

# **Brief review of elementary statistics: parameter estimation, confidence intervals, hypothesis testing**

Roger Levy

9.S916: Statistical data analysis for scientific inference in cognitive science

11 April 2024

# Running example

---

# Running example

---

- I'm about to join a game of betting on the heads/tails outcome of a potentially bent coin

# Running example

---

- I'm about to join a game of betting on the heads/tails outcome of a potentially bent coin
- I can't inspect the coin, but I can watch the coin being flipped "for a while" and record the outcomes

# Running example

---

- I'm about to join a game of betting on the heads/tails outcome of a potentially bent coin
- I can't inspect the coin, but I can watch the coin being flipped "for a while" and record the outcomes
- The coin flips constitute a sequence of **Bernoulli random variables** conditionally independent of each other given the coin weighting  $P(\text{heads}) = \pi$  with  $0 \leq \pi \leq 1$

# Running example

---

- I'm about to join a game of betting on the heads/tails outcome of a potentially bent coin
- I can't inspect the coin, but I can watch the coin being flipped "for a while" and record the outcomes
- The coin flips constitute a sequence of **Bernoulli random variables** conditionally independent of each other given the coin weighting  $P(\text{heads}) = \pi$  with  $0 \leq \pi \leq 1$
- Figuring out from observed data what the weighting is likely to be is **parameter estimation**

# Running example

---

- I'm about to join a game of betting on the heads/tails outcome of a potentially bent coin
- I can't inspect the coin, but I can watch the coin being flipped "for a while" and record the outcomes
- The coin flips constitute a sequence of **Bernoulli random variables** conditionally independent of each other given the coin weighting  $P(\text{heads}) = \pi$  with  $0 \leq \pi \leq 1$
- Figuring out from observed data what the weighting is likely to be is **parameter estimation**
- In general, here we will use  $\mathbf{y}$  to refer to observed-outcome **data** and  $\theta$  to refer to the model parameters to be estimated

# Characteristics of estimators

---



# Characteristics of estimators

---

- **Estimator:** a procedure for guessing a quantity of interest within a population from a sample from that population

# Characteristics of estimators

---

- **Estimator:** a procedure for guessing a quantity of interest within a population from a sample from that population
- For example, the **relative frequency estimator:** if we observe  $r$  instances of heads in  $n$  coin flips,

$$\hat{\pi} = \frac{r}{n}$$

# Characteristics of estimators

---

- **Estimator:** a procedure for guessing a quantity of interest within a population from a sample from that population
- For example, the **relative frequency estimator:** if we observe  $r$  instances of heads in  $n$  coin flips,

"this is an estimator"   $\hat{\pi} = \frac{r}{n}$

# Characteristics of estimators

---

- **Estimator:** a procedure for guessing a quantity of interest within a population from a sample from that population
- For example, the **relative frequency estimator:** if we observe  $r$  instances of heads in  $n$  coin flips,

"this is an estimator" 

$$\hat{\pi} = \frac{r}{n}$$

- Data are stochastic, so estimators give random variables!

# Characteristics of estimators

---

- **Estimator:** a procedure for guessing a quantity of interest within a population from a sample from that population
- For example, the **relative frequency estimator:** if we observe  $r$  instances of heads in  $n$  coin flips,

"this is an estimator"   $\hat{\pi} = \frac{r}{n}$

- Data are stochastic, so estimators give random variables!

so  $\hat{\pi}$  is **unbiased**

# Characteristics of estimators

---

- **Estimator:** a procedure for guessing a quantity of interest within a population from a sample from that population
- For example, the **relative frequency estimator:** if we observe  $r$  instances of heads in  $n$  coin flips,

"this is an estimator"   $\hat{\pi} = \frac{r}{n}$

- Data are stochastic, so estimators give random variables!

$$E[\hat{\pi}] = E\left[\frac{r}{n}\right] = \frac{1}{n}E[r] = \frac{r}{n} = \pi \quad \text{so } \hat{\pi} \text{ is unbiased}$$

# Characteristics of estimators

---

- **Estimator:** a procedure for guessing a quantity of interest within a population from a sample from that population
- For example, the **relative frequency estimator:** if we observe  $r$  instances of heads in  $n$  coin flips,

"this is an estimator" 

$$\hat{\pi} = \frac{r}{n}$$

- Data are stochastic, so estimators give random variables!
- **Bias** of an estimator is  $E[\hat{\theta}] - \theta$

$$E[\hat{\pi}] = E\left[\frac{r}{n}\right] = \frac{1}{n}E[r] = \frac{r}{n} = \pi \quad \text{so } \hat{\pi} \text{ is unbiased}$$

# Characteristics of estimators

---

- **Estimator:** a procedure for guessing a quantity of interest within a population from a sample from that population
- For example, the **relative frequency estimator:** if we observe  $r$  instances of heads in  $n$  coin flips,

"this is an estimator"   $\hat{\pi} = \frac{r}{n}$

- Data are stochastic, so estimators give random variables!

- **Bias** of an estimator is  $E[\hat{\theta}] - \theta$

$$E[\hat{\pi}] = E\left[\frac{r}{n}\right] = \frac{1}{n}E[r] = \frac{r}{n} = \pi \quad \text{so } \hat{\pi} \text{ is unbiased}$$

- **Variance** of an estimator is ordinary variance



# Characteristics of estimators

---

- **Estimator:** a procedure for guessing a quantity of interest within a population from a sample from that population
- For example, the **relative frequency estimator:** if we observe  $r$  instances of heads in  $n$  coin flips,

"this is an estimator"   $\hat{\pi} = \frac{r}{n}$

- Data are stochastic, so estimators give random variables!

- **Bias** of an estimator is  $E[\hat{\theta}] - \theta$

$$E[\hat{\pi}] = E\left[\frac{r}{n}\right] = \frac{1}{n}E[r] = \frac{r}{n} = \pi \quad \text{so } \hat{\pi} \text{ is unbiased}$$

- **Variance** of an estimator is ordinary variance

$$\text{Var}(X) \equiv E[(X - E[X])^2]$$

# Characteristics of estimators

---

- **Estimator:** a procedure for guessing a quantity of interest within a population from a sample from that population
- For example, the **relative frequency estimator:** if we observe  $r$  instances of heads in  $n$  coin flips,

"this is an estimator"   $\hat{\pi} = \frac{r}{n}$

- Data are stochastic, so estimators give random variables!

- **Bias** of an estimator is  $E[\hat{\theta}] - \theta$

$$E[\hat{\pi}] = E\left[\frac{r}{n}\right] = \frac{1}{n}E[r] = \frac{r}{n} = \pi \quad \text{so } \hat{\pi} \text{ is unbiased}$$

- **Variance** of an estimator is ordinary variance

$$\text{Var}(X) \equiv E[(X - E[X])^2] \quad \text{Var}(\hat{\pi}) = \frac{\pi(1 - \pi)}{n} \quad (\text{see reading materials})$$

# Characteristics of estimators

---

- **Estimator:** a procedure for guessing a quantity of interest within a population from a sample from that population
- For example, the **relative frequency estimator:** if we observe  $r$  instances of heads in  $n$  coin flips,

"this is an estimator"   $\hat{\pi} = \frac{r}{n}$

- Data are stochastic, so estimators give random variables!

- **Bias** of an estimator is  $E[\hat{\theta}] - \theta$

$$E[\hat{\pi}] = E\left[\frac{r}{n}\right] = \frac{1}{n}E[r] = \frac{r}{n} = \pi \quad \text{so } \hat{\pi} \text{ is unbiased}$$

- **Variance** of an estimator is ordinary variance

$$\text{Var}(X) \equiv E[(X - E[X])^2] \quad \text{Var}(\hat{\pi}) = \frac{\pi(1 - \pi)}{n} \quad (\text{see reading materials})$$

- Good estimators have favorable **bias–variance** tradeoff

# Maximum likelihood estimation

---

$$\text{Lik}(\boldsymbol{\theta}; \mathbf{y}) \equiv P(\mathbf{y}|\boldsymbol{\theta}) \quad \hat{\boldsymbol{\theta}}_{MLE} \stackrel{\text{def}}{=} \arg \max_{\boldsymbol{\theta}} \text{Lik}(\boldsymbol{\theta}; \mathbf{y})$$

$i$	$y_i$
1	T
2	T
3	H
4	T

*(repeat slide from lecture 3)*

# Maximum likelihood estimation

$$\text{Lik}(\boldsymbol{\theta}; \mathbf{y}) \equiv P(\mathbf{y}|\boldsymbol{\theta}) \quad \hat{\boldsymbol{\theta}}_{MLE} \stackrel{\text{def}}{=} \arg \max_{\boldsymbol{\theta}} \text{Lik}(\boldsymbol{\theta}; \mathbf{y})$$

- $p$  refers to the value of  $P(\text{coin toss}_i = \text{Heads})$
- Likelihood for the following dataset

