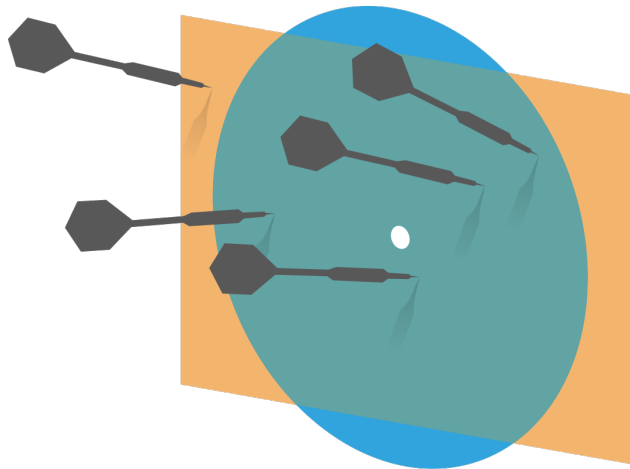# Probability and Statistics for Computer Science
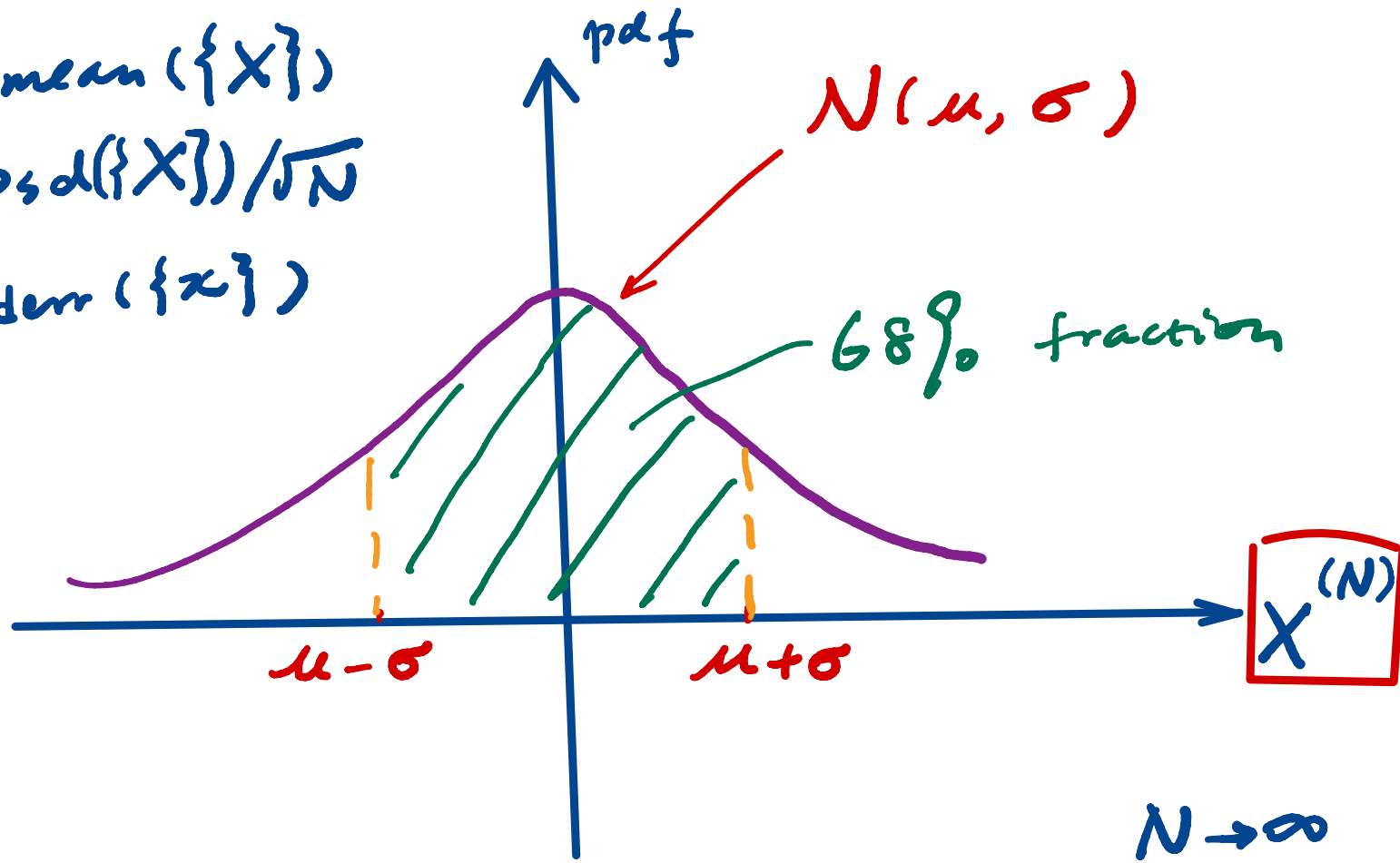
"Statistical thinking will one day be as necessary for efficient citizenship as the ability to read and write." H. G. Wells
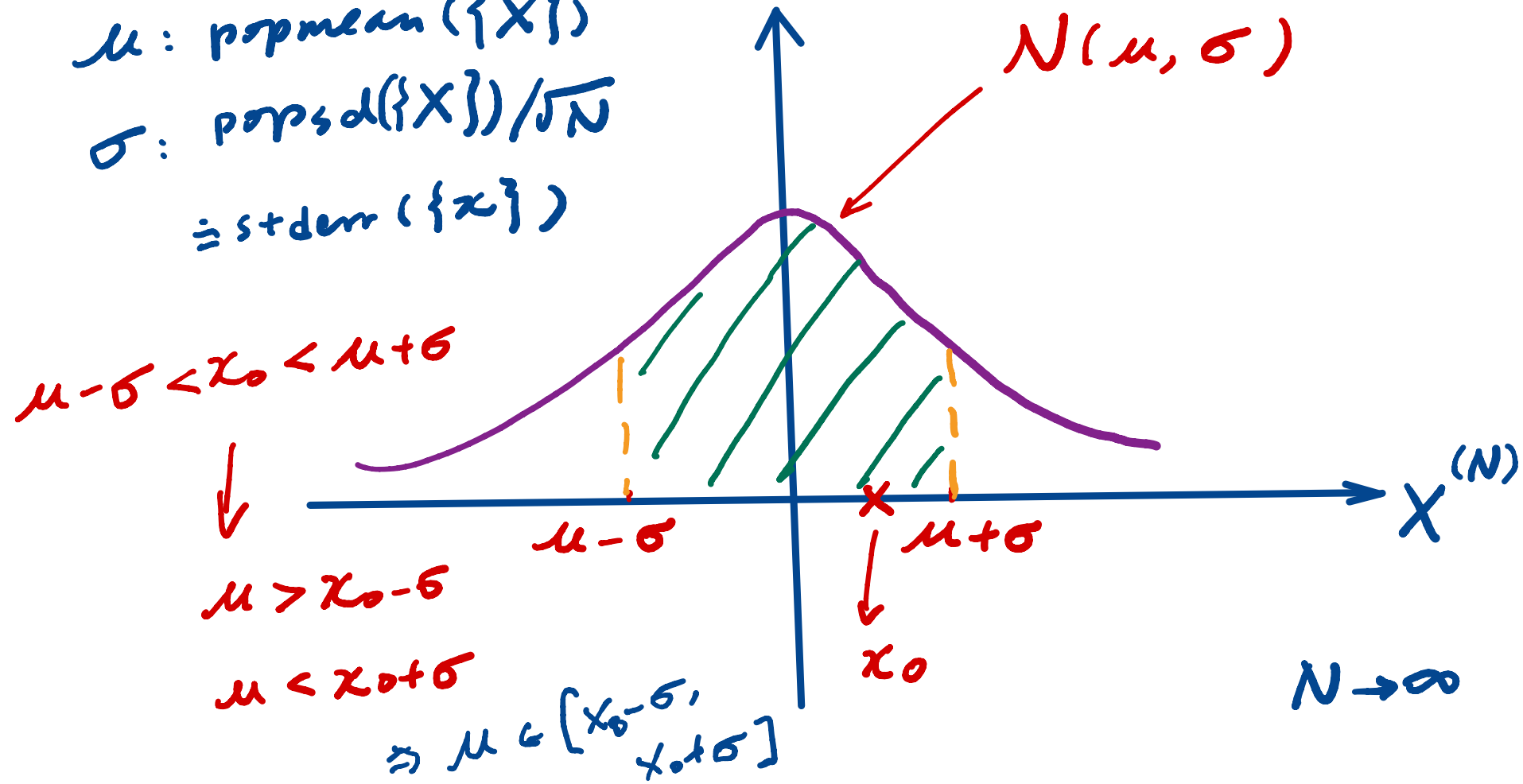
Credit: wikipedia

