

# **A Mini Introduction to Bayesian Statistics**

Yingbo Li

March 3, 2018

# Outline

Bayes rule: the conditional probability

Parameter estimation: dynamic belief updating

Hypothesis testing: no more  $p$ -values

# If the mammogram is positive, how likely s/he has breast cancer?



Source: <http://familyguy.wikia.com/wiki/File:Mammogram.png>

# If the mammogram is positive, how likely s/he has breast cancer?

- A mammogram is right 75% of the time, if a women has breast cancer

$$P(+ | \text{乳腺癌}) = 0.75$$

# If the mammogram is positive, how likely s/he has breast cancer?

- A mammogram is right 75% of the time, if a women has breast cancer

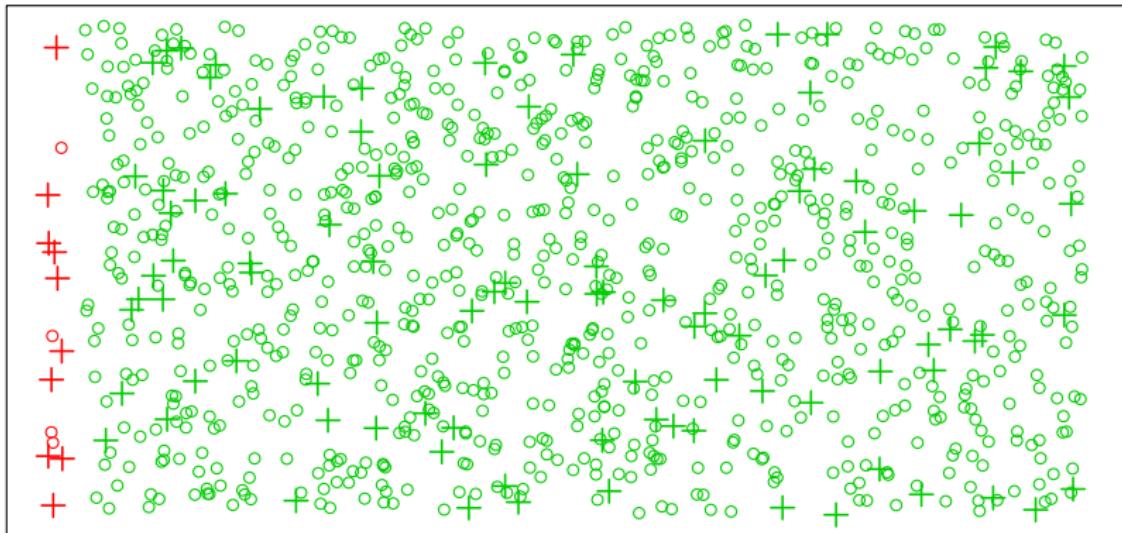
$$P(+ | \text{乳腺癌}) = 0.75$$

- A mammogram is right 90% of the time, if a women doesn't have breast cancer

$$P(+ | \text{无乳腺癌}) = 0.1$$

# If the mammogram is positive, how likely s/he has breast cancer?

- For a women in her 40's, the chance of having breast cancer is 1.4%.



# If the mammogram is positive, how likely s/he has breast cancer?

$$P(\text{😱} | +) = \frac{P(+ | \text{😱}) \cdot P(\text{😱})}{P(+ | \text{😱}) \cdot P(\text{😱}) + P(+ | \text{😊}) \cdot P(\text{😊})}$$

# If the mammogram is positive, how likely s/he has breast cancer?

$$P(\text{😱} | +) = \frac{P(+ | \text{😱}) \cdot P(\text{😱})}{P(+ | \text{😱}) \cdot P(\text{😱}) + P(+ | \text{😊}) \cdot P(\text{😊})}$$
$$= \frac{0.75 \times 0.014}{0.75 \times 0.014 + 0.1 \times 0.986}$$

# If the mammogram is positive, how likely s/he has breast cancer?

$$\begin{aligned} P(\text{😱} | +) &= \frac{P(+ | \text{😱}) \cdot P(\text{😱})}{P(+ | \text{😱}) \cdot P(\text{😱}) + P(+ | \text{😊}) \cdot P(\text{😊})} \\ &= \frac{0.75 \times 0.014}{0.75 \times 0.014 + 0.1 \times 0.986} \\ &= 0.096 \end{aligned}$$

# Prior vs posterior probabilities

- Event of interest: 🩺, having breast cancer

# Prior vs posterior probabilities

- Event of interest: 🙄, having breast cancer
- Prior probability: opinion on how likely the event will occur **before we collecting data**

$$P(\text{🙄}) = 1.4\%$$

# Prior vs posterior probabilities

- Event of interest: 😱, having breast cancer
- Prior probability: opinion on how likely the event will occur **before we collecting data**

$$P(\text{😱}) = 1.4\%$$

- Posterior probability: opinion on how likely the event will occur **after observing the data**

$$P(\text{😱} \mid \text{+}) = 9.6\%$$

# Bayes rule

$$P(A | D) = \frac{P(D | A) \cdot P(A)}{P(D | A) \cdot P(A) + P(D | A^c) \cdot P(A^c)}$$

A occurs

A doesn't occur

# **What if a 2nd mammogram is again +?**

- Assuming the two results are independent

# What if a 2nd mammogram is again +?

- Assuming the two results are independent

$$\begin{aligned} P(\text{😱} \mid +) &= \frac{P(+ \mid \text{😱}) \cdot P(\text{😱})}{P(+ \mid \text{😱}) \cdot P(\text{😱}) + P(+ \mid \text{😊}) \cdot P(\text{😊})} \\ &= \frac{0.75 \times 0.096}{0.75 \times 0.096 + 0.1 \times (1 - 0.096)} \\ &= 0.443 \end{aligned}$$

# What if a 2nd mammogram is again +?

- **Dynamic update of information:** she's no longer an average woman in her 40's, but who already has one positive mammogram.