$i$	$y_i$
1	T
2	T
3	H
4	T

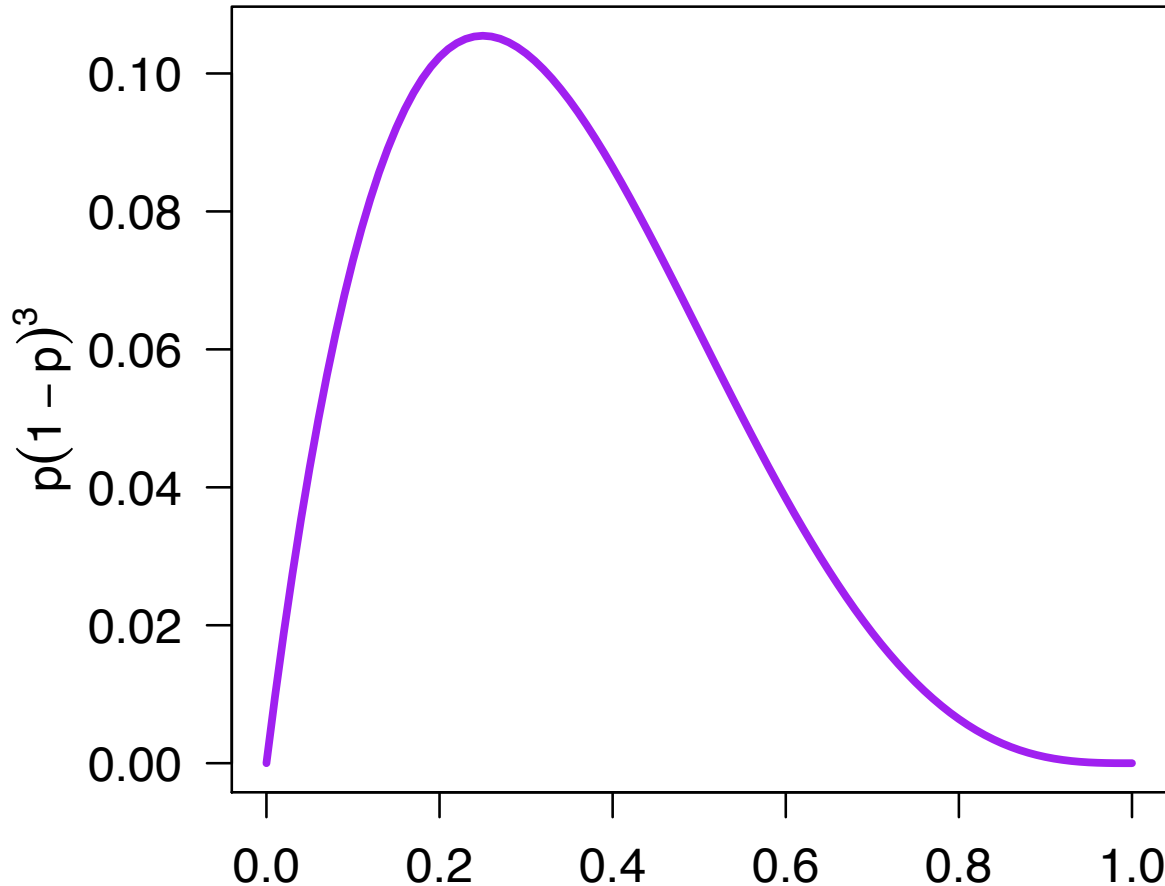
*(repeat slide from lecture 3)*

# Maximum likelihood estimation

$$\text{Lik}(\boldsymbol{\theta}; \mathbf{y}) \equiv P(\mathbf{y}|\boldsymbol{\theta}) \quad \hat{\boldsymbol{\theta}}_{MLE} \stackrel{\text{def}}{=} \arg \max_{\boldsymbol{\theta}} \text{Lik}(\boldsymbol{\theta}; \mathbf{y})$$

$i$	$y_i$
1	T
2	T
3	H
4	T

- $p$  refers to the value of  $P(\text{coin toss}_i = \text{Heads})$
- Likelihood for the following dataset



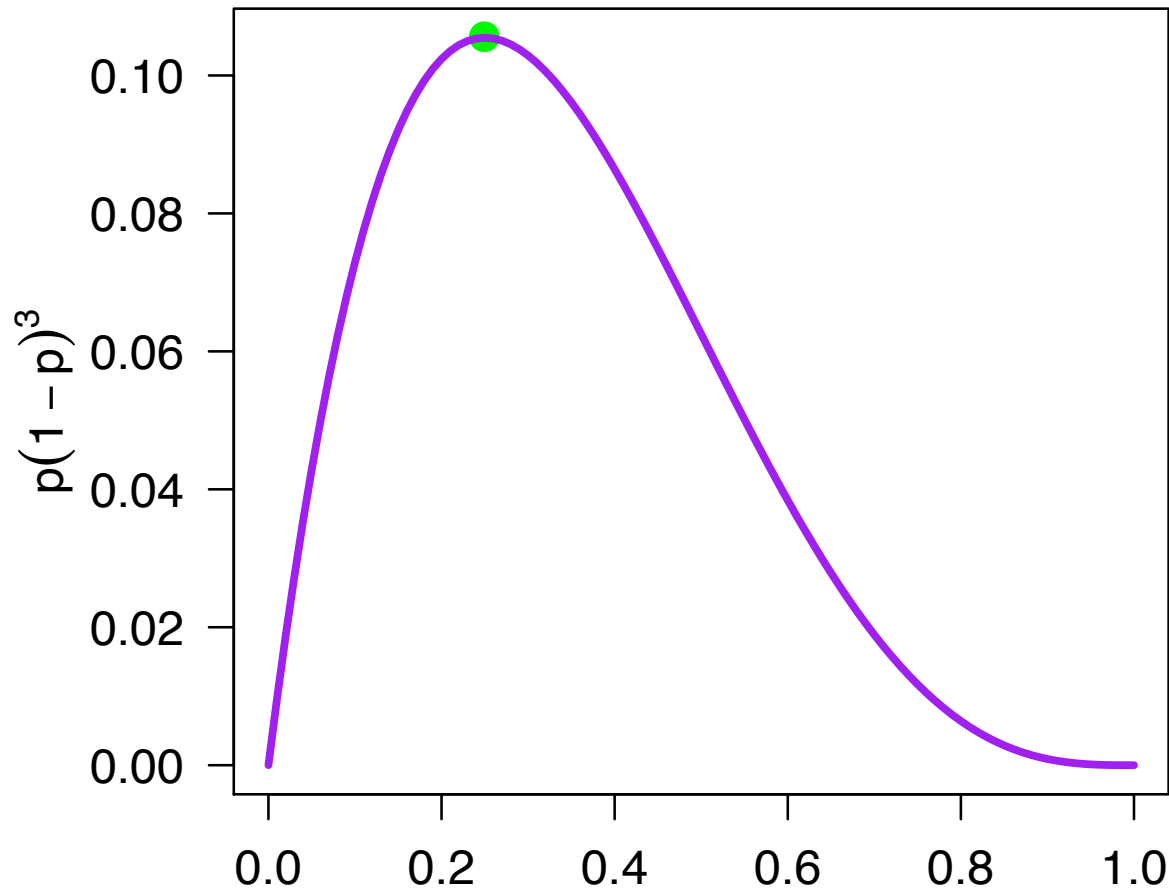
(repeat slide from lecture 3)

# Maximum likelihood estimation

$$\text{Lik}(\boldsymbol{\theta}; \mathbf{y}) \equiv P(\mathbf{y}|\boldsymbol{\theta}) \quad \hat{\boldsymbol{\theta}}_{MLE} \stackrel{\text{def}}{=} \arg \max_{\boldsymbol{\theta}} \text{Lik}(\boldsymbol{\theta}; \mathbf{y})$$

$i$	$y_i$
1	T
2	T
3	H
4	T

- $p$  refers to the value of  $P(\text{coin toss}_i = \text{Heads})$
- Likelihood for the following dataset



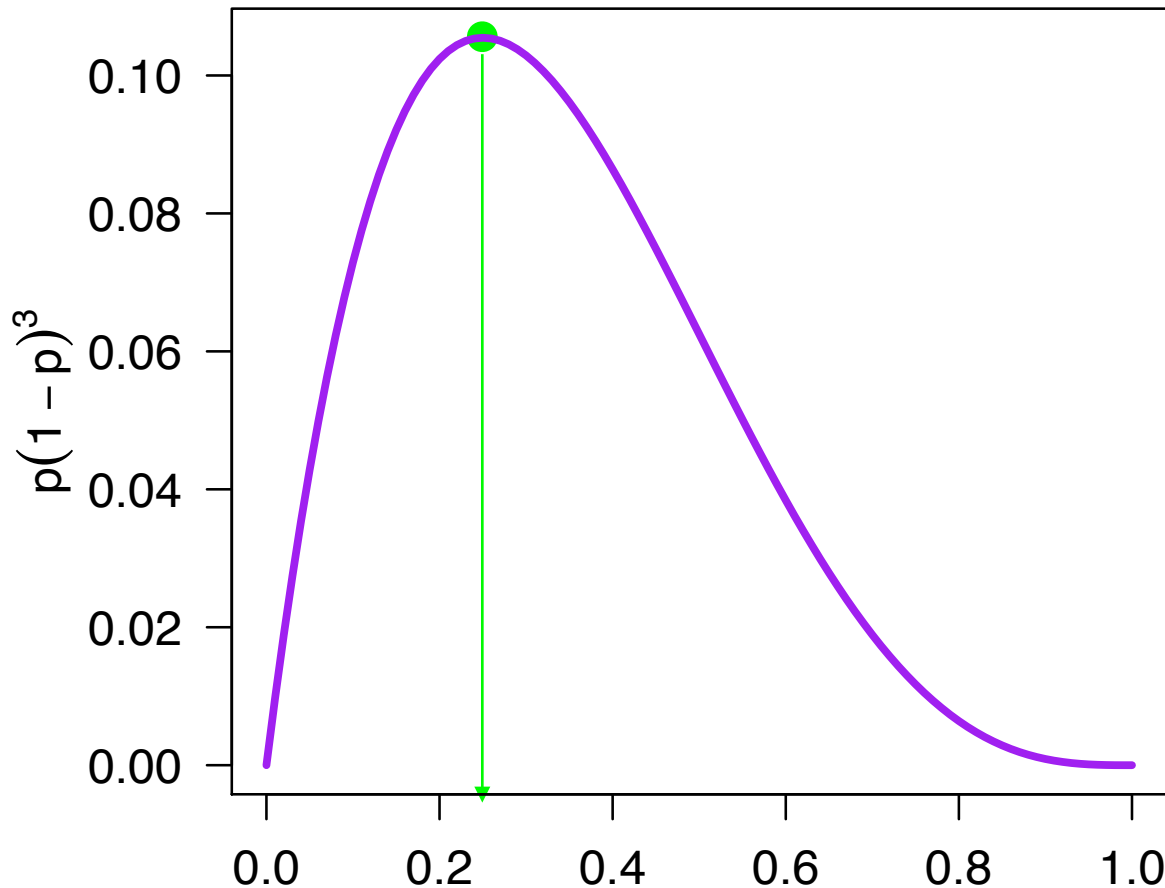
(repeat slide from lecture 3)

# Maximum likelihood estimation

$$\text{Lik}(\boldsymbol{\theta}; \mathbf{y}) \equiv P(\mathbf{y}|\boldsymbol{\theta}) \quad \hat{\boldsymbol{\theta}}_{MLE} \stackrel{\text{def}}{=} \arg \max_{\boldsymbol{\theta}} \text{Lik}(\boldsymbol{\theta}; \mathbf{y})$$

$i$	$y_i$
1	T
2	T
3	H
4	T

- $p$  refers to the value of  $P(\text{coin toss}_i = \text{Heads})$
- Likelihood for the following dataset



(repeat slide from lecture 3)

