



香港中文大學(深圳)

The Chinese University of Hong Kong, Shenzhen

# Introduction to Data Science

## Lecture 12 Statistics

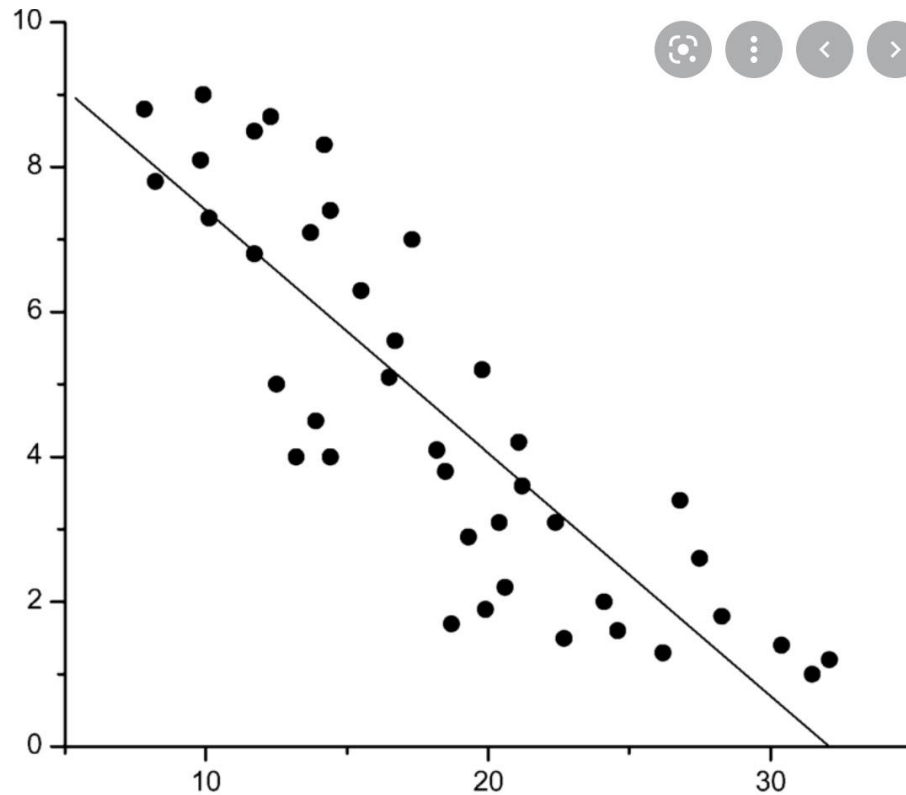
### Advanced Concepts: Confidence Interval

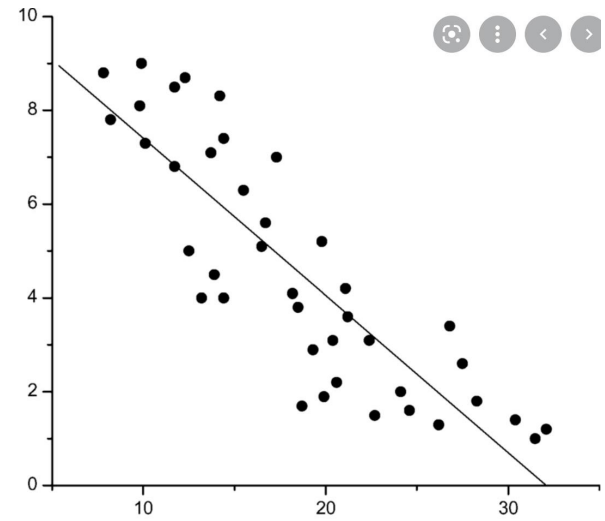
Zicheng Wang

# Recap

## Linear regression

Find the relationship  
between X and Y





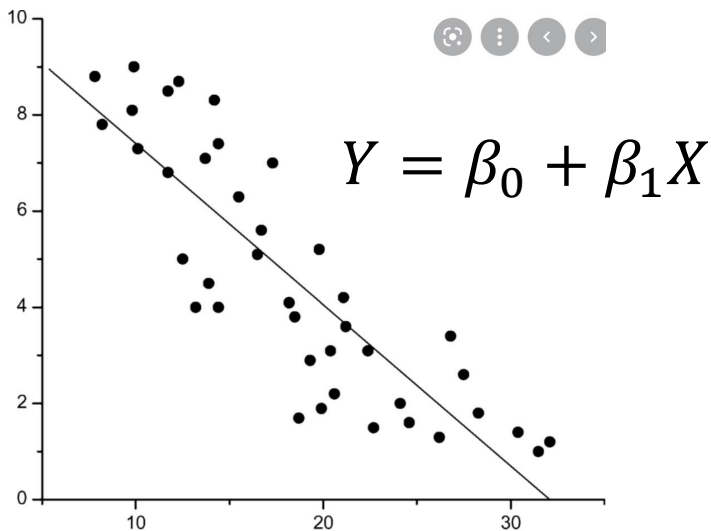
- Negative: larger  $x$  implies smaller  $y$ .
- **Question:** when  $x$  increases by a certain quantity, what's the reduction in  $y$ ?
- Use a line to approximate the relationship:
  - Regression analysis.