# Bayesian statistics: 250+ yrs history

Thomas Bayes  
(1701 - 1761)



Ronald Fisher  
(1890-1962)

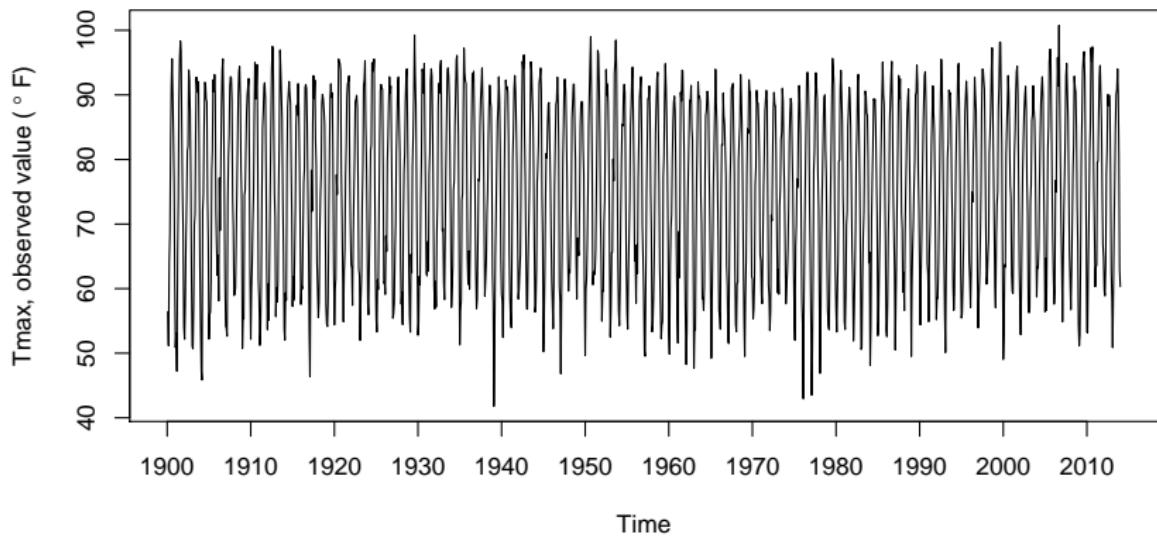


# Why Bayesian statistics?

- Can incorporate **subjective prior belief** or domain experts knowledge
- No large sample assumption needed
- Markov chain Monte Carlo (MCMC) enables easy inference for sophisticated models
- Revives with revolutionary computation power

