

Descriptive Statistics

Nguyen The Hien, DAH - Vietnam

CONTENT

I. Statistical Method

- 1. Descriptive statistics*
- 2. Inferential statistics*

II. Probability.

- 1. Definition of probability*
- 2. Normal distribution*
- 3. Binomial distribution*

I. Statistical Method

1. Descriptive statistics

- A method of collecting, organizing, displaying and drawing conclusion of data

2. Inferential statistics

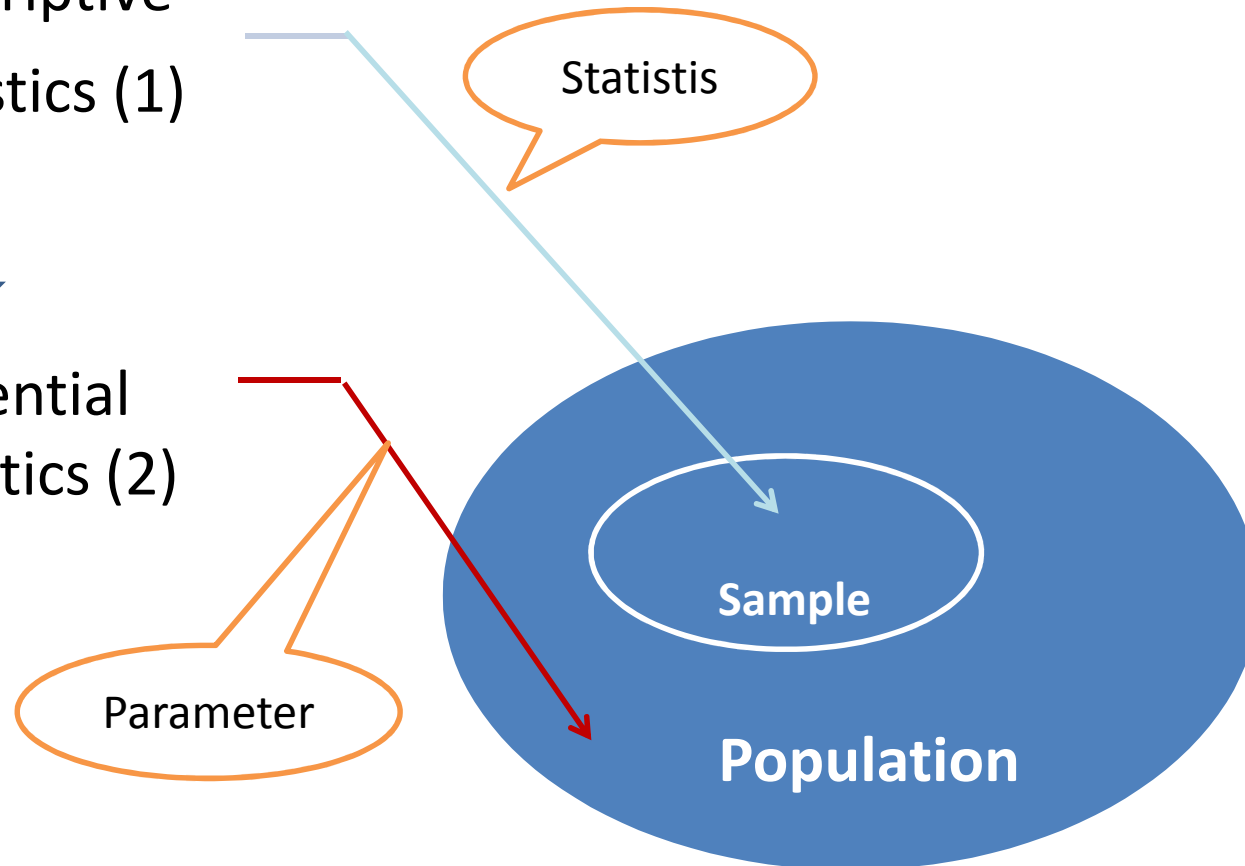
- Use the result of research from sample statistics to population parameter
- Type of inferential statistics methods:
 - + Parameter estimation
 - + Hypothesis testing

(1) Give a clearer picture of what we want to describe in a group (sample)

Descriptive
Statistics (1)



Inferential
Statistics (2)



(2) Give answer in the world of uncertainty (population)

