

# Interval estimation

# Agenda

- Sampling distribution
- Confidence interval
- Bootstrapping
- Bayesian credible interval

# Sampling distribution

# Statistic

Let  $T_n(X_1, \dots, X_n)$  be a transformation of a (either random or nonrandom) sample  $X_1, \dots, X_n$ .  $T_n(X_1, \dots, X_n)$  is called a statistic (統計量, statistics = 統計學).

# Sampling distribution

If  $X_1, \dots, X_n$  is a random sample,  $T_n(X_1, \dots, X_n)$  is also random. The distribution of  $T_n$  is called a sampling distribution.

# Example: sample mean

- If  $X_1, \dots, X_n \stackrel{iid}{\sim} N(\mu, \sigma^2)$ , then

$$\bar{X} = \frac{1}{n} \sum_{i=1}^n X_i \sim N\left(\mu, \frac{\sigma^2}{n}\right)$$

- If  $X_1, \dots, X_n \stackrel{iid}{\sim} f_X(x)$ , by CLT we have

$$\bar{X} \xrightarrow{d} N\left(\mu, \frac{\sigma^2}{n}\right) \text{ as } n \rightarrow \infty$$

# Example: sample variance

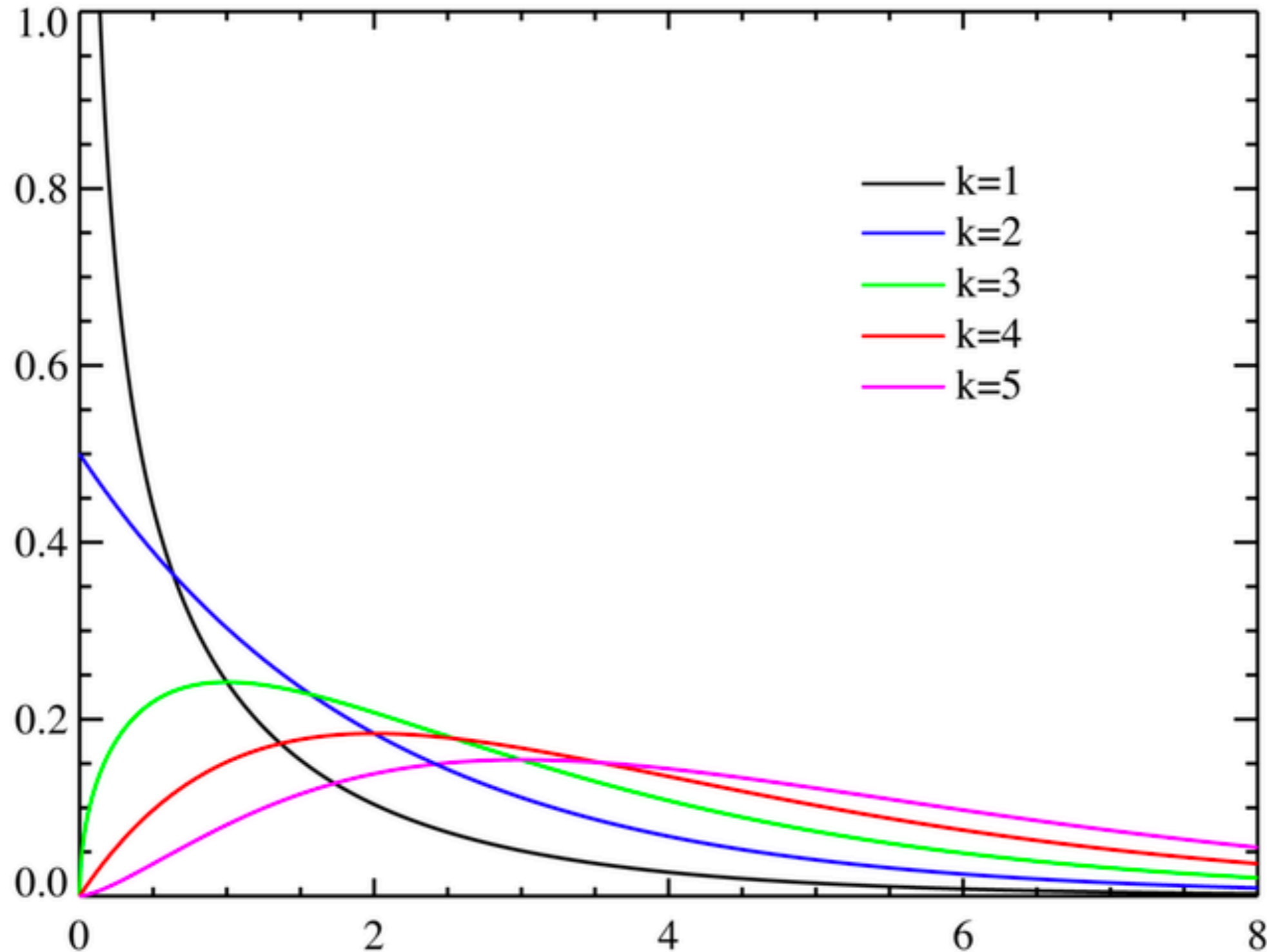
- If  $X_1, \dots, X_n \stackrel{iid}{\sim} N(\mu, \sigma^2)$ , then