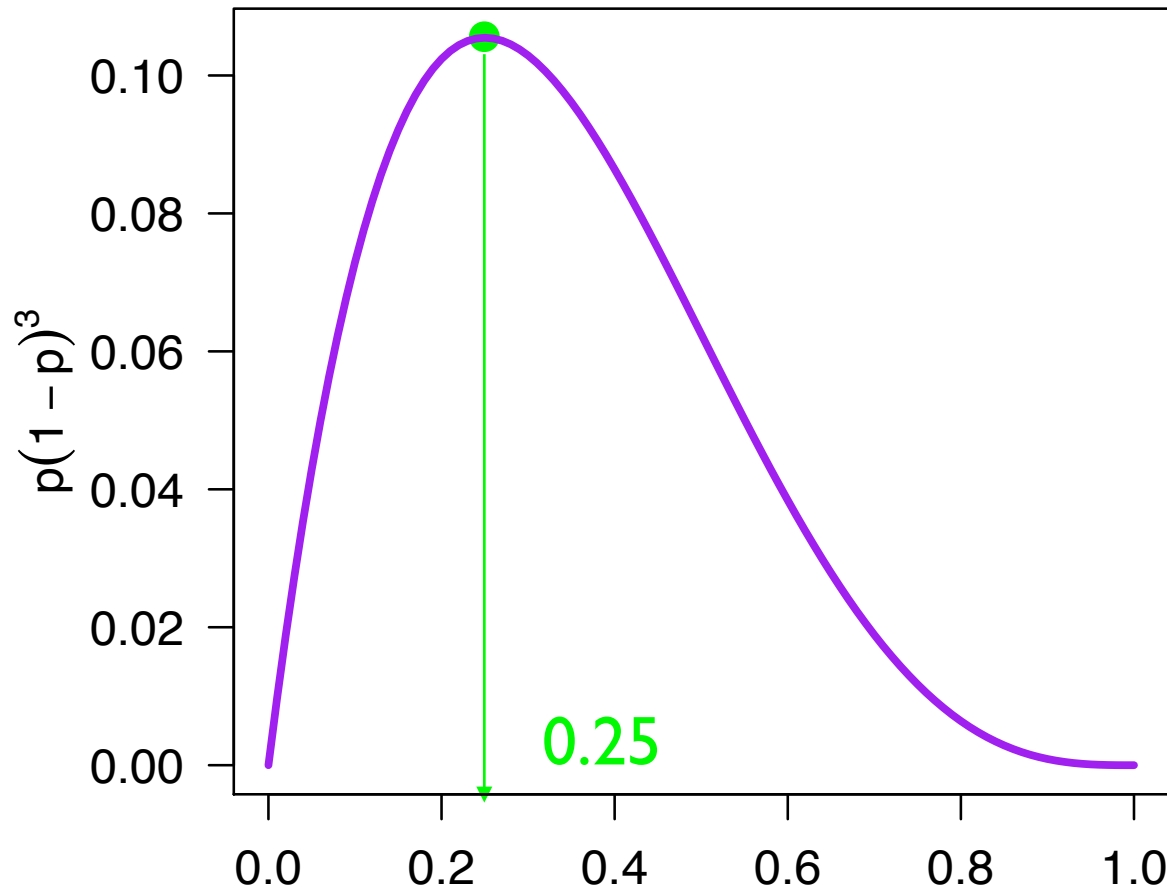


# Maximum likelihood estimation

$$\text{Lik}(\boldsymbol{\theta}; \mathbf{y}) \equiv P(\mathbf{y}|\boldsymbol{\theta}) \quad \hat{\boldsymbol{\theta}}_{MLE} \stackrel{\text{def}}{=} \arg \max_{\boldsymbol{\theta}} \text{Lik}(\boldsymbol{\theta}; \mathbf{y})$$

$i$	$y_i$
1	T
2	T
3	H
4	T

- $p$  refers to the value of  $P(\text{coin toss}_i = \text{Heads})$
- Likelihood for the following dataset



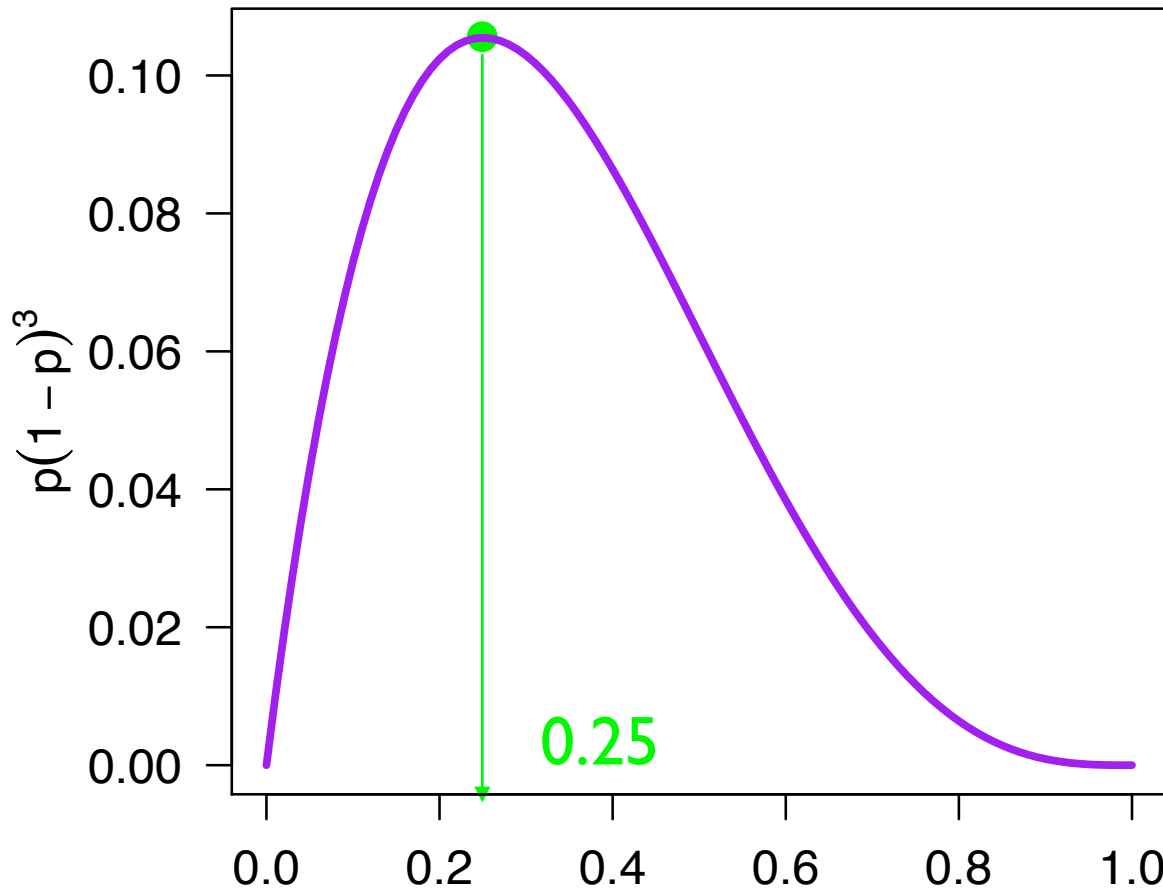
(repeat slide from lecture 3)

# Maximum likelihood estimation

$$\text{Lik}(\boldsymbol{\theta}; \mathbf{y}) \equiv P(\mathbf{y}|\boldsymbol{\theta}) \quad \hat{\boldsymbol{\theta}}_{MLE} \stackrel{\text{def}}{=} \arg \max_{\boldsymbol{\theta}} \text{Lik}(\boldsymbol{\theta}; \mathbf{y})$$

$i$	$y_i$
1	T
2	T
3	H
4	T

- $p$  refers to the value of  $P(\text{coin toss}_i = \text{Heads})$
- Likelihood for the following dataset



This is choosing the  
*maximum likelihood*  
estimate (**MLE**)

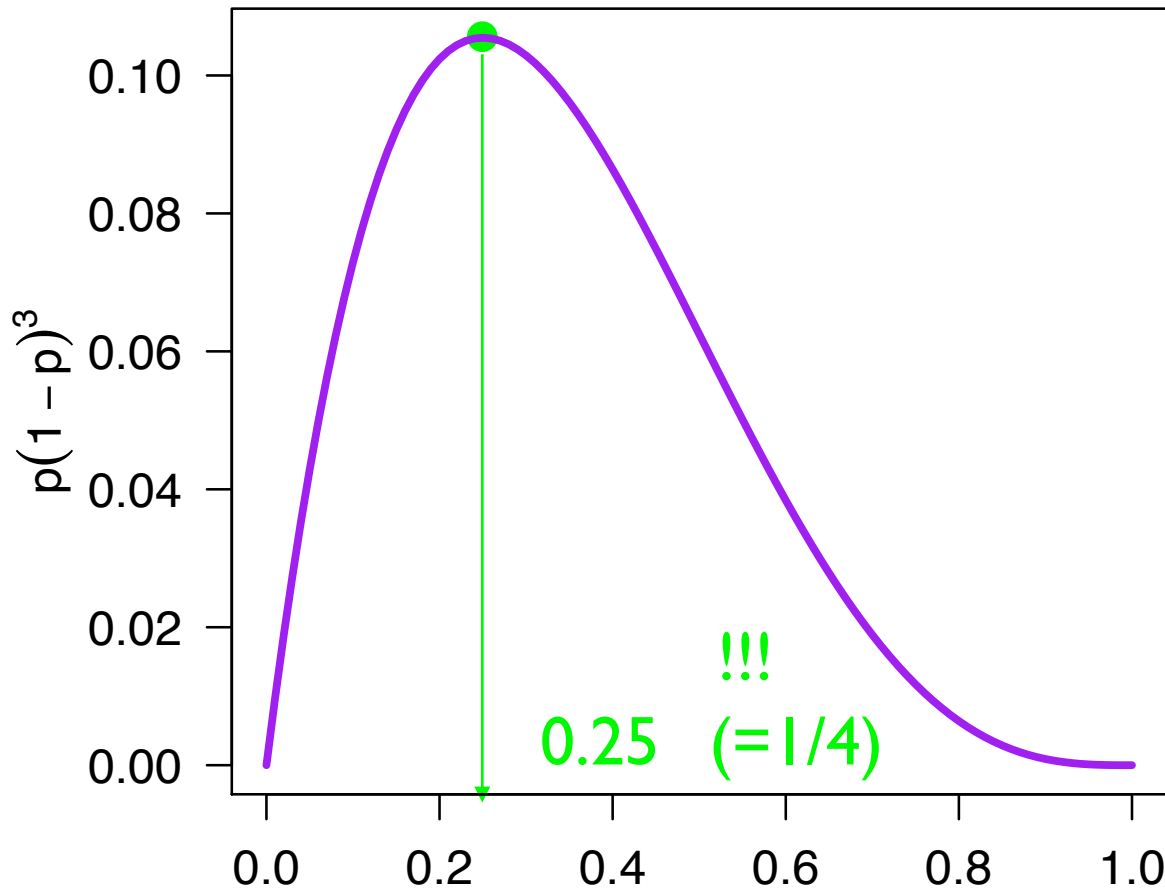
(repeat slide from lecture 3)