Hongye Liu, Teaching Assistant Prof, CS361, UIUC, 10.20.2020

# Interpretation of Confidence Interval

$\mu$ : popmean($\{X\}$)

$\sigma$ : popsd($\{X\}$)/$\sqrt{N}$

$\doteq$ stderr($\{x\}$)

pdf

$N(\mu, \sigma)$

68% fraction

$\mu - \sigma$

$\mu + \sigma$

$X^{(N)}$

$N \to \infty$

# Interpretation of Confidence Interval

$\mu$ : popmean $(\{X\})$

$\sigma$ : popsd $(\{X\})/\sqrt{N}$

$\doteq$ stderr $(\{x\})$

$N(\mu, \sigma)$

$\mu - \sigma < x_0 < \mu + \sigma$

$\mu > x_0 - \sigma$

$\mu < x_0 + \sigma$

$\Rightarrow \mu \in [x_0 - \sigma, x_0 + \sigma]$

$\mu - \sigma$

$\mu + \sigma$

$x_0$

$X^{(N)}$

$N \to \infty$

**Figure 8.5** A sample of one hundred observed 95% confidence intervals based on samples of size 26 from the normal distribution with mean $\mu = 5.1$ and standard deviation $\sigma = 1.6$. In this figure, 94% of the intervals contain the value of $\mu$.
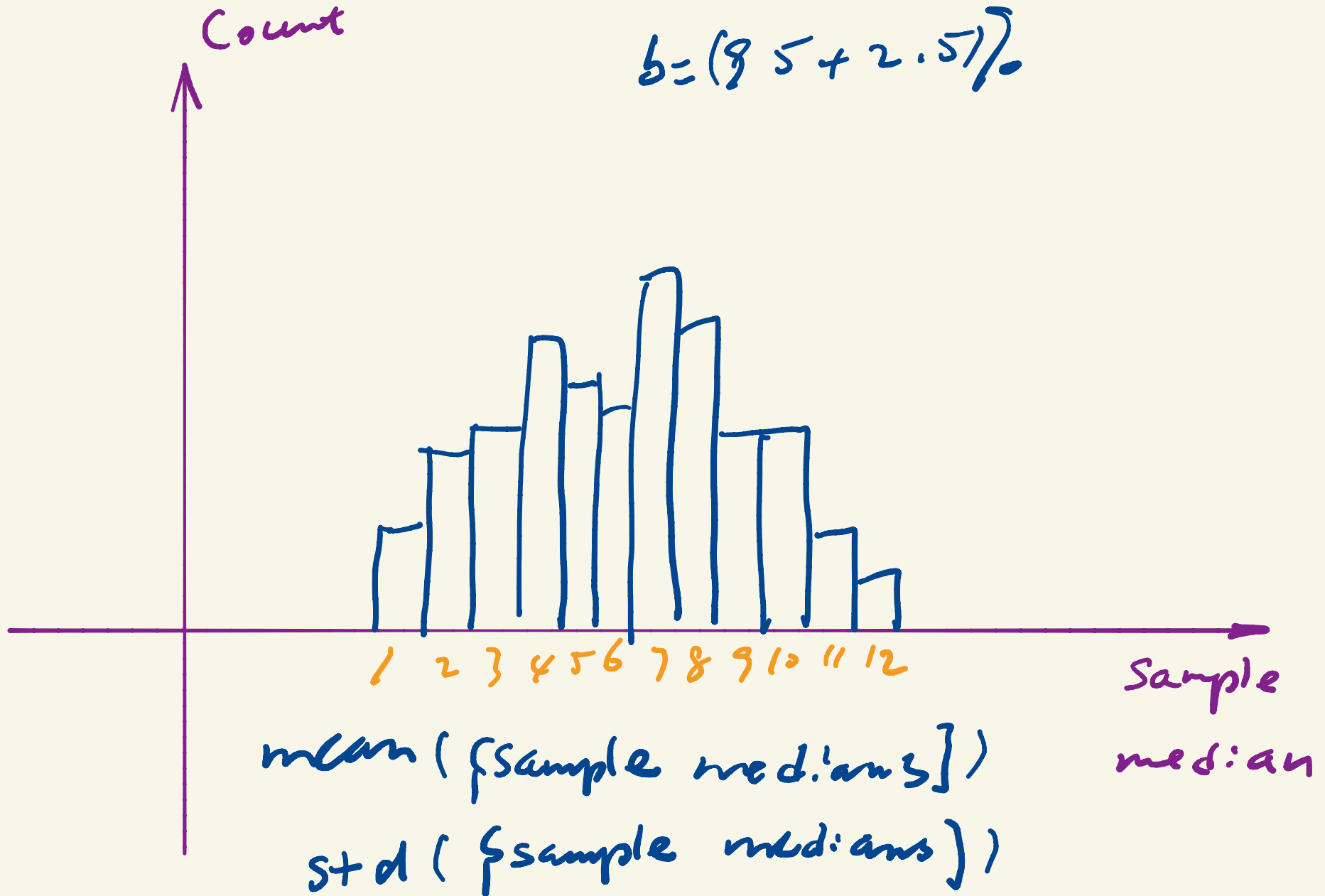


