# 22: MAP

Lisa Yan and Jerry Cain
November 2, 2020

# Quick slide reference

# Maximum a Posteriori Estimator

# Maximum Likelihood Estimator

Consider a sample of $n$ i.i.d. random variables $X_1, X_2, \ldots, X_n$ (data).

Maximum Likelihood Estimator (MLE)

What is the parameter $\theta$ that **maximizes the likelihood** of our observed data $(X_1, X_2, \ldots, X_n)$?

$$L(\theta) = f(X_1, X_2, \ldots, X_n | \theta)$$

$$= \prod_{i=1}^{n} f(X_i | \theta)$$

$$\theta_{MLE} = \arg \max_{\theta} f(X_1, X_2, \ldots, X_n | \theta)$$

likelihood of data

Observations:

- MLE maximizes probability of observing data given a parameter $\theta$.

- If we are estimating $\theta$, shouldn't we maximize the probability of $\theta$ directly?

Today: **Bayesian estimation** using the Bayesian definition of probability!

# Maximum A Posteriori (MAP) Estimator

Consider a sample of $n$ i.i.d. random variables $X_1, X_2, \dots, X_n$ (data).

Maximum Likelihood Estimator (MLE)

What is the parameter $\theta$ that **maximizes the likelihood** of our observed data $(X_1, X_2, \dots, X_n)$?

$$L(\theta) = f(X_1, X_2, \dots, X_n | \theta)$$
$$= \prod_{i=1}^{n} f(X_i | \theta)$$
$$\theta_{MLE} = \arg \max_{\theta} f(X_1, X_2, \dots, X_n | \theta)$$

likelihood of data

Maximum a Posteriori (MAP) Estimator

Given our observed data $(X_1, X_2, \dots, X_n)$, what is the **most likely parameter** $\theta$?

$$\theta_{MAP} = \arg \max_{\theta} f(\theta | X_1, X_2, \dots, X_n)$$

posterior distribution of $\theta$

# Maximum A Posteriori (MAP) Estimator

Consider a sample of $n$ i.i.d. random variables $X_1, X_2, \ldots, X_n$ (data).

<u>def</u> The Maximum a Posteriori (MAP) Estimator of $\theta$ is the value of $\theta$ that maximizes the posterior distribution of $\theta$.

$$\theta_{MAP} = \arg \max_{\theta} f(\theta | X_1, X_2, \ldots, X_n)$$

Intuition with Bayes' Theorem:

$L(\theta)$, probability of data given parameter $\theta$

likelihood    prior

After seeing data, posterior belief of $\theta$

posterior

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

Before seeing data, prior belief of $\theta$

# Solving for $\theta_{MAP}$

- Observe data: $X_1, X_2, \ldots, X_n$, all i.i.d.
- Let likelihood be same as MLE: $f(X_1, X_2, \ldots, X_n | \theta) = \prod_{i=1}^{n} f(X_i | \theta)$
- Let the prior distribution of $\theta$ be $g(\theta)$.

$$\theta_{MAP} = \arg\max_{\theta} f(\theta | X_1, X_2, \ldots, X_n) = \arg\max_{\theta} \frac{f(X_1, X_2, \ldots, X_n | \theta) g(\theta)}{h(X_1, X_2, \ldots, X_n)} \quad \text{(Bayes' Theorem)}$$

$$= \arg\max_{\theta} \frac{g(\theta) \prod_{i=1}^{n} f(X_i | \theta)}{h(X_1, X_2, \ldots, X_n)} \quad \text{(independence)}$$

$$= \arg\max_{\theta} g(\theta) \prod_{i=1}^{n} f(X_i | \theta) \quad (1/h(X_1, X_2, \ldots, X_n) \text{ is a positive constant w.r.t. } \theta)$$

$$= \arg\max_{\theta} \left( \log g(\theta) + \sum_{i=1}^{n} \log f(X_i | \theta) \right)$$

# $\theta_{MAP}$: Interpretation 1

- Observe data: $X_1, X_2, \ldots, X_n$, all i.i.d.
- Let likelihood be same as MLE: $f(X_1, X_2, \ldots, X_n | \theta) = \prod_{i=1}^{n} f(X_i | \theta)$
- Let the prior distribution of $\theta$ be $g(\theta)$.

$$\theta_{MAP} = \arg\max_{\theta} f(\theta | X_1, X_2, \ldots, X_n) = \arg\max_{\theta} \frac{f(X_1, X_2, \ldots, X_n | \theta) g(\theta)}{h(X_1, X_2, \ldots, X_n)} \quad \text{(Bayes' Theorem)}$$

$$= \arg\max_{\theta} \frac{g(\theta) \prod_{i=1}^{n} f(X_i | \theta)}{h(X_1, X_2, \ldots, X_n)} \quad \text{(independence)}$$

$$= \arg\max_{\theta} g(\theta) \prod_{i=1}^{n} f(X_i | \theta) \quad (1/h(X_1, X_2, \ldots, X_n) \text{ is a positive constant w.r.t. } \theta)$$

$$= \arg\max_{\theta} \left( \log g(\theta) + \sum_{i=1}^{n} \log f(X_i | \theta) \right)$$

$\theta_{MAP}$ maximizes
<span style="color:red">log prior</span> + <span style="color:orange">log-likelihood</span>

# $\theta_{MAP}$: Interpretation 2

- Observe data: $X_1, X_2, \ldots, X_n$, all i.i.d.
- Let likelihood be same as MLE: $f(X_1, X_2, \ldots, X_n | \theta) = \prod_{i=1}^{n} f(X_i | \theta)$
- Let the prior distribution of $\theta$ be $g(\theta)$.

$$\theta_{MAP} = \arg\max_{\theta} f(\theta | X_1, X_2, \ldots, X_n) = \arg \quad \text{(Bayes' Theorem)}$$