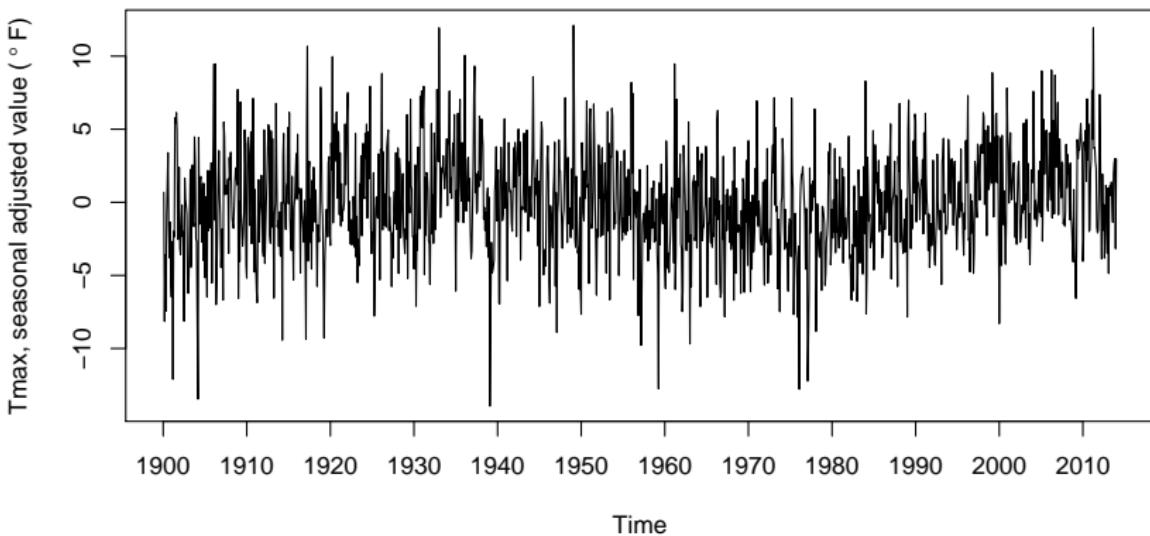
# Monthly temperature, Tuscaloosa AL

Observed data



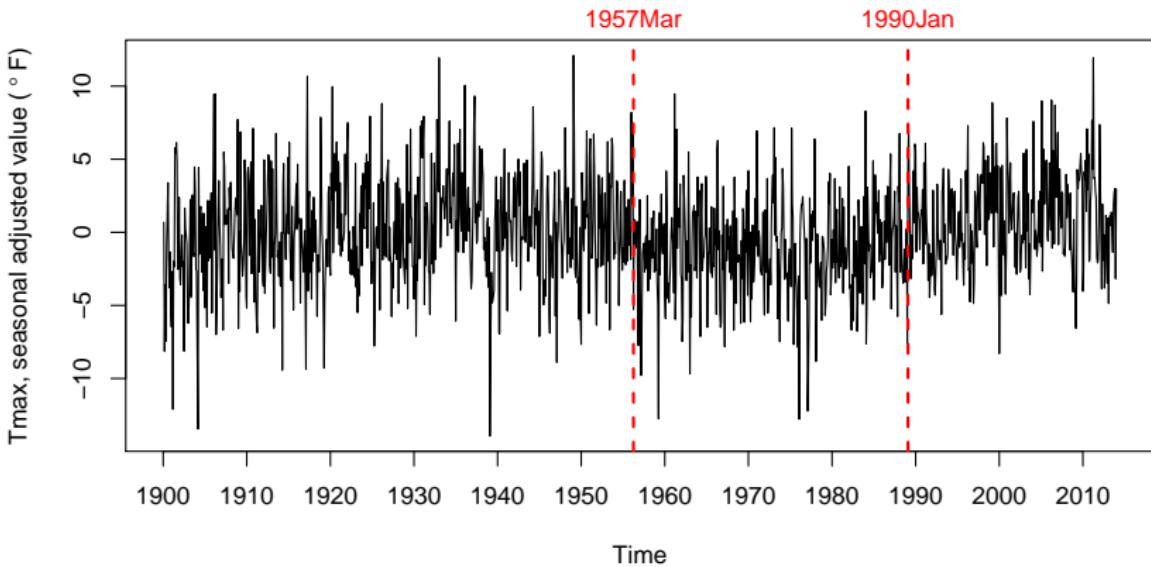
# Monthly temperature, Tuscaloosa AL

Observed data – sample seasonal mean



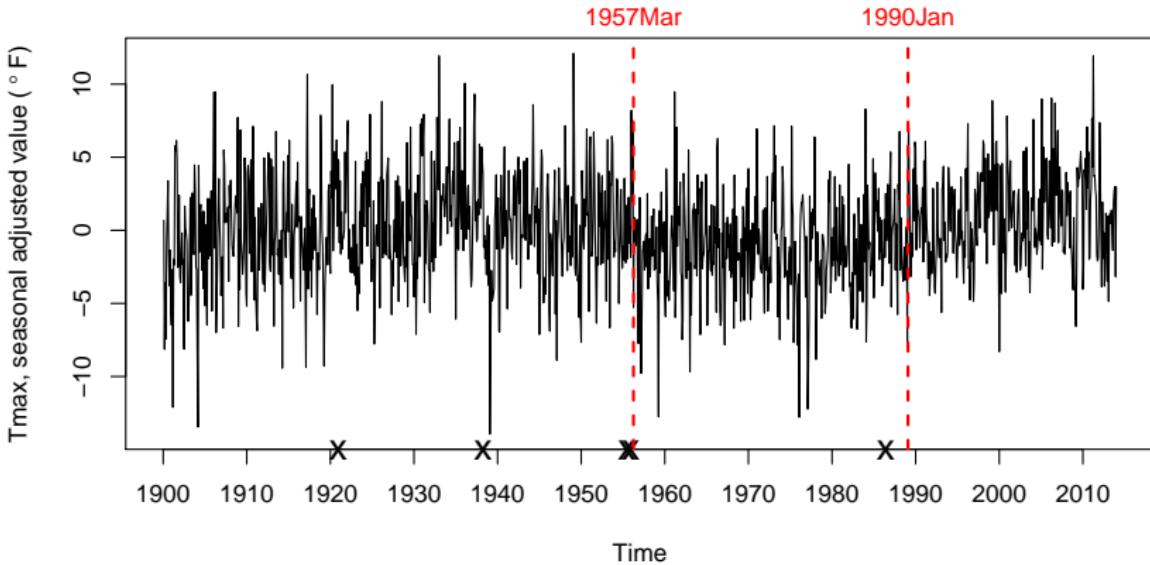
# Monthly temperature, Tuscaloosa AL

Observed data – sample seasonal mean



# Monthly temperature, Tuscaloosa AL

Observed data – sample seasonal mean

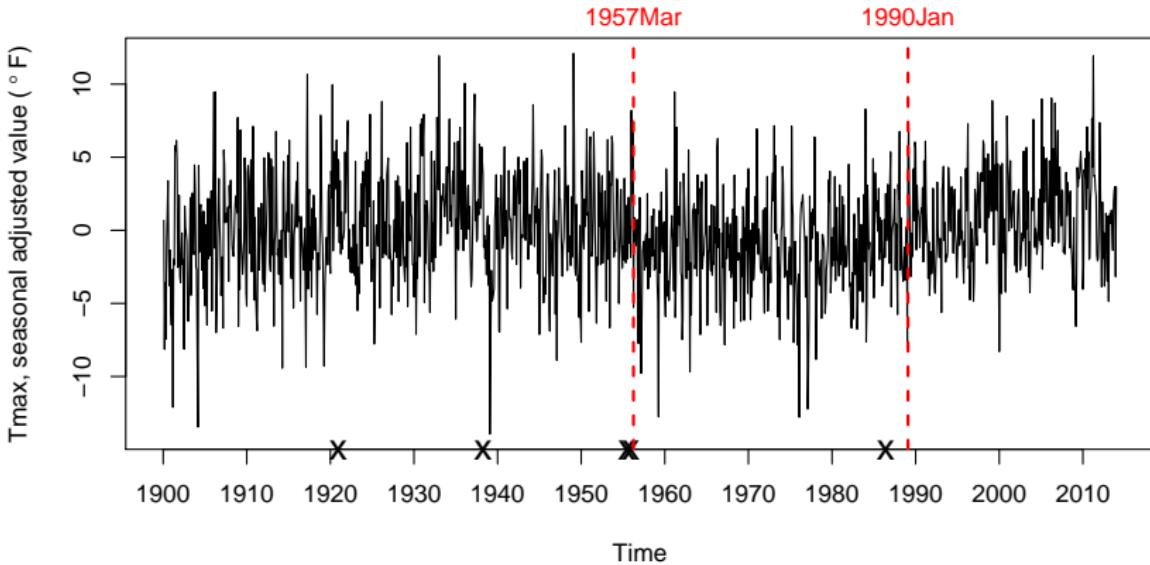


Metadata (station history logs)

- Station relocations: 1921/11, 1939/03, 1956/06
- Instrumentation changes: 1956/11, 1987/05

# Monthly temperature, Tuscaloosa AL

Observed data – sample seasonal mean

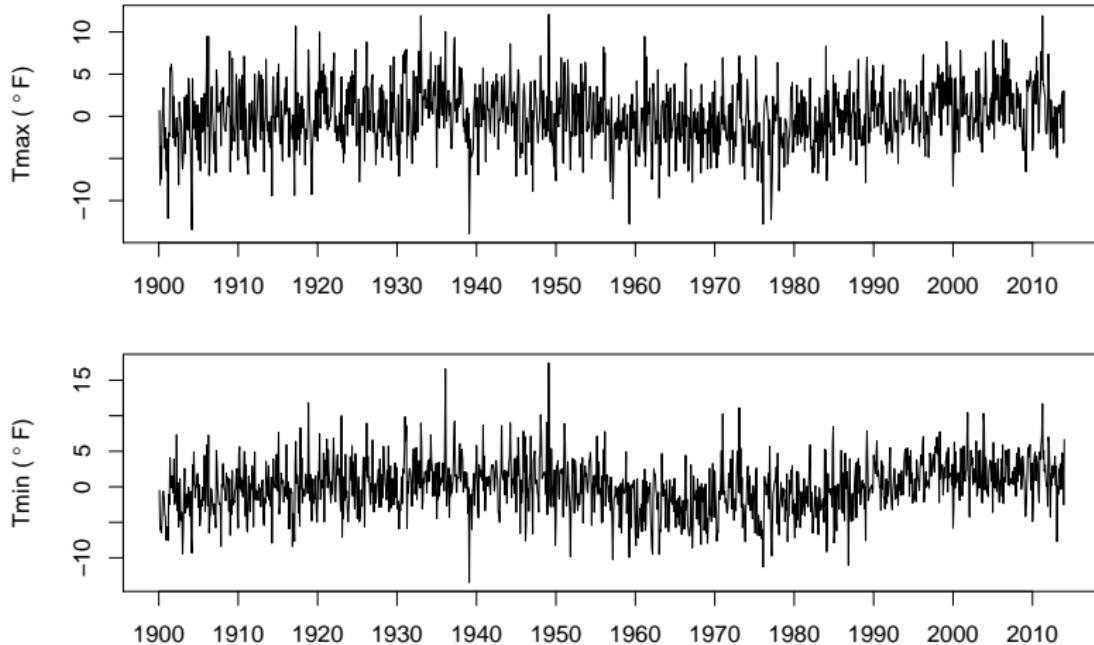


Metadata: more likely to induce mean shifts

- Station relocations: 1921/11, 1939/03, 1956/06
- Instrumentation changes: 1956/11, 1987/05