# Maximum likelihood estimation

$$\text{Lik}(\boldsymbol{\theta}; \mathbf{y}) \equiv P(\mathbf{y}|\boldsymbol{\theta}) \quad \hat{\boldsymbol{\theta}}_{MLE} \stackrel{\text{def}}{=} \arg \max_{\boldsymbol{\theta}} \text{Lik}(\boldsymbol{\theta}; \mathbf{y})$$

$i$	$y_i$
1	T
2	T
3	H
4	T

- $p$  refers to the value of  $P(\text{coin toss}_i = \text{Heads})$
- Likelihood for the following dataset



This is choosing the  
*maximum likelihood*  
estimate (**MLE**)

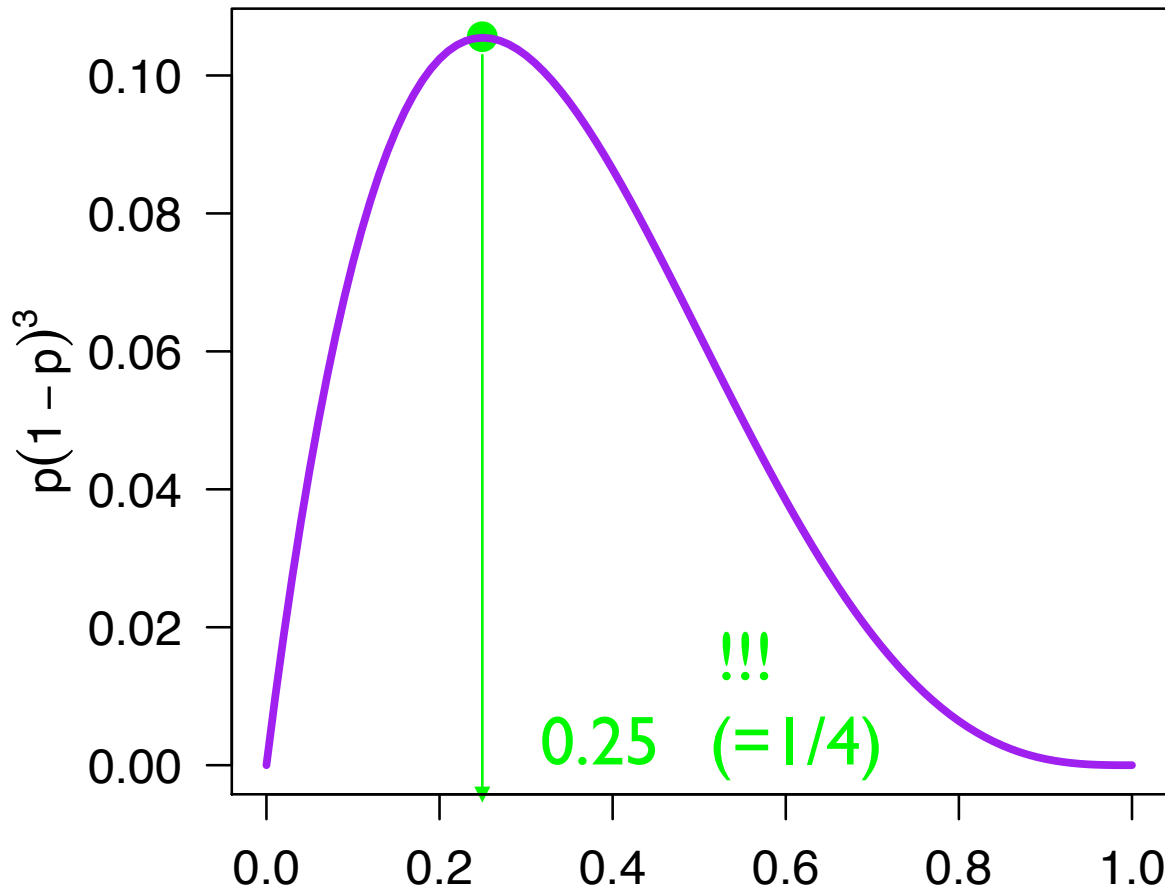
(repeat slide from lecture 3)

# Maximum likelihood estimation

$$\text{Lik}(\boldsymbol{\theta}; \mathbf{y}) \equiv P(\mathbf{y}|\boldsymbol{\theta}) \quad \hat{\boldsymbol{\theta}}_{MLE} \stackrel{\text{def}}{=} \arg \max_{\boldsymbol{\theta}} \text{Lik}(\boldsymbol{\theta}; \mathbf{y})$$

$i$	$y_i$
1	T
2	T
3	H
4	T

- $p$  refers to the value of  $P(\text{coin toss}_i = \text{Heads})$
- Likelihood for the following dataset



This is choosing the  
*maximum likelihood*  
estimate (**MLE**)

The **MLE** also turns  
out to be the *relative*  
*frequency estimate*  
(**RFE**)

(repeat slide from lecture 3)

# The binomial distribution

---

# The binomial distribution

---

- The **binomial distribution** is a two-parameter probability distribution over the number of **successes** in a number of independent, identically distributed (**iid**) **Bernoulli trials**

# The binomial distribution

---

- The **binomial distribution** is a two-parameter probability distribution over the number of **successes** in a number of independent, identically distributed (**iid**) **Bernoulli trials**
- Two parameters: Number of trials  $n$  & trial success parameter  $\pi$

# The binomial distribution

---

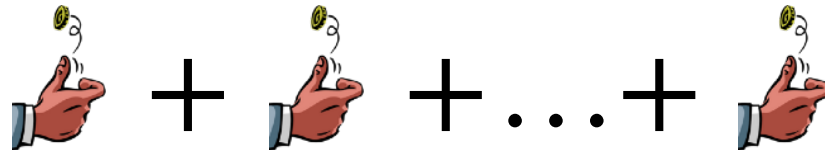
- The **binomial distribution** is a two-parameter probability distribution over the number of **successes** in a number of independent, identically distributed (**iid**) **Bernoulli trials**
- Two parameters: Number of trials  $n$  & trial success parameter  $\pi$
- A binomial-distributed random variable  $Y$  is simply the sum of  $n$  iid Bernoulli random variables with success parameter  $\pi$



# The binomial distribution

---

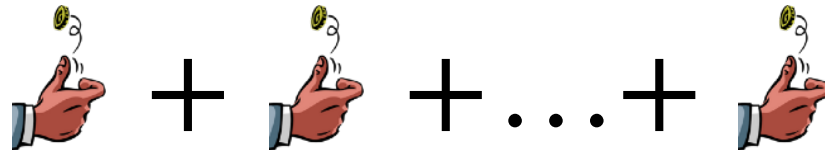
- The **binomial distribution** is a two-parameter probability distribution over the number of **successes** in a number of independent, identically distributed (**iid**) **Bernoulli trials**
- Two parameters: Number of trials  $n$  & trial success parameter  $\pi$
- A binomial-distributed random variable  $Y$  is simply the sum of  $n$  iid Bernoulli random variables with success parameter  $\pi$



# The binomial distribution

---

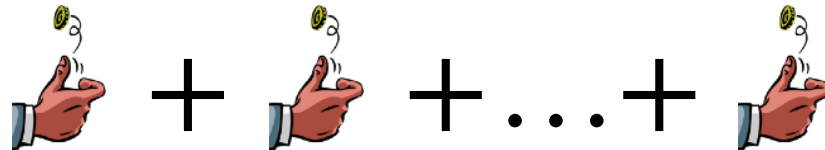
- The **binomial distribution** is a two-parameter probability distribution over the number of **successes** in a number of independent, identically distributed (**iid**) **Bernoulli trials**
- Two parameters: Number of trials  $n$  & trial success parameter  $\pi$
- A binomial-distributed random variable  $Y$  is simply the sum of  $n$  iid Bernoulli random variables with success parameter  $\pi$



- A binomial random variable has the following probability mass function:

# The binomial distribution

- The **binomial distribution** is a two-parameter probability distribution over the number of **successes** in a number of independent, identically distributed (**iid**) **Bernoulli trials**
- Two parameters: Number of trials  $n$  & trial success parameter  $\pi$
- A binomial-distributed random variable  $Y$  is simply the sum of  $n$  iid Bernoulli random variables with success parameter  $\pi$



- A binomial random variable has the following probability mass function:

$$P(Y = r) = \binom{n}{r} \pi^r (1 - \pi)^{n-r}$$

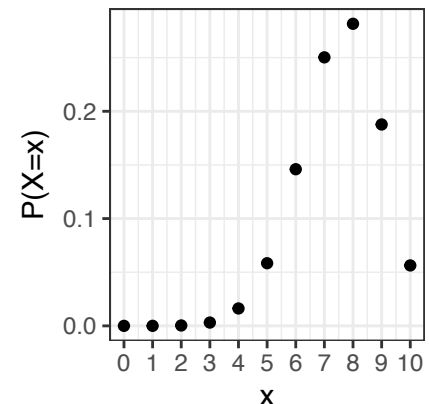
# The binomial distribution

- The **binomial distribution** is a two-parameter probability distribution over the number of **successes** in a number of independent, identically distributed (**iid**) **Bernoulli trials**
- Two parameters: Number of trials  $n$  & trial success parameter  $\pi$
- A binomial-distributed random variable  $Y$  is simply the sum of  $n$  iid Bernoulli random variables with success parameter  $\pi$



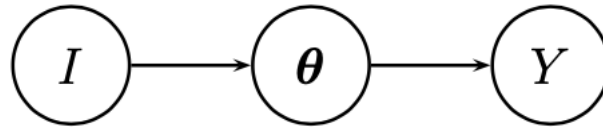
- A binomial random variable has the following probability mass function:

$$P(Y = r) = \binom{n}{r} \pi^r (1 - \pi)^{n-r}$$



# Bayesian parameter estimation

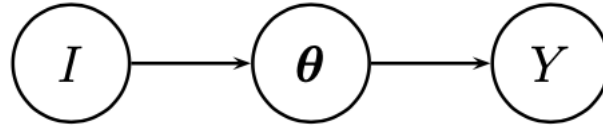
---



# Bayesian parameter estimation

---

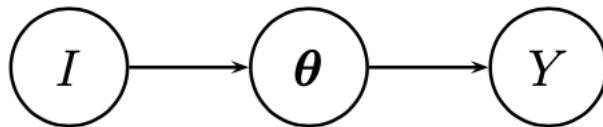
- Assume that the model parameters "intervene" between background knowledge  $I$  and data  $Y$ :



# Bayesian parameter estimation

---

- Assume that the model parameters "intervene" between background knowledge  $I$  and data  $Y$ :

