

# "Introduction to statistics"

①

Population: all people or things that you are studying.

Parameter: a numerical description measuring the variable of interest in the POPULATION.

Sample: A representative subset of the population.

Statistic: a numerical description measuring the variable of interest in the SAMPLE.

Variable: A value or characteristic that changes for each member of the population.

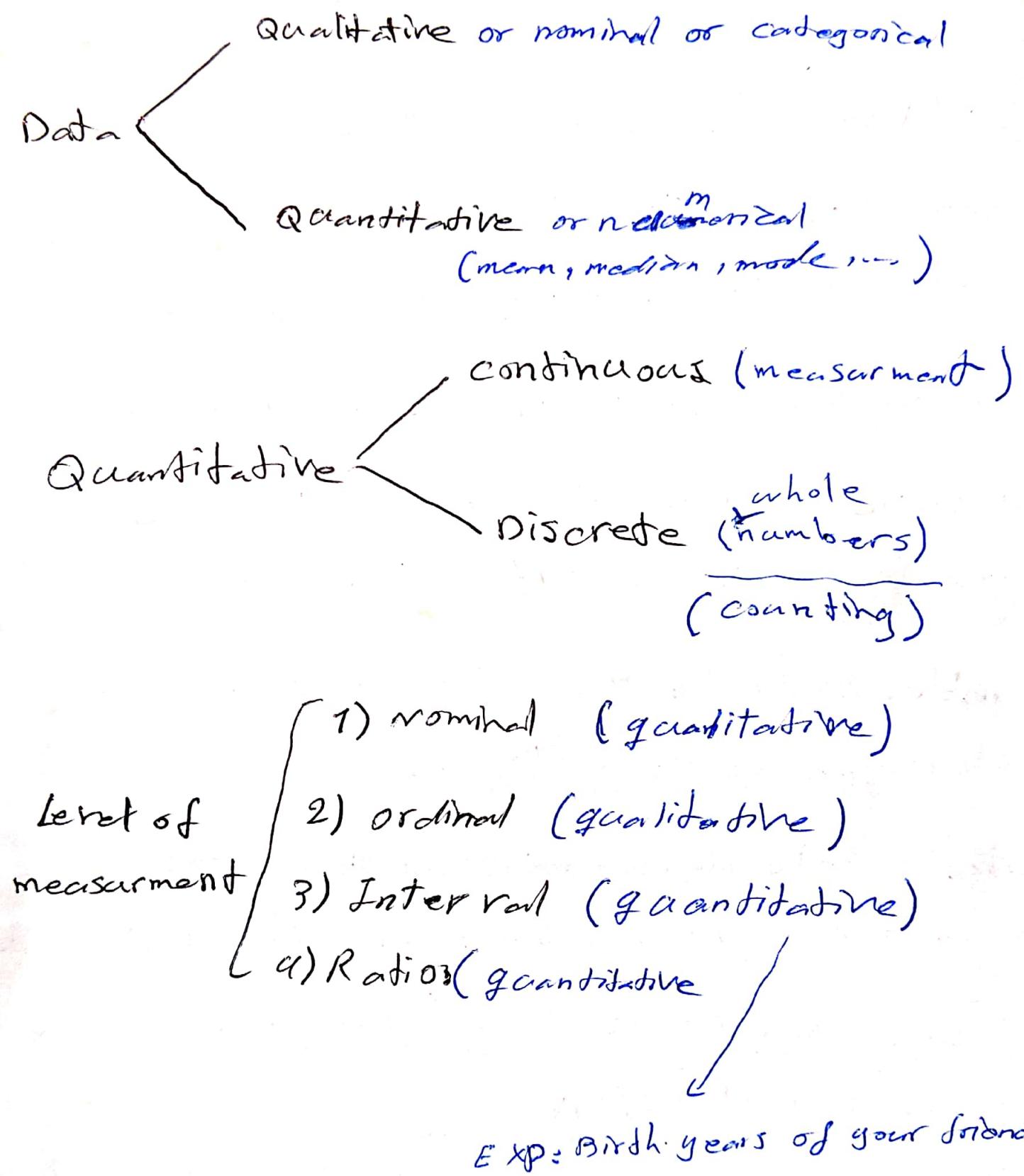
Data: counts, measurements or observations recorded about a specific variable to study in a population.

Census: when data is studied from every member of the POPULATION.

Descriptive statistics: Data that is collected, organized, summarized and displayed.

Inferential statistics: using data that is gathered to make predictions (or inferences) about the population.

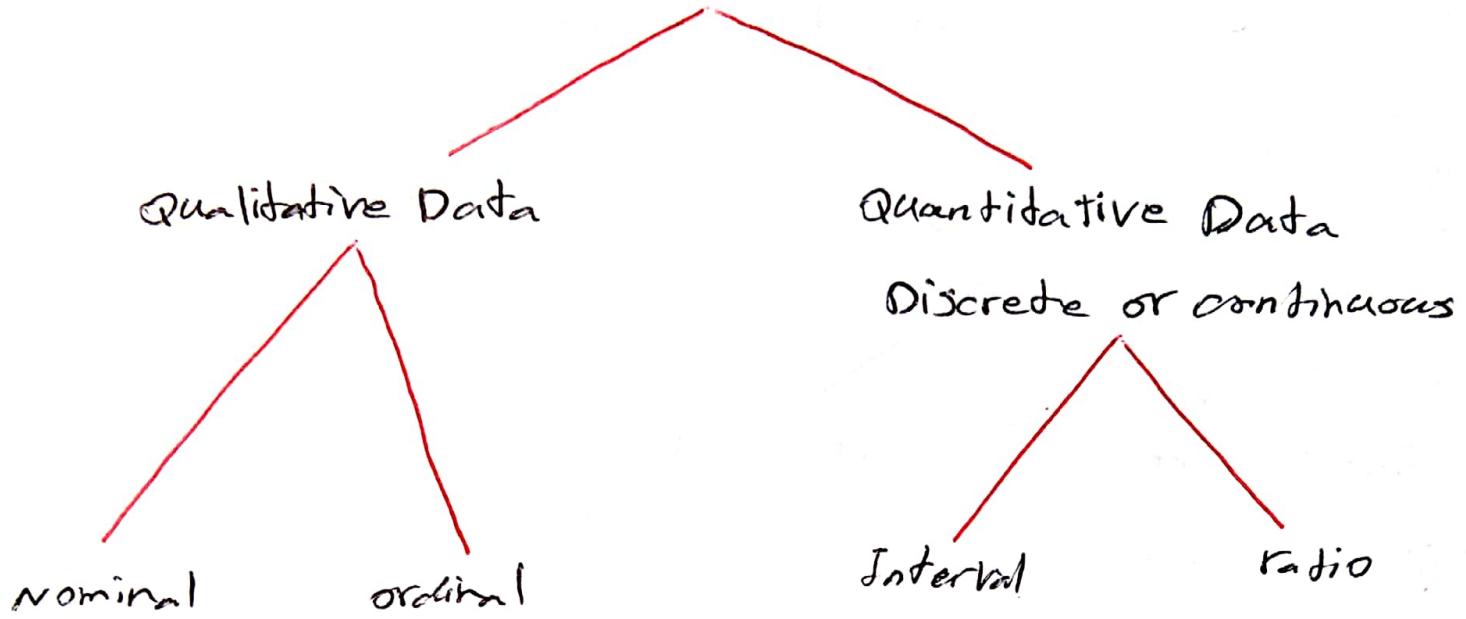
Parameters.



