat{\theta}_{MLE} = \arg \max_{\theta} L(\mathbf{x} | \theta) = \arg \max_{\theta} \prod_{i=1}^n f_X(x_i; \theta)$$

Shouldn't we maximize " $P(\theta | \mathbf{x})$ " instead? Doesn't make sense unless Θ is a rv.

Maximum a Posteriori (MAP) Estimation Idea: Actually, unknown parameter(s) is a random variable θ . We have a prior distribution $\pi_{\theta}(\theta)$ and posterior distribution (given data) $\pi_{\theta}(\theta | \mathbf{x})$. By Bayes' Theorem,

$$\pi_{\theta}(\theta | \mathbf{x}) = \frac{L(\mathbf{x} | \theta)\pi_{\theta}(\theta)}{P(\mathbf{x})} \propto L(\mathbf{x} | \theta)\pi_{\theta}(\theta)$$

$$\hat{\theta}_{MAP} = \arg \max_{\theta} \pi_{\theta}(\theta | \mathbf{x}) =$$



MAXIMUM A POSTERIORI (INTUITION)

In maximum likelihood estimation, we use iid samples $x = (x_1, \dots, x_n)$ from some distribution with unknown parameter(s) θ , in order to estimate θ .

$$\hat{\theta}_{MLE} = \arg \max_{\theta} L(x | \theta) = \arg \max_{\theta} \prod_{i=1}^n f_X(x_i; \theta)$$

Shouldn't we maximize " $P(\theta | x)$ " instead? Doesn't make sense unless Θ is a rv.

Maximum a Posteriori (MAP) Estimation Idea: Actually, unknown parameter(s) is a random variable θ . We have a prior distribution $\pi_{\theta}(\theta)$ and posterior distribution (given data) $\pi_{\theta}(\theta | x)$. By Bayes' Theorem,

$$\pi_{\theta}(\theta | x) = \frac{L(x | \theta)\pi_{\theta}(\theta)}{P(x)} \propto L(x | \theta)\pi_{\theta}(\theta)$$

$$\hat{\theta}_{MAP} = \arg \max_{\theta} \pi_{\theta}(\theta | x) = \arg \max_{\theta} L(x | \theta) \pi_{\theta}(\theta)$$

MAXIMUM A POSTERIORI (MAP) ESTIMATION

Maximum A Posteriori (MAP) Estimation: Let $\mathbf{x} = (x_1, \dots, x_n)$ be iid realizations from probability mass function $p_X(t; \Theta = \theta)$ (if X discrete), or from density $f_X(t; \Theta = \theta)$ (if X continuous), where Θ is the random variable representing the parameter (or vector of parameters). We define the maximum a posteriori (MAP) estimator $\hat{\theta}_{MAP}$ of Θ to be the parameter which maximizes the posterior distribution of Θ given the data.

$$\hat{\theta}_{MAP} = \arg \max_{\theta} \pi_{\Theta}(\theta | \mathbf{x}) = \arg \max_{\theta} \underbrace{L(\mathbf{x} | \theta)} \underbrace{\pi_{\Theta}(\theta)}$$

MAXIMUM A POSTERIORI (EXAMPLE)

$$\hat{\theta}_{MLE} = \frac{\sum x_i}{n}$$



- a. Suppose our samples are $x = (0, 0, \underline{1}, 1, 0)$, from $Bernoulli(\theta)$, where θ is unknown. Assume θ is unrestricted; that is, $\theta \in (0, 1)$. What is the MLE for θ ?

$$\frac{2}{5}$$



MAXIMUM A POSTERIORI (EXAMPLE)

- a. Suppose our samples are $x = (0,0,1,1,0)$, from $Bernoulli(\theta)$, where θ is unknown. Assume θ is unrestricted; that is, $\theta \in (0,1)$. What is the MLE for θ ?

$$L(x|\theta) = \theta^2 (1-\theta)^3$$



MAXIMUM A POSTERIORI (EXAMPLE)

- a. Suppose our samples are $x = (0,0,1,1,0)$, from $Bernoulli(\theta)$, where θ is unknown. Assume θ is unrestricted; that is, $\theta \in (0,1)$. What is the MLE for θ ?

$$L(x|\theta) = \theta^2(1-\theta)^3$$

$$\hat{\theta}_{MLE} = \arg \max_{\theta \in [0,1]} \theta^2(1-\theta)^3 = \frac{2}{5}$$



MAXIMUM A POSTERIORI (EXAMPLE)

- a. Suppose our samples are $x = (0,0,1,1,0)$, from $Bernoulli(\theta)$, where θ is unknown. Assume θ is unrestricted; that is, $\theta \in (0,1)$. What is the MLE for θ ?

$$L(x|\theta) = \theta^2(1 - \theta)^3$$

$$\hat{\theta}_{MLE} = \arg \max_{\theta \in [0,1]} \theta^2(1 - \theta)^3 = \boxed{\frac{2}{5}} = 0.4$$

- b. Suppose we impose $\theta \in \{0.2,0.5,0.7\}$. What is the MLE for θ ?