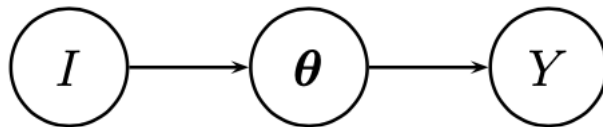


$$P(\boldsymbol{\theta}|\mathbf{y}, I) = \frac{P(\mathbf{y}|\boldsymbol{\theta}, I)P(\boldsymbol{\theta}|I)}{P(\mathbf{y}|I)}$$

# Bayesian parameter estimation

---

- Assume that the model parameters "intervene" between background knowledge  $I$  and data  $Y$ :



$$P(\boldsymbol{\theta}|\mathbf{y}, I) = \frac{P(\mathbf{y}|\boldsymbol{\theta}, I)P(\boldsymbol{\theta}|I)}{P(\mathbf{y}|I)}$$

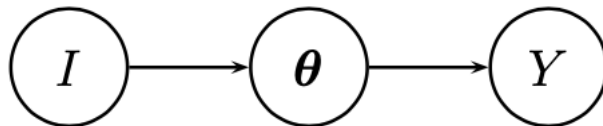
$$= \frac{P(\mathbf{y}|\boldsymbol{\theta}) P(\boldsymbol{\theta}|I)}{P(\mathbf{y}|I)} \quad (\text{because } \mathbf{y} \perp I \mid \boldsymbol{\theta})$$



# Bayesian parameter estimation

---

- Assume that the model parameters "intervene" between background knowledge  $I$  and data  $Y$ :



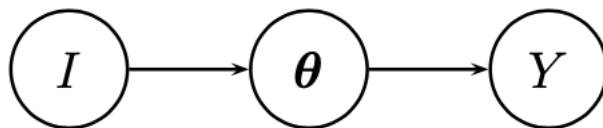
$$P(\theta|\mathbf{y}, I) = \frac{P(\mathbf{y}|\theta, I)P(\theta|I)}{P(\mathbf{y}|I)}$$

$$= \frac{\overbrace{P(\mathbf{y}|\theta)}^{\text{Likelihood for } \theta} P(\theta|I)}{P(\mathbf{y}|I)} \quad (\text{because } \mathbf{y} \perp I \mid \theta)$$

# Bayesian parameter estimation

---

- Assume that the model parameters "intervene" between background knowledge  $I$  and data  $Y$ :



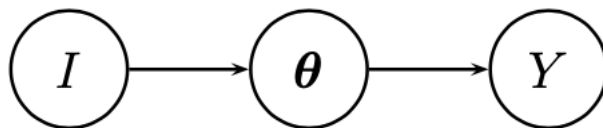
$$P(\boldsymbol{\theta}|\mathbf{y}, I) = \frac{P(\mathbf{y}|\boldsymbol{\theta}, I)P(\boldsymbol{\theta}|I)}{P(\mathbf{y}|I)}$$

$$= \frac{\overbrace{P(\mathbf{y}|\boldsymbol{\theta})}^{\text{Likelihood for } \boldsymbol{\theta}} \overbrace{P(\boldsymbol{\theta}|I)}^{\text{Prior over } \boldsymbol{\theta}}}{P(\mathbf{y}|I)} \quad (\text{because } \mathbf{y} \perp I \mid \boldsymbol{\theta})$$

# Bayesian parameter estimation

---

- Assume that the model parameters "intervene" between background knowledge  $I$  and data  $Y$ :

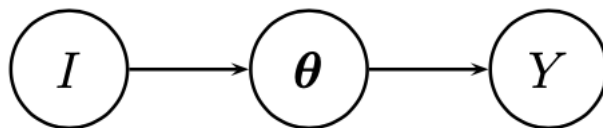


$$P(\boldsymbol{\theta}|\mathbf{y}, I) = \frac{P(\mathbf{y}|\boldsymbol{\theta}, I)P(\boldsymbol{\theta}|I)}{P(\mathbf{y}|I)}$$

$$= \frac{\overbrace{P(\mathbf{y}|\boldsymbol{\theta})}^{\text{Likelihood for } \boldsymbol{\theta}} \overbrace{P(\boldsymbol{\theta}|I)}^{\text{Prior over } \boldsymbol{\theta}}}{\underbrace{P(\mathbf{y}|I)}_{\text{Likelihood marginalized over } \boldsymbol{\theta}}} \quad (\text{because } \mathbf{y} \perp I \mid \boldsymbol{\theta})$$

# Bayesian parameter estimation

- Assume that the model parameters "intervene" between background knowledge  $I$  and data  $Y$ :



$$P(\theta | \mathbf{y}, I) = \frac{P(\mathbf{y} | \theta, I) P(\theta | I)}{P(\mathbf{y} | I)}$$

$$= \frac{\overbrace{P(\mathbf{y} | \theta)}^{\text{Likelihood for } \theta} \overbrace{P(\theta | I)}^{\text{Prior over } \theta}}{\underbrace{P(\mathbf{y} | I)}_{\text{Likelihood marginalized over } \theta}} \quad (\text{because } \mathbf{y} \perp I \mid \theta)$$

- Then, if we assume a parametric form for  $P(\mathbf{y} \mid \theta)$ , we just need the prior  $P(\theta \mid I)$

# Example for coin flips: the beta distribution

# Example for coin flips: the beta distribution

- Express background knowledge  $I$  as two "pseudo-count" parameters  $\alpha_1, \alpha_2$

# Example for coin flips: the beta distribution

- Express background knowledge  $I$  as two "pseudo-count" parameters  $\alpha_1, \alpha_2$
- The beta distribution has form

$$P(\pi|\alpha_1, \alpha_2) = \frac{1}{B(\alpha_1, \alpha_2)} \pi^{\alpha_1-1} (1 - \pi)^{\alpha_2-1}$$

# Example for coin flips: the beta distribution

- Express background knowledge  $I$  as two "pseudo-count" parameters  $\alpha_1, \alpha_2$
- The beta distribution has form

$$P(\pi|\alpha_1, \alpha_2) = \frac{1}{B(\alpha_1, \alpha_2)} \pi^{\alpha_1-1} (1-\pi)^{\alpha_2-1}$$

Where the action is!



# Example for coin flips: the beta distribution

- Express background knowledge  $I$  as two "pseudo-count" parameters  $\alpha_1, \alpha_2$
- The beta distribution has form

$$P(\pi|\alpha_1, \alpha_2) = \frac{1}{B(\alpha_1, \alpha_2)} \pi^{\alpha_1-1} (1-\pi)^{\alpha_2-1}$$

Normalizing constant,  
not of great interest for  
present purposes

**Where the action is!**

$$B(\alpha_1, \alpha_2) = \int_0^1 \pi^{\alpha_1-1} (1-\pi)^{\alpha_2-1} d\pi$$

# Example for coin flips: the beta distribution

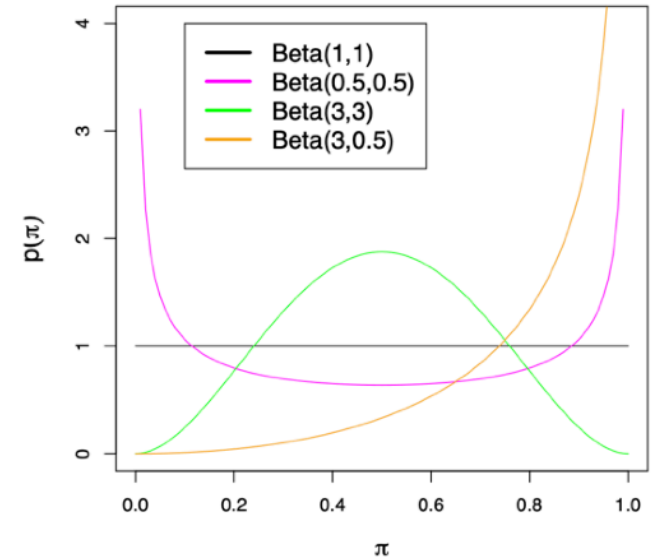
- Express background knowledge  $I$  as two "pseudo-count" parameters  $\alpha_1, \alpha_2$
- The beta distribution has form

$$P(\pi|\alpha_1, \alpha_2) = \frac{1}{B(\alpha_1, \alpha_2)} \pi^{\alpha_1-1} (1-\pi)^{\alpha_2-1}$$

Normalizing constant,  
not of great interest for  
present purposes

**Where the action is!**

$$B(\alpha_1, \alpha_2) = \int_0^1 \pi^{\alpha_1-1} (1-\pi)^{\alpha_2-1} d\pi$$



# Example for coin flips: the beta distribution

- Express background knowledge  $I$  as two "pseudo-count" parameters  $\alpha_1, \alpha_2$

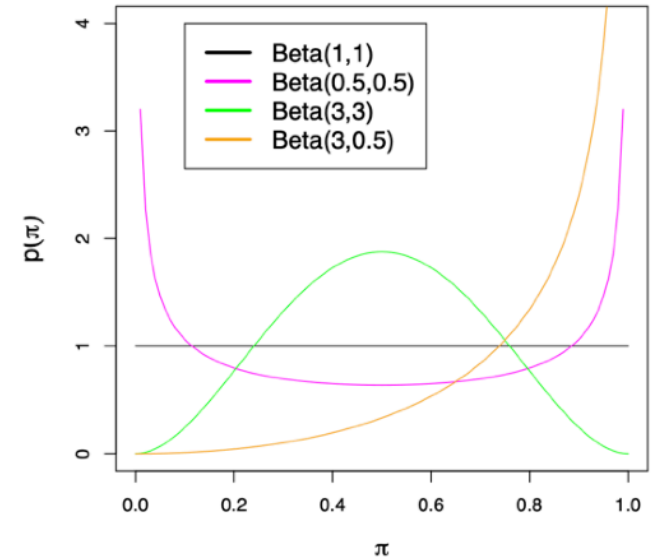
- The beta distribution has form

$$P(\pi|\alpha_1, \alpha_2) = \frac{1}{B(\alpha_1, \alpha_2)} \pi^{\alpha_1-1} (1-\pi)^{\alpha_2-1}$$

Normalizing constant,  
not of great interest for  
present purposes

Where the action is!

$$B(\alpha_1, \alpha_2) = \int_0^1 \pi^{\alpha_1-1} (1-\pi)^{\alpha_2-1} d\pi$$



- Cool thing about the beta distribution: the posterior is also beta distributed! For  $\mathbf{y} = m$  successes in  $n$  trials:

$$P(\pi|\mathbf{y}, \alpha_1, \alpha_2) \propto \overbrace{\pi^m (1-\pi)^{n-m}}^{\text{Likelihood}} \overbrace{\pi^{\alpha_1-1} (1-\pi)^{\alpha_2-1}}^{\text{Prior}}$$

# Example for coin flips: the beta distribution

- Express background knowledge  $I$  as two "pseudo-count" parameters  $\alpha_1, \alpha_2$