- Propose some models

- $Y \sim N(\beta_0 + \beta_1 X, \sigma^2)$

- Given the observation of  $X$
- $Y$  follows a normal distribution with mean  $\beta_0 + \beta_1 X$ , and variance  $\sigma^2$
- To simplify the analysis, we assume  $\sigma^2$  is known

- Regression analysis: knowing  $\beta_0, \beta_1, \sigma$ , you can predict  $X$  given  $Y$

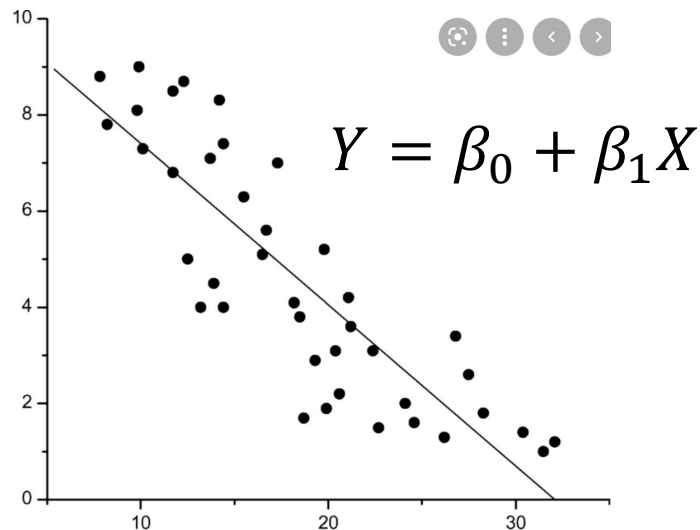


- Propose some models

- $Y \sim N(\beta_0 + \beta_1 X, \sigma^2)$

- Given the observation of  $X$
- $Y$  follows a normal distribution with mean  $\beta_0 + \beta_1 X$

- Regression analysis: knowing  $\beta_0, \beta_1, \sigma^2$ , you can predict  $X$  given  $Y$



MLE: choose the best  $\beta_0, \beta_1$

- PDF for normal

$$\frac{1}{\sqrt{2\pi}\sigma} \exp \left[ -\frac{1}{2} \left( \frac{x - \mu}{\sigma} \right)^2 \right].$$

- Samples:  $(X_1, Y_1), \dots, (X_N, Y_N)$

- For the model with  $\beta_0, \beta_1, \sigma^2$ , the likelihood is

$$Y \sim N(\beta_0 + \beta_1 X, \sigma^2)$$

$$\frac{1}{(\sqrt{2\pi})^n \sigma^n} \exp \left[ -\frac{1}{2} \frac{\sum_i (Y_i - \beta_1 X_i - \beta_0)^2}{\sigma^2} \right]$$

- Given  $\sigma^2$ , to maximize the likelihood, we only need to minimize

$$\sum_i (Y_i - \beta_1 X_i - \beta_0)^2$$

From High School:  
Least square regression  
最小二乘法

- Taking derivative over  $\beta_0$  and  $\beta_1$ , we have

$$\sum_i (Y_i - \beta_1 X_i - \beta_0) = 0$$

$$\sum_i (Y_i - \beta_1 X_i - \beta_0) X_i = 0$$

- PDF for normal

$$\frac{1}{\sqrt{2\pi}\sigma} \exp \left[ -\frac{1}{2} \left( \frac{x - \mu}{\sigma} \right)^2 \right].$$

- Samples:  $(X_1, Y_1), \dots, (X_N, Y_N)$

- For the model with  $\beta_0, \beta_1, \sigma^2$ , the likelihood is

**Step 1**

$$Y \sim N(\beta_0 + \beta_1 X, \sigma^2)$$

$$\frac{1}{(\sqrt{2\pi})^n \sigma^n} \exp \left[ -\frac{1}{2} \frac{\sum_i (Y_i - \beta_1 X_i - \beta_0)^2}{\sigma^2} \right]$$

- Given  $\sigma^2$ , to maximize the likelihood, we only need to minimize

$$\sum_i (Y_i - \beta_1 X_i - \beta_0)^2$$

- Taking derivative over  $\beta_0$  and  $\beta_1$ , we have

$$\sum_i (Y_i - \beta_1 X_i - \beta_0) = 0$$

$$\sum_i (Y_i - \beta_1 X_i - \beta_0) X_i = 0$$



## Step 1

$$f_{X_i}(Y_i) = \frac{1}{\sqrt{2\pi}\sigma} \exp \left[ -\frac{1}{2} \left( \frac{Y_i - (\beta_0 + \beta_1 X_i)}{\sigma} \right)^2 \right]$$



$$\begin{aligned} L(\beta_0, \beta_1, \sigma^2) &= f_{X_1}(Y_1) \times \cdots \times f_{X_N}(Y_N) \\ &= \frac{1}{(\sqrt{2\pi})^n \sigma^n} \exp \left[ -\frac{1}{2} \frac{\sum_i (Y_i - \beta_1 X_i - \beta_0)^2}{\sigma^2} \right] \end{aligned}$$

- PDF for normal

$$\frac{1}{\sqrt{2\pi}\sigma} \exp \left[ -\frac{1}{2} \left( \frac{x - \mu}{\sigma} \right)^2 \right].$$

- Samples:  $(X_1, Y_1), \dots, (X_N, Y_N)$

- For the model with  $\beta_0, \beta_1, \sigma^2$ , the likelihood is

$$Y \sim N(\beta_0 + \beta_1 X, \sigma^2)$$

$$\frac{1}{(\sqrt{2\pi})^n \sigma^n} \exp \left[ -\frac{1}{2} \frac{\sum_i (Y_i - \beta_1 X_i - \beta_0)^2}{\sigma^2} \right]$$