Degroot    Pg 487

# Bootstrap Histogram

$a = 2.5\%$

$b = (85 + 2.5)\%$

Count



1 2 3 4 5 6 7 8 9 10 11 12

Sample median

mean ( {sample medians} )

std ( {sample medians} )

# Last time

* Hypothesis test

* Chi-square test

* Maximum likelihood estimation

# Objectives

* More on Maximum likelihood Estimation (MLE)

* Bayesian Inference (MAP)

*Handwritten annotations:*

Likelihood

frequentist

Bayesian

Posterior distri.

$\theta$

# Maximum likelihood estimation (MLE)

✳ We write the probability of seeing the data D given parameter θ

$$L(\theta) = P(D|\theta)$$

✳ The **likelihood function** $L(\theta)$ is **not** a probability distribution

✳ The **maximum likelihood estimate (MLE)** of θ is

$$\hat{\theta} = arg\ \max_{\theta}\ L(\theta)$$

# Likelihood function: binomial example

❋ Suppose we have a coin with unknown probability of θ coming up heads

$$P(X=k) = \binom{N}{k} p^k (1-p)^{N-k}$$

❋ We toss it **10** times and observe **7** heads
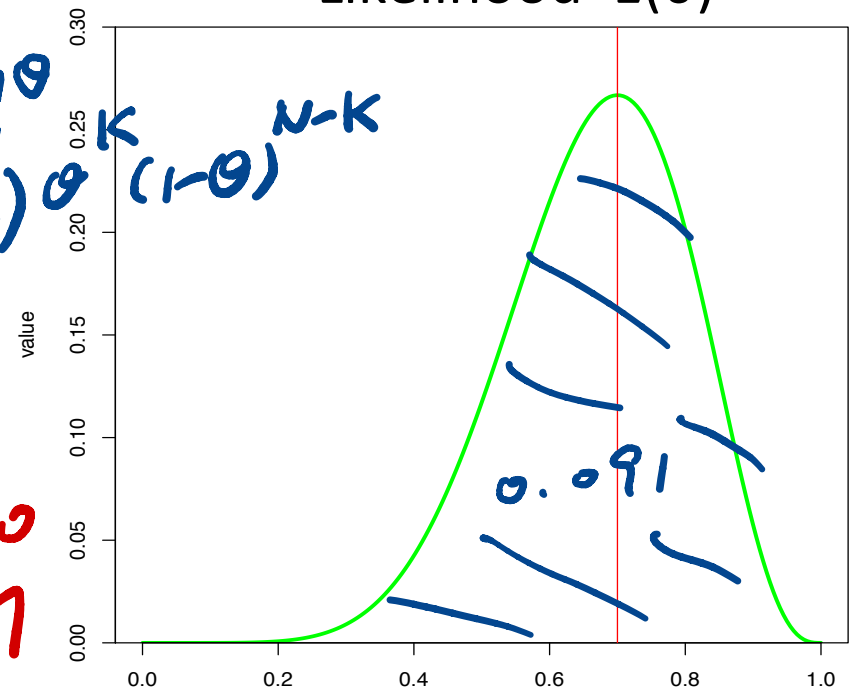
$$L(\theta) = \binom{N}{k} \theta^k (1-\theta)^{N-k}$$

$$p \to \theta$$

❋ The likelihood function is:

$$P(D|\theta) = \binom{10}{7} \theta^7 (1-\theta)^3$$

$$D: N = 10$$
$$k = 7$$

❋ The MLE is

$$\hat{\theta} = 0.7$$

Likelihood  **L**(θ)

0.091

$L(\theta)$ is not a distri. !!

# Q. What is the MLE of binomial N=12, k=7

A. 12!/7!/5!

B. 7/12

C. 5/12

D. 12/7

$$\hat{\theta} = \frac{k}{N}$$

for Bino. Likelihood

# Q. What is the MLE of Poisson k₁=5, k₂=7, n=2

A. 6

B. 35/2

C. 12

D. other

$\hat{\theta} = MLE \text{ (Poisson with } \lambda\text{)}$

$$= \frac{\sum K_i}{N}$$

$L(\theta) = \prod L(D_i | \theta)$

$\log \frac{L(\theta)}{\sum}$

# MLE Example

You find a 5-sided die and want to estimate its $\underset{\text{independently}}{\overset{\text{All}}{\underbrace{\phantom{xxxxx}}}}$ probability θ of coming up 5, you decided to roll it 12 times and then roll it until it comes up 5. You rolled 15 times altogether and found there were 3 times when the die came up 5. Write down the likelihood function L(θ).

$$L(\theta) = P(D \mid \theta) = P(D_1 \mid \theta) P(D_2 \mid \theta)$$

D ?

"5" → S

not "5" → F

# "S" = 3

12 rolls

3 rolls

2 "S"  10 "F"

F F S

# MLE Example

You find a 5-sided die and want to estimate its probability θ of coming up 5, you decided to roll it 12 times and then roll it until it comes up 5. You rolled 15 times altogether and found there were 3 times when the die came up 5. Write down the likelihood function L(θ).

$L_1(\theta)$  Exp-1     12 times to check
                              # of "5"    Bino.

$L_2(\theta)$  Exp-2     .... 1st "5"              Geom.