$$s^2 = \frac{1}{n-1} \sum_{i=1}^n (X_i - \bar{X})^2 \sim \chi_{n-1}^2$$

- If  $X_1, \dots, X_n \stackrel{iid}{\sim} f_X(x)$ , by CLT we have

$$s^2 \xrightarrow{d} \chi_{n-1}^2 \text{ as } n \rightarrow \infty$$

# Chi-square distribution





# Asymptotic distribution of MLE

- Let  $\hat{\theta}$  be the MLE of  $\theta$ . Then,

$$\hat{\theta} \xrightarrow{d} N(\theta_0, I(\theta_0)) \text{ as } n \rightarrow \infty$$

where  $I(\theta_0) = E \left[ -\frac{\partial^2}{\partial \theta \partial \theta^T} \ell(\theta_0) \right]$ ,  $\ell(\theta)$  denotes the

log-likelihood function of  $\theta$

Confidence interval

# Confidence interval (C.I.)

An interval  $[L_n(X_1, \dots, X_n), U_n(X_1, \dots, X_n)]$  is an  $(1 - \alpha)$  confidence interval of a parameter  $\theta$  if and only if

$$P\left(\theta \in [L_n(X_1, \dots, X_n), U_n(X_1, \dots, X_n)]\right) = 1 - \alpha$$

- the interval depends on random sample
- explanation

# Pivotal quantity method

1. start with a point estimation
2. find an appropriate pivotal quantity and derive its sampling distribution
3. derive the lower and upper bounds of the interval by the above sampling distribution

# Example: sample mean

$X_1, \dots, X_n \stackrel{iid}{\sim} N(\mu, \sigma^2) \leftarrow z \in \mathbb{R}$   
parameter of interest

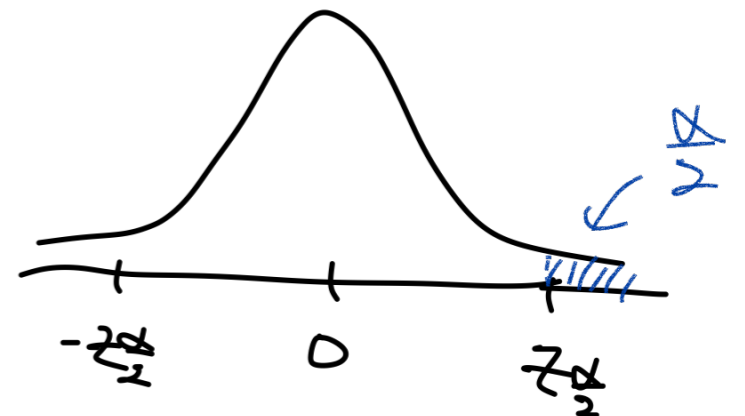
$\bar{X} = \frac{1}{n} \sum_{i=1}^n X_i$  : point estimation of  $\mu$