Step 2

- Given  $\sigma^2$ , to maximize the likelihood, we only need to minimize

$$\sum_i (Y_i - \beta_1 X_i - \beta_0)^2$$

- Taking derivative over  $\beta_0$  and  $\beta_1$ , we have

$$\sum_i (Y_i - \beta_1 X_i - \beta_0) = 0$$

$$\sum_i (Y_i - \beta_1 X_i - \beta_0) X_i = 0$$

- PDF for normal

$$\frac{1}{\sqrt{2\pi}\sigma} \exp \left[ -\frac{1}{2} \left( \frac{x - \mu}{\sigma} \right)^2 \right].$$

- Samples:  $(X_1, Y_1), \dots, (X_N, Y_N)$

- For the model with  $\beta_0, \beta_1, \sigma^2$ , the likelihood is

$$Y \sim N(\beta_0 + \beta_1 X, \sigma^2)$$

Constants  $\leftarrow$   $\frac{1}{(\sqrt{2\pi})^n \sigma^n} \exp \left[ -\frac{1}{2} \frac{\sum_i (Y_i - \beta_1 X_i - \beta_0)^2}{\sigma^2} \right]$

- Given  $\sigma^2$ , to maximize the likelihood, we only need to minimize

$$\sum_i (Y_i - \beta_1 X_i - \beta_0)^2$$

- Taking derivative over  $\beta_0$  and  $\beta_1$ , we have

$$\sum_i (Y_i - \beta_1 X_i - \beta_0) = 0$$

$$\sum_i (Y_i - \beta_1 X_i - \beta_0) X_i = 0$$

- PDF for normal

$$\frac{1}{\sqrt{2\pi}\sigma} \exp \left[ -\frac{1}{2} \left( \frac{x - \mu}{\sigma} \right)^2 \right].$$

- Samples:  $(X_1, Y_1), \dots, (X_N, Y_N)$

- For the model with  $\beta_0, \beta_1, \sigma^2$ , the likelihood is

$$Y \sim N(\beta_0 + \beta_1 X, \sigma^2)$$

$$\frac{1}{(\sqrt{2\pi})^n \sigma^n} \exp \left[ \boxed{-\frac{1}{2} \frac{\sum_i (Y_i - \beta_1 X_i - \beta_0)^2}{\sigma^2}} \right]$$

Negative Sign

- Given  $\sigma^2$ , to maximize the likelihood, we only need to minimize

$$\sum_i (Y_i - \beta_1 X_i - \beta_0)^2$$

- Taking derivative over  $\beta_0$  and  $\beta_1$ , we have

$$\sum_i (Y_i - \beta_1 X_i - \beta_0) = 0$$

$$\sum_i (Y_i - \beta_1 X_i - \beta_0) X_i = 0$$

- PDF for normal

$$\frac{1}{\sqrt{2\pi}\sigma} \exp \left[ -\frac{1}{2} \left( \frac{x - \mu}{\sigma} \right)^2 \right].$$

- Samples:  $(X_1, Y_1), \dots, (X_N, Y_N)$

- For the model with  $\beta_0, \beta_1, \sigma^2$ , the likelihood is

$$Y \sim N(\beta_0 + \beta_1 X, \sigma^2)$$

$$\frac{1}{(\sqrt{2\pi})^n \sigma^n} \exp \left[ -\frac{1}{2} \frac{\sum_i (Y_i - \beta_1 X_i - \beta_0)^2}{\sigma^2} \right]$$

- Given  $\sigma^2$ , to maximize the likelihood, we only need to minimize

$$\sum_i (Y_i - \beta_1 X_i - \beta_0)^2$$

- Taking derivative over  $\beta_0$  and  $\beta_1$ , we have

**Step 3**

$$\sum_i (Y_i - \beta_1 X_i - \beta_0) = 0$$

$$\sum_i (Y_i - \beta_1 X_i - \beta_0) X_i = 0$$

- **First order condition**
- **Take derivative with respect to  $\beta_0$  and  $\beta_1$  separately**
- **Set the derivative to be equal to zero**

$$\sum_i (Y_i - \beta_1 X_i - \beta_0) X_i = 0 \text{ AND } \sum_i (Y_i - \beta_1 X_i - \beta_0) = 0$$

Eliminate  $\beta_0$  first:

$$\sum_i (Y_i - \beta_1 X_i - \beta_0) = 0 \rightarrow \beta_0 = \frac{1}{N} \sum_i (Y_i - \beta_1 X_i) = \bar{Y} - \beta_1 \bar{X}$$

$$\text{MLE: } \left\{ \begin{array}{l} \widehat{\beta}_1 = \frac{\sum_i (X_i - \bar{X}) (Y_i - \bar{Y})}{\sum_i (X_i - \bar{X})^2} \\ \widehat{\beta}_0 = \bar{Y} - \widehat{\beta}_1 \bar{X} \end{array} \right.$$

# When simple regression is invalid?

- The model we propose is not correct.

$$Y \sim N(\beta_0 + \beta_1 X, \sigma^2) \text{ or } Y - \beta_0 - \beta_1 X \sim N(0, \sigma^2)$$

Linear regression assumes that...