$L(\theta) = L_1(\theta) \cdot L_2(\theta)$

$$L(\theta) = L_1(\theta) L_2(\theta)$$

$$= \binom{12}{2} \theta^2 (1-\theta)^{10} \cdot ((-\theta)^2 \theta^1$$

$$= \binom{12}{2} \theta^3 (1-\theta)^{12}$$

$$\hat{\theta} = \underset{\theta}{\arg\max} \; L(\theta)$$

$$\hat{\theta} = \cdots$$

$$\rightarrow L(\theta) = C \; \theta^3 (1-\theta)^{12}$$

$$\rightarrow \log L(\theta) = \log C + 3 \log \theta + 12 \log (1-\theta)$$

$$\frac{d \log L(\theta)}{d\theta} = 0 + \frac{3}{\theta} - \frac{12}{1-\theta} = 0$$

$$\frac{3}{\theta} = \frac{12}{1-\theta}$$

$$12\theta = 3 - 3\theta$$

$$\hat{\theta} = \frac{3}{15} = \frac{1}{5}$$

# Drawbacks of MLE

✳ Maximizing some likelihood or log-likelihood function is mathematically hard

✳ If there are few data items, the MLE estimate maybe very unreliable

  ✳ If we observe 3 heads in 10 coin tosses, should we accept that p(heads)= 0.3 ?

  ✳ If we observe 0 heads in 2 coin tosses, should we accept that p(heads)= 0 ?

# Bayesian inference

✳ In MLE, we maximized the likelihood function

$$L(\theta) = P(D|\theta)$$

✳ In Bayesian inference, we will maximize the **posterior**, which is the probability of the parameters **θ** given the observed data D.

$$P(\theta|D)$$

*θ is RV!*

✳ Unlike $L(\theta)$, the posterior is a probability distribution

✳ The value of **θ** that maximizes $P(\theta|D)$ is called the **maximum a posterior (MAP)** estimate $\hat{\theta}$

# The components of Bayesian Inference

* From Bayes rule

$$P(\theta|D) = \frac{P(D|\theta)\, P(\theta)}{P(D)}$$

# The components of Bayesian Inference

* From Bayes rule

$$P(\theta \mid D) = \frac{\overbrace{P(D \mid \theta)}^{L(\theta)}\; \boxed{P(\theta)}}{P(D)}$$

* Prior, assumed distribution of $\theta$ before seeing data **D**

* Likelihood function of $\theta$ seeing **D** : $L(\theta)$

* Total Probability seeing **D** --- P(**D**)

* Posterior, distribution of $\theta$ given **D**

# The usefulness of Bayesian inference

✳ From Bayes rule

$$P(\theta|D) = \frac{P(D|\theta)P(\theta)}{P(D)}$$

✳ Bayesian inference allows us to include prior beliefs about θ in the prior $P(\theta)$, which is useful

  ✳ When we have reasonable beliefs, such as a coin can not have P(heads) = 0

  ✳ When there isn't much data

  ✳ We get a distribution of the posterior, not just one maxima

# Bayesian Inference: a discrete prior

⁕ Suppose we have a coin of unknown probability θ of heads

⁕ We see 7 heads in 10 tosses (**D**)

⁕ We assume the prior about θ.

$$P(\theta) = \begin{cases} \frac{2}{3} & if \ \theta = 0.5 \\ \frac{1}{3} & if \ \theta = 0.6 \\ 0 & otherwise \end{cases}$$

⁕ We have this likelihood:

$$P(D|\theta) = \binom{10}{7} \theta^7 (1 - \theta)^3$$

⁕ What is the posterior $P(\theta|D)$ ?

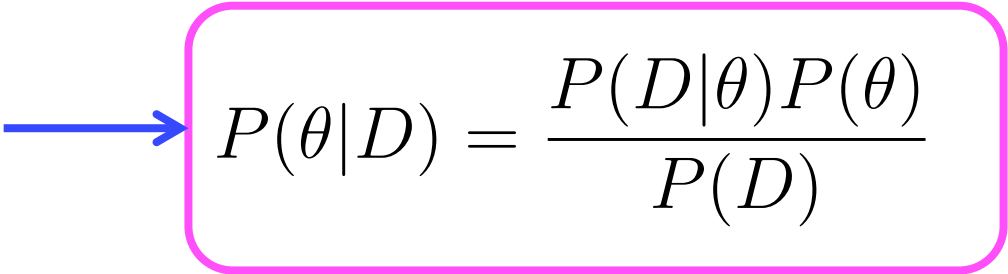# Bayesian Inference: a discrete prior

✳ We see 7 heads in 10 tosses (**D**)

✳ We assume the prior about θ.
$$P(\theta) = \begin{cases} \frac{2}{3} & if\ \theta = 0.5 \\ \frac{1}{3} & if\ \theta = 0.6 \\ 0 & otherwise \end{cases}$$

✳ We have this likelihood:

$$P(D|\theta) = \binom{10}{7} \theta^7 (1-\theta)^3$$

✳ What is the posterior $P(\theta|D)$ ?

$$P(\theta|D) = \frac{P(D|\theta)P(\theta)}{P(D)}$$

# Bayesian Inference: a discrete prior

✳ We see 7 heads in 10 tosses (**D**)

✳ We assume the prior about θ.

$$P(\theta) = \begin{cases} \frac{2}{3} & if \ \theta = 0.5 \\ \frac{1}{3} & if \ \theta = 0.6 \\ 0 & otherwise \end{cases}$$

✳ We have this likelihood:

$$P(D|\theta) = \binom{10}{7}\theta^7(1-\theta)^3$$

✳ What is the posterior $P(\theta|D)$ ?

$$P(\theta|D) = \frac{P(D|\theta)P(\theta)}{P(D)} \qquad P(D) = \sum_{\theta_i \in \theta} P(D|\theta_i)P(\theta_i)$$

$$P(\theta \mid D) = \frac{P(D \mid \theta) P(\theta)}{P(D)}$$

$$P(\theta) = \begin{cases} \frac{2}{3} & \theta = 0.5 \\ \frac{1}{3} & \theta = 0.6 \\ 0 & \text{other} \end{cases}$$

$$P(D \mid \theta) = \binom{10}{7} \theta^7 (1-\theta)^3$$

$$P(D) = \sum P(D \mid \theta_i) \cdot P(\theta_i) \qquad \underset{\theta = 0.5}{\swarrow} \qquad \underset{\theta = 0.6}{\swarrow}$$

$$= \underline{\binom{10}{7} 0.5^7 \cdot .5^3 \cdot \frac{2}{3}} + \underline{\binom{10}{7} 0.6^7 \cdot .4^3 \cdot \frac{1}{3}}$$

$$P(\theta \mid D) = \begin{cases} \boxed{0.52} & \theta = 0.5 \\ \underline{0.48} & \theta = 0.6 \\ 0 & \text{other} \end{cases}$$

which $\theta$ maximize $P(\theta \mid D)$?

$$\hat{\theta} = 0.5$$

# Bayesian Inference: a discrete prior

✳ We see 7 heads in 10 tosses (**D**)