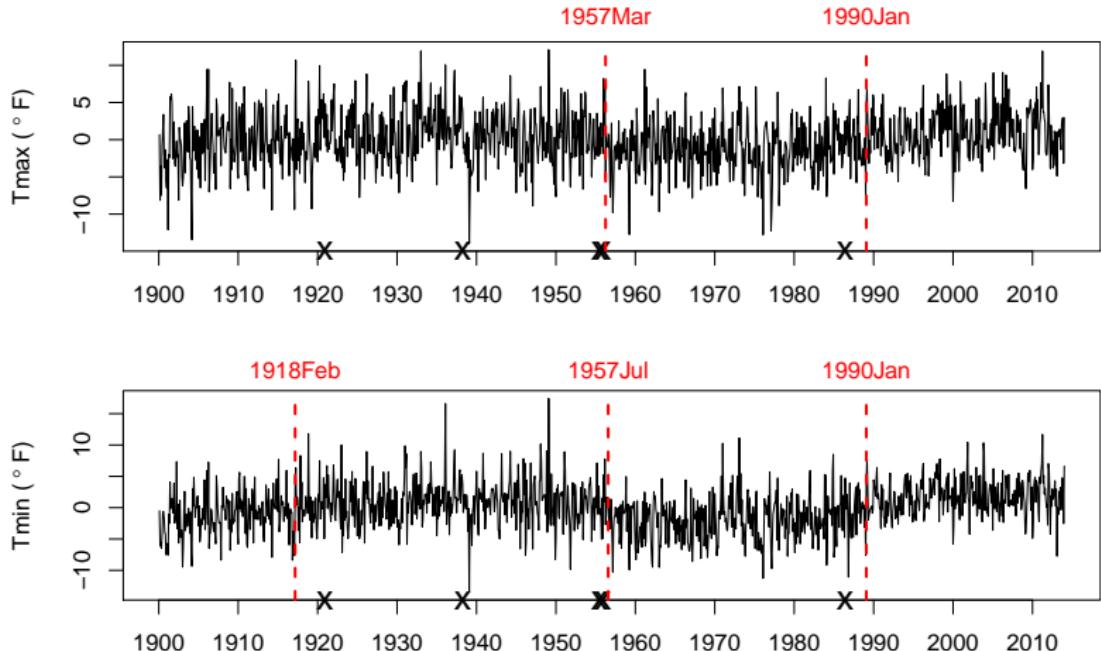
# Monthly max and min temperatures

Observed data – sample seasonal mean

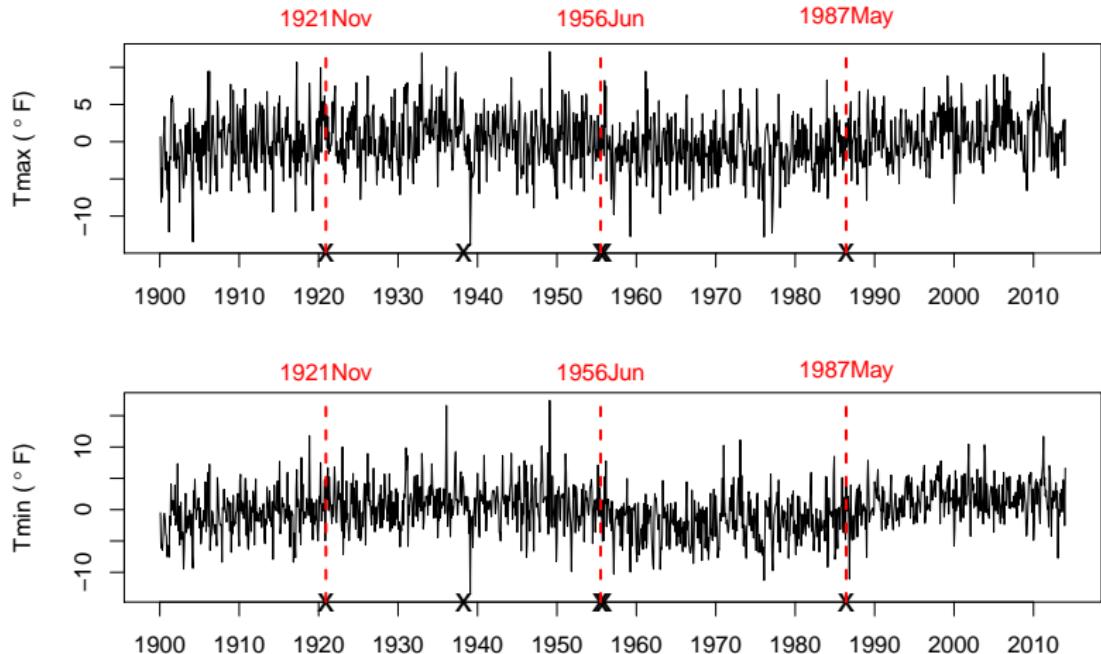


Tmax and Tmin: likely to shift at the same time

# Existing changepoint detection



# Bayesian changepoint detection



# The denominator is not important

$$P(A \mid D) = \frac{P(D \mid A) \cdot P(A)}{P(D \mid A) \cdot P(A) + P(D \mid A^c) \cdot P(A^c)}$$

$$P(A^c \mid D) = \frac{P(D \mid A^c) \cdot P(A^c)}{P(D \mid A) \cdot P(A) + P(D \mid A^c) \cdot P(A^c)}$$

# The denominator is not important

$$P(A \mid D) = \frac{P(D \mid A) \cdot P(A)}{P(D \mid A) \cdot P(A) + P(D \mid A^c) \cdot P(A^c)}$$

$$P(A^c \mid D) = \frac{P(D \mid A^c) \cdot P(A^c)}{P(D \mid A) \cdot P(A) + P(D \mid A^c) \cdot P(A^c)}$$