1. The relationship between X and Y is **linear**
2. The variance of  $Y - \beta_0 - \beta_1 X$  at every value of X is the **same** (homogeneity of variances)

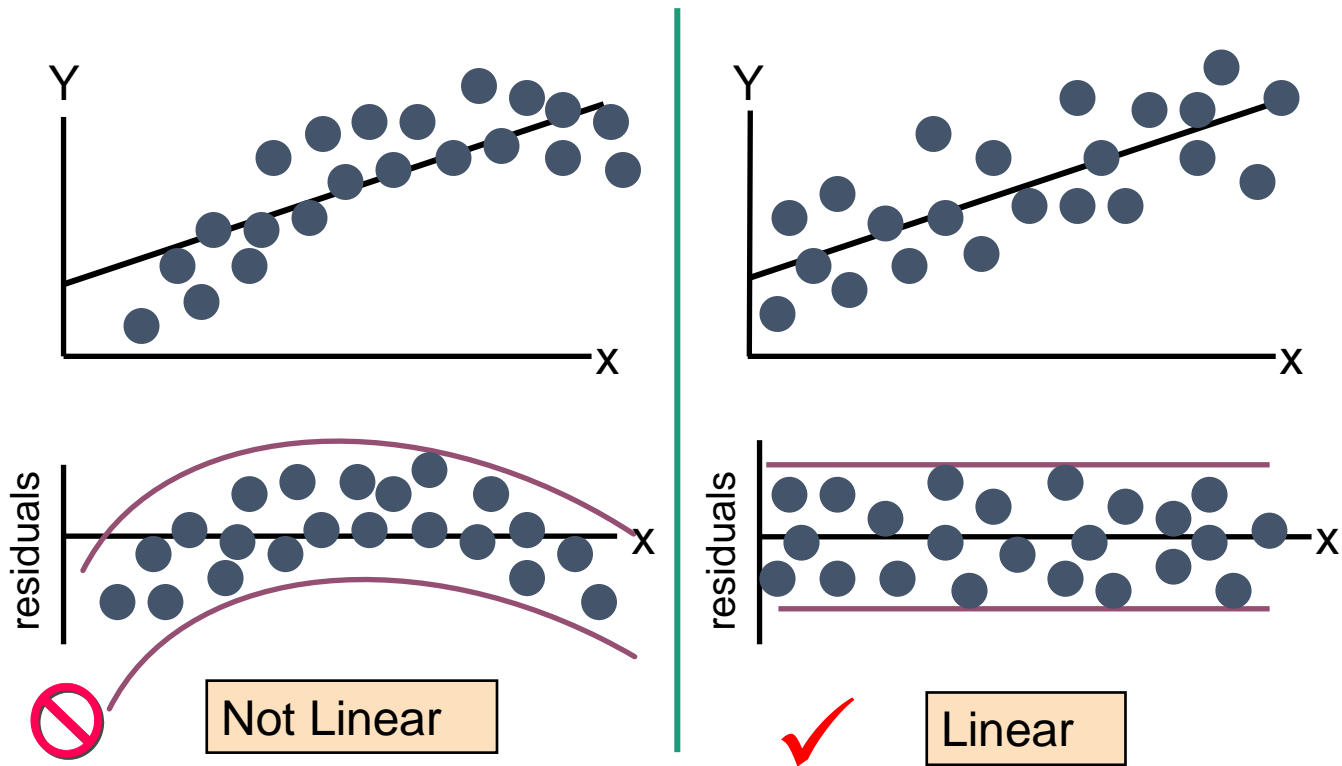
# Residual Analysis: check assumptions

$$\text{Residual: } e_i := Y_i - \widehat{\beta}_0 - \widehat{\beta}_1 X_i$$

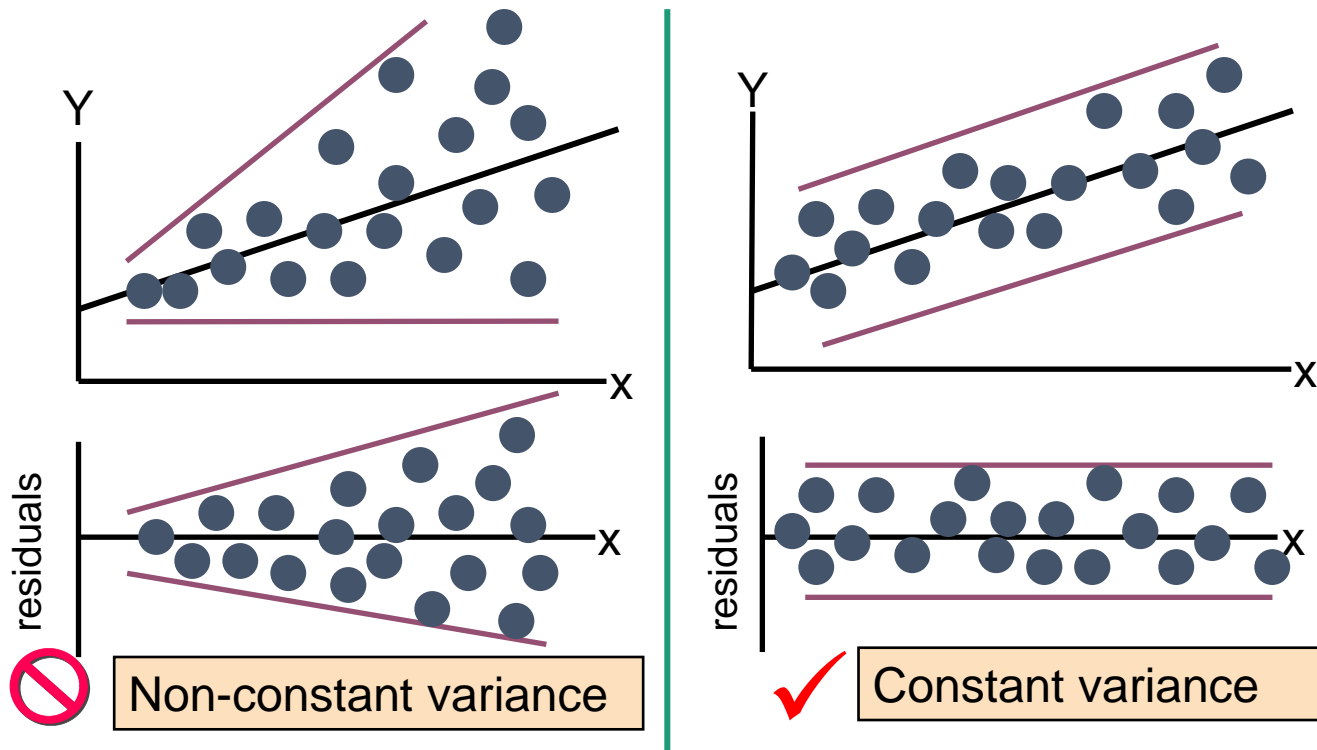
- Check the assumptions by examining the residuals
  - Examine for linearity assumption:
    - *$e_i$  does not depend on  $X_i$*
  - Evaluate constant-variance assumption:
    - *variance of  $e_i$  does not depend on  $X_i$*
- Graphical Analysis of Residuals: Can plot residuals vs. X