**Interval:** can be arranged in a meaningful order and where differences between data entries are meaningful, but not zero.

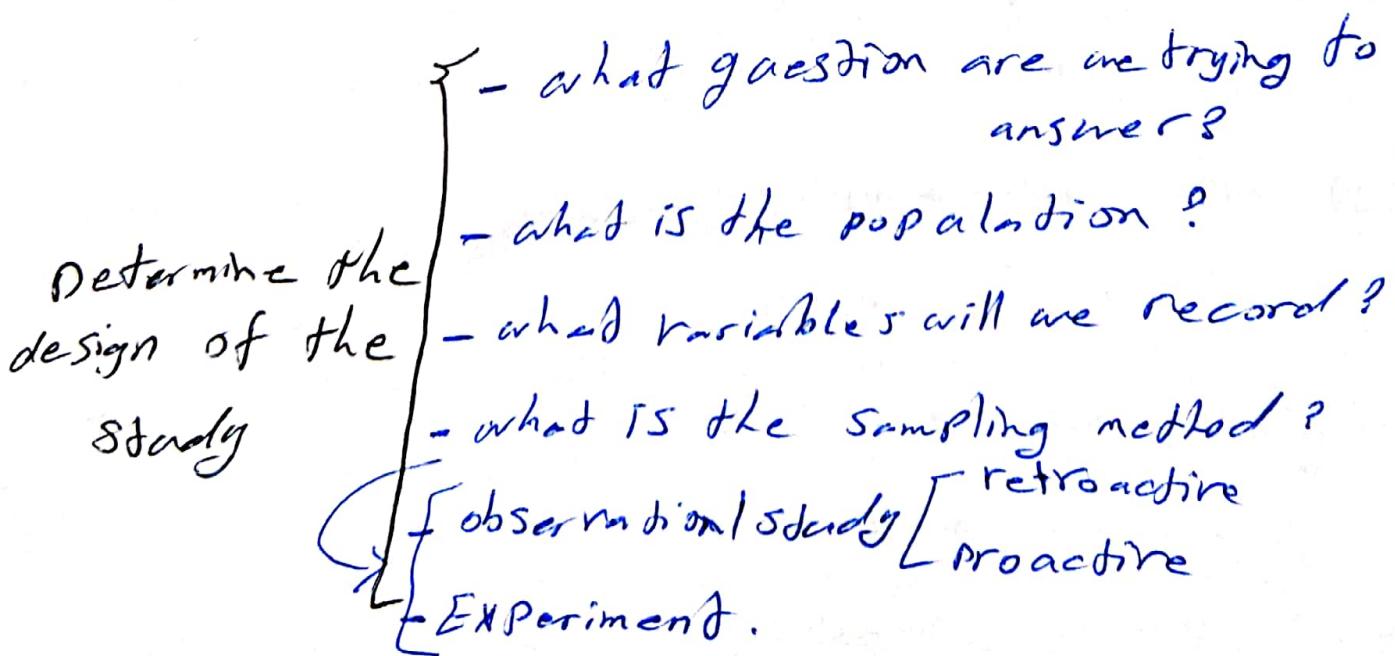
**Ratio:** can be arranged in a meaningful order, when differences between data entries are meaningful and the zero point indicates ~~of~~ the absence of something.

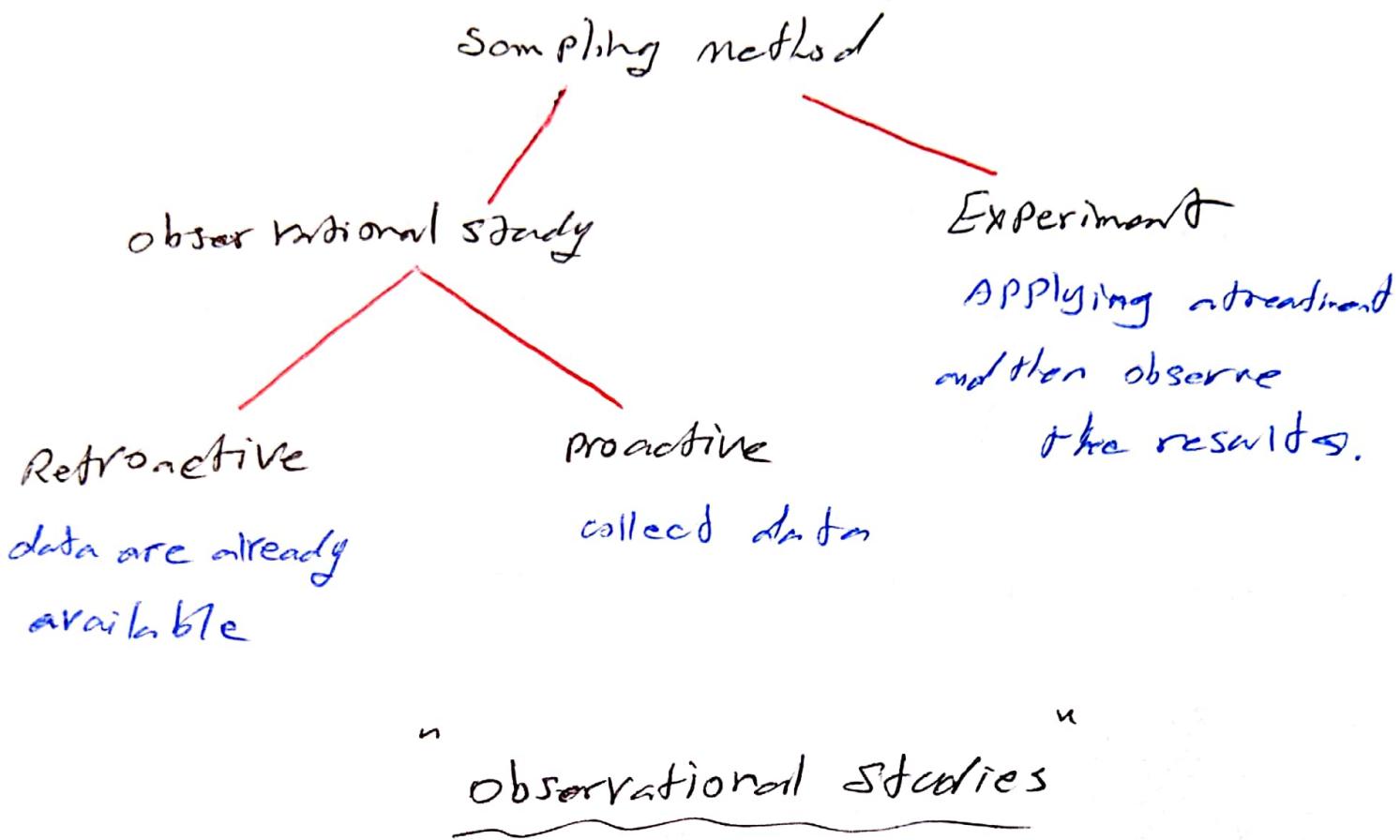
## Determine Data type



## The Process of a statistical study

- 1) Determine the design of the study
- 2) collect the data
- 3) organize the data
- a) Analyze the data to answer the question.