# The denominator is not important

$$P(A \mid D) \propto P(D \mid A) \cdot P(A)$$

$$P(A^c \mid D) \propto P(D \mid A^c) \cdot P(A^c)$$

# Bayes rule: distribution version

- Data  $D$  is randomly generated from a model
- With unknown parameter  $\theta$

$$p(\theta \mid D) \propto p(D \mid \theta) \times p(\theta)$$

posterior distribution      likelihood      prior distribution

# Bayes rule: distribution version

- Data  $D$  is randomly generated from a model
- With unknown parameter  $\theta$

$$p(\theta \mid D) \propto p(D \mid \theta) \times p(\theta)$$

posterior distribution      likelihood      prior distribution

- Estimate  $\theta$  using the posterior mean  $E(\theta \mid D)$

# Is it a fair coin?



Source: <https://www.li-lacchocolates.com/assets/images/Chanukah/Gold-Chocolate-Coins.jpg>

# Is it a fair coin?

- Parameter  $\theta$ : chance of Head
- Data  $D$ : coin flipping Head/Tail outcomes
- Likelihood function:

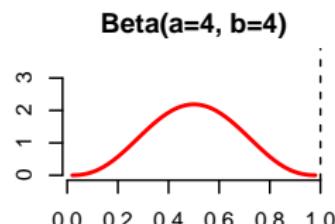
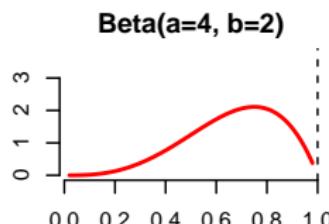
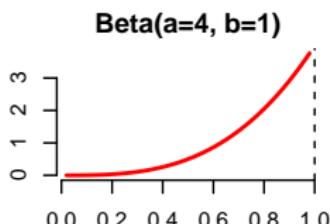
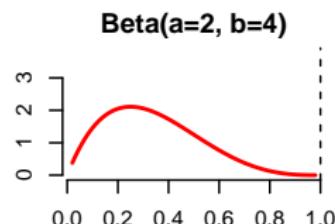
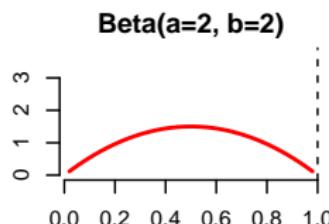
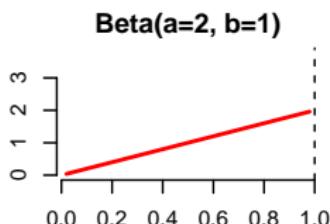
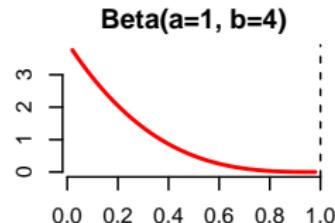
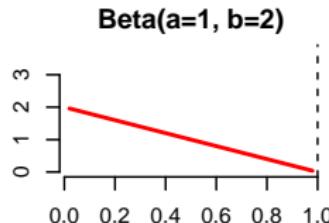
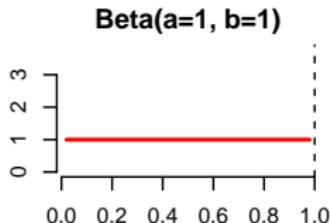
$$p(D \mid \theta) = \theta^{\#\text{Head}}(1 - \theta)^{\#\text{Tail}}$$

# Beta distribution: between 0 and 1

$$p(\theta) \propto \theta^{a-1} (1-\theta)^{b-1}$$

# Beta distribution: between 0 and 1

$$p(\theta) \propto \theta^{a-1}(1-\theta)^{b-1}$$



# Posterior distribution is again Beta

$$p(\theta \mid D) \propto p(D \mid \theta) \times p(\theta)$$

posterior              likelihood              prior

# Posterior distribution is again Beta

$$p(\theta | D) \propto [\theta^{\#Head} (1 - \theta)^{\#Tail}] \times p(\theta)$$

posterior                      likelihood                      prior

# Posterior distribution is again Beta

$$p(\theta | D) \propto [\theta^{\# \text{Head}} (1 - \theta)^{\# \text{Tail}}] \times [\theta^{a-1} (1 - \theta)^{b-1}]$$

posterior                      likelihood                      prior

# Posterior distribution is again Beta

$$p(\theta | D) \propto [\theta^{\#Head} (1 - \theta)^{\#Tail}] \times [\theta^{a-1} (1 - \theta)^{b-1}]$$

posterior                      likelihood                      prior

$$\theta | D \sim \text{Beta} (\#Head + a, \#Tail + b)$$

# Posterior distribution is again Beta

$$p(\theta | D) \propto [\theta^{\#Head} (1 - \theta)^{\#Tail}] \times [\theta^{a-1} (1 - \theta)^{b-1}]$$

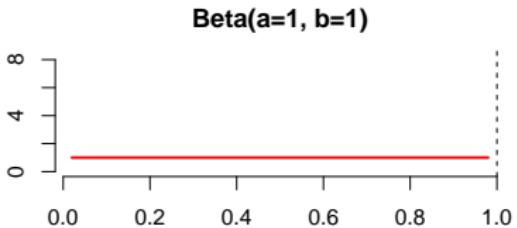
posterior                      likelihood                      prior

$$\theta | D \sim \text{Beta} (\#Head + a, \#Tail + b)$$

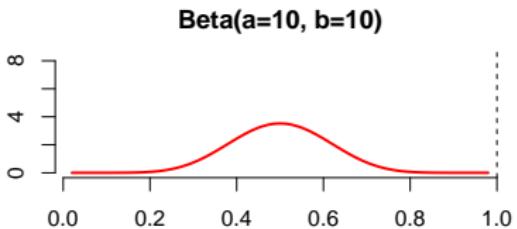
- Conjugacy: prior and posterior in the same distribution family

Person	Prior	Mean	SD
--------	-------	------	----

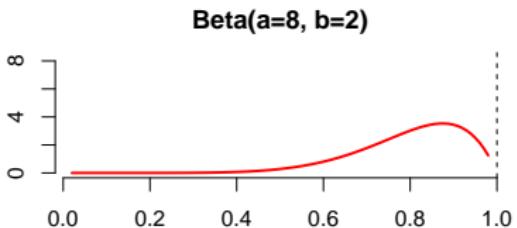
1

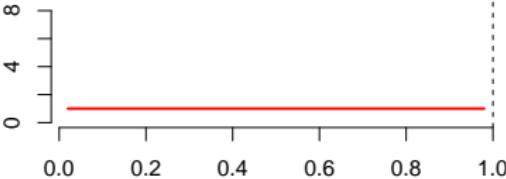
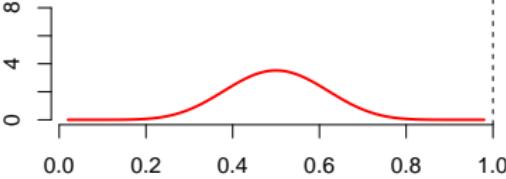
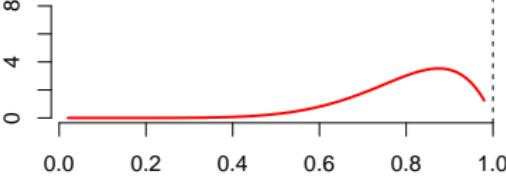