> The mode of the
> posterior distribution of $\theta$

$$= \arg\max_{\theta} \frac{g(\theta) \prod_{i=1}^{n} f(X_i | \theta)}{h(X_1, X_2, \ldots, X_n)} \qquad \text{(independence)}$$

$$= \arg\max_{\theta} g(\theta) \prod_{i=1}^{n} f(X_i | \theta) \qquad (1/h(X_1, X_2, \ldots, X_n) \text{ is a positive constant w.r.t. } \theta)$$

$$= \arg\max_{\theta} \left( \log g(\theta) + \sum_{i=1}^{n} \log f(X_i | \theta) \right)$$

> $\theta_{MAP}$ maximizes
> log prior + log-likelihood

# Mode: A statistic of a random variable

The **mode** of a random variable $X$ is defined as:

(*X* discrete, PMF $p(x)$)
$$\arg \max_{x} p(x)$$

$$\arg \max_{x} f(x)$$
(*X* continuous, PDF $f(x)$)

- Intuitively: The value of $X$ that is "most likely."
- Note that some distributions may not have a unique mode (e.g., Uniform distribution, or Bernoulli(0.5))

$$\theta_{MAP} = \arg \max_{\theta} f(\theta | X_1, X_2, \ldots, X_n)$$

$\theta_{MAP}$ is the most likely $\theta$ given the data $X_1, X_2, \ldots, X_n$.

# Bernoulli MAP: Choosing a prior

# How does MAP work? (for Bernoulli)

Observe data                    $n$ heads, $m$ tails

Choose model                    Bernoulli($p$)

Choose prior on $\theta$

(some $g(\theta)$)

Find $\theta_{MAP} =$
$\underset{\theta}{\arg\max} \ f(\theta | X_1, X_2, \ldots, X_n)$

maximize
log prior + log-likelihood

$$\log g(\theta) + \sum_{i=1}^{n} \log f(X_i | \theta)$$

- Differentiate, set to 0
- Solve

A lot of our effort in MAP depends on the $g(\theta)$ we choose.

# MAP for Bernoulli

- Flip a coin 8 times. Observe $n = 7$ heads and $m = 1$ tail.
- Choose a prior on $\theta$. What is $\theta_{MAP}$?

Suppose we pick a prior $\theta \sim \mathcal{N}(0.5, 1^2)$. $g(\theta) = \frac{1}{\sqrt{2\pi}} e^{-(p-0.5)^2/2}$

1. Determine log prior + log likelihood

$$\log g(\theta) + \log f(X_1, X_2, \dots, X_n | \theta)$$

$$= \log \left( \frac{1}{\sqrt{2\pi}} e^{-(p-0.5)^2/2} \right) + \log \left( \binom{n+m}{n} p^n (1-p)^m \right)$$

$$= -\log(\sqrt{2\pi}) - (p-0.5)^2/2 + \log \binom{n+m}{n} + n \log p + m \log(1-p)$$

2. Differentiate w.r.t. (each) $\theta$, set to 0

$$-(p-0.5) + \frac{n}{p} - \frac{m}{1-p} = 0$$

3. Solve resulting equations

cubic equations why

We should choose an "easier" prior. This one is hard!

# A better approach: Use conjugate distributions

Observe data $\qquad\qquad n$ heads, $m$ tails

Choose model $\qquad\qquad$ Bernoulli$(p)$ $\qquad\qquad\qquad$ (choose conjugate distribution)

Choose prior on $\theta$ $\qquad$ (some $g(\theta)$)

Find $\theta_{MAP} =$ $\underset{\theta}{\arg\max} \, f(\theta|X_1, X_2, \dots, X_n)$

maximize
log prior + log-likelihood

$$\log g(\theta) + \sum_{i=1}^{n} \log f(X_i|\theta)$$

- Differentiate, set to 0
- Solve

⭐

Up next: Conjugate priors are great for MAP!

# Bernoulli MAP: Conjugate prior

# Beta is a conjugate distribution for Bernoulli

Beta is a **conjugate distribution** for Bernoulli, meaning:

- Prior and posterior parametric forms are the same

- Practically, conjugate means easy update:
  Add numbers of "successes" and "failures" seen to Beta parameters.

- You can set the prior to reflect how fair/biased you think the experiment is apriori.

| | |
|---|---|
| **Prior** | $\text{Beta}(a = n_{imag} + 1, b = m_{imag} + 1)$ |
| **Experiment** | Observe $n$ successes and $m$ failures |
| **Posterior** | $\text{Beta}(a = n_{imag} + n + 1, b = m_{imag} + m + 1)$ |

Mode of $\text{Beta}(a, b)$: $\dfrac{a - 1}{a + b - 2}$

(we'll prove this in a few minutes)