MAXIMUM A POSTERIORI (EXAMPLE)

- a. Suppose our samples are $x = (0,0,1,1,0)$, from $Bernoulli(\theta)$, where θ is unknown. Assume θ is unrestricted; that is, $\theta \in (0,1)$. What is the MLE for θ ?

$$L(x|\theta) = \theta^2(1 - \theta)^3$$

$$\hat{\theta}_{MLE} = \arg \max_{\theta \in [0,1]} \theta^2(1 - \theta)^3 = \boxed{\frac{2}{5}}$$

- b. Suppose we impose $\theta \in \{0.2,0.5,0.7\}$. What is the MLE for θ ?

$$\hat{\theta}_{MLE} = \arg \max_{\theta \in \{0.2,0.5,0.7\}} L(x|\theta)$$



MAXIMUM A POSTERIORI (EXAMPLE)

- a. Suppose our samples are $x = (0,0,1,1,0)$, from $Bernoulli(\theta)$, where θ is unknown. Assume θ is unrestricted; that is, $\theta \in (0,1)$. What is the MLE for θ ?

$$L(x|\theta) = \theta^2(1-\theta)^3$$

$$0.2^2(1-0.2)^3$$

$$\hat{\theta}_{MLE} = \arg \max_{\theta \in [0,1]} \theta^2(1-\theta)^3 = \boxed{\frac{2}{5}}$$

- b. Suppose we impose $\theta \in \{0.2,0.5,0.7\}$. What is the MLE for θ ?

$$\hat{\theta}_{MLE} = \arg \max_{\theta \in \{0.2,0.5,0.7\}} L(x|\theta)$$

$$L(x|0.2) = \underline{(0.2^2 0.8^3)}$$

$$L(x|0.5) = \underline{(0.5^2 0.5^3)}$$

$$L(x|0.7) = \underline{(0.7^2 0.3^3)}$$



MAXIMUM A POSTERIORI (EXAMPLE)

- a. Suppose our samples are $x = (0,0,1,1,0)$, from $Bernoulli(\theta)$, where θ is unknown. Assume θ is unrestricted; that is, $\theta \in (0,1)$. What is the MLE for θ ?

$$L(x|\theta) = \theta^2(1 - \theta)^3$$

$$\hat{\theta}_{MLE} = \arg \max_{\theta \in [0,1]} \theta^2(1 - \theta)^3 = \boxed{\frac{2}{5}}$$

- b. Suppose we impose $\theta \in \{0.2,0.5,0.7\}$. What is the MLE for θ ?

$$\hat{\theta}_{MLE} = \arg \max_{\theta \in \{0.2,0.5,0.7\}} L(x|\theta)$$

$$L(x|0.2) = (0.2^2 0.8^3) = 0.02048$$

$$L(x|0.5) = (0.5^2 0.5^3) = 0.03125$$

$$L(x|0.7) = (0.7^2 0.3^3) = 0.01323$$



MAXIMUM A POSTERIORI (EXAMPLE)

- a. Suppose our samples are $x = (0,0,1,1,0)$, from $Bernoulli(\theta)$, where θ is unknown. Assume θ is unrestricted; that is, $\theta \in (0,1)$. What is the MLE for θ ?

$$L(x|\theta) = \theta^2(1 - \theta)^3$$

$$\hat{\theta}_{MLE} = \arg \max_{\theta \in [0,1]} \theta^2(1 - \theta)^3 = \frac{2}{5} = 0.4$$

- b. Suppose we impose $\theta \in \{0.2,0.5,0.7\}$. What is the MLE for θ ?

$$\hat{\theta}_{MLE} = \arg \max_{\theta \in \{0.2,0.5,0.7\}} L(x|\theta) = 0.5$$

$$L(x|0.2) = (0.2^2 0.8^3) = 0.02048 \leftarrow$$

$$L(x|0.5) = (0.5^2 0.5^3) = 0.03125 \leftarrow$$

$$L(x|0.7) = (0.7^2 0.3^3) = 0.01323 \leftarrow$$

MAXIMUM A POSTERIORI (EXAMPLE)

0.2

0.5

0.7



- c. Assume Θ is restricted as in part b (now a random variable for MAP). Suppose we have a (discrete) prior $\pi_{\Theta}(0.2) = 0.1$, $\pi_{\Theta}(0.5) = 0.01$, and $\pi_{\Theta}(0.7) = 0.89$. What is the MAP for θ ?