# Residual Analysis for Linearity



# Residual Analysis for constant-variance



# Advanced Concepts: Confidence Interval

## Reading Materials

- Applied Statistics and Probability for Engineers, Third Edition, Douglas C. Montgomery and George C. Runger.
- Chap 8-2.1, ..., 8-2.5

# Experiments

Whether a drug can cure a disease:  $\hat{p} = \frac{\sum_i X_i}{n}$  (MLE)

- Drug 1:  $\hat{p}_1 = 90\%$ .
- Drug 2:  $\hat{p}_2 = 80\%$ .

Which drug do you think is more effective?

# Experiments

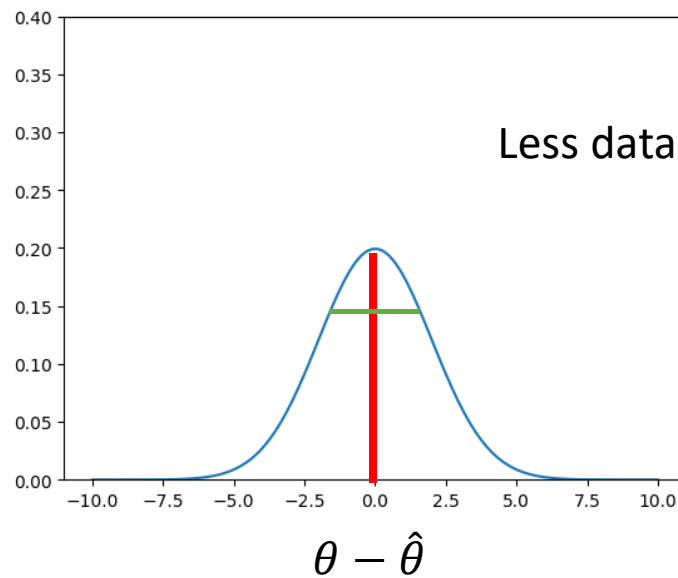
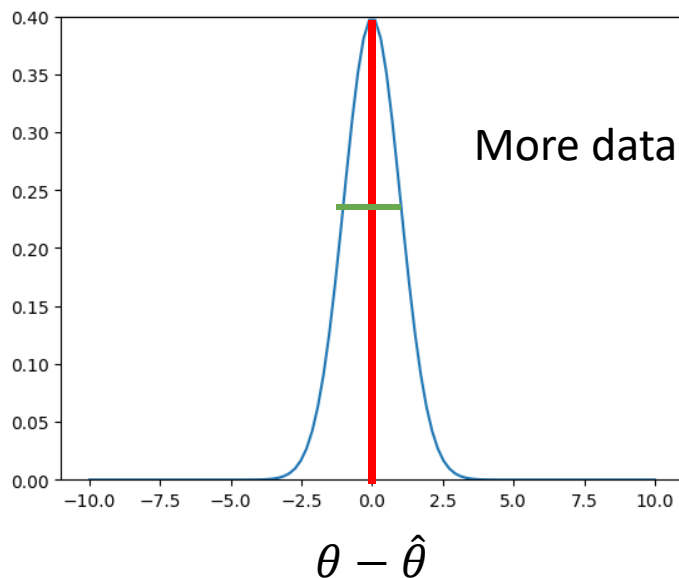
Whether a drug can cure a disease:  $\hat{p} = \frac{\sum_i X_i}{n}$

- Drug 1:  $\hat{p}_1 = 90\%$ . 10 experiments.
- Drug 2:  $\hat{p}_2 = 80\%$ . 10000 experiments.