$\bar{X} \sim N(\mu, \frac{\sigma^2}{n})$  sampling distribution of  $\bar{X}$

$Z = \frac{\bar{X} - \mu}{\sigma/\sqrt{n}} \sim N(0, 1)$  pivotal quantity

$$\Rightarrow P(-z_{\frac{\alpha}{2}} \leq Z \leq z_{\frac{\alpha}{2}}) = 1 - \alpha$$

$$\Rightarrow \dots \Rightarrow P(\mu \in \bar{X} \pm \frac{\sqrt{n} z_{\frac{\alpha}{2}}}{n}) = 1 - \alpha$$



If  $\sigma$  is unknown,

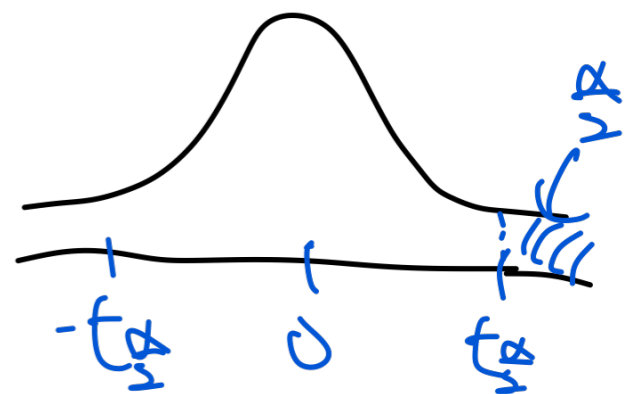
$Z = \frac{\bar{X} - \mu}{\sigma/\sqrt{n}}$  is not an appropriate pivotal quantity  
 $\uparrow$   
estimate  $\sigma$

$\hat{\sigma}^2 = \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{X})^2$  is an estimator of  $\sigma^2$