✳ We assume the prior about θ.
$$P(\theta) = \begin{cases} \frac{2}{3} & if \ \theta = 0.5 \\ \frac{1}{3} & if \ \theta = 0.6 \\ 0 & otherwise \end{cases}$$

✳ We have this likelihood:
$$P(D|\theta) = \binom{10}{7} \theta^7 (1-\theta)^3$$

✳ What is the posterior $P(\theta|D)$ ?

$$P(\theta|D) = \begin{cases} 0.52 & if \ \theta = 0.5 \\ 0.48 & if \ \theta = 0.6 \\ 0 & otherwise \end{cases}$$

**MAP** $\hat{\theta}$ **=0.5**
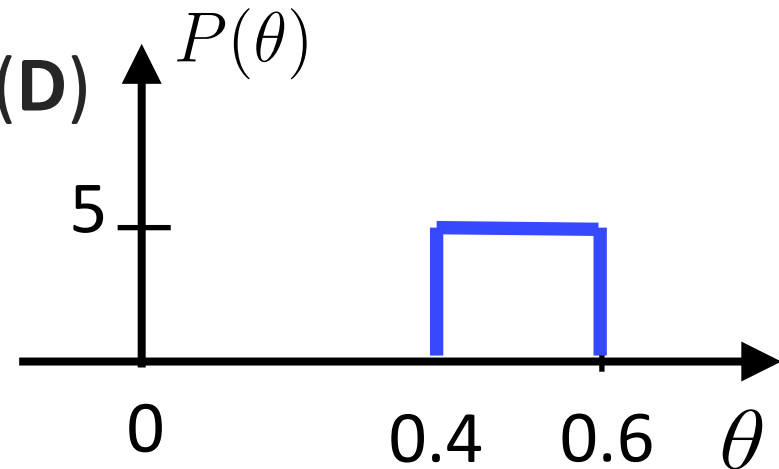
Biased by the prior

# Bayesian Inference: a continuous prior

✳ Suppose we have a coin of unknown probability θ of heads

✳ We see 7 heads in 10 tosses (**D**)

✳ We assume
$$P(\theta) = \begin{cases} 5 & if \ \theta \in [0.4, 0.6] \\ 0 & if \ \theta \notin [0.4, 0.6] \end{cases}$$
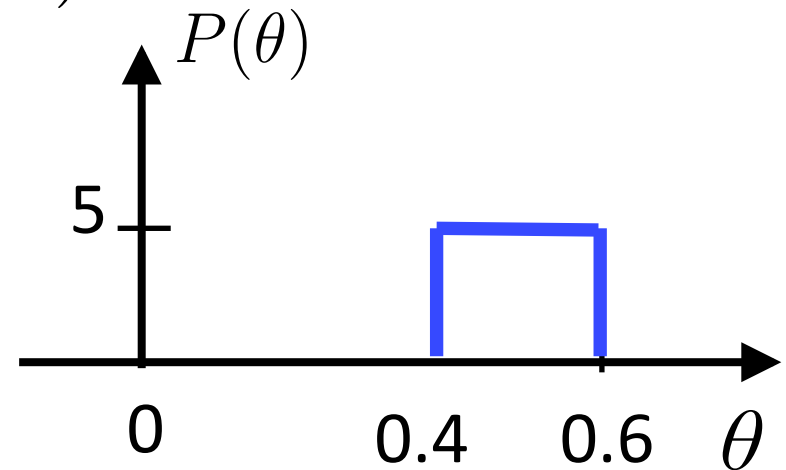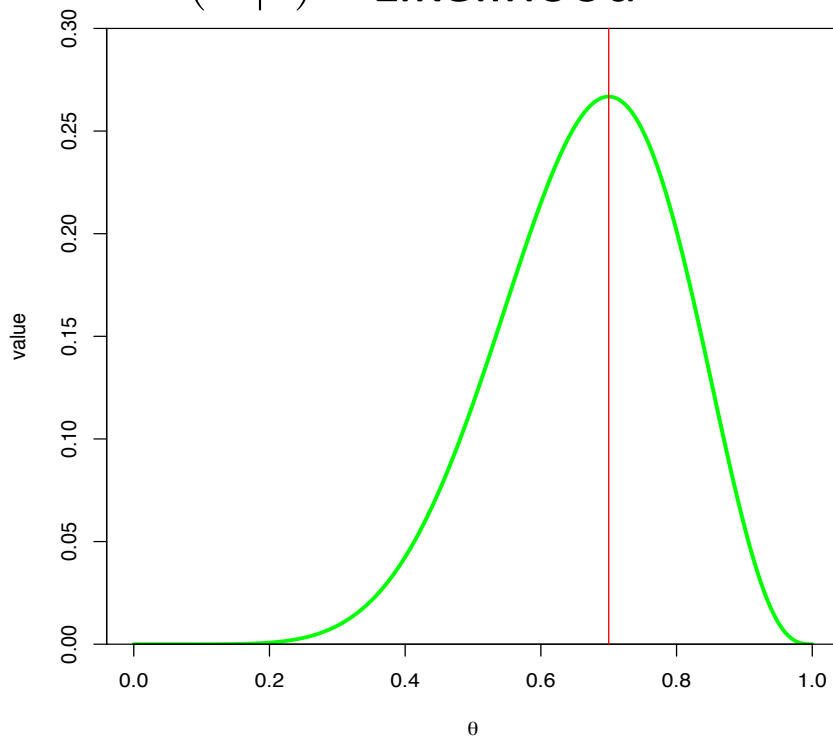
✳ What is the posterior $P(\theta|D)$ ?