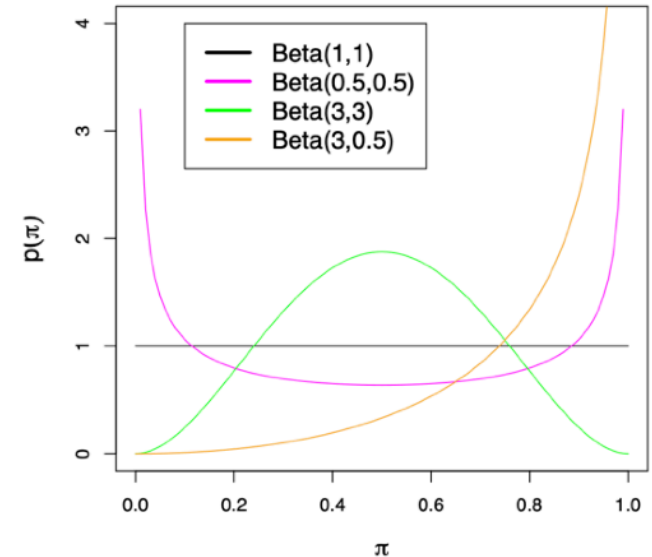
- The beta distribution has form

$$P(\pi|\alpha_1, \alpha_2) = \frac{1}{B(\alpha_1, \alpha_2)} \pi^{\alpha_1-1} (1-\pi)^{\alpha_2-1}$$

Normalizing constant,  
not of great interest for  
present purposes

Where the action is!

$$B(\alpha_1, \alpha_2) = \int_0^1 \pi^{\alpha_1-1} (1-\pi)^{\alpha_2-1} d\pi$$



- Cool thing about the beta distribution: the posterior is also beta distributed! For  $\mathbf{y} = m$  successes in  $n$  trials:

$$\begin{aligned} P(\pi|\mathbf{y}, \alpha_1, \alpha_2) &\propto \overbrace{\pi^m (1-\pi)^{n-m}}^{\text{Likelihood}} \overbrace{\pi^{\alpha_1-1} (1-\pi)^{\alpha_2-1}}^{\text{Prior}} \\ &\propto \pi^{m+\alpha_1-1} (1-\pi)^{n-m+\alpha_2-1} \end{aligned}$$

# Example for coin flips: the beta distribution

- Express background knowledge  $I$  as two "pseudo-count" parameters  $\alpha_1, \alpha_2$

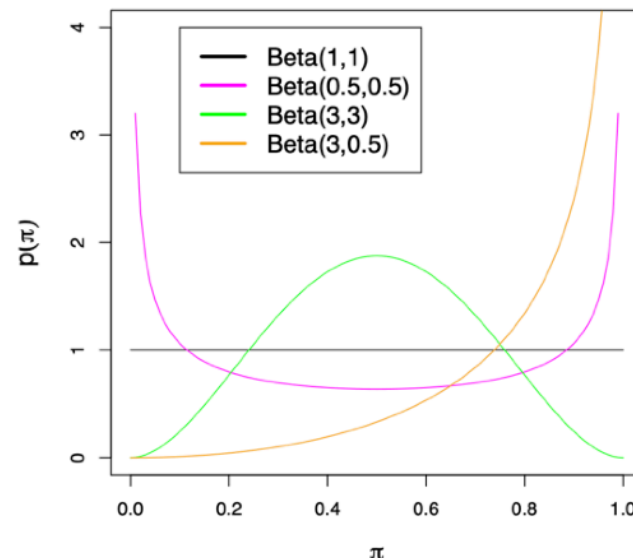
- The beta distribution has form

$$P(\pi|\alpha_1, \alpha_2) = \frac{1}{B(\alpha_1, \alpha_2)} \pi^{\alpha_1-1} (1-\pi)^{\alpha_2-1}$$

Normalizing constant,  
not of great interest for  
present purposes

Where the action is!

$$B(\alpha_1, \alpha_2) = \int_0^1 \pi^{\alpha_1-1} (1-\pi)^{\alpha_2-1} d\pi$$



- Cool thing about the beta distribution: the posterior is also beta distributed! For  $\mathbf{y} = m$  successes in  $n$  trials:

$$\begin{aligned} P(\pi|\mathbf{y}, \alpha_1, \alpha_2) &\propto \overbrace{\pi^m (1-\pi)^{n-m}}^{\text{Likelihood}} \overbrace{\pi^{\alpha_1-1} (1-\pi)^{\alpha_2-1}}^{\text{Prior}} \\ &\propto \pi^{m+\alpha_1-1} (1-\pi)^{n-m+\alpha_2-1} \end{aligned}$$

- This property is called **conjugacy** and is convenient where available!

# Example of Bayesian parameter estimation



# Example of Bayesian parameter estimation

---

- I inspect my coin and notice serious irregularities!



# Example of Bayesian parameter estimation

---

- I inspect my coin and notice serious irregularities!



- My prior for  $P(\text{heads})$ : a  
 $\alpha_1 = 3, \alpha_2 = 24$   
Beta prior



# Example of Bayesian parameter estimation

---

- I inspect my coin and notice serious irregularities!



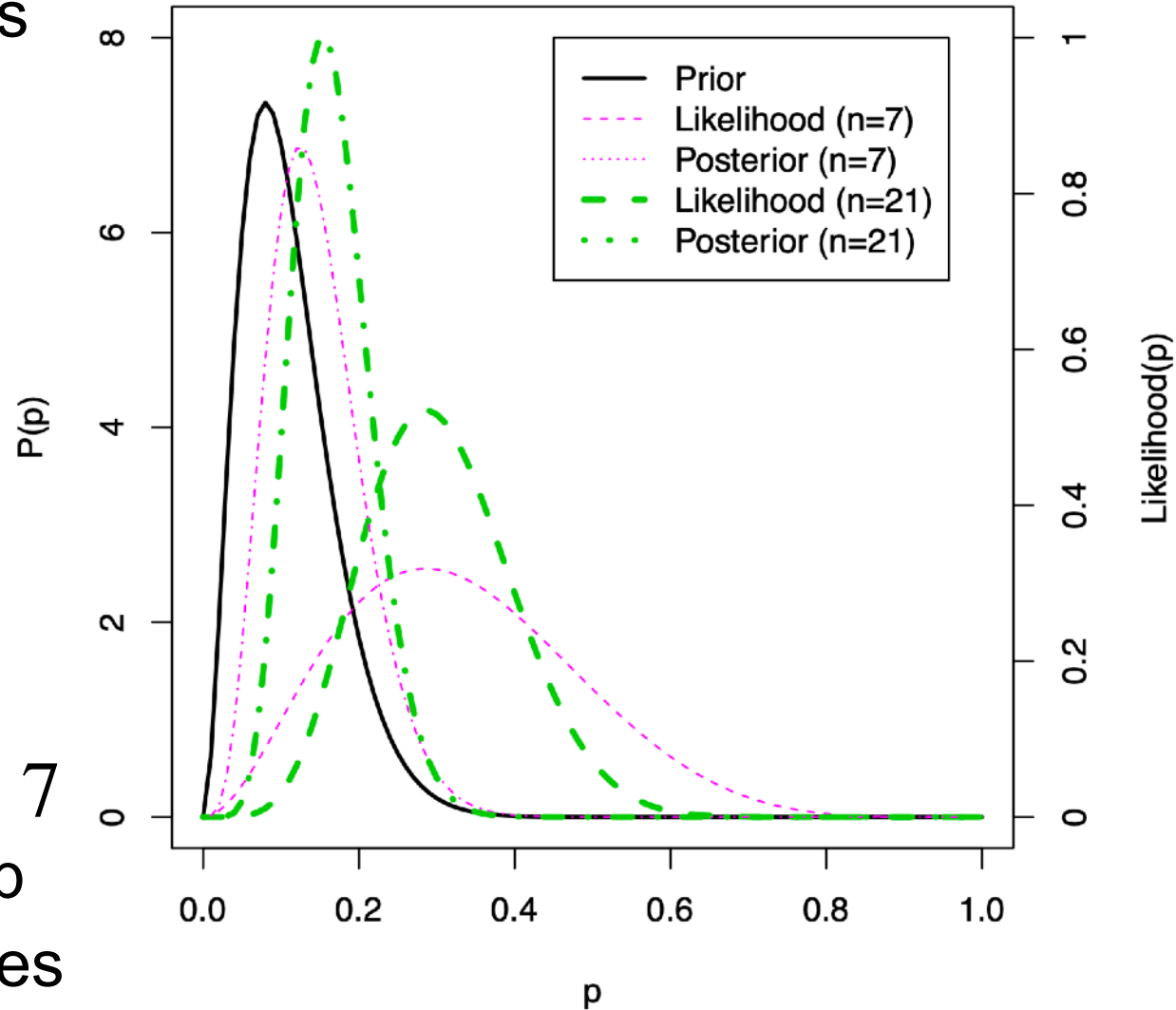
- My prior for  $P(\text{heads})$ : a  
 $\alpha_1 = 3, \alpha_2 = 24$   
Beta prior
- I flip the coin  $n = 7$   
times, it comes up  
heads  $m = 2$  times

# Example of Bayesian parameter estimation

- I inspect my coin and notice serious irregularities!