Summary: Descriptive Statistics

Type of variable	Measurement scale	Statistics	Presentation
Numerical discrete	<ul style="list-style-type: none"> - Interval - Ratio 	Mean, or Median, or Mode (rare), Range, SD, IQR	Histogram scatter plot line graph
Numerical continuous	<ul style="list-style-type: none"> - Interval - Ratio 	Mean, or Median, or Mode (rare), Range, SD, IQR	Histogram scatter plot line graph
Categorical dichotomous	<ul style="list-style-type: none"> - Nominal - Ordinal 	Proportion (%)	Table, Pie, Bar chart
Categorical polychotomous	<ul style="list-style-type: none"> - Nominal - Ordinal 	Proportion (%)	Table, Pie, Bar chart

Notation

	Statistics	Parameter
Mean	\bar{X}	μ
Variance	S^2	σ^2
Standard deviation	S	σ
Standard error	$S_{\bar{X}}$	$\sigma_{\bar{X}}$
Proportion	P	π
Correlation coefficient	r	ρ

II. Probability

1. Definition and characteristic

Definition: The chance, or likelihood that an event occurs.

Characteristics:

- Value of probabilities to happen an event will be around 0 – 1
- The total of Pr of all events are 1.0

Ex.

– $Pr(x = n) = X_1$ -----> probability of an value/point

– $Pr(x < n_a) = \text{Sum of } Pr(x = n_i) \text{ with } i = 0, 1, \dots, n_a$

– $Pr(n_a < x < n_b) = \text{Sum of } Pr(x = n_i) \text{ with } i = a, a+1, \dots, b$

} -----> *pace/interval*

2. Probability distribution

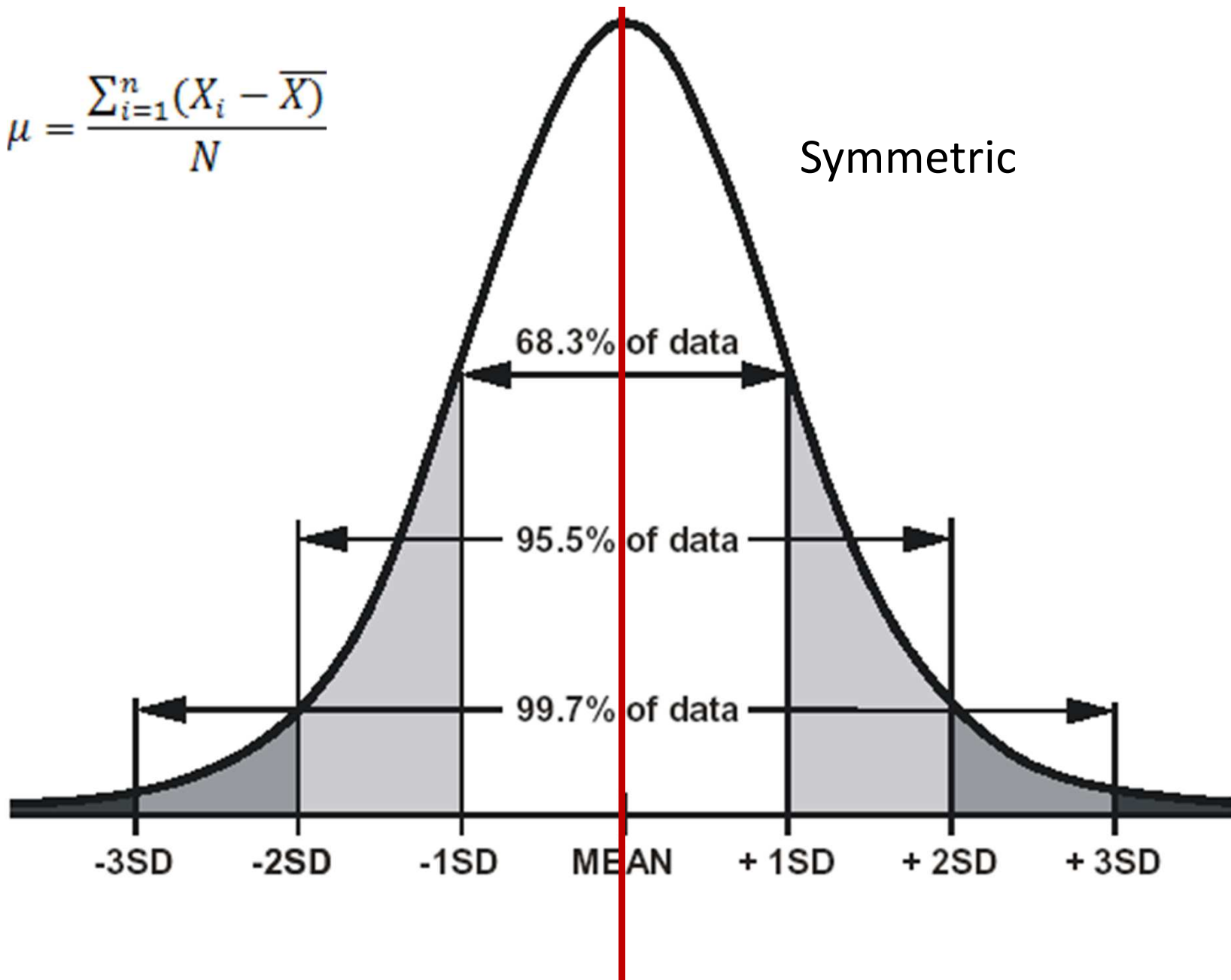
- Continuous probability distribution
 - Normal distribution
- Discrete probability distribution
 - Binomial distribution
 - Poisson distribution