MAXIMUM A POSTERIORI (EXAMPLE)



- c. Assume Θ is restricted as in part b (now a random variable for MAP). Suppose we have a (discrete) prior $\pi_{\Theta}(0.2) = 0.1$, $\pi_{\Theta}(0.5) = 0.01$, and $\pi_{\Theta}(0.7) = 0.89$. What is the MAP for θ ?

$$\hat{\theta}_{MAP} = \arg \max_{\theta \in \{0.2, 0.5, 0.7\}} L(x|\theta)\pi_{\Theta}(\theta)$$



MAXIMUM A POSTERIORI (EXAMPLE)

- c. Assume Θ is restricted as in part b (now a random variable for MAP). Suppose we have a (discrete) prior $\pi_{\Theta}(0.2) = 0.1$, $\pi_{\Theta}(0.5) = 0.01$, and $\pi_{\Theta}(0.7) = 0.89$. What is the MAP for θ ?

$$\hat{\theta}_{MAP} = \arg \max_{\theta \in \{0.2, 0.5, 0.7\}} L(x|\theta)\pi_{\Theta}(\theta)$$

$$\begin{aligned}\pi_{\Theta}(0.2|x) &= L(\underline{x|0.2})\pi_{\Theta}(\underline{0.2}) \\ \pi_{\Theta}(0.5|x) &= L(x|0.5)\pi_{\Theta}(0.5) \\ \pi_{\Theta}(0.7|x) &= L(x|0.7)\pi_{\Theta}(0.7)\end{aligned}$$



MAXIMUM A POSTERIORI (EXAMPLE)

- c. Assume Θ is restricted as in part b (now a random variable for MAP). Suppose we have a (discrete) prior $\pi_{\Theta}(0.2) = 0.1$, $\pi_{\Theta}(0.5) = 0.01$, and $\pi_{\Theta}(0.7) = 0.89$. What is the MAP for θ ?

$$\hat{\theta}_{MAP} = \arg \max_{\theta \in \{0.2, 0.5, 0.7\}} L(x|\theta)\pi_{\Theta}(\theta)$$

$$\begin{aligned} \pi_{\Theta}(0.2|x) &= L(x|0.2)\pi_{\Theta}(0.2) = (0.2^2 0.8^3)(0.1) = 0.0020480 \\ \pi_{\Theta}(0.5|x) &= L(x|0.5)\pi_{\Theta}(0.5) = (0.5^2 0.5^3)(0.01) = 0.0003125 \\ \pi_{\Theta}(0.7|x) &= L(x|0.7)\pi_{\Theta}(0.7) = (0.7^2 0.3^3)(0.89) = 0.0117747 \end{aligned}$$



MAXIMUM A POSTERIORI (EXAMPLE)

- c. Assume Θ is restricted as in part b (now a random variable for MAP). Suppose we have a (discrete) prior $\pi_{\Theta}(0.2) = 0.1$, $\pi_{\Theta}(0.5) = 0.01$, and $\pi_{\Theta}(0.7) = 0.89$. What is the MAP for θ ?

$$\hat{\theta}_{MAP} = \arg \max_{\theta \in \{0.2, 0.5, 0.7\}} L(x|\theta)\pi_{\Theta}(\theta) = \boxed{0.7}$$

$$\begin{aligned}\pi_{\Theta}(0.2|x) &= L(x|0.2)\pi_{\Theta}(0.2) = (0.2^2 0.8^3)(0.1) = 0.0020480 \\ \pi_{\Theta}(0.5|x) &= L(x|0.5)\pi_{\Theta}(0.5) = (0.5^2 0.5^3)(0.01) = 0.0003125 \\ \pi_{\Theta}(0.7|x) &= L(x|0.7)\pi_{\Theta}(0.7) = (0.7^2 0.3^3)(0.89) = 0.0117747\end{aligned}$$



MAXIMUM A POSTERIORI (EXAMPLE)

- c. Assume Θ is restricted as in part b (now a random variable for MAP). Suppose we have a (discrete) prior $\pi_{\Theta}(0.2) = 0.1$, $\pi_{\Theta}(0.5) = 0.01$, and $\pi_{\Theta}(0.7) = 0.89$. What is the MAP for θ ?

$$\hat{\theta}_{MAP} = \arg \max_{\theta \in \{0.2, 0.5, 0.7\}} L(x|\theta)\pi_{\Theta}(\theta) = \boxed{0.7}$$

$$\pi_{\Theta}(0.2|x) = L(x|0.2)\pi_{\Theta}(0.2) = (0.2^2 0.8^3)(0.1) = 0.0020480$$

$$\pi_{\Theta}(0.5|x) = L(x|0.5)\pi_{\Theta}(0.5) = (0.5^2 0.5^3)(0.01) = 0.0003125$$

$$\pi_{\Theta}(0.7|x) = L(x|0.7)\pi_{\Theta}(0.7) = (0.7^2 0.3^3)(0.89) = 0.0117747$$

- d. Show that we can make the MAP whatever we like, by finding a prior over $\{0.2, 0.5, 0.7\}$ so that the MAP is 0.2, another so that it is 0.5, and another so that it is 0.7.