# Bayesian Inference: a continuous prior

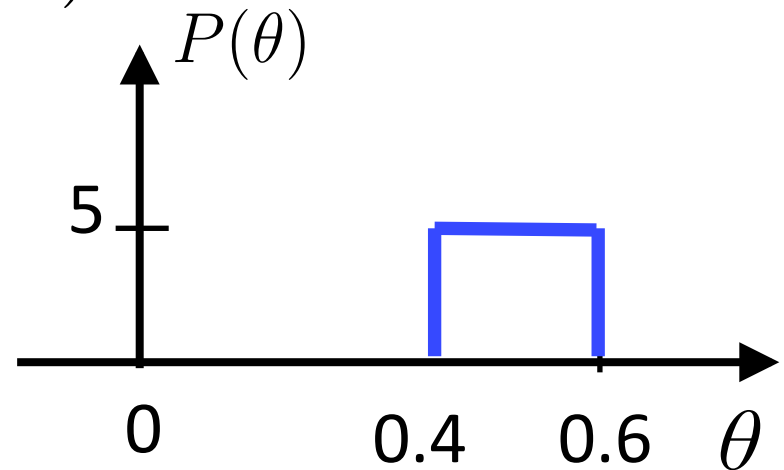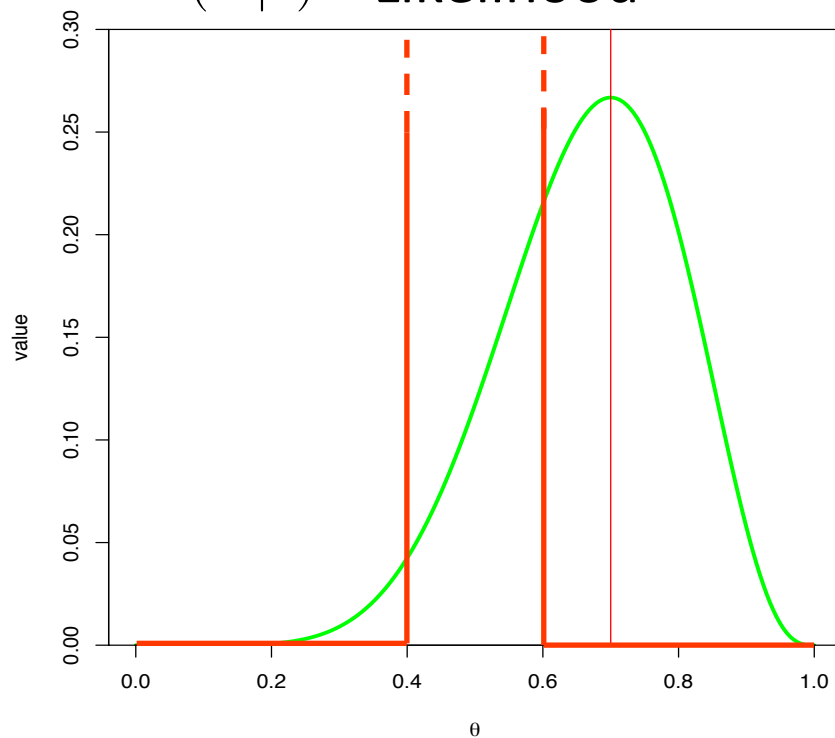✳ What is the posterior $P(\theta|D)$ ?

$P(D|\theta)$ = Likelihood



$$P(\theta) = \begin{cases} 5 & if\ \theta \in [0.4, 0.6] \\ 0 & if\ \theta \notin [0.4, 0.6] \end{cases}$$

$$P(\theta|D) \propto P(D|\theta)P(\theta)$$

# Bayesian Inference: a continuous prior

✳ What is the posterior $P(\theta|D)$ ?

$P(D|\theta)$ = Likelihood



$P(\theta)$

5

0    0.4    0.6    $\theta$

$$P(\theta) = \begin{cases} 5 & if\ \theta \in [0.4, 0.6] \\ 0 & if\ \theta \notin [0.4, 0.6] \end{cases}$$
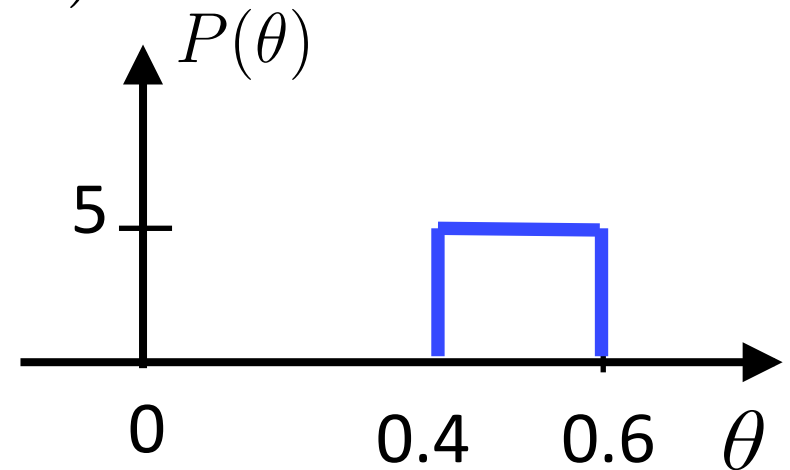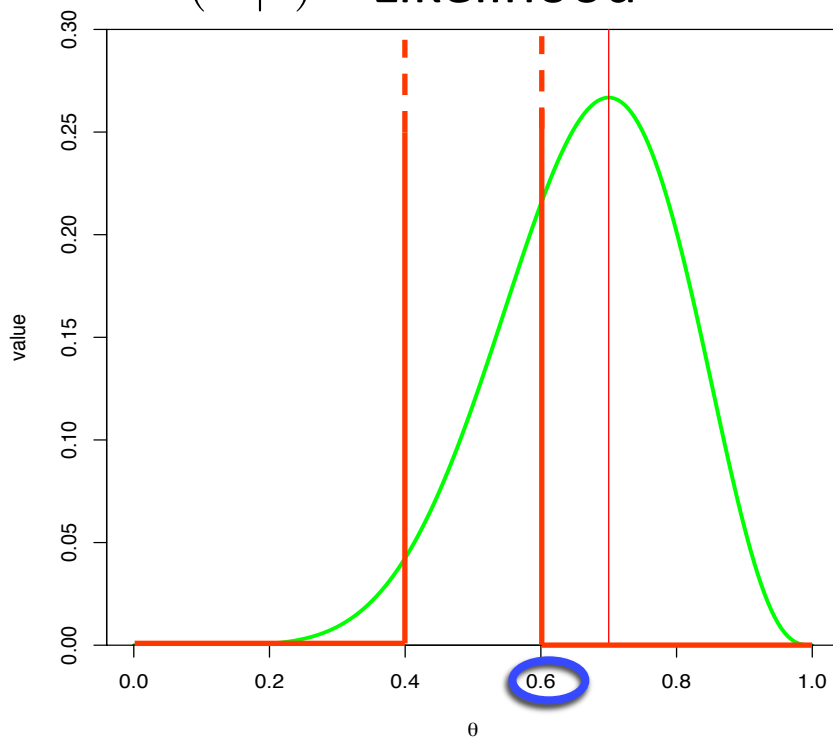
without $P(D)$

$$P(\theta|D) \propto P(D|\theta)P(\theta)$$

$\hat{\theta} = 0.6$

# Bayesian Inference: a continuous prior

❋ What is the posterior $P(\theta|D)$ ?