Beta parameters $a, b$ are called **hyperparameters**.
Interpret $\text{Beta}(a, b)$: $a + b - 2$ trials,
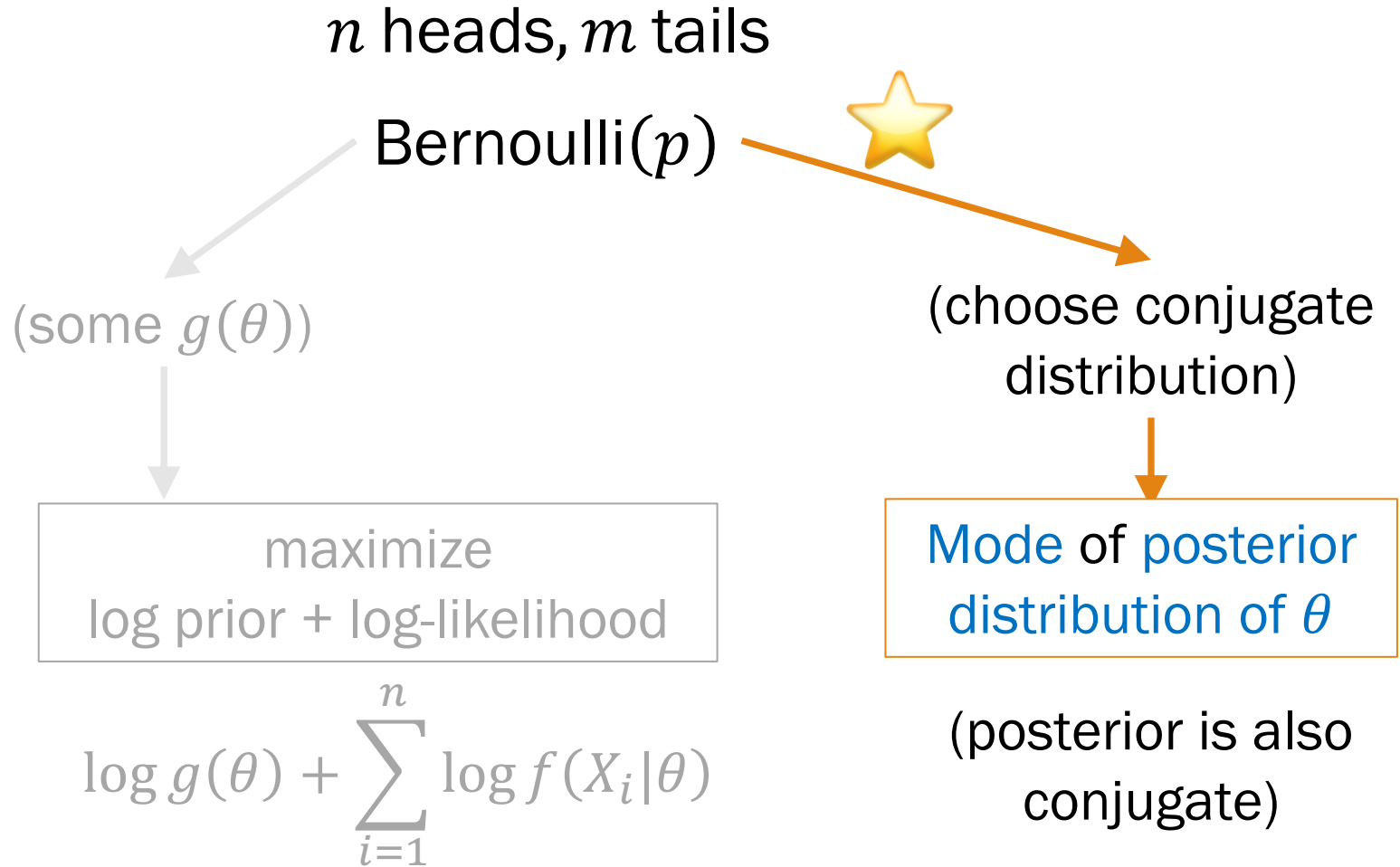of which $a - 1$ are successes

# How does MAP work? (for Bernoulli)

Observe data                     $n$ heads, $m$ tails

Choose model                     Bernoulli($p$) ⭐

Choose prior on $\theta$         (some $g(\theta)$)         (choose conjugate distribution)

Find $\theta_{MAP} =$
$\underset{\theta}{\arg\max} \, f(\theta|X_1, X_2, \ldots, X_n)$

$$\begin{array}{c} \text{maximize} \\ \text{log prior + log-likelihood} \end{array}$$

$$\log g(\theta) + \sum_{i=1}^{n} \log f(X_i|\theta)$$

- Differentiate, set to 0
- Solve

Mode of posterior distribution of $\theta$

(posterior is also conjugate)

# Conjugate strategy: MAP for Bernoulli

- Flip a coin 8 times. Observe $n = 7$ heads and $m = 1$ tail. ⎫ Define as data, $D$
- Choose a prior on $\theta$. What is $\theta_{MAP}$?

1. Choose a prior

   Suppose we pick a prior $\theta \sim \text{Beta}(a, b)$.

2. Determine posterior

   Because Beta is a conjugate distribution for Bernoulli, the posterior distribution is $\theta|D \sim \text{Beta}(a + n, b + m)$

3. Compute MAP

   $$\theta_{MAP} = \frac{a + n - 1}{a + n + b + m - 2}$$   (mode of $\text{Beta}(a + n, b + m)$)
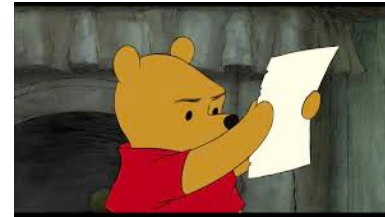
# MAP in practice

- Flip a coin 8 times. Observe $n = 7$ heads and $m = 1$ tail.
- What is the MAP estimator of the Bernoulli parameter $p$, if we assume a prior on $p$ of $\text{Beta}(2, 2)$?

# MAP in practice

- Flip a coin 8 times. Observe $n = 7$ heads and $m = 1$ tail.
- What is the MAP estimator of the Bernoulli parameter $p$, if we assume a prior on $p$ of $\text{Beta}(2, 2)$?