MAXIMUM A POSTERIORI (EXAMPLE)

- c. Assume Θ is restricted as in part b (now a random variable for MAP). Suppose we have a (discrete) prior $\pi_{\Theta}(0.2) = 0.1$, $\pi_{\Theta}(0.5) = 0.01$, and $\pi_{\Theta}(0.7) = 0.89$. What is the MAP for θ ?

$$\hat{\theta}_{MAP} = \arg \max_{\theta \in \{0.2, 0.5, 0.7\}} L(x|\theta)\pi_{\Theta}(\theta) = \boxed{0.7}$$

$$\begin{aligned}\pi_{\Theta}(0.2|x) &= L(x|0.2)\pi_{\Theta}(0.2) = (0.2^2 0.8^3)(0.1) = 0.0020480 \\ \pi_{\Theta}(0.5|x) &= L(x|0.5)\pi_{\Theta}(0.5) = (0.5^2 0.5^3)(0.01) = 0.0003125 \\ \pi_{\Theta}(0.7|x) &= L(x|0.7)\pi_{\Theta}(0.7) = (0.7^2 0.3^3)(0.89) = 0.0117747\end{aligned}$$

- d. Show that we can make the MAP whatever we like, by finding a prior over $\{0.2, 0.5, 0.7\}$ so that the MAP is 0.2, another so that it is 0.5, and another so that it is 0.7.

Choose $\pi_{\Theta}(\theta) = 1$ for the θ you want.

MAXIMUM A POSTERIORI (EXAMPLE)



- e. Typically, for the Bernoulli/Binomial distribution, we want to be able to get any value $\in (0,1)$, (not just ones in a finite set such as $\{0.2,0.5,0.7\}$). So we assign $\Theta \sim \text{Beta}(\alpha, \beta)$ with parameters $\alpha, \beta > 0$ and density $\pi_{\Theta}(\theta) = \frac{1}{B(\alpha, \beta)} \theta^{\alpha-1} (1-\theta)^{\beta-1}$ for $\theta \in (0,1)$. Recall the mode of a $W \sim \text{Beta}(\alpha, \beta)$ rv is $\frac{\alpha-1}{(\alpha-1)+(\beta-1)}$ (the value with highest density $\arg \max_w f_W(w)$).



MAXIMUM A POSTERIORI (EXAMPLE)

- e. Typically, for the Bernoulli/Binomial distribution, we want to be able to get any value $\in (0,1)$, (not just ones in a finite set such as $\{0.2, 0.5, 0.7\}$). So we assign $\Theta \sim \text{Beta}(\alpha, \beta)$ with parameters $\alpha, \beta > 0$ and density $\pi_{\Theta}(\theta) = \frac{1}{B(\alpha, \beta)} \theta^{\alpha-1} (1-\theta)^{\beta-1}$ for $\theta \in (0,1)$. Recall the mode of a $W \sim \text{Beta}(\alpha, \beta)$ rv is $\frac{\alpha-1}{(\alpha-1)+(\beta-1)}$ (the value with highest density $\arg \max_w f_W(w)$).

$k = \text{heads}$
 $n = \text{total flips}$

Suppose x_1, \dots, x_n are iid from a Bernoulli distribution with unknown parameter. Recall the MLE is k/n , where $k = \sum x_i$. Show that the posterior $\pi_{\Theta}(\theta|x)$ is $\text{Beta}(k + \alpha, n - k + \beta)$, and find the MAP. (Hint: use the mode given).

MAXIMUM A POSTERIORI (EXAMPLE)



- e. Typically, for the Bernoulli/Binomial distribution, we want to be able to get any value $\in (0,1)$, (not just ones in a finite set such as $\{0.2, 0.5, 0.7\}$). So we assign $\Theta \sim \text{Beta}(\alpha, \beta)$ with parameters $\alpha, \beta > 0$ and density $\pi_{\Theta}(\theta) = \frac{1}{B(\alpha, \beta)} \theta^{\alpha-1} (1-\theta)^{\beta-1}$ for $\theta \in (0,1)$. Recall the mode of a $W \sim \text{Beta}(\alpha, \beta)$ rv is $\frac{\alpha-1}{(\alpha-1)+(\beta-1)}$ (the value with highest density $\arg \max_w f_W(w)$).

Suppose x_1, \dots, x_n are iid from a Bernoulli distribution with unknown parameter. Recall the MLE is k/n , where $k = \sum x_i$. Show that the posterior $\pi_{\Theta}(\theta|x)$ is $\text{Beta}(k + \alpha, n - k + \beta)$, and find the MAP. (Hint: use the mode given).

$$\pi_{\Theta}(\theta|x) \propto L(x|\theta) \cdot \pi_{\Theta}(\theta)$$



MAXIMUM A POSTERIORI (EXAMPLE)

- e. Typically, for the Bernoulli/Binomial distribution, we want to be able to get any value $\in (0,1)$, (not just ones in a finite set such as $\{0.2,0.5,0.7\}$). So we assign $\Theta \sim \text{Beta}(\alpha, \beta)$ with parameters $\alpha, \beta > 0$ and density $\pi_{\Theta}(\theta) = \frac{1}{B(\alpha, \beta)} \theta^{\alpha-1} (1-\theta)^{\beta-1}$ for $\theta \in (0,1)$. Recall the mode of a $W \sim \text{Beta}(\alpha, \beta)$ rv is $\frac{\alpha-1}{(\alpha-1)+(\beta-1)}$ (the value with highest density $\arg \max_w f_W(w)$).