- our goal is to collect data about our population.
- A census may be too expensive
- so we use sample!

types of

collect Sample's ~~types~~



selected.

- 1) random sample: every member has same chance to be
- 2) simple random samples: every sample has same chance to be selected
- 3) stratified samples: the population is divided in subgroups and a proportional random sample from each
- 4) Cluster Sample: same with same with one stratified but without subgroups called clusters.
- 5) systematic sample,
- 6) convenience sample,

## Types of observational studies

- 1) cross-sectional study: data are collected at a single point in time.
  - 2) longitudinal study: data are gathered over a period of time.
  - 3) meta-analysis: information from previous studies.
- a) case study: look at multiple variables that affect a single event.

## "Experiments"

- Experiments are the only way to show a cause-and-effect relationship.
- Treatment: some condition applied to a group of subjects in an experiment
- subjects: people (also known as participants) or things being studied in an experiment.
- Response variable: The variable in an experiment that responds to the treatment.
- Explanatory variable: the variable or other in an experiment that caused the change in the response variable.

- control group: a group of subjects to which no treatment or a placebo is applied.
- placebo effect: a response to the power of suggestion, rather than the treatment itself.
- treatment groups: a group of subjects to which researchers apply a treatment.

~~single~~

- confounding variables: factors other than the treatment that cause an effect on the subjects of an experiment.

### Experiment's steps:

- 1) Randomize the control and treatment groups.
- 2) Control for outside effects on the response variable.
- 3) Replicate the experiment a significant number of times to see meaningful patterns.

### "critiquing a published study"

- A) consider the study:
- |                                                                                                                                                                                                          |
|----------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------|
| <ul style="list-style-type: none"> <li>- who paid for the study?</li> <li>- where were the data collected?</li> <li>- when was the information collected?</li> <li>- who published the study?</li> </ul> |
|----------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------|

- B) consider the variable
- EXP: (7)
- Bottle of ashampoo "makes hair goi smoother"
- smoother than what?
  - under what condition?
  - How did they mean "smoothness"?
  - in what type of weather was it tested?
  - under what condition?
  - what types of hair were tested?

C) consider the set-up:

Bias: is favoring of a certain outcome in a study.

Sampling bias: occurs when the sample does not accurately represent the population.

- Dropouts: are participants who begin a study but fail to complete it.
- Processing errors: errors that occur from the data processing
- nonadherents: a participant who remains in the study but stray from the directions
- Researcher bias: when a researcher influences the result of a study.
- Response bias: when a researcher's behavior causes a participant to alter their response or when a participant gives an inaccurate response.
- Participation bias: occurs when there is a lack of participation

-nonresponse bias: when a person refuses to participate in a survey or when a respondent omits questions when answering a survey.

D) Consider the conclusion

- Do the data support the conclusion?
- Do the results represent the whole picture or just a part?
- Could there be other conclusions drawn?
- Could there be other reasons for the same conclusion drawn?
- Does the study have any practical applications?

### "Measures of center"

Mean:  $\bar{x} = \frac{\sum_{i=1}^n x_i}{n}$

Median?

Mode?

### "measures of spread"

Range: max - min

- Population:  $\sigma^2 = \frac{\sum_{i=1}^N (x_i - \mu)^2}{N}$

Variance

- Sample:  $s^2 = \frac{\sum_{i=1}^n (x_i - \bar{x})^2}{n-1}$

(9)

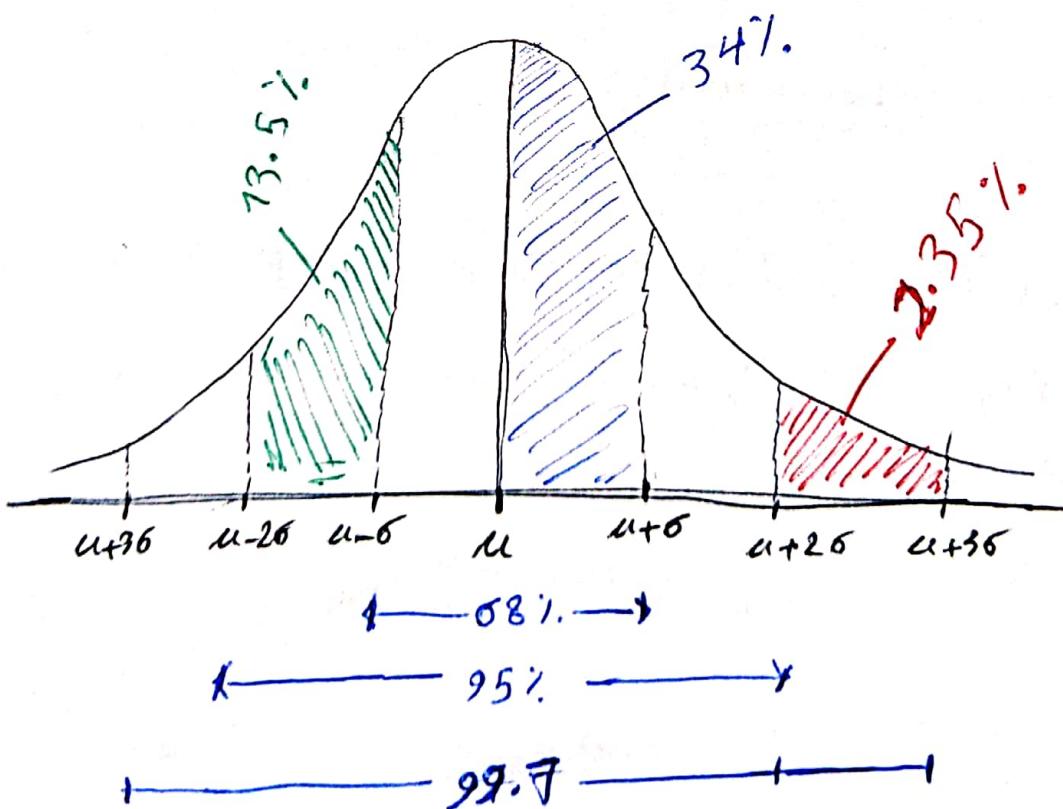
Standard deviation

- population  $\sigma = \sqrt{\frac{\sum (x_i - \mu)^2}{n}}$

- Sample  $s = \sqrt{\frac{\sum (x_i - \bar{x})^2}{n-1}}$

"Empirical Rule" (for Just for normal distribution)

empirical rule for normal distributions it gives us an estimate for how much data falls within 1, 2 or 3 standard deviations of our mean.



## "Chebychev's theorem" [for other distribution]

(10)

Chebychev's theorem: estimate of spread of data based on standard deviation for non-normal data.