2.1. Normal distribution

- Symmetric about the mean: μ
- Mean = median = mode = μ , $SD = \sigma$
- Area under the curve = 1 implies 50% of the area to the left of the mean (or median)
- Notation: $X \sim N(\mu, \sigma^2)$
- *About 95% of the area is between $\mu - 1.96\sigma$ and $\mu + 1.96\sigma$*

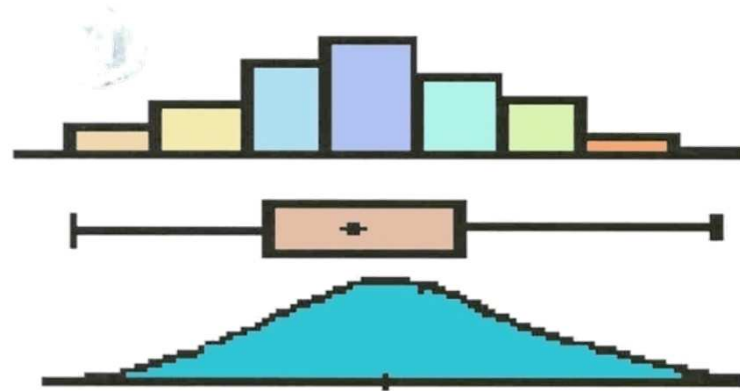
Areas under the normal curve that lie between 1, 2, and 3 standard deviations on each side of the mean