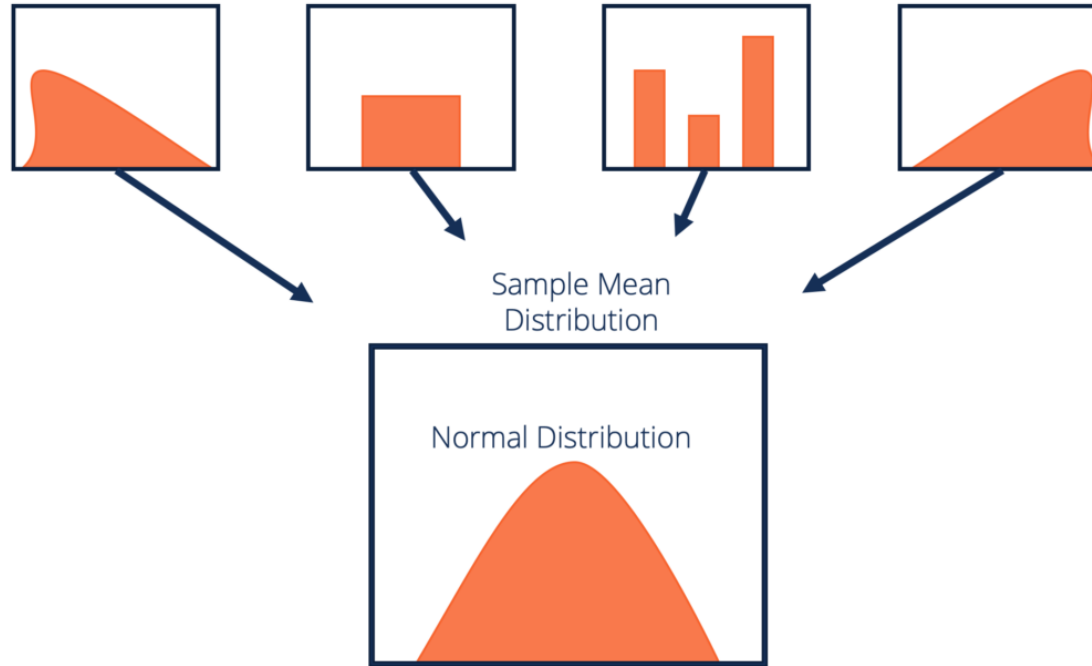
Which drug do you think is more effective?

Which estimation is more reliable?

- Number of samples can affect the accuracy!!!
- With **more** data, we **believe** the estimator is **closer** to the true parameter.



# Central limit theorem



No matter what the true distribution is, the **sample mean** will be very close to the **normal distribution**, as long as the sample size is **large**.

# Central limit theorem

$X_1, \dots, X_n$  can be non-normal    **Mean:  $\mu$ ; Variance:  $\sigma^2$**

**Sample mean**     $\bar{X} = \frac{X_1 + X_2 + \dots + X_n}{n}$

$$\bar{X} \sim N\left(\mu, \frac{\sigma^2}{n}\right)$$

Or write as:

$$\frac{\sqrt{n}(\bar{X} - \mu)}{\sigma} \sim N(0,1)$$

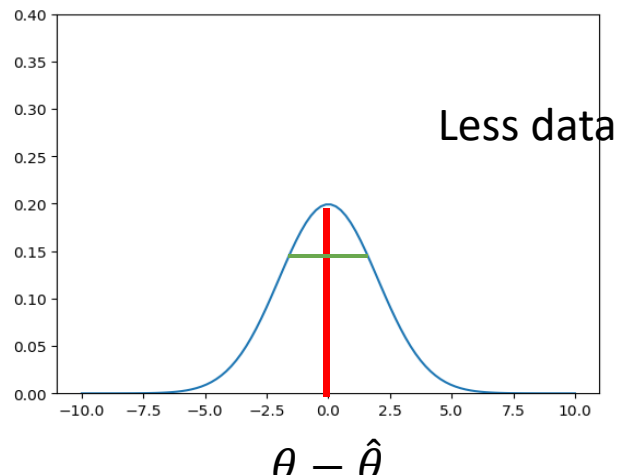
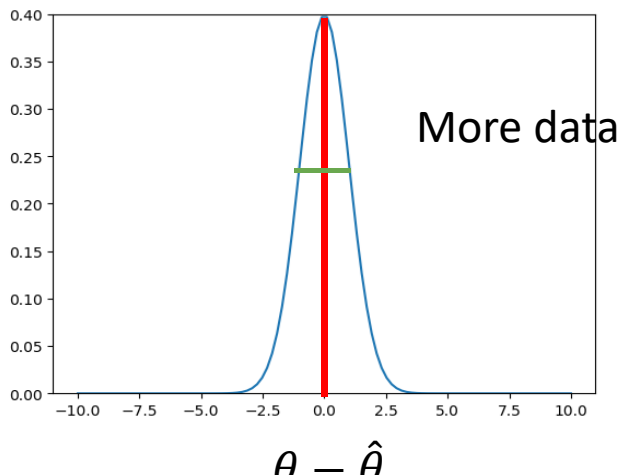
**Standard Normal**



# Target

We will use normal distribution to show:

- with different size of data, **how close** the estimator is to the true parameter.
- With what probability, the true parameter falls in a region.



# Interval Estimation – example 1

- We have data  $X_1, X_2, \dots, X_n$  that are sampled from some distribution with a **known** variance  $\sigma^2$
- Their mean is  $\mu$ , which we want to estimate
- We can easily give a point estimate:  $\bar{X}$  (sample mean)
- How to get an interval estimate??
  - Use Central Limit Theorem!

# Interval Estimation – example 1

- $\frac{\sqrt{n}(\bar{X}-\mu)}{\sigma} \sim \mathcal{N}(0,1)$
- $P(a \leq \frac{\sqrt{n}(\bar{X}-\mu)}{\sigma} \leq b) = \Phi(b) - \Phi(a)$ 
  - $\Phi(x)$ : CDF of a standard normal distribution  $\mathcal{N}(0,1)$ .
- $P(\bar{X} - \frac{b\sigma}{\sqrt{n}} \leq \mu \leq \bar{X} - \frac{a\sigma}{\sqrt{n}}) = \Phi(b) - \Phi(a)$
- W.P.  $\Phi(b) - \Phi(a)$ ,  $\mu$  is within  $[\bar{X} - \frac{b\sigma}{\sqrt{n}}, \bar{X} - \frac{a\sigma}{\sqrt{n}}]$