Suppose x_1, \dots, x_n are iid from a Bernoulli distribution with unknown parameter. Recall the MLE is k/n , where $k = \sum x_i$. Show that the posterior $\pi_{\Theta}(\theta|x)$ is $\text{Beta}(k + \alpha, n - k + \beta)$, and find the MAP. (Hint: use the mode given).

$$\begin{aligned} \pi_{\Theta}(\theta|x) &\propto L(x|\theta) \cdot \pi_{\Theta}(\theta) \\ &= \binom{n}{k} \theta^k (1-\theta)^{n-k}. \end{aligned}$$



MAXIMUM A POSTERIORI (EXAMPLE)

- e. Typically, for the Bernoulli/Binomial distribution, we want to be able to get any value $\in (0,1)$, (not just ones in a finite set such as $\{0.2,0.5,0.7\}$). So we assign $\Theta \sim \text{Beta}(\alpha, \beta)$ with parameters $\alpha, \beta > 0$ and density $\pi_{\Theta}(\theta) = \frac{1}{B(\alpha, \beta)} \theta^{\alpha-1} (1-\theta)^{\beta-1}$ for $\theta \in (0,1)$. Recall the mode of a $W \sim \text{Beta}(\alpha, \beta)$ rv is $\frac{\alpha-1}{(\alpha-1)+(\beta-1)}$ (the value with highest density $\arg \max_w f_W(w)$).

Suppose x_1, \dots, x_n are iid from a Bernoulli distribution with unknown parameter. Recall the MLE is k/n , where $k = \sum x_i$. Show that the posterior $\pi_{\Theta}(\theta|x)$ is $\text{Beta}(k + \alpha, n - k + \beta)$, and find the MAP. (Hint: use the mode given).

$$\begin{aligned} \pi_{\Theta}(\theta|x) &\propto L(x|\theta) \cdot \pi_{\Theta}(\theta) \\ &= \binom{n}{k} \theta^k (1-\theta)^{n-k} \cdot \left(\frac{1}{B(\alpha, \beta)} \theta^{\alpha-1} (1-\theta)^{\beta-1} \right) \end{aligned}$$

MAXIMUM A POSTERIORI (EXAMPLE)



- e. Typically, for the Bernoulli/Binomial distribution, we want to be able to get any value $\in (0,1)$, (not just ones in a finite set such as $\{0.2,0.5,0.7\}$). So we assign $\Theta \sim \text{Beta}(\alpha, \beta)$ with parameters $\alpha, \beta > 0$ and density $\pi_{\Theta}(\theta) = \frac{1}{B(\alpha, \beta)} \theta^{\alpha-1} (1-\theta)^{\beta-1}$ for $\theta \in (0,1)$. Recall the mode of a $W \sim \text{Beta}(\alpha, \beta)$ rv is $\frac{\alpha-1}{(\alpha-1)+(\beta-1)}$ (the value with highest density $\arg \max_w f_W(w)$).

Suppose x_1, \dots, x_n are iid from a Bernoulli distribution with unknown parameter. Recall the MLE is k/n , where $k = \sum x_i$. Show that the posterior $\pi_{\Theta}(\theta|x)$ is $\text{Beta}(k + \alpha, n - k + \beta)$, and find the MAP. (Hint: use the mode given).

$$\begin{aligned} \pi_{\Theta}(\theta|x) &\propto L(x|\theta) \cdot \pi_{\Theta}(\theta) \\ &= \left(\binom{n}{k} \theta^k (1-\theta)^{n-k} \right) \cdot \left(\frac{1}{B(\alpha, \beta)} \theta^{\alpha-1} (1-\theta)^{\beta-1} \right) \\ &\propto \theta^{(k+\alpha)-1} (1-\theta)^{(n-k+\beta)-1} \end{aligned}$$



MAXIMUM A POSTERIORI (EXAMPLE)

- e. Typically, for the Bernoulli/Binomial distribution, we want to be able to get any value $\in (0,1)$, (not just ones in a finite set such as $\{0.2,0.5,0.7\}$). So we assign $\Theta \sim \text{Beta}(\alpha, \beta)$ with parameters $\alpha, \beta > 0$ and density $\pi_{\Theta}(\theta) = \frac{1}{B(\alpha, \beta)} \theta^{\alpha-1} (1-\theta)^{\beta-1}$ for $\theta \in (0,1)$. Recall the mode of a $W \sim \text{Beta}(\alpha, \beta)$ rv is $\frac{\alpha-1}{(\alpha-1)+(\beta-1)}$ (the value with highest density $\arg \max_w f_W(w)$).