$$\mu = \frac{\sum_{i=1}^n (X_i - \bar{X})}{N}$$



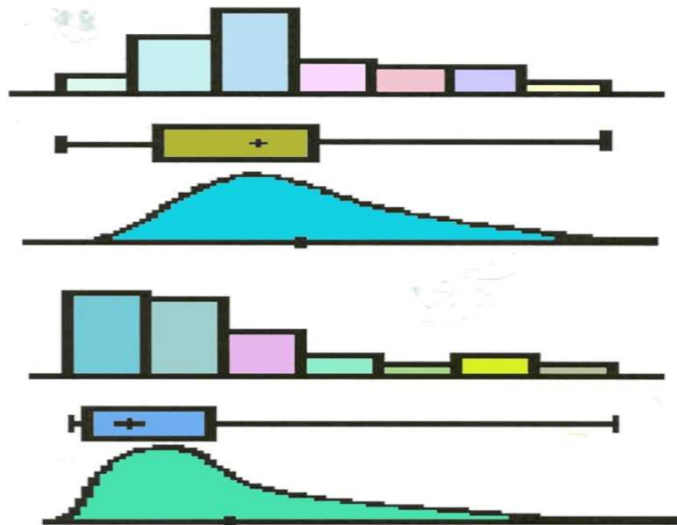
Shape of data distribution

- Normal Distribution

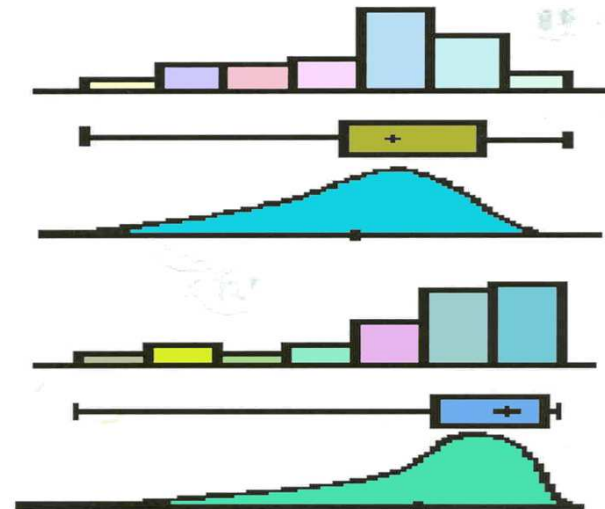


- Skewed Distribution

Positively skewed



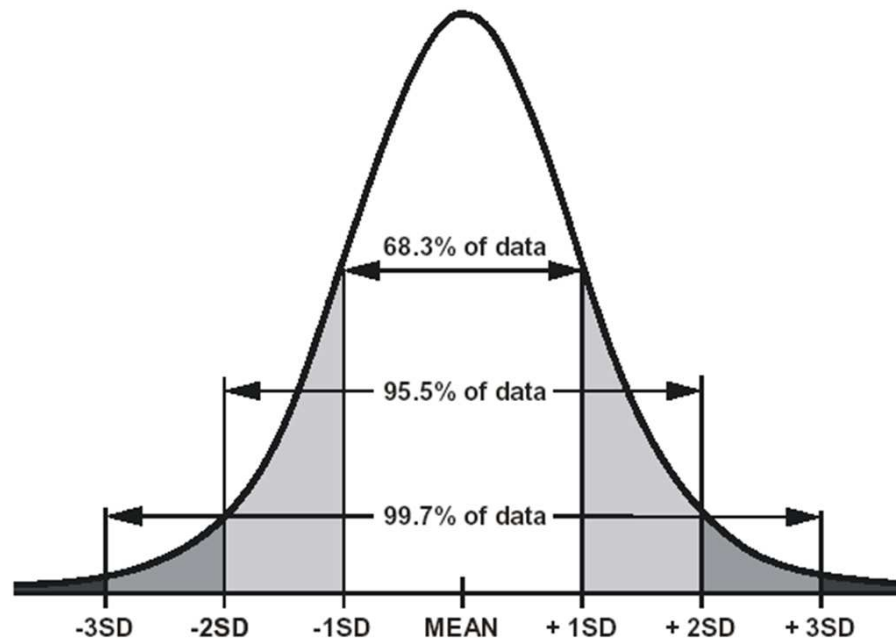
Negatively skewed



Standard normal distribution

- Normal distribution
- Notation: $X \sim N(0,1)$
- Symmetric about mean = 0
- Standard deviation = 1

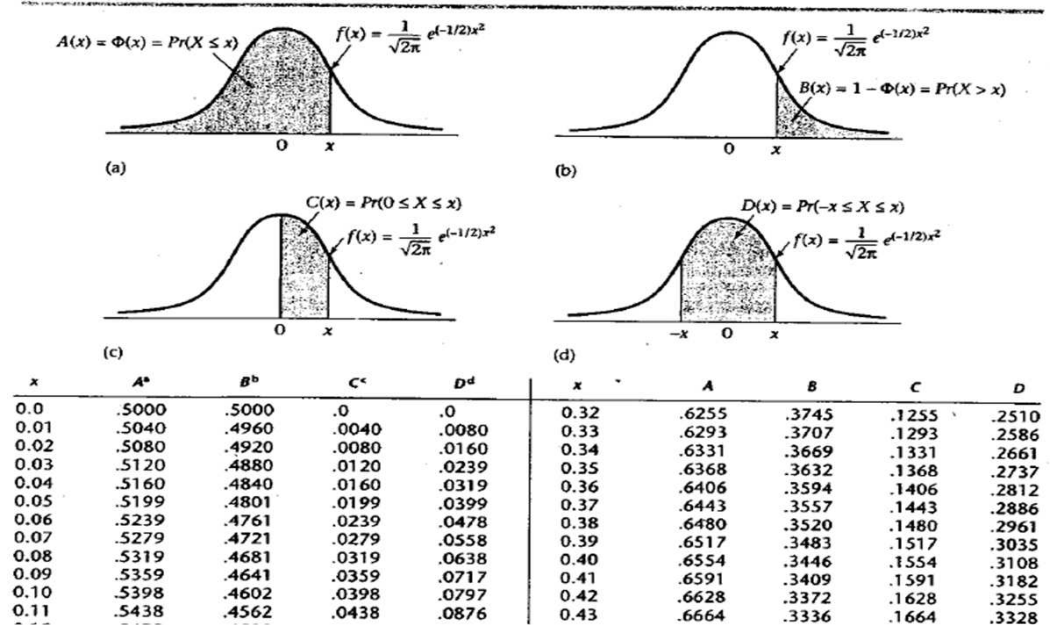
Areas under the normal curve that lie between 1, 2, and 3 standard deviations on each side of the mean