The proportion of data that lie within  $K$  standard deviations of the mean is at least  $1 - \frac{1}{K^2}$  for  $K > 1$ .

$$\text{Exp: } K=2 = 1 - \frac{1}{2^2} = 1 - \frac{1}{4} = \frac{3}{4} = 75\% \text{ falls within } 2 \text{ std.}$$

$$K=3 = 1 - \frac{1}{3^2} = \frac{8}{9} \approx 88.9\%$$

$$K=4 = 1 - \frac{1}{4^2} = \frac{15}{16} = 93.75\%$$

## "Percentiles"

- To find the data value for the  $P^{\text{th}}$  percentile, we use the location ( $L$ ) of the data value ( $n$ ) in the data set is given by,

$$L = n \cdot \frac{P}{100}$$

$$\text{Exp: } \begin{cases} n=735 \\ P=10\% \end{cases}$$

$$L = n \cdot \frac{P}{100} = 735 \times \frac{10}{100} = 73.5$$

Find the value of the 10<sup>th</sup> percentile. that is the position of the value of our data set in

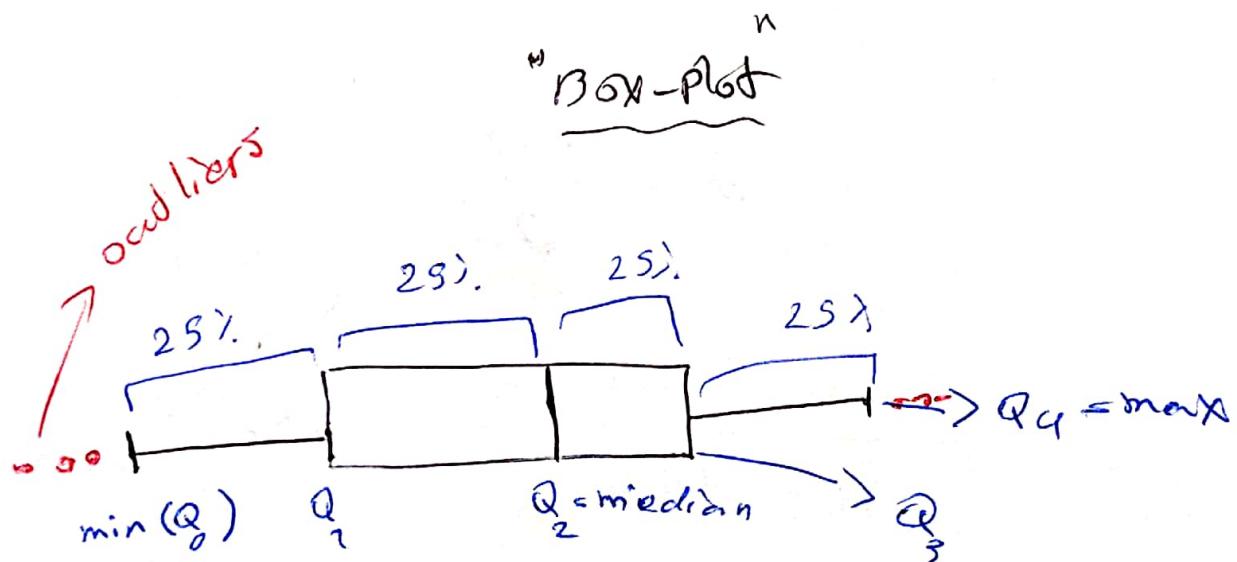
## ~~Statistik~~ Quartiles ( $Q_i$ )

→ a special type of percentile.

$Q_1 = 25^{\text{th}}$  percentile

$Q_2 = 50^{\text{th}}$  percentile or median

$Q_3 = 75^{\text{th}}$  percentile



$$IQR = Q_3 - Q_1$$

- To determine if our data has outliers, we can move  $\pm 1.5 \cdot IQR$ 's outside of our quartiles.

$$\begin{aligned} \text{lower fence} &= Q_1 - 1.5(IQR) \\ \text{upper fence} &= Q_3 + 1.5(IQR) \end{aligned} \quad ] \begin{array}{l} \text{every data outside} \\ \text{these bounds is} \\ \text{an outlier.} \end{array}$$

## "Z-scores" or Standard Scores

- it is a standard score, tells us how far a value is from the mean by computing how many std. it is away from the mean.

Population Z-score

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

Sample Z-score

$$Z = \frac{x - \bar{x}}{s}$$

$$Z = \frac{\text{observed} - \text{expected}}{\text{std. } (\sigma \text{ or } s)}$$

$$\begin{aligned} & \rightarrow (x - \bar{x}) \\ & \text{or} \\ & (x - \mu) \end{aligned}$$

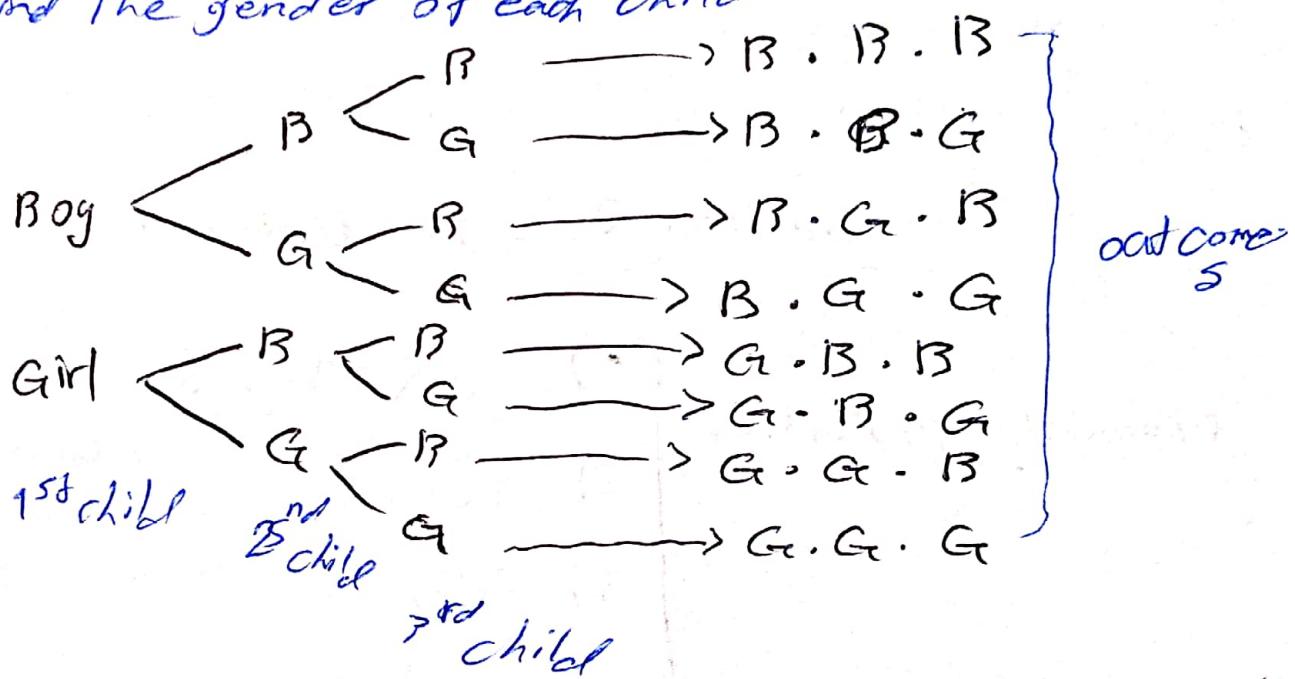
## "Probability"

- A probability experiment or trial is any process with a result determined by chance.
- An outcome is each individual result that is possible for a probability experiment.
- The sample space is the set of all possible outcomes for a given probability ~~exp~~ experiment.

an Event is a subset of outcomes from the sample space. (13)

A Tree diagram.