# The best interval?

- W.P.  $\Phi(b) - \Phi(a)$ ,  $\mu$  is within  $[\bar{X} - \frac{b\sigma}{\sqrt{n}}, \bar{X} - \frac{a\sigma}{\sqrt{n}}]$
- Fix  $\Phi(b) - \Phi(a)$ , there are too many **a** and **b** to choose from.
- At least, we want  $\bar{X}$  to be within the interval
  - $a < 0$
  - $b > 0$

# The best interval?

- W.P.  $\Phi(b) - \Phi(a)$ ,  $\mu$  is within  $[\bar{X} - \frac{b\sigma}{\sqrt{n}}, \bar{X} - \frac{a\sigma}{\sqrt{n}}]$
- Fix  $\Phi(b) - \Phi(a)$ , there are too many **a** and **b** to choose from.
- $\mu$  has an **upper bound**, say U.
  - If  $\bar{X}$  is too close to U, choose **a** such that  $\bar{X} - \frac{a\sigma}{\sqrt{n}} = U$

# The best interval?

- W.P.  $\Phi(b) - \Phi(a)$ ,  $\mu$  is within  $[\bar{X} - \frac{b\sigma}{\sqrt{n}}, \bar{X} - \frac{a\sigma}{\sqrt{n}}]$
- Fix  $\Phi(b) - \Phi(a)$ , there are too many **a** and **b** to choose from.
- $\mu$  has a **lower bound**, say L.
  - If  $\bar{X}$  is too close to L, choose **b** such that  $\bar{X} - \frac{b\sigma}{\sqrt{n}} = L$ .

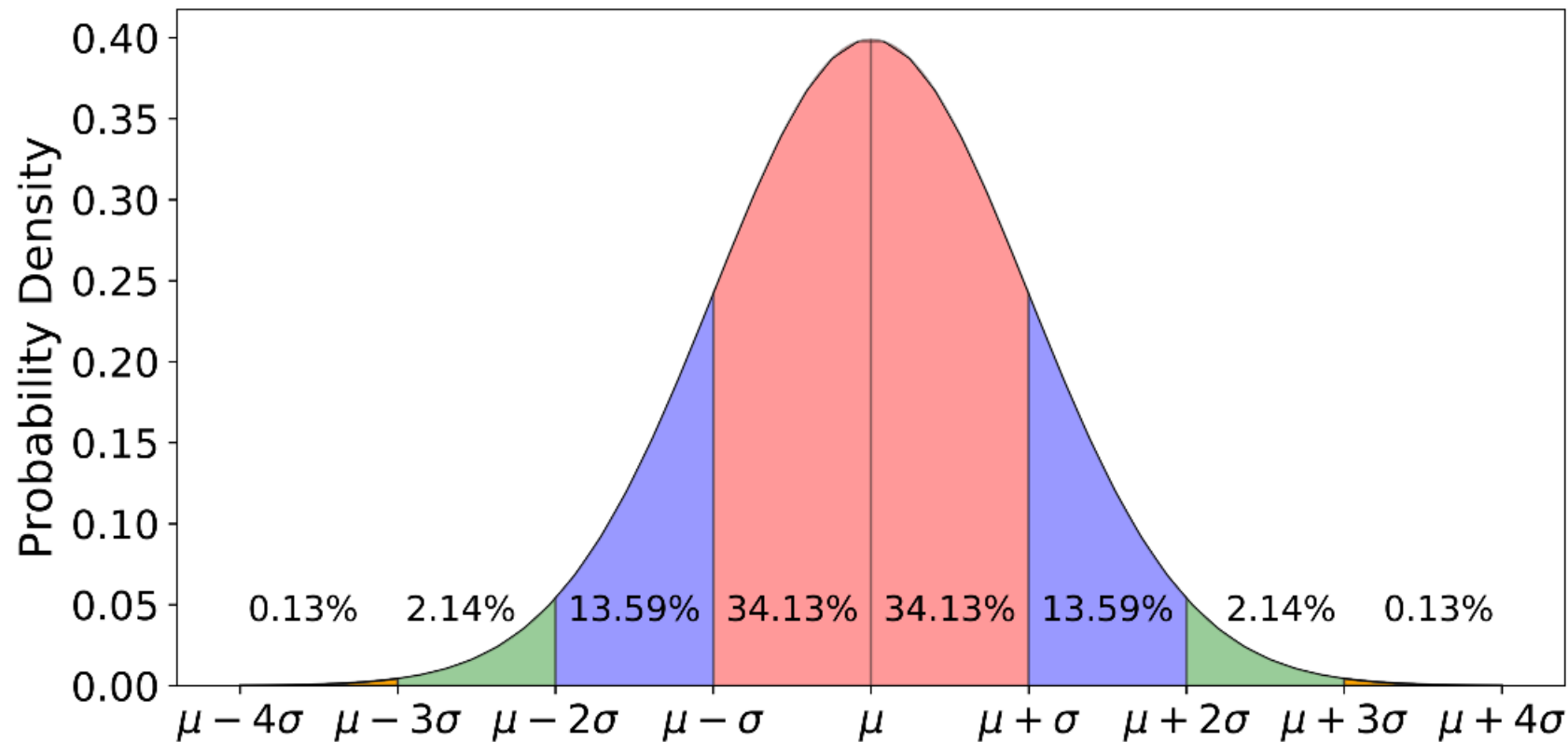
# The best interval?

- W.P.  $\Phi(b) - \Phi(a)$ ,  $\mu$  is within  $[\bar{X} - \frac{b\sigma}{\sqrt{n}}, \bar{X} - \frac{a\sigma}{\sqrt{n}}]$
- Fix  $\Phi(b) - \Phi(a)$ , there are too many **a** and **b** to choose from.
- $\mu$  has no **bound**.
  - Choose a and b such that b-a is minimized.