1. Choose a prior          $\theta \sim \text{Beta}(2,2).$

<span style="color:red">Before flipping the coin, we imagined 2 trials: 1 imaginary head, 1 imaginary tail.</span>

2. Determine posterior    Posterior distribution of $\theta$ given observed data is $\text{Beta}(9, 3)$

3. Compute MAP            $\theta_{MAP} = \dfrac{8}{10}$

<span style="color:blue">After the coin, we saw 10 trials: 8 heads (imaginary and real), 2 tails (imaginary and real).</span>

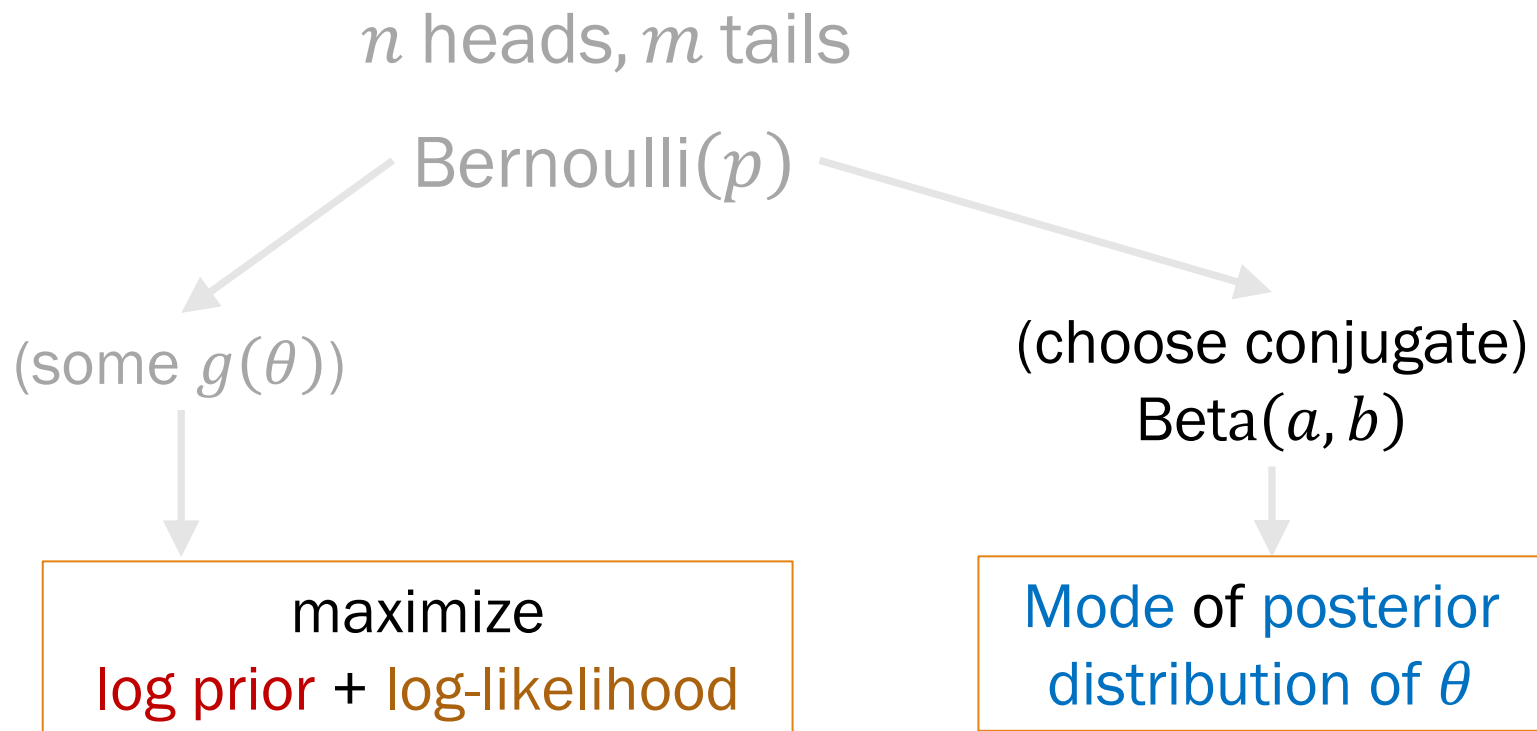# Proving the mode of Beta

Observe data                                  $n$ heads, $m$ tails

Choose model                              Bernoulli($p$)

Choose prior on $\theta$          (some $g(\theta)$)                (choose conjugate)
                                                                                              Beta($a, b$)

Find  $\theta_{MAP} =$
  $\arg\max\limits_{\theta} f(\theta | X_1, X_2, \ldots, X_n)$

$$\boxed{\begin{array}{c} \text{maximize} \\ \text{log prior} + \text{log-likelihood} \end{array}}$$

$$\boxed{\begin{array}{c} \text{Mode of posterior} \\ \text{distribution of } \theta \end{array}}$$

These are equivalent
interpretations of $\theta_{MAP}$.
We'll use this equivalence
to prove the mode of Beta.

$$\log g(\theta) + \sum_{i=1}^{n} \log f(X_i | \theta)$$

(posterior is also
conjugate)

- Differentiate, set to 0
- Solve

# From first principles: MAP for Bernoulli, conjugate prior

- Flip a coin 8 times. Observe $n = 7$ heads and $m = 1$ tail.
- Choose a prior on $\theta$. What is $\theta_{MAP}$?

Suppose we pick a prior $\theta \sim \text{Beta}(a, b)$. $g(\theta = p) = \frac{1}{\beta} p^{a-1}(1-p)^{b-1}$

normalizing constant, $\beta$

1. Determine log prior + log likelihood

$$\log g(\theta) + \log f(X_1, X_2, \ldots, X_n | \theta) = \log\left(\frac{1}{\beta} p^{a-1}(1-p)^{b-1}\right) + \log\left(\binom{n+m}{n} p^n (1-p)^m\right)$$

$$= \log \frac{1}{\beta} + (a-1)\log(p) + (b-1)\log(1-p) + \log\binom{n+m}{n} + n\log p + m\log(1-p)$$

2. Differentiate w.r.t. (each) $\theta$, set to 0

$$\frac{a-1}{p} + \frac{n}{p} - \frac{b-1}{1-p} - \frac{m}{1-p} = 0$$

3. Solve

(next slide)

# From first principles: MAP for Bernoulli, conjugate prior

- Flip a coin 8 times. Observe $n = 7$ heads and $m = 1$ tail.
- Choose a prior $\theta$. What is $\theta_{MAP}$?