- My prior for  $P(\text{heads})$ : a  $\alpha_1 = 3, \alpha_2 = 24$  Beta prior
- I flip the coin  $n = 7$  times, it comes up heads  $m = 2$  times



# Posterior prediction

---

$$P(\text{heads}) = \pi$$

$$P(\pi) = \text{Beta}(\alpha_1, \alpha_2)$$

Observe  $m$  heads out of  $n$  flips

# Posterior prediction

---

$$P(\text{heads}) = \pi$$

$$P(\pi) = \text{Beta}(\alpha_1, \alpha_2)$$

Observe  $m$  heads out of  $n$  flips

**Posterior mean**

**Posterior mode**  
(when it exists)

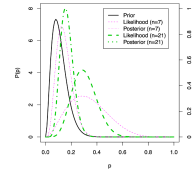
# Posterior prediction

$$P(\text{heads}) = \pi$$

$$P(\pi) = \text{Beta}(\alpha_1, \alpha_2)$$

Observe  $m$  heads out of  $n$  flips

**Beta distribution**



**Posterior mean**

**Posterior mode**  
(when it exists)

# Posterior prediction

$$P(\text{heads}) = \pi$$

$$P(\pi) = \text{Beta}(\alpha_1, \alpha_2)$$

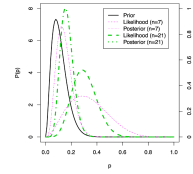
Observe  $m$  heads out of  $n$  flips

**Posterior mean**

$$E[\pi | I] = \frac{\alpha_1}{\alpha_1 + \alpha_2}$$

**Posterior mode**  
(when it exists)

**Beta distribution**



# Posterior prediction

$$P(\text{heads}) = \pi$$

$$P(\pi) = \text{Beta}(\alpha_1, \alpha_2)$$

Observe  $m$  heads out of  $n$  flips

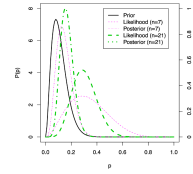
**Posterior mean**

$$E[\pi | I] = \frac{\alpha_1}{\alpha_1 + \alpha_2}$$

**Posterior mode**  
(when it exists)

$$\frac{\alpha_1 - 1}{\alpha_1 + \alpha_2 - 2}$$

**Beta distribution**



# Posterior prediction

$$P(\text{heads}) = \pi$$

$$P(\pi) = \text{Beta}(\alpha_1, \alpha_2)$$

Observe  $m$  heads out of  $n$  flips

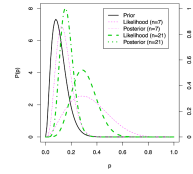
**Posterior mean**

$$E[\pi | I] = \frac{\alpha_1}{\alpha_1 + \alpha_2}$$

**Posterior mode**  
(when it exists)

$$\frac{\alpha_1 - 1}{\alpha_1 + \alpha_2 - 2}$$

**Beta distribution**



**Our example**



# Posterior prediction

$$P(\text{heads}) = \pi$$

$$P(\pi) = \text{Beta}(\alpha_1, \alpha_2)$$

Observe  $m$  heads out of  $n$  flips

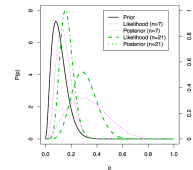
**Posterior mean**

$$E[\pi | I] = \frac{\alpha_1}{\alpha_1 + \alpha_2}$$

**Posterior mode**  
(when it exists)

$$\frac{\alpha_1 - 1}{\alpha_1 + \alpha_2 - 2}$$

**Beta distribution**



**Our example**

$$E[\pi | y, I] = \frac{\alpha_1 + m}{\alpha_1 + \alpha_2 + n}$$

# Posterior prediction

$$P(\text{heads}) = \pi$$

$$P(\pi) = \text{Beta}(\alpha_1, \alpha_2)$$

Observe  $m$  heads out of  $n$  flips

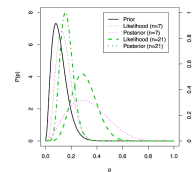
**Posterior mean**

**Posterior mode**  
(when it exists)

**Beta distribution**

$$E[\pi | I] = \frac{\alpha_1}{\alpha_1 + \alpha_2}$$

$$\frac{\alpha_1 - 1}{\alpha_1 + \alpha_2 - 2}$$



**Our example**

$$E[\pi | y, I] = \frac{\alpha_1 + m}{\alpha_1 + \alpha_2 + n}$$

$$\frac{\alpha_1 + m - 1}{\alpha_1 + \alpha_2 + n - 2}$$

# Posterior prediction

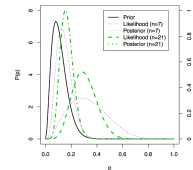
$$P(\text{heads}) = \pi$$

$$P(\pi) = \text{Beta}(\alpha_1, \alpha_2)$$

Observe  $m$  heads out of  $n$  flips

**Posterior mean**

$$E[\pi | I] = \frac{\alpha_1}{\alpha_1 + \alpha_2}$$



**Our example**

$$E[\pi | y, I] = \frac{\alpha_1 + m}{\alpha_1 + \alpha_2 + n}$$

**Posterior mode**  
(when it exists)

$$\frac{\alpha_1 - 1}{\alpha_1 + \alpha_2 - 2}$$

$$\frac{\alpha_1 + m - 1}{\alpha_1 + \alpha_2 + n - 2}$$

**Posterior predictive distribution**

# Posterior prediction

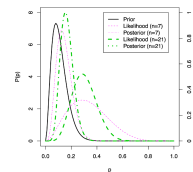
$$P(\text{heads}) = \pi$$

$$P(\pi) = \text{Beta}(\alpha_1, \alpha_2)$$

Observe  $m$  heads out of  $n$  flips

**Posterior mean**

$$E[\pi | I] = \frac{\alpha_1}{\alpha_1 + \alpha_2}$$



**Our example**

$$E[\pi | y, I] = \frac{\alpha_1 + m}{\alpha_1 + \alpha_2 + n}$$

**Posterior mode**  
(when it exists)

$$\frac{\alpha_1 - 1}{\alpha_1 + \alpha_2 - 2}$$

$$\frac{\alpha_1 + m - 1}{\alpha_1 + \alpha_2 + n - 2}$$

## Posterior predictive distribution

If I flip the same coin  $k$  more times, what is the distribution on the resulting # heads  $r$ ?

# Posterior prediction

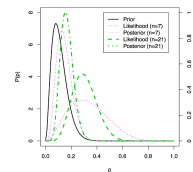
$$P(\text{heads}) = \pi$$

$$P(\pi) = \text{Beta}(\alpha_1, \alpha_2)$$

Observe  $m$  heads out of  $n$  flips

**Posterior mean**

$$E[\pi | I] = \frac{\alpha_1}{\alpha_1 + \alpha_2}$$



**Our example**

$$E[\pi | y, I] = \frac{\alpha_1 + m}{\alpha_1 + \alpha_2 + n}$$

**Posterior mode**  
(when it exists)

$$\frac{\alpha_1 - 1}{\alpha_1 + \alpha_2 - 2}$$

$$\frac{\alpha_1 + m - 1}{\alpha_1 + \alpha_2 + n - 2}$$

## Posterior predictive distribution

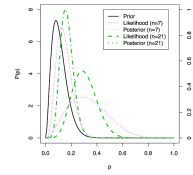
If I flip the same coin  $k$  more times, what is the distribution on the resulting # heads  $r$ ?

$$P(\mathbf{y}_{new} | \mathbf{y}, I)$$

# Posterior prediction

$P(\text{heads}) = \pi$   
 $P(\pi) = \text{Beta}(\alpha_1, \alpha_2)$   
 Observe  $m$  heads out of  $n$  flips

**Beta distribution**



**Our example**

**Posterior mean**

$$E[\pi | I] = \frac{\alpha_1}{\alpha_1 + \alpha_2}$$

$$E[\pi | y, I] = \frac{\alpha_1 + m}{\alpha_1 + \alpha_2 + n}$$

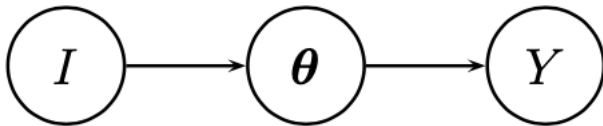
**Posterior mode**  
(when it exists)

$$\frac{\alpha_1 - 1}{\alpha_1 + \alpha_2 - 2}$$

$$\frac{\alpha_1 + m - 1}{\alpha_1 + \alpha_2 + n - 2}$$

## Posterior predictive distribution

If I flip the same coin  $k$  more times, what is the distribution on the resulting # heads  $r$ ?

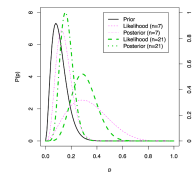


$$P(\mathbf{y}_{new} | \mathbf{y}, I)$$

# Posterior prediction

$P(\text{heads}) = \pi$   
 $P(\pi) = \text{Beta}(\alpha_1, \alpha_2)$   
 Observe  $m$  heads out of  $n$  flips

**Beta distribution**



**Our example**

**Posterior mean**

$$E[\pi | I] = \frac{\alpha_1}{\alpha_1 + \alpha_2}$$

$$E[\pi | y, I] = \frac{\alpha_1 + m}{\alpha_1 + \alpha_2 + n}$$

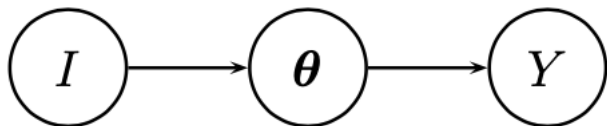
**Posterior mode**  
(when it exists)

$$\frac{\alpha_1 - 1}{\alpha_1 + \alpha_2 - 2}$$

$$\frac{\alpha_1 + m - 1}{\alpha_1 + \alpha_2 + n - 2}$$

## Posterior predictive distribution

If I flip the same coin  $k$  more times, what is the distribution on the resulting # heads  $r$ ?

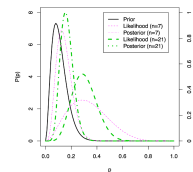


$$P(\mathbf{y}_{new} | \mathbf{y}, I) = \int_{\theta} P(\mathbf{y}_{new} | \theta) P(\theta | \mathbf{y}, I) d\theta$$

# Posterior prediction

$P(\text{heads}) = \pi$   
 $P(\pi) = \text{Beta}(\alpha_1, \alpha_2)$   
 Observe  $m$  heads out of  $n$  flips

**Beta distribution**



**Our example**

**Posterior mean**

$$E[\pi | I] = \frac{\alpha_1}{\alpha_1 + \alpha_2}$$

$$E[\pi | y, I] = \frac{\alpha_1 + m}{\alpha_1 + \alpha_2 + n}$$

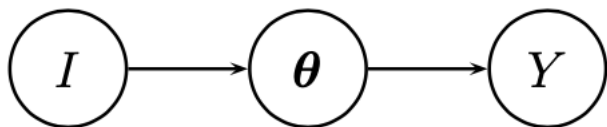
**Posterior mode**  
(when it exists)

$$\frac{\alpha_1 - 1}{\alpha_1 + \alpha_2 - 2}$$

$$\frac{\alpha_1 + m - 1}{\alpha_1 + \alpha_2 + n - 2}$$

## Posterior predictive distribution

If I flip the same coin  $k$  more times, what is the distribution on the resulting # heads  $r$ ?



$$P(\mathbf{y}_{new} | \mathbf{y}, I) = \int_{\theta} P(\mathbf{y}_{new} | \theta) P(\theta | \mathbf{y}, I) d\theta$$

→ The **Beta-Binomial** model:  $P(r | k, I, \mathbf{y}) = \binom{k}{r} \frac{B(\alpha_1 + m + r, \alpha_2 + n - m + k - r)}{B(\alpha_1 + m, \alpha_2 + n - m)}$