- W.P.  $\Phi(b) - \Phi(a)$ ,  $\mu$  is within  $[\bar{X} - \frac{b\sigma}{\sqrt{n}}, \bar{X} - \frac{a\sigma}{\sqrt{n}}]$
- For **the  $1 - \alpha$  confidence interval**, we **need**
  - $\Phi(b) - \Phi(a) = 1 - \alpha$
  - $b - a$  is minimized.
- As pdf of normal is symmetric and has a single peak, for **the  $1 - \alpha$  confidence interval**,
  - $b + a = 0$ .

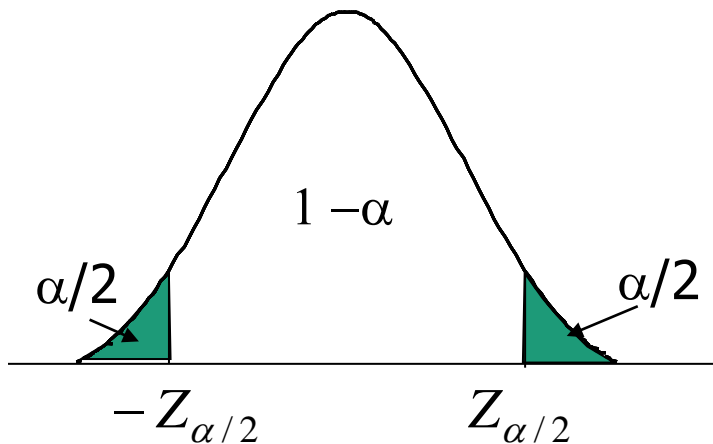


# Normal Distribution



# Notation

$N(0, 1)$ : **Standard Normal Distribution**



let  $z_{\alpha/2}$  be the number such that the area under the **standard normal density function** to the right of  $z_{\alpha/2}$  is  $\alpha/2$ .

Then if  $Z \sim N(0, 1)$

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

# Interval Estimation – example 1

- W.P.  $\Phi(b) - \Phi(a)$ ,  $\mu$  is within  $[\bar{X} - \frac{b\sigma}{\sqrt{n}}, \bar{X} + \frac{a\sigma}{\sqrt{n}}]$
- $P(-z_{\alpha/2} \leq Z \leq z_{\alpha/2}) = 1 - \alpha$
- Let  $b = z_{\alpha/2}$ ,  $a = -z_{\alpha/2}$ . Then  $\Phi(b) - \Phi(a) = 1 - \alpha$ .
- The  $1 - \alpha$  confidence interval for  $\mu$  :  $[\bar{X} - z_{\alpha/2} \frac{\sigma}{\sqrt{n}}, \bar{X} + z_{\alpha/2} \frac{\sigma}{\sqrt{n}}]$

# Some observations

**$1-\alpha$  Confidence Interval (CI):**

$$[\bar{X} - z_{\alpha/2} \frac{\sigma}{\sqrt{n}}, \bar{X} + z_{\alpha/2} \frac{\sigma}{\sqrt{n}}]$$

We typically call  $1 - \alpha$  as the **confidence level**.

The length of the confidence interval is affected by several factors

- As the sample size  **$n$  increases**, the length of CI **decreases**
- As the variance  **$\sigma^2$  increases**, the length of CI **increases**
- As the confidence level increases ( **$\alpha$  decreases**), the length of CI **increases**.

# Interval Estimation - example

- $n$  patients use the new drug, whether the drug can cure the disease is a Bernoulli RV
- We have data  $X_1, X_2, \dots, X_n$  that are sampled from this Bernoulli distribution with unknown cure rate  $p$  to be estimated
- Clearly, the mean of  $\text{Bernoulli}(p)$  is  $p$
- We can easily give a point estimate:  $\hat{p} = \bar{X}$  (sample cure rate)
- How to get an interval estimate??
  - Use Central Limit Theorem!

# Interval Estimation - example

$$Z = \frac{\sqrt{n}(\hat{p} - p)}{\sqrt{p(1-p)}} \sim N(0,1)$$

Mean:  $p$

Variance:  $p(1-p)$

$\hat{p} = \bar{X}$  (sample cure rate)

- $1 - \alpha$  confidence interval:

$$[\hat{p} - z_{\alpha/2} \sqrt{\frac{p(1-p)}{n}}, \hat{p} + z_{\alpha/2} \sqrt{\frac{p(1-p)}{n}}]$$

- As  $p$  is unknown, replace  $p$  above by  $\hat{p}$ ,  $1 - \alpha$  confidence interval:

$$[\hat{p} - z_{\alpha/2} \sqrt{\frac{\hat{p}(1-\hat{p})}{n}}, \hat{p} + z_{\alpha/2} \sqrt{\frac{\hat{p}(1-\hat{p})}{n}}]$$

# Experiments

$$\hat{p} \pm z_{\alpha/2} \sqrt{\frac{\hat{p}(1 - \hat{p})}{n}}$$

Whether a drug can cure a disease:  $\hat{p} = \frac{\sum_i X_i}{n}$

95% Confidence Interval

- Drug 1:  $\hat{p}_1 = 90\%$ . 10 experiments.

[71.41%, 100%]

- Drug 2:  $\hat{p}_2 = 80\%$ . 10000 experiments.

[79.22%, 80.78%]

# Confidence Statements

- Fortune Teller



**“I believe the cure  
rate is 80%”**

point

- Scientist



**“I believe the cure  
rate is 80% plus or  
minus 5%”**

**interval**