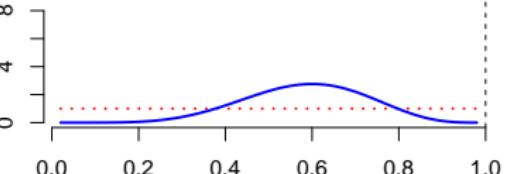
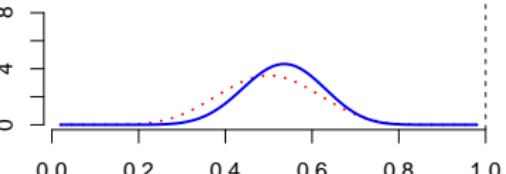
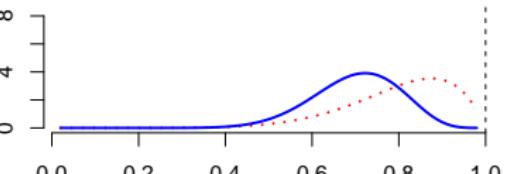
2

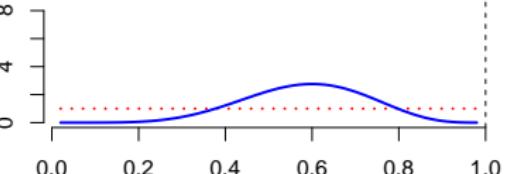
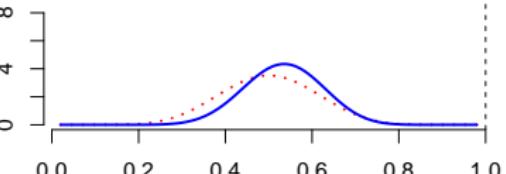
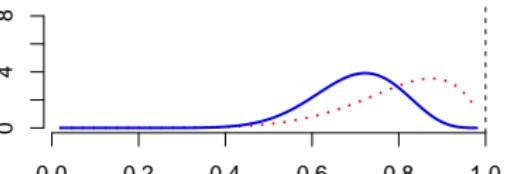


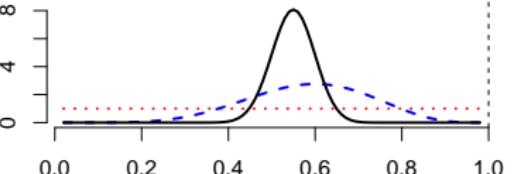
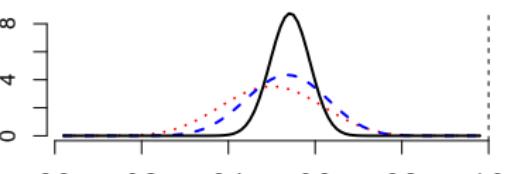
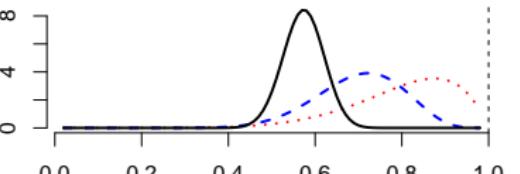
3



Person	Prior	Mean	SD
1	Beta( $a=1, b=1$ ) 	0.50	0.29
2	Beta( $a=10, b=10$ ) 	0.50	0.11
3	Beta( $a=8, b=2$ ) 	0.80	0.12

Person	Posterior, 6 head / 4 tails	Mean	SD
1	<p style="text-align: center;"><b>Beta(a=7, b=5)</b></p> 	0.50	0.29
2	<p style="text-align: center;"><b>Beta(a=16, b=14)</b></p> 	0.50	0.11
3	<p style="text-align: center;"><b>Beta(a=14, b=6)</b></p> 	0.80	0.12

Person	Posterior, 6 head / 4 tails	Mean	SD
1	<p style="text-align: center;"><b>Beta(a=7, b=5)</b></p> 	0.50	0.29
2	<p style="text-align: center;"><b>Beta(a=16, b=14)</b></p> 	0.50	0.11
3	<p style="text-align: center;"><b>Beta(a=14, b=6)</b></p> 	0.80	0.12

Person	Posterior, 55 head / 45 tails	Mean	SD
1	<p style="text-align: center;"><math>\text{Beta}(a=56, b=46)</math></p>  <p>A plot showing a posterior distribution for Person 1. The x-axis represents probability from 0.0 to 1.0, and the y-axis represents density from 0 to 8. A solid black curve represents the beta(a=56, b=46) distribution, which is centered at approximately 0.58. A dashed blue curve shows the prior distribution, and a dotted red curve shows the likelihood. A vertical dashed line is at 1.0.</p>	0.50	0.29
2	<p style="text-align: center;"><math>\text{Beta}(a=65, b=55)</math></p>  <p>A plot showing a posterior distribution for Person 2. The x-axis represents probability from 0.0 to 1.0, and the y-axis represents density from 0 to 8. A solid black curve represents the beta(a=65, b=55) distribution, which is centered at approximately 0.53. A dashed blue curve shows the prior distribution, and a dotted red curve shows the likelihood. A vertical dashed line is at 1.0.</p>	0.50	0.11
3	<p style="text-align: center;"><math>\text{Beta}(a=63, b=47)</math></p>  <p>A plot showing a posterior distribution for Person 3. The x-axis represents probability from 0.0 to 1.0, and the y-axis represents density from 0 to 8. A solid black curve represents the beta(a=63, b=47) distribution, which is centered at approximately 0.70. A dashed blue curve shows the prior distribution, and a dotted red curve shows the likelihood. A vertical dashed line is at 1.0.</p>	0.80	0.12

Person	Posterior, 55 head / 45 tails	Mean	SD
1	<p style="text-align: center;"><b>Beta(a=56, b=46)</b></p>	0.50	0.29
2	<p style="text-align: center;"><b>Beta(a=65, b=55)</b></p>	0.50	0.11
3	<p style="text-align: center;"><b>Beta(a=63, b=47)</b></p>	0.80	0.12

# Movie ratings on Rotten Tomatoes

Rotten Tomatoes staffs

- collect reviews from certified movie critics, and
- determine whether each review is

positive “fresh” 

vs

negative “rotten” 