$P(D|\theta)$ = Likelihood



$P(\theta)$

$$P(\theta) = \begin{cases} 5 & if \ \theta \in [0.4, 0.6] \\ 0 & if \ \theta \notin [0.4, 0.6] \end{cases}$$
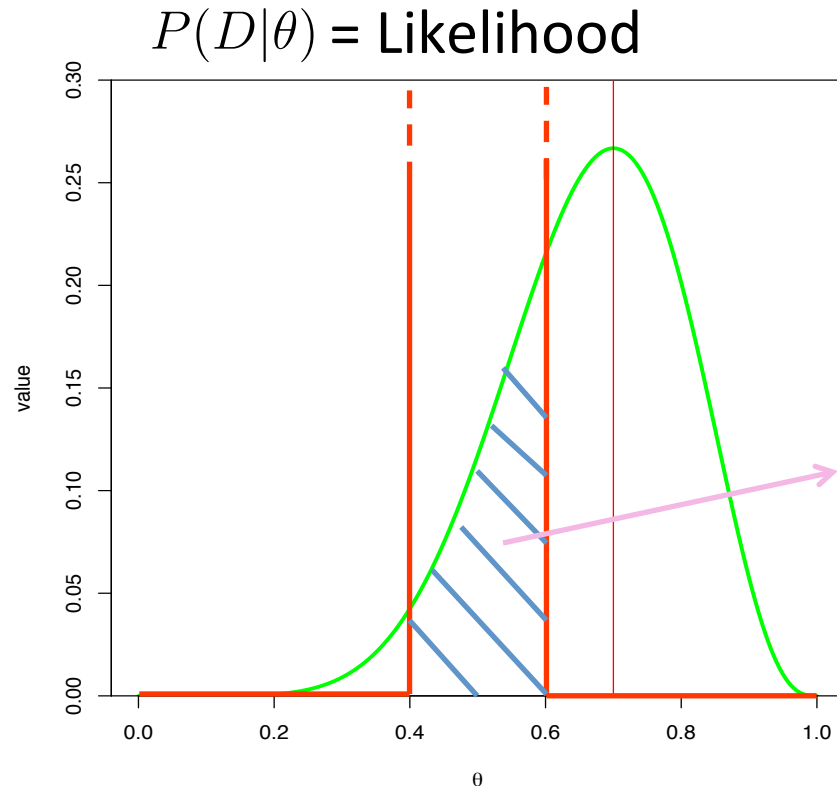
$$P(\theta|D) \propto P(D|\theta)P(\theta)$$

**MAP** $\hat{\theta}$ **=0.6**

# The constant in the Bayesian inference

$$P(D) = \int_\theta P(D|\theta)P(\theta)d\theta$$

$P(\theta)$

✳ It's not always possible to calculating P(D) in closed form.

✳ There are a lot of approximation methods.

$P(D|\theta)$ = Likelihood



Scale by 5 for this example

# Drawbacks of Bayesian inference

* Maximizing some posteriors $P(\theta|D)$ is difficult

* Some choices of prior $P(\theta)$ can overwhelm any data observed.

* It's hard to justify a choice of prior

# The concept of conjugacy

✳ For a given likelihood function $P(D|\theta)$, a prior $P(\theta)$ is its conjugate prior if it has the following properties:

  ✳ $P(\theta)$ belongs to a family of distributions that are expressive

  ✳ The posterior $P(\theta|D) \propto P(D|\theta)P(\theta)$ belongs to the same family of distribution as the prior $P(\theta)$

  ✳ The posterior $P(\theta|D)$ is easy to maximize

✳ For example, a conjugate prior for binomial likelihood function is Beta distribution

# Beta distribution

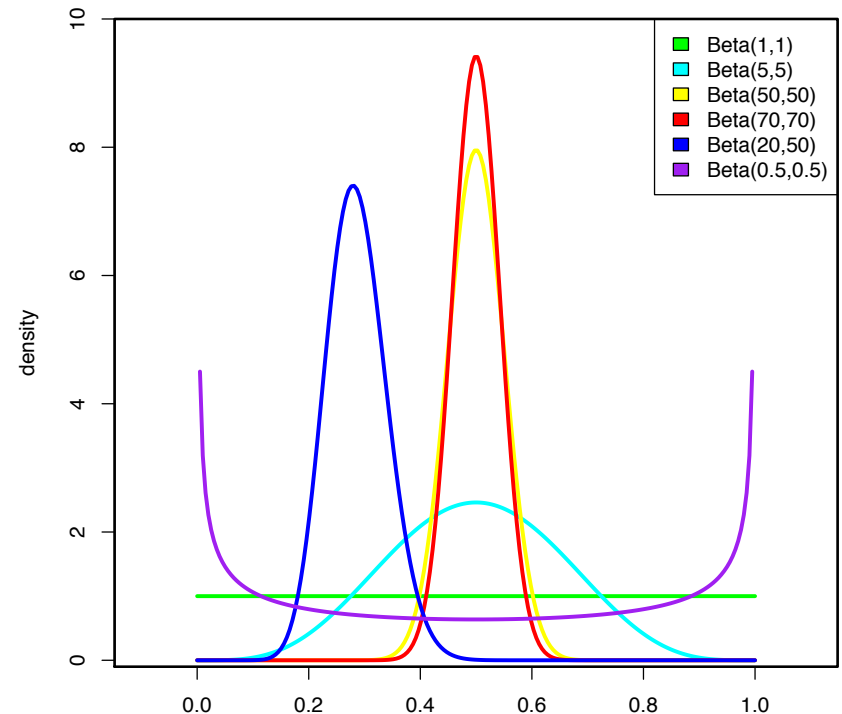* A distribution is Beta distribution if it has the following pdf:

$$P(\theta) = \begin{cases} K(\alpha, \beta)\theta^{\alpha-1}(1-\theta)^{\beta-1} \\ 0 \qquad\qquad\qquad D.W. \end{cases}$$

$$0 \leq \theta \leq 1$$
$$\alpha > 0, \beta > 0$$

$$K(\alpha, \beta) = \frac{\Gamma(\alpha + \beta)}{\Gamma(\alpha)\Gamma(\beta)}$$