-a family with three children. use the tree diagram to find the gender of each child.



-There are three methods for calculating the probability of an outcome:

1) subjective probability: regarding of chance, the least precise type of probability. Exp: predicting the weather.

2) Experimental probability (aka Empirical): based on actual data. Exp:  $P(E) = \frac{f}{n} = \frac{48}{100} = 0.48$

f = frequency of event E

n = the total number of time the experiment is performed

3) Classical probability (aka Theoretical): based on what should happen.

$$P(E) = \frac{n(E)}{n(S)}, \quad \text{for example in a coin-flipping}$$

74

$$\text{then } P(\text{Heads}) = \frac{1}{2} = 0.5$$

- The law of large numbers: The more time an experiment performed, the closer the probability will get to the classical probability.

$$P(A \cup B) = P(A) + P(B) - P(A \cap B) \quad (\text{experimental probability})$$

"Addition Rule for Probability"

probability properties

- 1)  $0 \leq P(E) \leq 1$
- 2)  $P(S) = 1$  same probability of the sample space
- 3)  $P(\emptyset) = 0$  even not in the sample space
- a)  $P(\bar{E}) = 1 - P(E)$  the complement

- The addition rule for probability - "OR" rule.

6-sided die

$$P(3 \text{ or } 2) = \frac{1}{6} + \frac{1}{6} = \frac{2}{6} = \frac{1}{3}$$

$$P(\text{less than } 4) = \frac{1}{6} + \frac{1}{6} + \frac{1}{6} = \frac{3}{6} = \frac{1}{2}$$

$$P(\text{not less than } 4) = \frac{1}{6} + \frac{1}{6} + \frac{1}{6} + \frac{1}{6} = \frac{4}{6} = \frac{2}{3}$$

$$P(\text{less than } 4 \text{ or even}) = ? \quad P(A \cup B) = P(A) + P(B) - P(A \cap B)$$

$$P(\text{less than } 4) = P(1) + P(2) + P(3) = \frac{3}{6} = \frac{1}{2}$$

$$P(\text{even}) = P(2) + P(4) + P(6) = \frac{3}{6} = \frac{1}{2}$$

$$P(\text{less than } 4 \cap \text{even}) = P(2) = \frac{1}{6}$$

$$\text{so } \frac{1}{2} + \frac{1}{2} - \frac{1}{6} = \frac{5}{6}$$

## "multiplication Rule for probability" 75

multiplication Rules  $P(A \text{ and } B) = P(A) \times P(B)$  → with repl  
General Rule  $= P(A \text{ and } B) = P(A) \times P(B|A)$  ↓  
acement

$P(B|A) \rightarrow$  P of given A

without replacement

conditional probability - "Given"

$$P(B|A) = \frac{P(A \text{ and } B)}{P(A)} = \frac{P(A \cap B)}{P(A)}$$

- Fundamental counting principle: if event A can happen in  $n$  ways and B can happen in  $m$  ways, then A and B can happen in  $m \times n$  ways.

$\nearrow$  order matter       $\searrow$  order does not matter  
Permutation and Combination

- factorials

$$\frac{7!}{0!} = \frac{7 \times 6 \times 5 \times 4 \times 3 \times 2 \times 1}{1}$$

$$\frac{95!}{93!} = \frac{95 \times 94 \times 93!}{93!} = 95 \times 94$$

- combination (order does not matter):

$$n \underset{\substack{\text{number of items} \\ \downarrow \\ \text{number of choices}}}{C_r} = \binom{n}{r} = \frac{n!}{r!(n-r)!}$$

- Permutation (order matters):

$$n \underset{\substack{\text{number of items} \\ \downarrow \\ \text{number of choices}}}{P_r} = P(n, r) = \frac{n!}{(n-r)!}$$

- Permutation with repetition:

$$\frac{n!}{k_1! k_2! \dots k_p!}$$

Exps: How many different ways can you arrange the letters in the word STATISTICS?

$$S = 3$$

$$T = 3$$

$$A = 1$$

$$I = 2$$

$$C = 1$$

$$n = 10$$

$$\frac{n!}{k_1! k_2! \dots k_p!} = \frac{10!}{S! T! A! I! C!}$$

## expected value of Discrete probability Distributions

Expected value = weighted mean =  $\sum_{i=1}^n (x_i \cdot P(x=x_i))$

Exp: rolling a six-sided die:

$$\cancel{P(3) = \frac{1}{6}}$$

$x$	1	2	3	4	5	6
$P(x=x)$	$\frac{1}{6}$	$\frac{1}{6}$	$\frac{1}{6}$	$\frac{1}{6}$	$\frac{1}{6}$	$\frac{1}{6}$

$$\sum_{i=1}^n (x_i \cdot P(x=x_i)) = (\frac{1}{6} \cdot 1) + (2 \cdot \frac{1}{6}) + (3 \cdot \frac{1}{6}) + (4 \cdot \frac{1}{6}) + (5 \cdot \frac{1}{6}) + (6 \cdot \frac{1}{6}) =$$

Exp: let's play a game called Hearts - It costs \$5 to play - from a deck of cards, you randomly draw 1 card - If you draw:

- Queen or King of Hearts you get \$ 150
- Any other Queen or King you get \$ 10
- any other Heart you get \$ 5

Is this a good game to play?

outcome	Q/K hearts	Other Q/K	other Hearts	other card's
$x$	\$ 145	\$ 5	0	-5
$P(x=x)$	$\frac{2}{52}$	$\frac{6}{52}$	$\frac{11}{52}$	$\frac{33}{52}$

$$\sum_{i=1}^n \left( x_i \cdot P(X=x_i) \right) = \left( 145 \times \frac{2}{52} \right) + \left( 5 \times \frac{5}{52} \right) + \left( 0 \times \frac{11}{52} \right) - \left( 5 \times \frac{73}{52} \right) \approx 2.98$$

~~We~~ we pay \$5 to play the game but just can win  $\approx \$3$  per game so that is not a good game to play.