Suppose we pick a prior $\theta \sim \text{Beta}(a, b)$. $g(\theta) = \frac{1}{\beta} p^{a-1}(1-p)^{b-1}$

normalizing constant, $\beta$

3. Solve for $p$
$$\frac{a-1}{p} + \frac{n}{p} - \frac{b-1}{1-p} - \frac{m}{1-p} = 0 \quad \text{(from previous slide)}$$

$$\implies \frac{a+n-1}{p} - \frac{b+m-1}{1-p} = 0$$

$$\theta_{MAP} = \frac{a+n-1}{a+n+b+m-2} \quad ✅$$

The mode of the posterior, $\text{Beta}(a+n, b+m)$!

If we choose a conjugate prior, we avoid calculus with MAP: just report mode of posterior.

# (live)

# 22: MAP

Lisa Yan and Jerry Cain

November 2, 2020

# Maximum A Posteriori (MAP) Estimator

Consider a sample of $n$ i.i.d. random variables $X_1, X_2, \ldots, X_n$ (data).

**Maximum Likelihood Estimator (MLE)**

What is the parameter $\theta$ that **maximizes the likelihood** of our observed data $(X_1, X_2, \ldots, X_n)$?

$$L(\theta) = f(X_1, X_2, \ldots, X_n | \theta)$$
$$= \prod_{i=1}^{n} f(X_i | \theta)$$

$$\theta_{MLE} = \arg \max_{\theta} f(X_1, X_2, \ldots, X_n | \theta)$$

likelihood of data

---

**Maximum a Posteriori (MAP) Estimator**

Given our observed data $(X_1, X_2, \ldots, X_n)$, what is the **most likely parameter** $\theta$?

Bayes rule

$$\theta_{MAP} = \arg \max_{\theta} f(\theta | X_1, X_2, \ldots, X_n)$$

posterior distribution of $\theta$

# How does MAP work?

Observe data

$Uni(0,1) \rightsquigarrow Beta(1,1)$

Choose model with parameter $\theta$

$Bernoulli / Binomial$

Choose prior on $\theta$

$g(\theta)$

Two valid approaches to computing $\theta_{MAP}$

Find $\theta_{MAP} = \arg\max_{\theta} f(\theta | X_1, X_2, \ldots, X_n)$

$$= \arg\max_{\theta} \left( \log g(\theta) + \sum_{i=1}^{n} \log f(X_i | \theta) \right)$$

Mode of posterior distribution of $\theta$

or

maximize
log prior + log-likelihood

If we choose a conjugate prior, we avoid calculus with MAP: just report mode of posterior.

# Conjugate distributions

# Quick MAP for Bernoulli and Binomial

Beta$(a, b)$ is a conjugate prior for the probability of success in Bernoulli and Binomial distributions.

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

**Prior**     Beta$(a, b)$
Saw $a + b - 2$ imaginary trials: $a - 1$ successes, $b - 1$ failures

**Experiment**  Observe $n + m$ new trials: $n$ successes, $m$ failures

**Posterior**  Beta$(a + n, b + m)$

MAP:    $$p = \frac{a + n - 1}{a + b + n + m - 2}$$

# Conjugate distributions

MAP estimator:

$$\theta_{MAP} = \arg \max_{\theta} f(\theta | X_1, X_2, \ldots, X_n)$$

The mode of the posterior distribution of $\theta$

| Distribution parameter | Conjugate distribution |
|---|---|
| Bernoulli $p$ | Beta |
| Binomial $p$ | Beta |
| Multinomial $p_i$ | Dirichlet |
| Poisson $\lambda$ | Gamma |
| Exponential $\lambda$ | Gamma |
| Normal $\mu$ | Normal |
| Normal $\sigma^2$ | Inverse Gamma |

Don't need to know Inverse Gamma... but it will know you ☺

CS109: We'll only focus on MAP for Bernoulli/Binomial $p$, Multinomial $p_i$, and Poisson $\lambda$.

# Multinomial is Multiple times the fun

$x_1^a x_2^b$
$x_2 = 1 - x_1$

Dirichlet$(a_1, a_2, \ldots, a_m)$ is a conjugate for Multinomial.

- Generalizes Beta in the same way Multinomial generalizes Bernoulli/Binomial:

$$f(x_1, x_2, \ldots, x_m) = \frac{1}{B(a_1, a_2, \ldots, a_m)} \prod_{i=1}^{m} x_i^{a_i - 1}$$

$(a_1 - 1) + (a_2 - 1) + (\ldots) + (a_m - 1)$

**Prior**
Dirichlet$(a_1, a_2, \ldots, a_m)$
Saw $\left(\sum_{i=1}^{m} a_i\right) - m$ imaginary trials, with $a_i - 1$ of outcome $i$

**Experiment** Observe $n_1 + n_2 + \cdots + n_m$ new trials, with $n_i$ of outcome $i$

**Posterior** Dirichlet$(a_1 + n_1, a_2 + n_2, \ldots, a_m + n_m)$

$m$

MAP:
$$p_i = \frac{a_i + n_i - 1}{\left(\sum_{i=1}^{m} a_i\right) + \left(\sum_{i=1}^{m} n_i\right) - m}$$

# Good times with Gamma



Gamma$(\alpha, \beta)$ is a conjugate for Poisson.