[https://en.wikipedia.org/wiki/Rotten\\_Tomatoes](https://en.wikipedia.org/wiki/Rotten_Tomatoes)

# Harry Potter and the Philosopher's Stone (movie 1)



# Harry Potter and the Philosopher's Stone (movie 1)



- Prior distribution

mean = 50%, sd = 28.8%

# Harry Potter and the Philosopher's Stone (movie 1)



- Prior distribution

mean = 50%, sd = 28.8%

- Data

🍅 155, 🌿 38

# Harry Potter and the Philosopher's Stone (movie 1)



- Prior distribution

mean = 50%, sd = 28.8%

- Data

🍅 155, 🌿 38

- Posterior distribution

mean = 80%, sd = 2.9%

# Harry Potter and the Chamber of Secrets (movie 2)



# Harry Potter and the Chamber of Secrets (movie 2)



- Prior distribution

mean = 80%, sd = 2.9%

# Harry Potter and the Chamber of Secrets (movie 2)



- Prior distribution  
 $\text{mean} = 80\%$ ,  $\text{sd} = 2.9\%$
- Data

🍅 191, 🌿 41

# Harry Potter and the Chamber of Secrets (movie 2)



- Prior distribution  
 $\text{mean} = 80\%$ ,  $\text{sd} = 2.9\%$
- Data  
    ● 191, ✖ 41
- Posterior distribution  
 $\text{mean} = 81.3\%$ ,  $\text{sd} = 1.9\%$

# Harry Potter and the Prisoner of Azkaban (movie 3)



# Harry Potter and the Prisoner of Azkaban (movie 3)



- Prior distribution

mean = 81.3%, sd = 1.9%

# Harry Potter and the Prisoner of Azkaban (movie 3)



- Prior distribution  
mean = 81.3%, sd = 1.9%
- Data

🍅 229, 🌿 24

# Harry Potter and the Prisoner of Azkaban (movie 3)



- Prior distribution  
mean = 81.3%, sd = 1.9%
- Data  
 229,  24
- Posterior distribution  
mean = 84.7%, sd = 1.4%

# Harry Potter and the Goblet of Fire (movie 4)



# Harry Potter and the Goblet of Fire (movie 4)



- Prior distribution

mean = 84.7%, sd = 1.4%

# Harry Potter and the Goblet of Fire (movie 4)



- Prior distribution  
mean = 84.7%, sd = 1.4%
- Data

🍅 218, 🌿 30

# Harry Potter and the Goblet of Fire (movie 4)



- Prior distribution  
 $\text{mean} = 84.7\%$ ,  $\text{sd} = 1.4\%$
- Data  
    🍅 218, ✳️ 30
- Posterior distribution  
 $\text{mean} = 85.6\%$ ,  $\text{sd} = 1.2\%$

# Harry Potter and the Order of the Phoenix (movie 5)



# Harry Potter and the Order of the Phoenix (movie 5)



- Prior distribution

mean = 85.6%, sd = 1.2%

# Harry Potter and the Order of the Phoenix (movie 5)



- Prior distribution  
mean = 85.6%, sd = 1.2%
- Data

🍅 192, 🌿 53

# Harry Potter and the Order of the Phoenix (movie 5)



- Prior distribution  
mean = 85.6%, sd = 1.2%
- Data  
 192,  53
- Posterior distribution  
mean = 84.1%, sd = 1.1%

# Harry Potter and the Half-Blood Prince (movie 6)



# Harry Potter and the Half-Blood Prince (movie 6)

- Prior distribution

mean = 84.1%, sd = 1.1%



# Harry Potter and the Half-Blood Prince (movie 6)



- Prior distribution  
mean = 84.1%, sd = 1.1%
- Data

🍅 226, 🍃 43

# Harry Potter and the Half-Blood Prince (movie 6)



- Prior distribution  
mean = 84.1%, sd = 1.1%
- Data  
 226,  43
- Posterior distribution  
mean = 84.0%, sd = 1.0%

# Harry Potter and the Deathly Hallows Part 1 (movie 7)



# Harry Potter and the Deathly Hallows

## Part 1 (movie 7)



- Prior distribution

mean = 84.0%, sd = 1.0%

# Harry Potter and the Deathly Hallows Part 1 (movie 7)



- Prior distribution

mean = 84.0%, sd = 1.0%

- Data

🍅 208, 🌿 57

# Harry Potter and the Deathly Hallows Part 1 (movie 7)



- Prior distribution  
 $\text{mean} = 84.0\%$ ,  $\text{sd} = 1.0\%$
- Data  
 208,  57
- Posterior distribution  
 $\text{mean} = 83.2\%$ ,  $\text{sd} = 0.9\%$