- Variance and Standard Deviation of Discrete Prob.

$$\sigma^2 = \sum \left[ (x_i - \mu)^2 \cdot P(X=x_i) \right]$$

$$\sigma = \sqrt{\sigma^2} = \sqrt{\sum \left[ (x_i - \mu)^2 \cdot P(X=x_i) \right]}$$

### The Binomial Distribution

- in binomial distribution, the experiment must have the following characteristics:

- A) The experiment consists of a fixed number,  $n$ , of identical and independent trials.
- B) There are only two possible outcomes - success or failure.

- ~~(a)~~) probability of success =  $P$
- ~~(b)~~) probability of failure =  $1 - P = q$
- number of trials =  $n$

The random variable ( $X$ ) representing the number of successes in  $n$  trials.

$$P(X=x) = {}_n C_x \cdot p^x (1-p)^{n-x}$$

Exps what is the probability of getting 6 heads in 10 tosses of a coin?

$$\underline{p=0.5}$$

$$P(X=6) = {}_{10} C_6 \cdot (0.5)^6 (1-0.5)^{10-6}$$

$$\underline{q=0.5}$$

$$P(X=6) = {}_{10} C_6 (0.5)^6 (0.5)^9$$

$$\underline{n=10}$$

$$P(X=6) = \frac{10!}{6!(10-6)!} \cdot (0.5)^6 \cdot (0.5)^9$$

$$P(X=6) = \frac{10 \times 9 \times 8 \times 7}{4 \times 3 \times 2 \times 1} \cdot (0.5)^6 \cdot (0.5)^9$$

$$P(X=6) \approx 0.2057$$

### "Cumulative Binomial Distribution"

Exps: in an experiment of flipping a coin 10 times, what is the probability of flipping 7 or more heads?

$$P(X \geq 7) = P(X=7) + P(X=8) + P(X=9) + P(X=10)$$

or  $(1 - P(X < 7))$

## "The Poisson Distribution"

- in the Poisson distribution, the experiment must have the following characteristics:

- A) Each success must be independent of other successes
- B) The mean number of successes in a given interval must remain constant

- mean number of success in a given interval =  $\lambda$
- The random variable =  $X \rightarrow$  it represents the number of   
      success in the given interval.

$$P(X=x) = \frac{e^{-\lambda} \lambda^x}{x!}$$

Exp = A barber usually finishes one haircut every 75 mins.  
calculate the probability that our barber is feeling extra  
speedy one day and finishes 6 haircuts in one hour.

$$P(X=x) = \frac{e^{-\lambda} \lambda^x}{x!} = \frac{e^{-9} 9^6}{6!} = ?$$

$$\lambda = \frac{60}{75} = 4$$

$$x = 6$$

## "The Hypergeometric Distribution"

- in this case, the experiment must have the following characteristics

- A) each trial consists of an outcome of  $P$  or  $q$ .
- B) The population is of a known size.
- C) The trials are dependent (that is, selections are made without replacement).

- number of successes:  $X$
- number of dependent trials:  $n$
- number of successes in the population:  $K$
- number of items in the population:  $N$

$$P(X=x) = \frac{\binom{K}{x} \binom{N-K}{n-x}}{\binom{N}{n}}$$

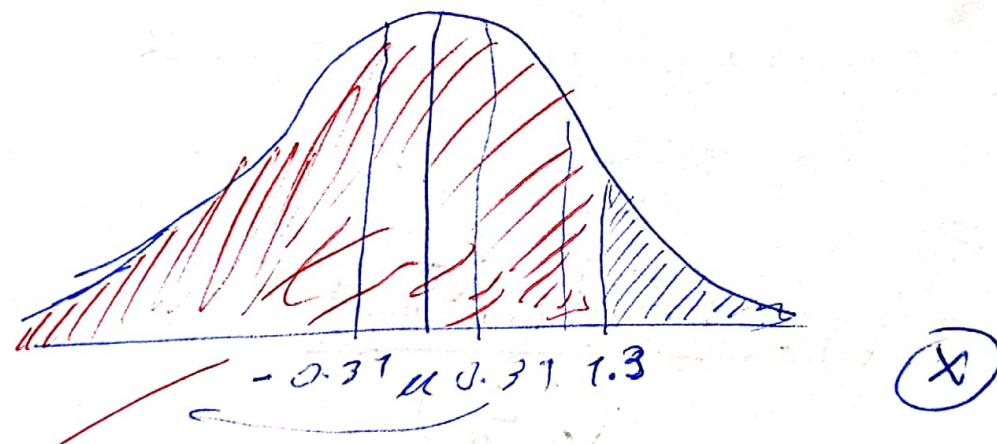
Exps At a local grocery store there are twenty boxes of cereal on one shelf, half of which contain a prize, suppose that you buy three boxes of cereal. what is the probability that all three boxes contain a prize?

$$\begin{array}{c} \text{population} \\ N=20 \\ K=10 \end{array} \quad \left. \begin{array}{c} \text{sample} \\ n=3 \\ x=3 \end{array} \right\}$$

$$P(x=3) = \frac{\binom{10}{3} \binom{20-10}{3-3}}{\binom{20}{3}} = \frac{\binom{10}{3} \binom{10}{0}}{\binom{20}{3}}$$

$$= \frac{(120)(1)}{(7920)} = \frac{2}{79} \approx 0.1053$$

### "The normal Distribution and Z-Score"



$$P(Z < 0.37) = \text{using Z-table} = 0.6277$$

$$P(Z < 1.37) = 0.9116$$

$$P(Z < -2.03) = 0.0277$$

Exps Body temperatures for adults are normally distributed with a mean of  $98.6^{\circ}\text{F}$  and a standard deviation of  $0.73^{\circ}\text{F}$ . (23)

A) what is the probability of a healthy adult having a body temperature less than  $97^{\circ}\text{F}$ ?

$$Z = \frac{x - \mu}{\sigma} = \frac{97 - 98.6}{0.73} = -2.19178 \xrightarrow{\text{use Z-table}}$$

$$P(Z < 97^{\circ}\text{F}) \xrightarrow{\text{Z-table}} 0.019197 \leftarrow \text{probability of less than } 97^{\circ}\text{F}$$

B) what is the probability of a healthy adult having a body temperature more than  $99.2^{\circ}\text{F}$ ?

$$Z = \frac{x - \mu}{\sigma} = \frac{99.2 - 98.6}{0.73} = 0.82$$

$$P(Z > 0.82) = 1 - P(Z < 0.82) = 0.2056$$

Exp: Compute the Z-scores for each situations

A) an area to the left of 0.239:

$$0.239 \xrightarrow{\text{use Z-table}} -0.7095 \rightarrow \text{Z-score}$$

if  $\begin{cases} \text{mean} = 52 \\ \text{std} = 3 \end{cases} \Rightarrow Z = \frac{x - \mu}{\sigma} \Rightarrow x = Z\sigma + \mu$

$$\Rightarrow x = (-0.7095)(3) + 52$$

$$x = 49.871$$

B) The 80th percentile:  $\bar{x} = 52, \sigma = 3$

(29)

$$0.8 \xrightarrow[\text{using Z-table}]{\text{near}} \boxed{Z = 0.8916}$$

in reverse

$$x = Z\sigma + \mu$$

$$= (0.8916)(3) + 52 = 59.525$$

C) an area to the right of 0.963.

$$0.963 \longrightarrow \boxed{Z = -1.79}$$

~~1 - 0.963~~

Approximating a Binomial Dist. with a normal Dist.