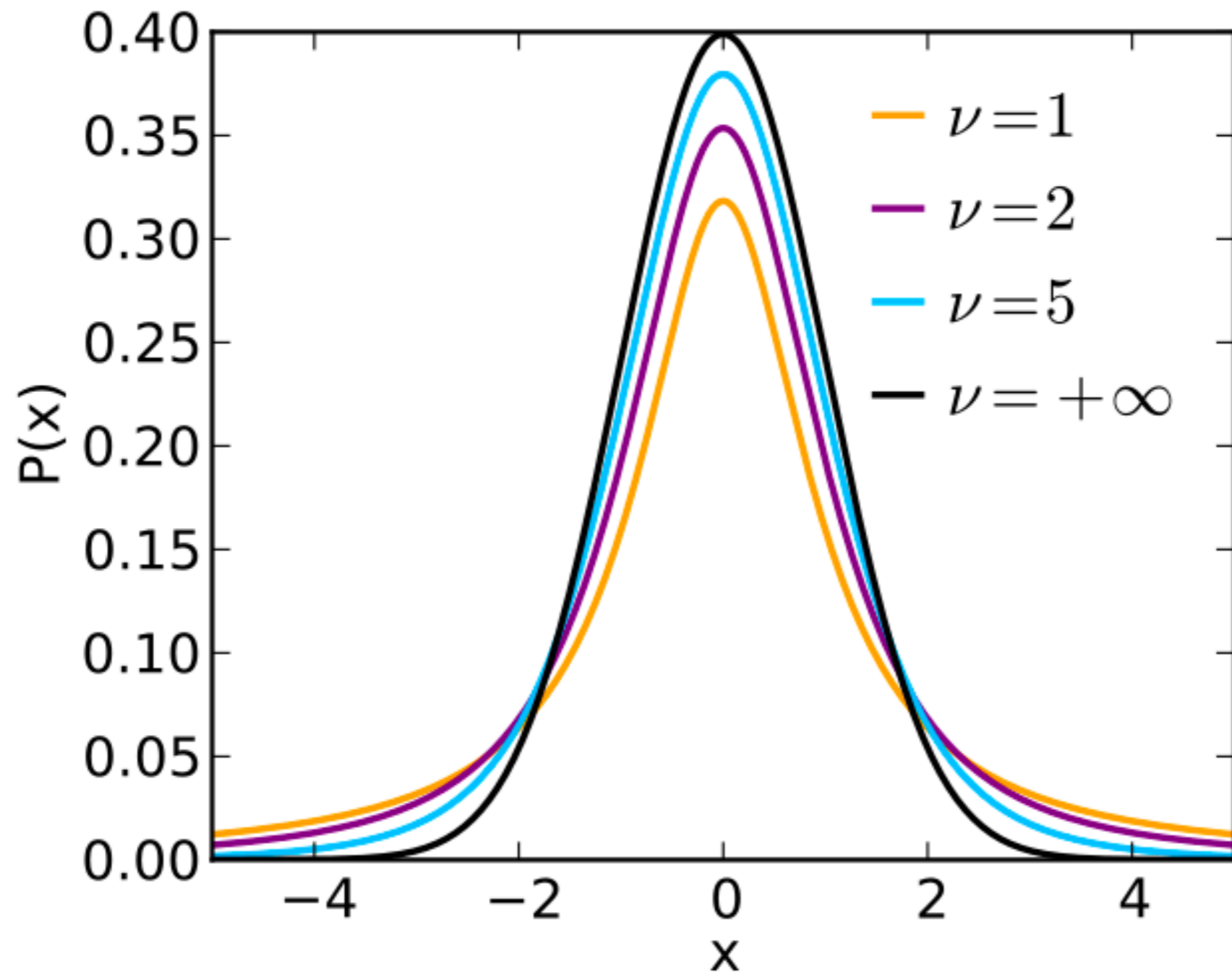
$\Rightarrow T = \frac{\bar{X} - \mu}{\hat{\sigma}/\sqrt{n}} \sim t_{n-1}$  pivotal quantity

$\Rightarrow \dots$

$\Rightarrow P(\mu \in \bar{X} \pm t_{\frac{\alpha}{2}} \cdot \hat{\sigma}/\sqrt{n}) = 1 - \alpha$



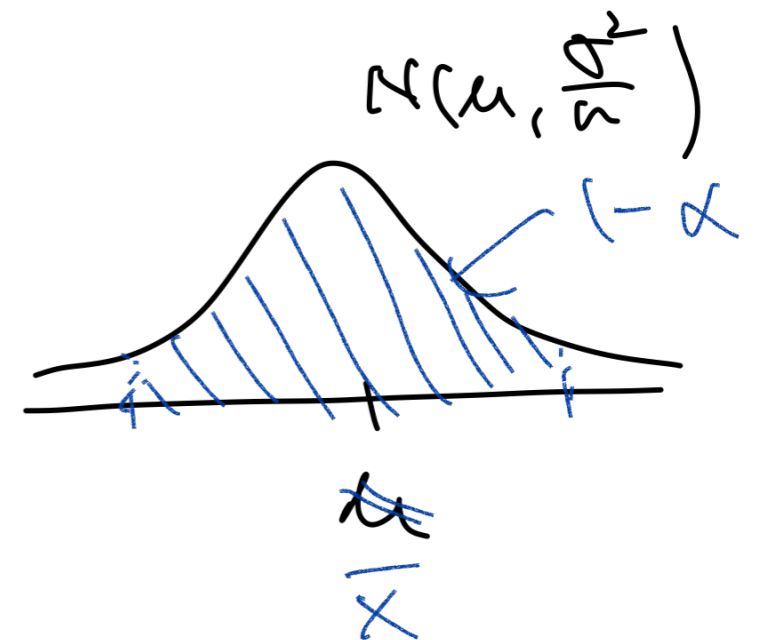
# Student's-t distribution



# Percentile interval

- If  $\hat{\theta} \xrightarrow{d} N(\theta_0, \text{Var}(\hat{\theta}))$ , the interval  $[\hat{\theta}_{\alpha/2}, \hat{\theta}_{1-\alpha/2}]$  is naturally a  $1 - \alpha$  confidence interval for  $\theta$ , where  $\hat{\theta}_{\alpha}$  is the  $\alpha$ th percentile of the sampling distribution of  $\hat{\theta}$