Suppose x_1, \dots, x_n are iid from a Bernoulli distribution with unknown parameter. Recall the MLE is k/n , where $k = \sum x_i$. Show that the posterior $\pi_{\Theta}(\theta|x)$ is $\text{Beta}(k + \alpha, n - k + \beta)$, and find the MAP. (Hint: use the mode given).

$$\begin{aligned} \pi_{\Theta}(\theta|x) &\propto L(x|\theta) \cdot \pi_{\Theta}(\theta) \\ &= \binom{n}{k} \theta^k (1-\theta)^{n-k} \cdot \left(\frac{1}{B(\alpha, \beta)} \theta^{\alpha-1} (1-\theta)^{\beta-1} \right) \\ &\propto \theta^{(k+\alpha)-1} (1-\theta)^{(n-k+\beta)-1} \end{aligned}$$

Hence our posterior is $\text{Beta}(k + \alpha, n - k + \beta)$. The mode of this beta is given:

MAXIMUM A POSTERIORI (EXAMPLE)



- e. Typically, for the Bernoulli/Binomial distribution, we want to be able to get any value $\in (0,1)$, (not just ones in a finite set such as $\{0.2, 0.5, 0.7\}$). So we assign $\Theta \sim \text{Beta}(\alpha, \beta)$ with parameters $\alpha, \beta > 0$ and density $\pi_{\Theta}(\theta) = \frac{1}{B(\alpha, \beta)} \theta^{\alpha-1} (1-\theta)^{\beta-1}$ for $\theta \in (0,1)$. Recall the mode of a $W \sim \text{Beta}(\alpha, \beta)$ rv is $\frac{\alpha-1}{(\alpha-1)+(\beta-1)}$ (the value with highest density $\arg \max_w f_W(w)$).

Suppose x_1, \dots, x_n are iid from a Bernoulli distribution with unknown parameter. Recall the MLE is k/n , where $k = \sum x_i$. Show that the posterior $\pi_{\Theta}(\theta|x)$ is $\text{Beta}(k + \alpha, n - k + \beta)$, and find the MAP. (Hint: use the mode given).

$$\begin{aligned} \pi_{\Theta}(\theta|x) &\propto L(x|\theta) \cdot \pi_{\Theta}(\theta) \\ &= \binom{n}{k} \theta^k (1-\theta)^{n-k} \cdot \left(\frac{1}{B(\alpha, \beta)} \theta^{\alpha-1} (1-\theta)^{\beta-1} \right) \\ &\propto \theta^{(k+\alpha)-1} (1-\theta)^{(n-k+\beta)-1} \end{aligned}$$

Hence our posterior is $\text{Beta}(k + \alpha, n - k + \beta)$. The mode of this beta is given:

$$\frac{k + \alpha - 1}{(k + \alpha - 1) + (n - k + \beta - 1)}$$

MAXIMUM A POSTERIORI (EXAMPLE)



- e. Typically, for the Bernoulli/Binomial distribution, we want to be able to get any value $\in (0,1)$, (not just ones in a finite set such as $\{0.2,0.5,0.7\}$). So we assign $\Theta \sim \text{Beta}(\alpha, \beta)$ with parameters $\alpha, \beta > 0$ and density $\pi_{\Theta}(\theta) = \frac{1}{B(\alpha, \beta)} \theta^{\alpha-1} (1-\theta)^{\beta-1}$ for $\theta \in (0,1)$. Recall the mode of a $W \sim \text{Beta}(\alpha, \beta)$ rv is $\frac{\alpha-1}{(\alpha-1)+(\beta-1)}$ (the value with highest density $\arg \max_w f_W(w)$).

Suppose x_1, \dots, x_n are iid from a Bernoulli distribution with unknown parameter. Recall the MLE is k/n , where $k = \sum x_i$. Show that the posterior $\pi_{\Theta}(\theta|x)$ is $\text{Beta}(k + \alpha, n - k + \beta)$, and find the MAP. (Hint: use the mode given).

$$\begin{aligned} \pi_{\Theta}(\theta|x) &\propto L(x|\theta) \cdot \pi_{\Theta}(\theta) \\ &= \binom{n}{k} \theta^k (1-\theta)^{n-k} \cdot \left(\frac{1}{B(\alpha, \beta)} \theta^{\alpha-1} (1-\theta)^{\beta-1} \right) \\ &\propto \theta^{(k+\alpha)-1} (1-\theta)^{(n-k+\beta)-1} \end{aligned}$$

Hence our posterior is $\text{Beta}(k + \alpha, n - k + \beta)$. The mode of this beta is given:

$$\frac{k + \alpha - 1}{(k + \alpha - 1) + (n - k + \beta - 1)} = \frac{k + (\alpha - 1)}{n + (\alpha - 1) + (\beta - 1)}$$



MAXIMUM A POSTERIORI (EXAMPLE)

f. Notice that $Beta(1,1) \equiv Unif(0,1)$. If we used this as the prior, how would the MLE and MAP compare?

$$MLE = \left[\frac{k}{n} \right] \quad MAP = \frac{k + \underline{0} - 1}{n + (\underline{\alpha} - 1) + (\underline{\beta} - 1)} = \left[\frac{k}{n} \right]$$



MAXIMUM A POSTERIORI (EXAMPLE)

f. Notice that $Beta(1,1) \equiv Unif(0,1)$. If we used this as the prior, how would the MLE and MAP compare?