- If the conditions that  $np \geq 5$  and  $n(\frac{1-p}{p}) \geq 5$  are met for a given binomial distribution, and the distribution is very large, then a normal distribution can be used to approximate the binomial probability distribution with the mean and std. given by.

$$\bar{x} = np \quad \sigma = \sqrt{np(1-p)}$$

one problem:

(25)

- we are trying to use the normal distribution, which is continuous, to approximate the binomial distribution, which is discrete. Therefore, a continuity correction must be used to convert the whole number value of the discrete binomial random variable to an interval range of continuous normal random variable.
- To make this correction, determine the value  $x$  of the binomial random variable and convert it to the interval from  $(x - 0.5)$  to  $(x + 0.5)$ .

- Ex: use a normal distribution to estimate the probability of more than 55 girls being born in 100 births. Assume that the probability of a girl being born in an individual birth is 50%.

first check the conditions

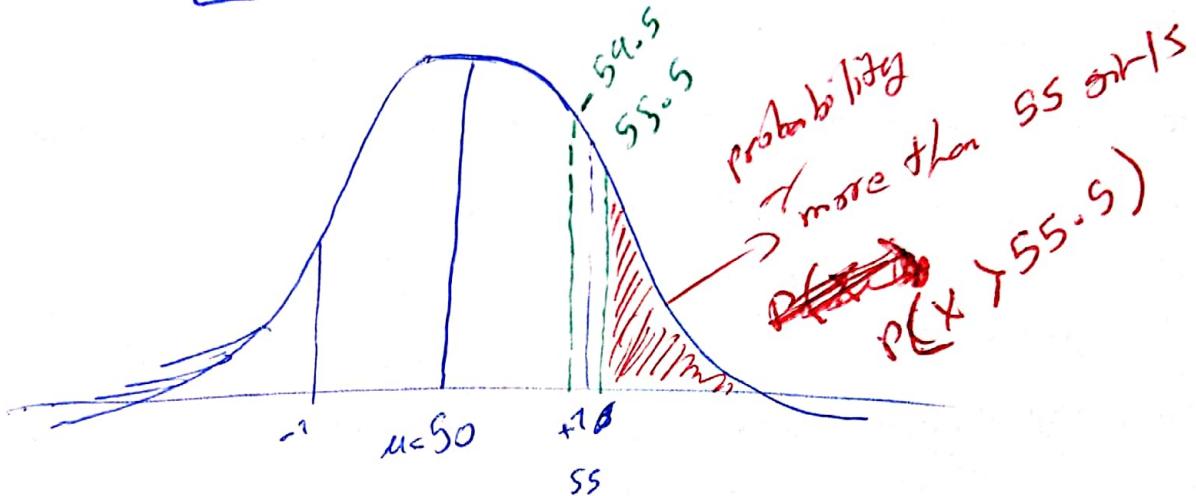
$$\left\{ \begin{array}{l} np = 100(0.5) = 50 \geq 5 \\ n(1-p) = 100(0.5) = 50 \geq 5 \end{array} \right.$$

second

$$\left\{ \begin{array}{l} \mu = np = 100(0.5) = 50 \\ \sigma = \sqrt{np(1-p)} = \sqrt{100(0.5)(0.5)} = \sqrt{25} = 5 \end{array} \right.$$

third: calculate the continuity corrections

$$\begin{cases} 55 + 0.5 = 55.5 \\ 55 - 0.5 = 54.5 \end{cases}$$



$$0.555 \xrightarrow[\text{Z-score table}]{\text{using Z-table}} 1.1$$

$$p(x > 55.5) = 1 - p(x < 55.5) = \underline{\underline{0.1356}}$$

### "The Central Limit Theorem"

- it states that as the sample size increases, the distribution of the sample means approaches a normal distribution regardless of the underlying distribution of the population from which the sample is drawn.

The ~~sample~~ variables of the sample must be

- 1) independent
- 2) identically distributed random variables
- 3) the sample size must be large enough ( $30 < n$ )

- if the requirement is met, the sampling distribution of sample means will have the following three characteristics:

$$1) \quad \mu_{\bar{x}} = \mu$$

$\downarrow$   
the mean of a sampling distribution

$$2) \quad \text{the standard deviation of a sampling distribution} \rightarrow \sigma_{\bar{x}} = \frac{\sigma}{\sqrt{n}}$$

of sample mean

### "The (CLT) for means"

$$Z = \frac{\bar{x} - \mu_{\bar{x}}}{\sigma_{\bar{x}}} = \frac{\bar{x} - \mu}{\frac{\sigma}{\sqrt{n}}}$$

$\frac{\sigma}{\sqrt{n}}$   $\rightarrow$  Standard error

### "The (CLT) for proportions"

$$\begin{cases} p = \frac{x}{n} \rightarrow \text{population proportion (success)} \\ \hat{p} = \frac{x}{n} \rightarrow \text{sample proportion} \end{cases}$$

The requirements:

A)  $np \geq 5 \text{ & } n(1-p) \geq 5$

B) the population is ~~not~~ binomially distributed given by proportional data.

The mean of the sampling distribution of sample proportion

$$\text{mean} \rightarrow m\hat{p} = p$$

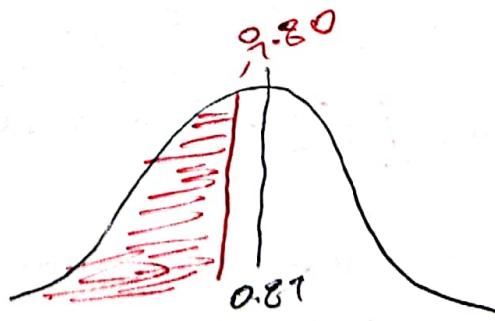
$$\text{std.} \rightarrow \sigma_{\hat{p}} = \sqrt{\frac{p(1-p)}{n}}$$

$$\text{Z-score} \rightarrow Z = \frac{\hat{p} - np}{\sigma_{\hat{p}}} = \frac{\hat{p} - p}{\sqrt{\frac{p(1-p)}{n}}}$$

Ex: In a certain liberal precinct across town, 81% of the voters are registered Democrats. What is the probability that, in a random sample of 100 voters from this precinct, no more than 80 of the voters would be registered Democrats?