- Also conjugate for Exponential, but we won't delve into that

- Mode of gamma: $(\alpha - 1)/\beta$

  *(handwritten: $\frac{a-1}{2}$)*

**Prior**  $\theta \sim \text{Gamma}(\alpha, \beta) = \dfrac{\beta^{\alpha} x^{\alpha-1} e^{-\beta x}}{\Gamma(\alpha)}$

  *(handwritten: $\Gamma(\alpha) = (\alpha-1)!$)*

Saw $\alpha - 1$ total imaginary events during $\beta$ prior time periods

**Experiment**  Observe $n$ events during next $k$ time periods

**Posterior**  $(\theta | n \text{ events in } k \text{ periods}) \sim \text{Gamma}(\alpha + n, \beta + k)$

MAP:  $\theta_{MAP} = \dfrac{a + n - 1}{\beta + k}$

*(handwritten on graph: $\text{Gamma}(\alpha, \beta) = \dfrac{\beta^{\alpha} x^{\alpha-1} e^{-\beta x}}{\Gamma(\alpha)}$  ignore)*

# MAP for Poisson

Let $\lambda$ be the average # of successes in a time period.

1. What does it mean to have a prior of $\theta \sim$ Gamma(11,5)?

Observe 10 imaginary events in 5 time periods, i.e., observe at Poisson rate = 2

Now perform the experiment and see 11 events in next 2 time periods.

2. Given your prior, what is the posterior distribution?

3. What is $\theta_{MAP}$?

# MAP for Poisson

Let $\lambda$ be the average # of successes in a time period.

1. What does it mean to have a prior of $\theta \sim$ Gamma(11,5)?

Observe 10 imaginary events in 5 time periods, i.e., observe at Poisson rate = 2

Now perform the experiment and see 11 events in next 2 time periods.

2. Given your prior, what is the posterior distribution?

$(\theta | n$ events in $k$ periods$) \sim$ Gamma($\cancel{22}$, 7)

$21-1$

3. What is $\theta_{MAP}$?

$\frac{22 - 1}{7} = 3$

$\theta_{MAP} = 3$, the updated Poisson rate

$R$

Gauss, Dirichlet, Laplace

group theory

Abel

$|\epsilon| < 0$

# Interlude for jokes/announcements

# Announcements

Quiz 2 Grades Released Soon

Wednesday's Lecture: Optional

No Discussion Section This Week!

Lisa and I Still Have Wednesday OH!



https://en.wikipedia.org/wiki/Quicksort

# Choosing hyperparameters for conjugate prior
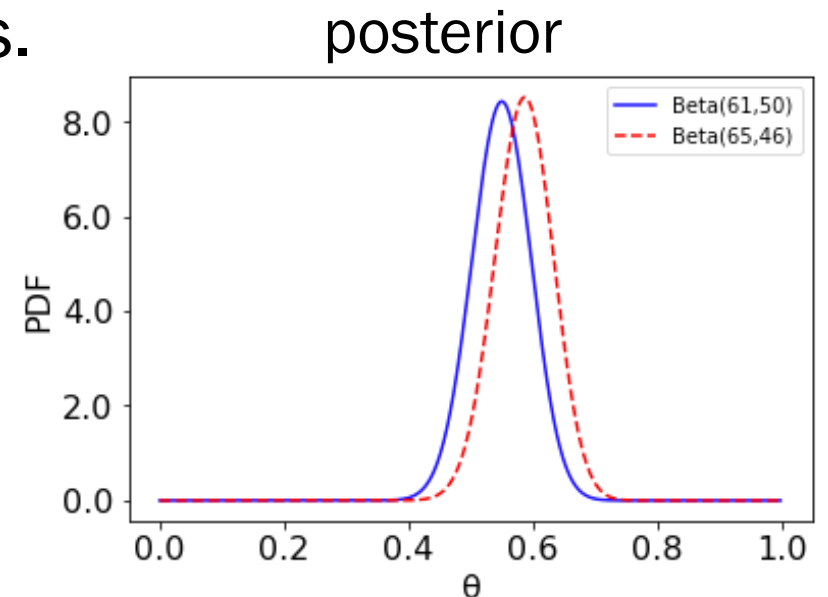
# Where'd you get them priors?

- Let $\theta$ be the probability a coin turns up heads.
- Model $\theta$ with 2 different priors:
  - Prior 1: Beta(3,8): 2 imaginary heads, 7 imaginary tails    mode: $\frac{2}{9}$
  - Prior 2: Beta(7,4): 6 imaginary heads, 3 imaginary tails    mode: $\frac{6}{9}$



prior

Now flip 100 coins and get 58 heads and 42 tails.
1. What are the two posterior distributions?
2. What are the modes of the two posterior distributions?

# Where'd you get them priors?

- Let $\theta$ be the probability a coin turns up heads.
- Model $\theta$ with 2 different priors:
  - Prior 1: Beta(3,8): 2 imaginary heads, 7 imaginary tails    mode: $\frac{2}{9}$
  - Prior 2: Beta(7,4): 6 imaginary heads, 3 imaginary tails    mode: $\frac{6}{9}$


prior

Now flip 100 coins and get 58 heads and 42 tails.

Posterior 1: Beta(61,50)    mode: $\frac{60}{109}$

Posterior 2: Beta(65,46)    mode: $\frac{64}{109}$

Provided we collect enough data, posteriors will converge to the true value.


posterior

Lisa Yan and Jerry Cain, CS109, 2020

# Laplace smoothing

MAP with **Laplace smoothing**:    a prior which represents $k$ imagined observations of each outcome.

- Categorical data (i.e., Multinomial, Bernoulli/Binomial)

- Also known as additive smoothing

**Laplace estimate**        Imagine $k = 1$ of each outcome
(follows from Laplace's "law of succession")

Example:        Laplace estimate for coin probabilities from aforementioned experiment (100 coins: 58 heads, 42 tails)

heads    $\dfrac{59}{102}$        tails    $\dfrac{43}{102}$

Laplace smoothing:
- Easy to implement/remember

# Back to our happy Laplace

Consider our previous 6-sided die.

- Roll the dice $n = 12$ times.
- Observe: 3 ones, 2 twos, 0 threes, 3 fours, 1 fives, 3 sixes

Recall $\theta_{MLE}$: $\quad p_1 = 3/12, p_2 = 2/12, \textcolor{red}{p_3 = 0/12}, \quad \triangle$
$\qquad\qquad\qquad p_4 = 3/12, p_5 = 1/12, p_6 = 3/12$

What are your Laplace estimates for each roll outcome?