From previous slide, if $\alpha = \beta = 1$, then our MAP estimate is the same as our ML estimate!

$$\hat{\theta}_{MAP} = \frac{k + (\alpha - 1)}{n + (\alpha - 1) + (\beta - 1)}$$



MAXIMUM A POSTERIORI (EXAMPLE)

f. Notice that $Beta(1,1) \equiv Unif(0,1)$. If we used this as the prior, how would the MLE and MAP compare?

From previous slide, if $\alpha = \beta = 1$, then our MAP estimate is the same as our ML estimate!

$$\hat{\theta}_{MAP} = \frac{k + (\alpha - 1)}{n + (\alpha - 1) + (\beta - 1)} = \frac{k}{n} =$$



MAXIMUM A POSTERIORI (EXAMPLE)

f. Notice that $Beta(1,1) \equiv Unif(0,1)$. If we used this as the prior, how would the MLE and MAP compare?

From previous slide, if $\alpha = \beta = 1$, then our MAP estimate is the same as our ML estimate!

$$\hat{\theta}_{MAP} = \frac{k + (\alpha - 1)}{n + (\alpha - 1) + (\beta - 1)} = \frac{k}{n} = \hat{\theta}_{MLE}$$



MAXIMUM A POSTERIORI (EXAMPLE)

- f. Notice that $Beta(1,1) \equiv Unif(0,1)$. If we used this as the prior, how would the MLE and MAP compare?

From previous slide, if $\alpha = \beta = 1$, then our MAP estimate is the same as our ML estimate!

$$\hat{\theta}_{MAP} = \frac{k + (\alpha - 1)}{n + (\alpha - 1) + (\beta - 1)} = \frac{k}{n} = \hat{\theta}_{MLE}$$

- g. Since the posterior is also a Beta distribution, we call Beta the **conjugate prior** to the Bernoulli/Binomial distribution. Interpret α, β .



MAXIMUM A POSTERIORI (EXAMPLE)

- f. Notice that $Beta(1,1) \equiv Unif(0,1)$. If we used this as the prior, how would the MLE and MAP compare?

From previous slide, if $\alpha = \beta = 1$, then our MAP estimate is the same as our ML estimate!

TRUE HEADS

$$\hat{\theta}_{MAP} = \frac{k + (\alpha - 1)}{n + (\alpha - 1) + (\beta - 1)} = \frac{k}{n} = \hat{\theta}_{MLE}$$

- g. Since the posterior is also a Beta distribution, we call Beta the **conjugate prior** to the Bernoulli/Binomial distribution. Interpret α, β .

TRUE TRIALS



MAXIMUM A POSTERIORI (EXAMPLE)

- f. Notice that $Beta(1,1) \equiv Unif(0,1)$. If we used this as the prior, how would the MLE and MAP compare?

From previous slide, if $\alpha = \beta = 1$, then our MAP estimate is the same as our ML estimate!

FAKE HEADS

TRUE HEADS

$$\hat{\theta}_{MAP} = \frac{k + (\alpha - 1)}{n + (\alpha - 1) + (\beta - 1)} = \frac{k}{n} = \hat{\theta}_{MLE}$$

- g. Since the posterior is also a Beta distribution, we call Beta the **conjugate prior** to the Bernoulli/Binomial distribution. Interpret α, β .

TRUE TRIALS



MAXIMUM A POSTERIORI (EXAMPLE)

f. Notice that $Beta(1,1) \equiv Unif(0,1)$. If we used this as the prior, how would the MLE and MAP compare?

From previous slide, if $\alpha = \beta = 1$, then our MAP estimate is the same as our ML estimate!

FAKE HEADS

TRUE HEADS

$$\hat{\theta}_{MAP} = \frac{k + (\alpha - 1)}{n + (\alpha - 1) + (\beta - 1)} = \frac{k}{n} = \hat{\theta}_{MLE}$$

g. Since the posterior is also a Beta distribution, we call Beta the **conjugate prior** to the Bernoulli/Binomial distribution. Interpret α, β .

TRUE TRIALS

FAKE TRIALS



MAXIMUM A POSTERIORI (EXAMPLE)

- f. Notice that $Beta(1,1) \equiv Unif(0,1)$. If we used this as the prior, how would the MLE and MAP compare?

From previous slide, if $\alpha = \beta = 1$, then our MAP estimate is the same as our ML estimate!

$$\hat{\theta}_{MAP} = \frac{k + (\alpha - 1)}{n + (\alpha - 1) + (\beta - 1)} = \frac{k}{n} = \hat{\theta}_{MLE}$$

- g. Since the posterior is also a Beta distribution, we call Beta the **conjugate prior** to the Bernoulli/Binomial distribution. Interpret α, β .

It means: pretend we saw $\alpha - 1$ heads ahead of time, and $\beta - 1$ tails ahead of time. Then our total heads is $k + \alpha - 1$ and our total trials is $n + \alpha + \beta - 2$.



MAXIMUM A POSTERIORI (EXAMPLE)

f. Notice that $Beta(1,1) \equiv Unif(0,1)$. If we used this as the prior, how would the MLE and MAP compare?