# Point estimation vs Bayesian prediction

---

- Say we'll flip the coin  $k = 50$  more times

# Point estimation vs Bayesian prediction

---

- Say we'll flip the coin  $k = 50$  more times

$$\alpha_1 + m = 5$$

# Point estimation vs Bayesian prediction

---

- Say we'll flip the coin  $k = 50$  more times

$$\alpha_1 + m = 5$$

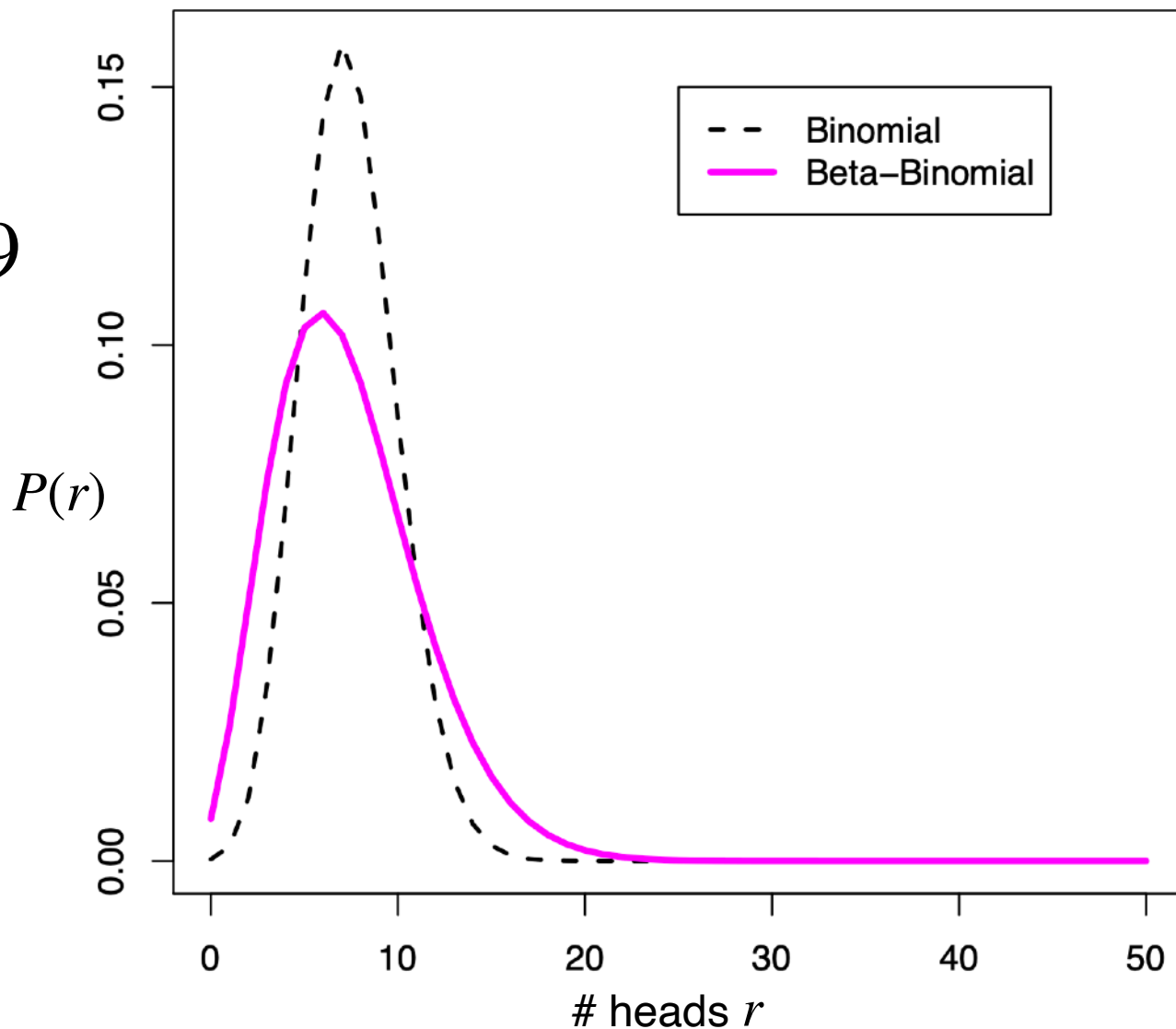
$$\alpha_2 + n - m = 29$$

# Point estimation vs Bayesian prediction

- Say we'll flip the coin  $k = 50$  more times

$$\alpha_1 + m = 5$$

$$\alpha_2 + n - m = 29$$



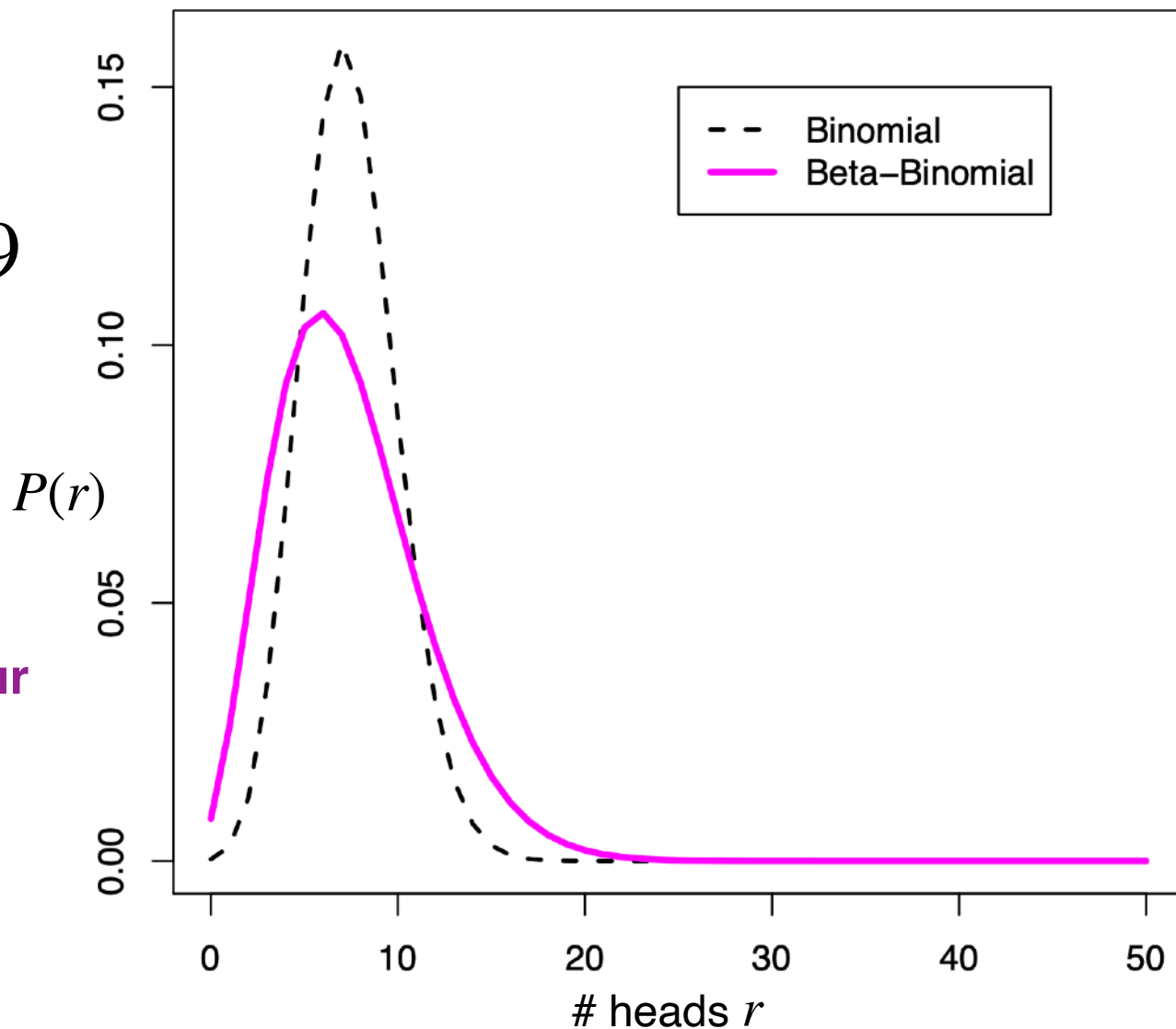
# Point estimation vs Bayesian prediction

- Say we'll flip the coin  $k = 50$  more times

$$\alpha_1 + m = 5$$

$$\alpha_2 + n - m = 29$$

Bayesian inference takes into account our uncertainty about model parameters  $\theta$ , leading to more hedged predictions



# A note on Bayesian priors

---

# A note on Bayesian priors

---

- The Bayesian prior is a double-edged sword

# A note on Bayesian priors

---

- The Bayesian prior is a double-edged sword
  - We get to specify it



# A note on Bayesian priors

---

- The Bayesian prior is a double-edged sword
  - We get to specify it
  - We have to specify it

# A note on Bayesian priors

---

- The Bayesian prior is a double-edged sword
  - We get to specify it
  - We have to specify it
- In the above example, we used an **informative prior**

# A note on Bayesian priors

---

- The Bayesian prior is a double-edged sword
  - We get to specify it
  - We have to specify it
- In the above example, we used an **informative prior**
  - When we have strong domain knowledge, this can potentially be useful in various ways

# A note on Bayesian priors

---

- The Bayesian prior is a double-edged sword
  - We get to specify it
  - We have to specify it
- In the above example, we used an **informative prior**
  - When we have strong domain knowledge, this can potentially be useful in various ways
- In scientific data analysis, however, our general goal is to *allow the data to speak to us about what we care about*

# A note on Bayesian priors

---

- The Bayesian prior is a double-edged sword
  - We get to specify it
  - We have to specify it
- In the above example, we used an **informative prior**
  - When we have strong domain knowledge, this can potentially be useful in various ways
- In scientific data analysis, however, our general goal is to *allow the data to speak to us about what we care about*
- For this reason, I generally advocate **vague priors for any part of the model whose posterior we care about**

# A note on Bayesian priors

---

- The Bayesian prior is a double-edged sword
  - We get to specify it
  - We have to specify it
- In the above example, we used an **informative prior**
  - When we have strong domain knowledge, this can potentially be useful in various ways
- In scientific data analysis, however, our general goal is to *allow the data to speak to us about what we care about*
- For this reason, I generally advocate **vague priors for any part of the model whose posterior we care about**
- If your qualitative conclusions depend on choice of prior, it is a reason to be wary of the robustness of your analysis!

# A note on Bayesian priors

---

- The Bayesian prior is a double-edged sword
  - We get to specify it
  - We have to specify it
- In the above example, we used an **informative prior**
  - When we have strong domain knowledge, this can potentially be useful in various ways
- In scientific data analysis, however, our general goal is to *allow the data to speak to us about what we care about*
- For this reason, I generally advocate **vague priors for any part of the model whose posterior we care about**
- If your qualitative conclusions depend on choice of prior, it is a reason to be wary of the robustness of your analysis!
- As data become plentiful\*, choice of prior *often but not always* recedes in importance

\*What counts as "plentiful" depends on size of the model and structure of the data