Standard normal distribution table

- See additional handout
- Column A : Area from $-\infty$ (infinite) to x . This is the $\Pr(X \leq x)$
- Column B : Area from x to $+\infty$. This is the $\Pr(X \geq x)$
- Column C : Area from 0 to a point $\Pr(0 \leq X \leq x)$
- Column D : Area symmetric about 0

TABLE 3 The normal distribution



Converting from a $N(\mu, \sigma^2)$ Distribution to a $N(0, 1)$ Distribution

- The formula is:

$$Z = \frac{X - \mu}{\sigma}$$

- $Z \sim N(0, 1)$

Therefore, $\Pr(a < X < b) =$

$$\Pr\left(\frac{a - \mu}{\sigma} < Z < \frac{b - \mu}{\sigma}\right)$$

Example

- $\mu = 200$, $\sigma^2 = 400$ implies $\sigma = 20$, $a = 220$, $b = 260$

$$Z_1 = \frac{a - \mu}{\sigma} = \frac{220 - 200}{20} = 1.00$$

$$Z_2 = \frac{b - \mu}{\sigma} = \frac{260 - 200}{20} = 3.00$$

- We have converted the probability statement about $N(200, 400)$ to one about $N(0, 1)$. We can evaluate by using the standard normal tables.

$$\begin{aligned} \Pr(220 < X < 260) &= \Pr\left(\frac{220 - 200}{20} < Z < \frac{260 - 200}{20}\right) \\ &= \Pr(1.00 < Z < 3.00) \end{aligned}$$

2.2. Binomial Distribution

- n independent trials, each with only two possible outcomes.
 - Example: Yes/No, Sick/not sick
- Probability of success at each trial is some constant, p
- *Probability of failure = $1 - p = q$*

2.2. Binomial Distribution

- Probability of k successes in n trials is

$$\Pr(X = k) = \frac{n!}{(n-k)!k!} p^k q^{n-k}$$

– $k = 0, 1, 2, \dots, n$

– $q = 1 - p$

$(Pr = p)$

– $P(x \leq n_a) = \text{Sum of } Pr(x = n_i) \text{ with } i = 0, 1, \dots, n_a$

– $P(n_a \leq x \leq n_b) = \text{Sum of } Pr(x = n_i) \text{ with } i = a, a+1, \dots, b$

Normal approximation to the binomial

- When **n is large**, there are no tables and using the formula is difficult. The Normal Distribution can be used as an approximation to the **Binomial Distribution**
- The Normal Distribution is a good approximation **if n is large and p not too near 0 or 1**. The binomial then will be nearly symmetric when **$n * p * q \geq 5$**
- Binomial: mean = np , variance = npq . Let's create a normal distribution that is $N(np, npq)$
- **If n is large enough, $npq \geq 5 \rightarrow X \sim B(np, npq)$, $\mu = n * p$; $\sigma^2 = n * p * q$**

Normal approximation to the binomial

- We want $\Pr(a \leq X \leq b)$ where X is binomially distributed with parameters n and p . The first choice for an approximation is to examine the $N(np, npq)$ normal distribution by getting the area under this curve from a to b
- Continuity Correction: Better to get area under the normal curve
 - For X is between a and b : $\Pr(a - \frac{1}{2} \leq X \leq b + \frac{1}{2})$
 - For $X > a$: $\Pr(X \geq a - \frac{1}{2})$
 - For $X < a$: $\Pr(X \leq b + \frac{1}{2})$

Example: operative complication

- $\Pr(\text{at most 5 complications}) = \Pr(X \leq 5)$
 - $p = 0.20, n = 50$
 - Normal approximation:
 - mean = $np = 10$
 - variance = $npq = 8 > 5$
 - SD = Square root or SQRT (8) = 2.8284
- $$\begin{aligned}\Pr(X \leq 5.5) &= \Pr\left(Z \leq \frac{5.5 - 10}{2.8284}\right) \\ &= \Pr(Z \leq -1.59) \\ &= 1 - \Pr(Z \leq 1.59) \\ &= 1 - 0.9441 \\ &= 0.0559\end{aligned}$$

Thank you for your attention!