# **Harry Potter and the Deathly Hallows Part 2 (movie 8)**



# Harry Potter and the Deathly Hallows Part 2 (movie 8)



- Prior distribution

mean = 83.2%, sd = 0.9%

# Harry Potter and the Deathly Hallows Part 2 (movie 8)



- Prior distribution

mean = 83.2%, sd = 0.9%

- Data

🍅 303, 🍃 13

# Harry Potter and the Deathly Hallows Part 2 (movie 8)



- Prior distribution

mean = 83.2%, sd = 0.9%

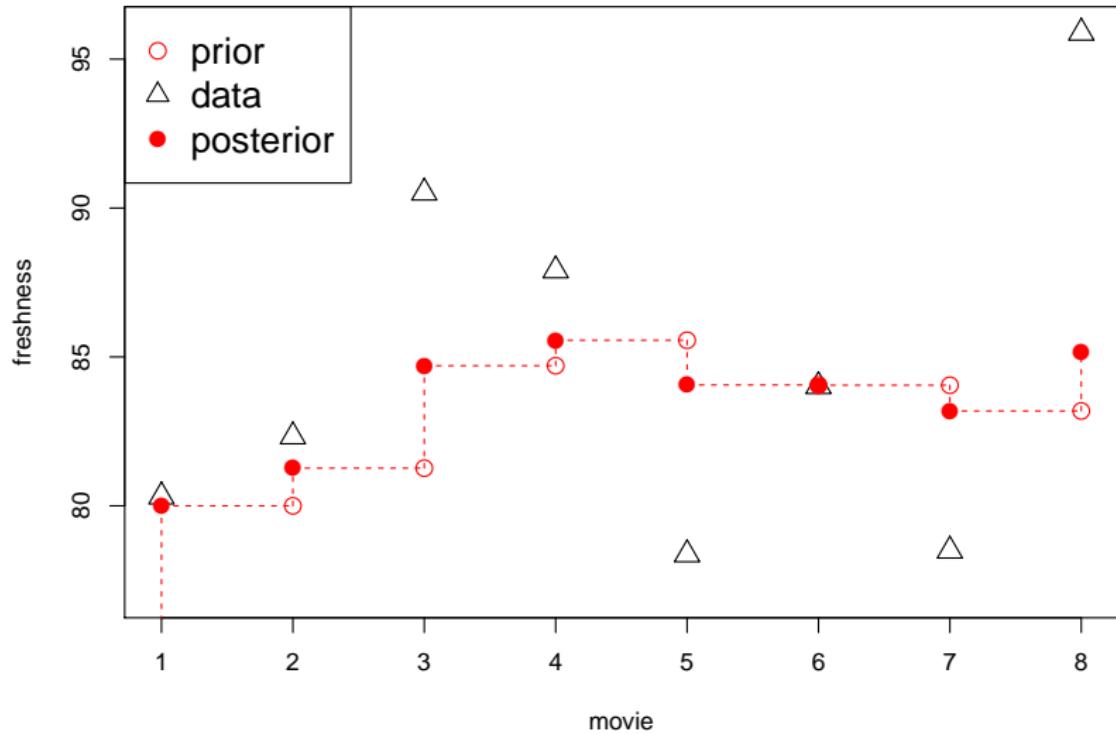
- Data

🍅 303, 🍃 13

- Posterior distribution

mean = 85.2%, sd = 0.8%

# The overall Harry Potter movie series



# Monte Carlo approximation

- For complicate problems, posterior  $p(\theta \mid D)$  doesn't look like a distribution we know.
- Draw many samples from the posterior

$$\theta^{(1)}, \dots, \theta^{(N)} \stackrel{\text{iid}}{\sim} p(\theta \mid D)$$

# Monte Carlo approximation

- For complicate problems, posterior  $p(\theta | D)$  doesn't look like a distribution we know.
- Draw many samples from the posterior

$$\theta^{(1)}, \dots, \theta^{(N)} \stackrel{\text{iid}}{\sim} p(\theta | D)$$

- Use sample mean to approximate posterior mean

$$E(\theta | Y) \approx \frac{1}{N} \left[ \theta^{(1)} + \theta^{(2)} + \dots + \theta^{(N)} \right]$$

# Markov chain Monte Carlo (MCMC)

- Drawing independent samples can be hard...
- We can draw samples sequentially:  
next sample's value depends on the current  
sample's value (Markov property)

$$\theta^{(1)} \longrightarrow \theta^{(2)} \longrightarrow \dots \longrightarrow \theta^{(N)}$$

# Markov chain Monte Carlo (MCMC)

- Drawing independent samples can be hard...
- We can draw samples sequentially:  
next sample's value depends on the current sample's value (Markov property)

$$\theta^{(1)} \longrightarrow \theta^{(2)} \longrightarrow \dots \longrightarrow \theta^{(N)}$$

- If the number of iteration  $N$  large enough,
  - Use MCMC samples' mean to approximate the posterior mean  $E(\theta | D)$
  - Starting point  $\theta^{(1)}$  doesn't matter

# Hypothesis testing

- Suppose we want to test whether the coin is fair or not. What are the hypotheses?

# Hypothesis testing

- Suppose we want to test whether the coin is fair or not. What are the hypotheses?
- Null hypothesis

$$H_0 : \theta = 0.5$$

Alternative hypothesis

$$H_1 : \theta \neq 0.5$$

# Hypothesis testing

- Suppose we want to test whether the coin is fair or not. What are the hypotheses?
- Null hypothesis

$$H_0 : \theta = 0.5$$

Alternative hypothesis

$$H_1 : \theta \neq 0.5$$

- Data: in 100 flips, we got 55 heads and 45 tails

# P-value: the frequentist (classic) way

*p*-value = 0.31

# P-value: the frequentist (classic) way

$p\text{-value} = 0.31$

- Reject the null if the p-value is small ( $< 0.05$ ).

# P-value: the frequentist (classic) way

$p\text{-value} = 0.31$

- Reject the null if the p-value is small ( $< 0.05$ ).
- What's the meaning of the  $p$ -value?

# P-value: the frequentist (classic) way

$p\text{-value} = 0.31$

- Reject the null if the p-value is small ( $< 0.05$ ).
- What's the meaning of the  $p$ -value?

$P(\text{current data or more extreme} \mid H_0) = 0.31$

# Bayesian hypothesis testing

- Equal prior probabilities

$$P(H_0) = P(H_1) = 0.5$$

# Bayesian hypothesis testing

- Equal prior probabilities

$$P(H_0) = P(H_1) = 0.5$$

- How likely to get current data

$$P(D \mid H_0) = 0.048, \quad P(D \mid H_1) = 0.010$$

# Bayesian hypothesis testing

- Equal prior probabilities

$$P(H_0) = P(H_1) = 0.5$$

- How likely to get current data

$$P(D \mid H_0) = 0.048, \quad P(D \mid H_1) = 0.010$$

- Use the Bayes rule to get posterior probabilities

$$P(H_0 \mid D) = \frac{p(D \mid H_0)P(H_0)}{p(D \mid H_0)P(H_0) + p(D \mid H_1)P(H_1)}$$

# Bayesian hypothesis testing

- Equal prior probabilities

$$P(H_0) = P(H_1) = 0.5$$

- How likely to get current data

$$P(D \mid H_0) = 0.048, \quad P(D \mid H_1) = 0.010$$

- Use the Bayes rule to get posterior probabilities

$$P(H_0 \mid D) = 0.83$$

$$P(H_1 \mid D) = 0.17$$

# Bayesian hypothesis testing

- Equal prior probabilities

$$P(H_0) = P(H_1) = 0.5$$

- How likely to get current data

$$P(D \mid H_0) = 0.048, \quad P(D \mid H_1) = 0.010$$

- Use the Bayes rule to get posterior probabilities

$$P(H_0 \mid D) = 0.83$$

$$P(H_1 \mid D) = 0.17$$

- Favor  $H_0$  if  $P(H_0 \mid D) > P(H_1 \mid D)$

# References and resources

- Textbook of my Bayesian intro class:  
Hoff (2010) A First Course in Bayesian Statistical Methods
- Another good reference book:  
Gelman, *et al* (2013) Bayesian Data Analysis, 3rd edition
- A textbook on MCMC:  
Robert and Casella (2004) Monte Carlo Statistical Methods, 2nd edition
- Feel free to check out my teaching slides  
<https://github.com/yingboli/MyTeachingSlides>