* Is an expressive family of distributions

* $Beta(\alpha = 1, \beta = 1)$ is underline{uniform}

**pdf of Beta – distribution**



Legend:
- Beta(1,1)
- Beta(5,5)
- Beta(50,50)
- Beta(70,70)
- Beta(20,50)
- Beta(0.5,0.5)

density vs $\theta$

# Q. Beta distribution is a continuous probability distribution

A. TRUE

B. FALSE

# Beta distribution as the conjugate prior for Binomial likelihood

* The likelihood is Binomial ($N$, $k$)

$$P(D|\theta) = \binom{N}{k} \theta^k (1-\theta)^{N-k}$$

* The Beta distribution is used as the prior

$$P(\theta) = K(\alpha, \beta)\theta^{\alpha-1}(1-\theta)^{\beta-1}$$

$$\alpha = 1$$
$$\beta = 1$$

* So $\quad P(\theta|D) \propto \theta^{\alpha+k-1}(1-\theta)^{\beta+N-k-1}$

$$\theta^{\hat{\alpha}-1} \quad (1-\theta)^{\hat{\beta}-1}$$

$$\hat{\alpha} = \alpha + k$$

* Then the posterior is $Beta(\alpha+k, \beta+N-k)$

$$\hat{\beta} = \beta + N - k$$

$$P(\theta|D) = K(\underbrace{\alpha+k}_{\hat{\alpha}}, \underbrace{\beta+N-k}_{\hat{\beta}})\theta^{\alpha+k-1}(1-\theta)^{\beta+N-k-1}$$
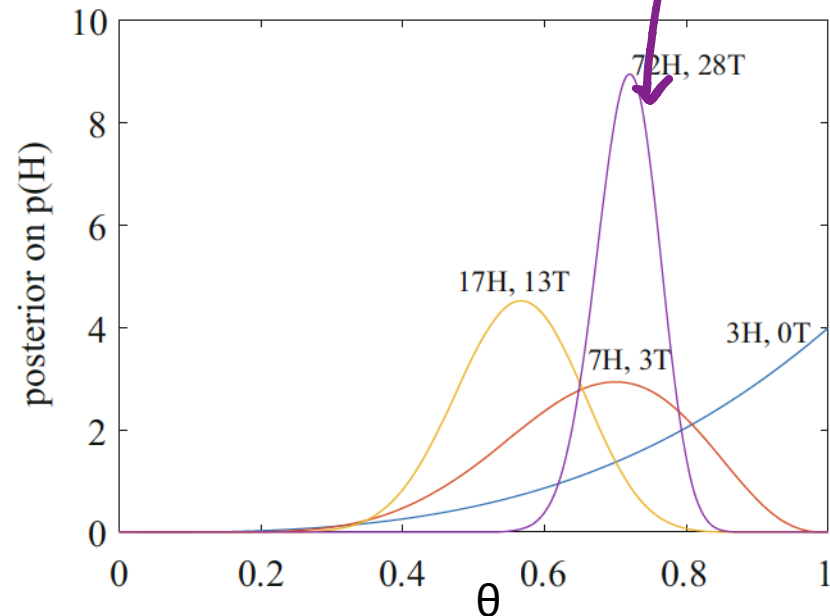
# The update of Bayesian posterior

✳ Since the posterior is in the same family as the conjugate prior, the posterior can be used as a new prior if more data is observed.

$\alpha = 1 \qquad \beta = 1 \qquad\qquad p(\theta | D)$

✳ Suppose we start with a uniform prior on the probability θ of heads

  ✳ Then we see 3H 0T

  ✳ Then we see 4H 3T for 7H 3T in total

  ✳ Then we see 10H 10T for 17H 13T in total

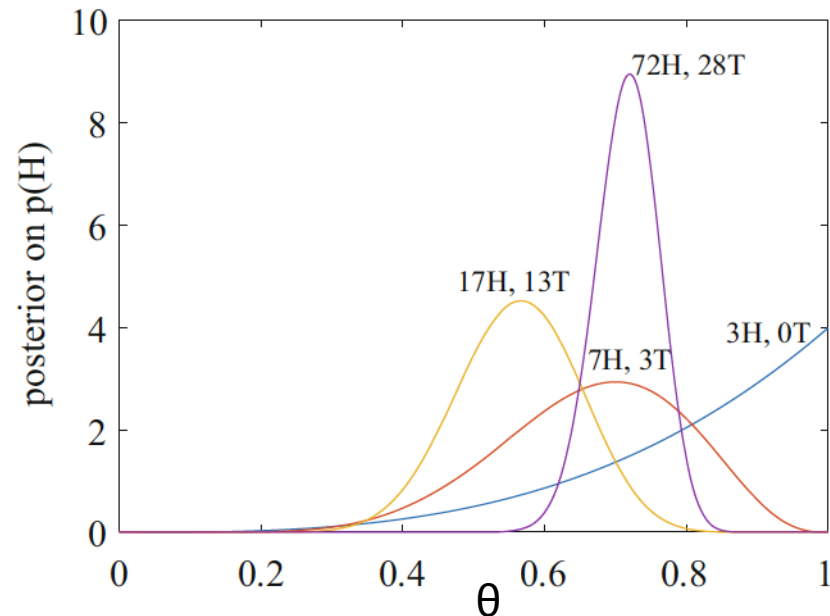  ✳ Then we see 55H 15T for 72H 28T in total

# The update of Bayesian posterior

✳ Since the posterior is in the same family as the conjugate prior, the posterior can be used as a new prior if more data is observed.

✳ Suppose we start with a uniform prior on the probability θ of heads

| N | k | $\hat{\alpha}$ | $\hat{\beta}$ |
|---|---|---|---|
| | | 1 | 1 |
| 3 | 0 | 1 | 4 |
| 10 | 7 | 8 | 7 |
| 30 | 17 | 25 | 20 |
| 100 | 72 | 97 | 48 |

# Simulation of the update of Bayesian posterior

https://seeing-theory.brown.edu/bayesian-inference/index.html

# Maximize the Bayesian posterior (MAP)

✳ The posterior of the previous example is

$$P(\theta|D) = K(\alpha + k, \beta + N - k)\theta^{\alpha+k-1}(1-\theta)^{\beta+N-k-1}$$

✳ Differentiating and setting to 0 gives the MAP estimate

$$\hat{\theta} = \frac{\alpha - 1 + k}{\alpha + \beta - 2 + N}$$

if $\alpha = 1$

$\beta = 1$

$\hat{\theta} = \frac{k}{N}$

θ

# Conjugate prior for other likelihood functions

✳  If the likelihood is Bernoulli or geometric, the conjugate prior is Beta

✳  If the likelihood is Poisson or Exponential, the conjugate prior is Gamma

✳  If the likelihood is normal with known variance, the conjugate prior is normal

$\theta$

# Assignments

* Finish Chapter 9 of the textbook

* Next time: Covariance matrix, PCA

# Additional References

✳ Robert V. Hogg, Elliot A. Tanis and Dale L. Zimmerman. "Probability and Statistical Inference"

✳ Morris H. Degroot and Mark J. Schervish "Probability and Statistics"

# See you next time

*See You!*