From previous slide, if $\alpha = \beta = 1$, then our MAP estimate is the same as our ML estimate!

$$\hat{\theta}_{MAP} = \frac{k + (\alpha - 1)}{n + (\alpha - 1) + (\beta - 1)} = \frac{k}{n} = \hat{\theta}_{MLE}$$

g. Since the posterior is also a Beta distribution, we call Beta the **conjugate prior** to the Bernoulli/Binomial distribution. Interpret α, β .

It means: pretend we saw $\alpha - 1$ heads ahead of time, and $\beta - 1$ tails ahead of time. Then our total heads is $k + \alpha - 1$ and our total trials is $n + \alpha + \beta - 2$.

$$\begin{aligned} \pi_{\theta}(\theta|x) &\propto L(x|\theta) \cdot \pi_{\theta}(\theta) \\ &= \binom{n}{k} \theta^k (1-\theta)^{n-k} \cdot \left(\frac{1}{B(\alpha, \beta)} \theta^{\alpha-1} (1-\theta)^{\beta-1} \right) \end{aligned}$$

“LUCKY” US!!

$$\propto \theta^{(k+\alpha)-1} (1-\theta)^{(n-k+\beta)-1}$$

“LUCKY” US!!

MAXIMUM A POSTERIORI (EXAMPLE)

- h. As the number of samples goes to infinity, what is the relationship between the MLE and MAP? What does this say about our prior when n is small, or n is large?



MAXIMUM A POSTERIORI (EXAMPLE)



- h. As the number of samples goes to infinity, what is the relationship between the MLE and MAP? What does this say about our prior when n is small, or n is large?

They become equal! The prior is important if we don't have much data, but as we get more, the evidence overwhelms the prior.

MAXIMUM A POSTERIORI (EXAMPLE)



- h. As the number of samples goes to infinity, what is the relationship between the MLE and MAP? What does this say about our prior when n is small, or n is large?

They become equal! The prior is important if we don't have much data, but as we get more, the evidence overwhelms the prior.

- i. Which do you think is "better", MLE or MAP?

MAXIMUM A POSTERIORI (EXAMPLE)



- h. As the number of samples goes to infinity, what is the relationship between the MLE and MAP? What does this say about our prior when n is small, or n is large?

They become equal! The prior is important if we don't have much data, but as we get more, the evidence overwhelms the prior.

- i. Which do you think is "better", MLE or MAP?
- There is no right answer. There are two main schools in statistics: Bayesians and Frequentists.



MAXIMUM A POSTERIORI (EXAMPLE)

- h. As the number of samples goes to infinity, what is the relationship between the MLE and MAP? What does this say about our prior when n is small, or n is large?

They become equal! The prior is important if we don't have much data, but as we get more, the evidence overwhelms the prior.

- i. Which do you think is "better", MLE or MAP?
- There is no right answer. There are two main schools in statistics: Bayesians and Frequentists.
 - Frequentists prefer MLE since they don't believe you should be putting a prior belief on anything, and you should only make judgment based on what you've seen. They believe the parameter being estimated is a fixed quantity.

MAXIMUM A POSTERIORI (EXAMPLE)



- h. As the number of samples goes to infinity, what is the relationship between the MLE and MAP? What does this say about our prior when n is small, or n is large?

They become equal! The prior is important if we don't have much data, but as we get more, the evidence overwhelms the prior.

- i. Which do you think is "better", MLE or MAP?
- There is no right answer. There are two main schools in statistics: Bayesians and Frequentists.
 - Frequentists prefer MLE since they don't believe you should be putting a prior belief on anything, and you should only make judgment based on what you've seen. They believe the parameter being estimated is a fixed quantity.
 - On the other hand, Bayesians prefer MAP, since they can incorporate their prior knowledge into the estimation. Hence the parameter being estimated is a random variable, and we seek the mode - the value with the highest probability or density. An example would be estimating the probability of heads of a coin - is it reasonable to assume it is more likely fair than not? If so, what distribution should we put on the parameter space?

MAXIMUM A POSTERIORI (EXAMPLE)



- h. As the number of samples goes to infinity, what is the relationship between the MLE and MAP? What does this say about our prior when n is small, or n is large?

They become equal! The prior is important if we don't have much data, but as we get more, the evidence overwhelms the prior.

- i. Which do you think is "better", MLE or MAP?
- There is no right answer. There are two main schools in statistics: Bayesians and Frequentists.
 - Frequentists prefer MLE since they don't believe you should be putting a prior belief on anything, and you should only make judgment based on what you've seen. They believe the parameter being estimated is a fixed quantity.
 - On the other hand, Bayesians prefer MAP, since they can incorporate their prior knowledge into the estimation. Hence the parameter being estimated is a random variable, and we seek the mode - the value with the highest probability or density. An example would be estimating the probability of heads of a coin - is it reasonable to assume it is more likely fair than not? If so, what distribution should we put on the parameter space?
 - Anyway, in the long run, the prior "washes out", and the only thing that matters is the likelihood; the observed data. For small sample sizes like this, the prior significantly influences the MAP estimate. However, as the number of samples goes to infinity, the MAP and MLE are equal.