# Back to our happy Laplace

Consider our previous 6-sided die.

- Roll the dice $n = 12$ times.
- Observe: 3 ones, 2 twos, 0 threes, 3 fours, 1 fives, 3 sixes

Recall $\theta_{MLE}$:
$p_1 = 3/12, p_2 = 2/12, {\color{red}p_3 = 0/12}, \quad \triangle !$
$p_4 = 3/12, p_5 = 1/12, p_6 = 3/12$

What are your Laplace estimates for each roll outcome?

$$p_i = \frac{X_i + 1}{n + m}$$

$p_1 = 4/18, p_2 = 3/18, {\color{green}p_3 = 1/18}, \quad ✅$
$p_4 = 4/18, p_5 = 2/18, p_6 = 4/18$

Laplace smoothing:
- Easy to implement/remember
- **Avoids estimating a parameter of 0**

# Bayesian Envelope Demo

# Two envelopes

Two envelopes: One contains $\$X$, the other contains $\$2X$.  $3X$

- Select an envelope and <u>open it</u>.
- Before opening the envelope, think either <u>equally</u> good. $\rightarrow +X$
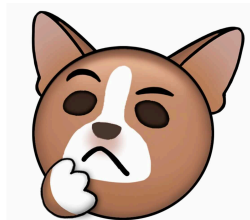
$-X \leftarrow$

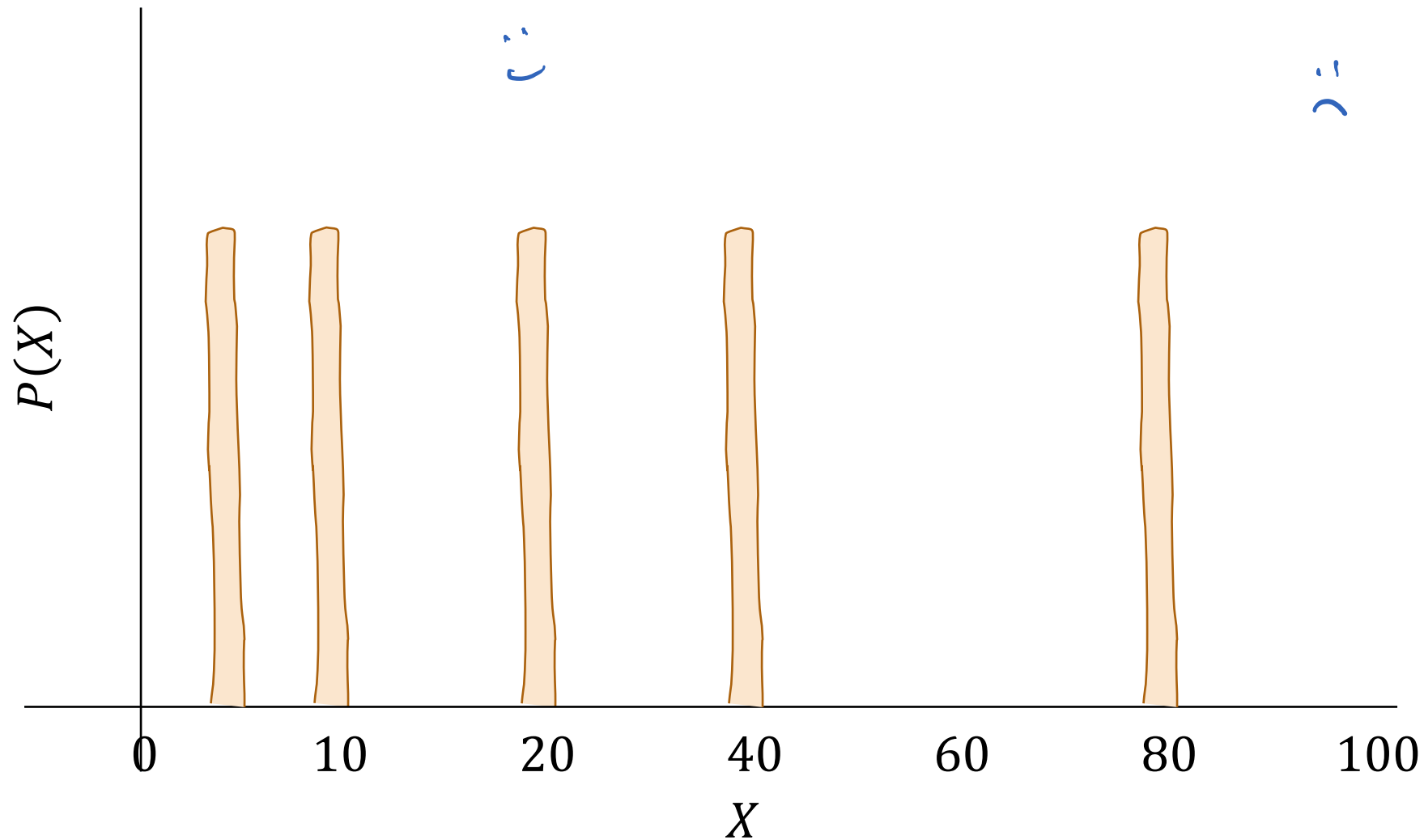Is the following reasoning valid?

- Let $Y = \$$ in envelope you selected.
- Let $Z = \$$ in other envelope.

$$E[Z|Y] = \frac{1}{2} \cdot \frac{Y}{2} + \frac{1}{2} \cdot 2Y = \frac{5}{4}Y$$

$$\frac{1}{2}(-x) + \frac{1}{2}(x) = 0$$

Follow-up: What happened by opening the envelope?