e.g.  $\bar{X} \sim N(\mu, \frac{\sigma^2}{n})$



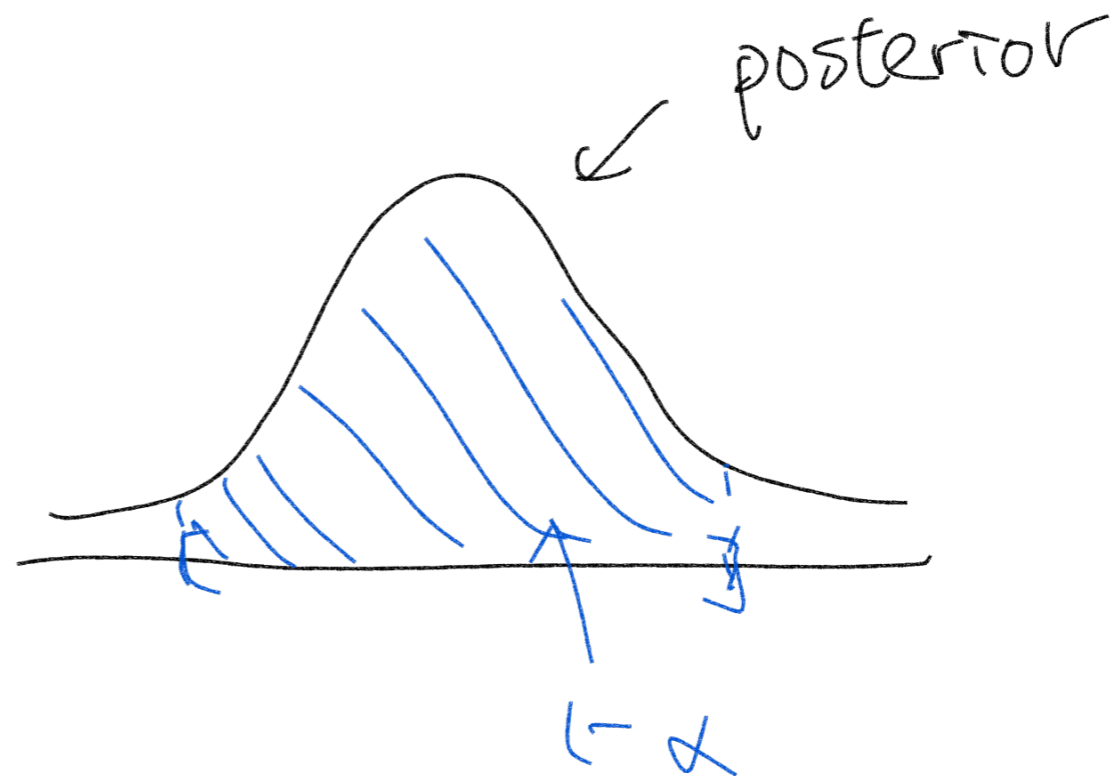


# Bootstrap C.I.

Bayesian credible  
interval

# Bayesian interval estimation

- In Bayesian perspective, the interval estimation of a parameter  $\theta$  can be directly derived from the posterior probability  $P(\theta | X_1, \dots, X_n)$



# Readings

- Chapters 10 and 13 of Computational and Inferential Thinking
- Chapter 2.3–2.5 of Practical Statistics for Data Scientists

# Homework: bootstrap C.I.

- Find a 95% bootstrap confidence interval for the logistic regression coefficients with  $B = 10000$ .
- Dataset: banknote authentication data set
- Compare the bootstrap C.I. with the C.I. given by the StatsModels package (example).
- Deadline: 11/6