# Two envelopes

Two envelopes: One contains $\$X$, the other contains $\$2X$.
- Select an envelope and <u>open it</u>.
- Before opening the envelope, think either <u>equally</u> good.

**Is the following reasoning valid?**
- Let $Y = \$$ in envelope you selected.
- Let $Z = \$$ in other envelope.

$$E[Z|Y] = \frac{1}{2} \cdot \frac{Y}{2} + \frac{1}{2} \cdot 2Y = \frac{5}{4} Y$$

- Assumes all values of $X$ (where $0 < X < \infty$) equally likely
- Infinitely many values of $X$
- So not a true probability distribution over $X$ (does not integrate to 1)

**Follow-up:** What happened by opening the envelope?

# Are all values equally likely?



Infinite powers of two times 10

# Two envelopes: The subjectivity of probability

Your belief about the content of envelopes:
- Since implied distribution over $X$ is not a true probability distribution, what *is* our distribution over $X$?

## Frequentist

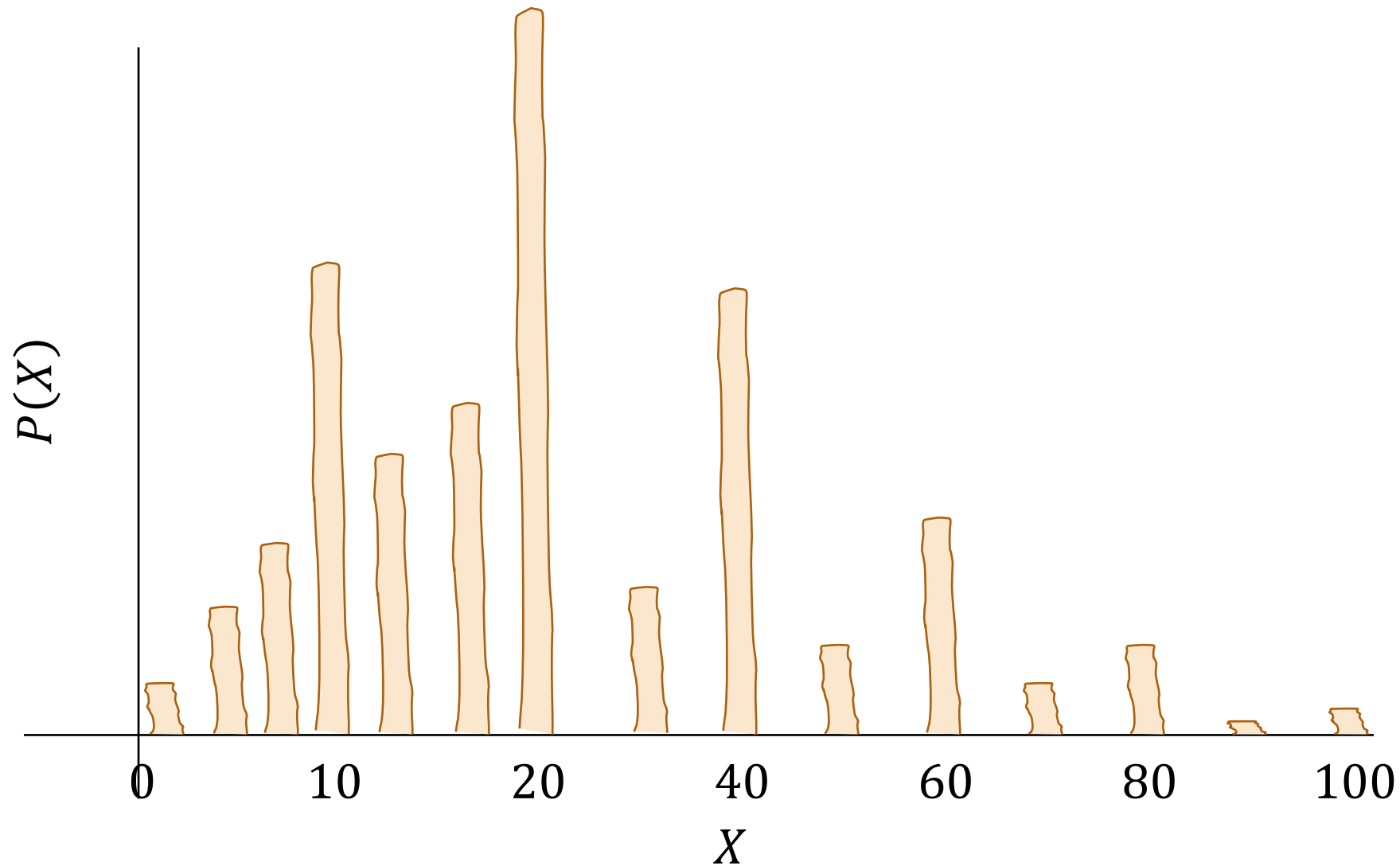Play game infinitely many times, see how often different values come up

Dilemma: You can only play the game once!

## Bayesian

Have <u>prior</u> belief of distribution of $X$
- Prior belief is a *subjective* probability
- Allows us to answer questions with limited data, or even no data at all
- As we run more experiments, all prior beliefs are eclipsed by data

# Two envelopes: The subjectivity of probability

# The envelope, please

Bayesian: Have a prior distribution over $X$, $P(X)$
- Let $Y = \$$ in envelope you selected. Open envelope to determine $Y$.
- Let $Z = \$$ in other envelope.

If $Y > E[Z|Y]$, keep your envelope, otherwise switch.   No inconsistency!!
- Opening envelope provides data to compute $P(X|Y)$
- …which allows you to compute $E[Z|Y]$

Of course, need to think about your prior distribution over $X$, but…

Bayesian probability: It doesn't matter how you construct your prior, but you **must** have one (whatever it is)

Imagine if envelope you opened contained $20.01. Should you switch?

# How much is a half cent?

# Have a wonderful Monday!