Step 1: Identify  $\hat{p}$ ,  $p$  and  $n$

$$\begin{cases} n = 100 \\ p = 0.81 \\ \hat{p} = 0.8 \end{cases}$$



Step 2: Draw the curve

$$\text{Step 3: Find the Z-score } Z = \frac{0.80 - 0.81}{\sqrt{\frac{0.81(0.19)}{100}}} \approx -0.2599$$

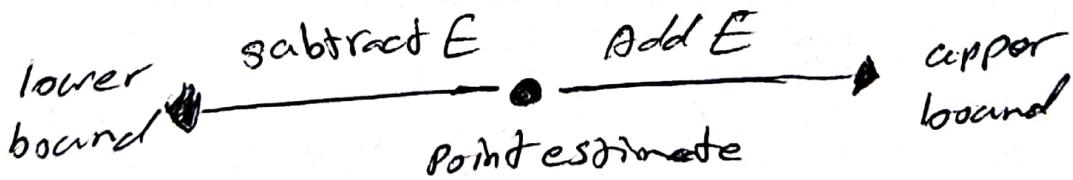
## "Confidence Intervals"

- The level of confidence is represented by the letter C, thus, for a 95% confidence interval,  $C = 0.95$
- It tells us if there is a 95% level of confidence that given, that we are confident that using this statistical method will result in 95% of all possible interval estimates containing the true population parameter.

do

- How are find a confidence interval for a population parameter?

- The margin of Error or maximum error of estimate, E, is the largest possible distance from the point estimate that a confidence interval will cover.



- Exp: A college student researching study habits collects data from a random sample of 250 college students on her campus and calculates that the sample mean is  $\bar{x} = 23.7$  hours per week. If the margin of error for her data using a 95% level of confidence

is  $E = 0.6$  hours, construct a 95% confidence interval for her data. Interpret your ~~data~~ results.

point estimate:  $\bar{x} = 15.7$

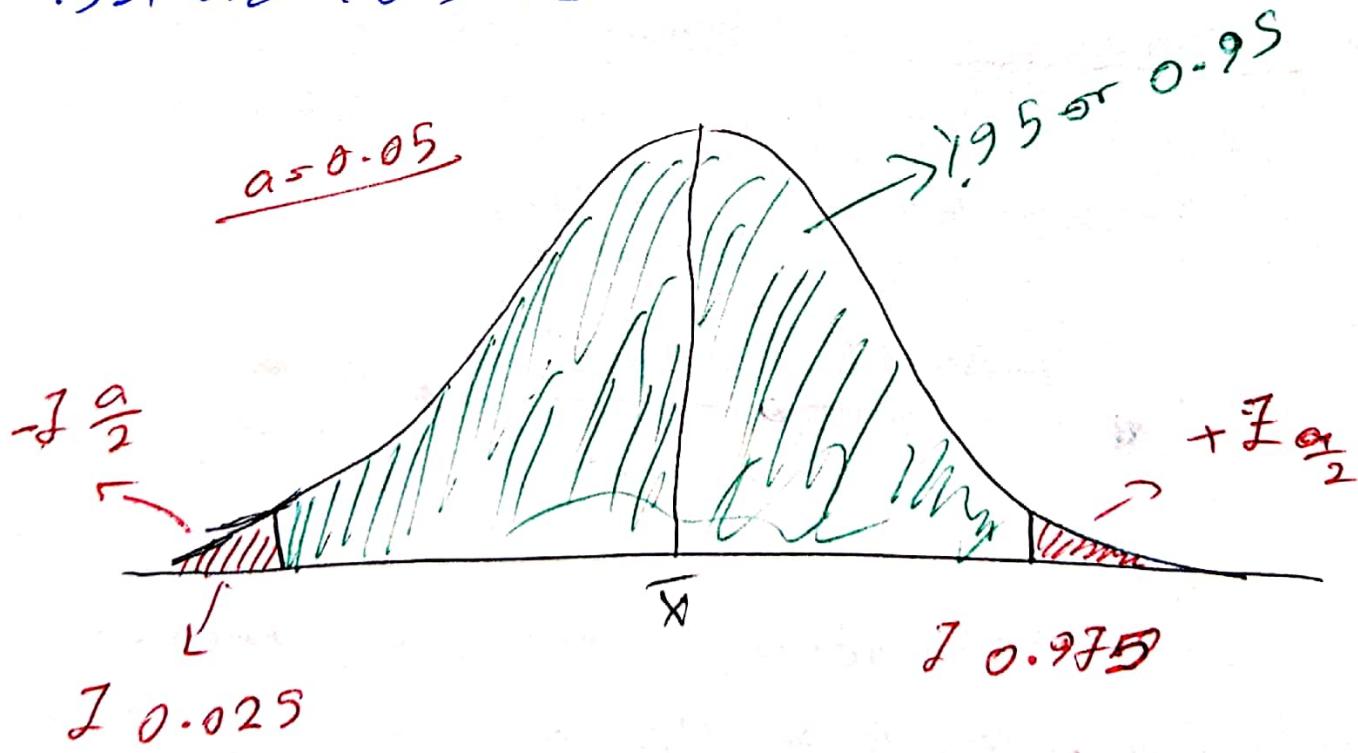
margin of error:  $E = 0.6$

$$\text{lower bound: } \bar{x} - E = 15.7 - 0.6 = 15.1 \text{ (hrs/week)}$$

$$\text{upper bound: } \bar{x} + E = 15.7 + 0.6 = 16.3 \text{ (hrs/week)}$$

$$\text{confidence interval: } 15.1 < \mu < 16.3 \text{ or } (15.1, 16.3)$$

Interpretation: we are 95% confident that the true population mean for the number of hrs/week that students on this campus spend studying is between 15.1 and 16.3 hrs.



$$E = (Z_{\frac{\alpha}{2}}) \left( \frac{\sigma}{\sqrt{n}} \right)$$

Expt: In order to estimate the number of calls to expect at a new suicide hotline, volunteers contact a random sample of 35 similar hotlines across the nation and find that the sample mean is 92.0 calls per month. Construct a 95% confidence interval for the mean number of calls per month. Assume that the population standard deviation is known to be 6.5 calls per month. (31)

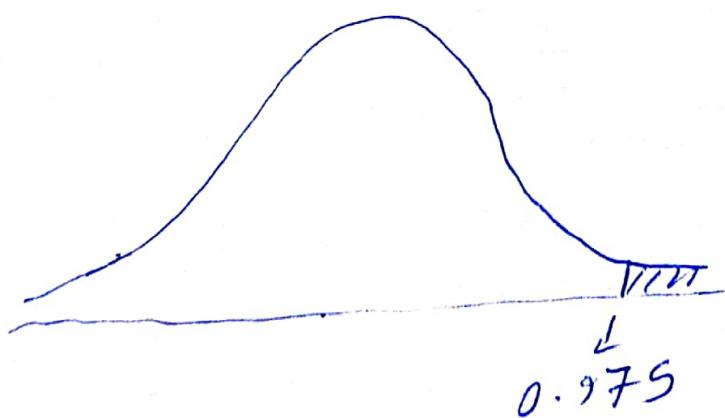
1) check the conditions

- random sample ✓
- $\sigma$  is known ✓
- $n > 30$  ✓

2) find  $\bar{x}$  (point estimate) -  $\bar{x} = 92$

3) find  $Z_{\frac{\alpha}{2}}$ ,  $n$  and  $\sigma$

- $n = 35$
- $\sigma = 6.5$
- $Z_{\frac{\alpha}{2}} = 1.96$



$$0.975 \xrightarrow{\text{Z-score table}} 1.96$$

4) calculate the margin of error,  $E = \left( Z_{\frac{\alpha}{2}} \right) \left( \frac{\sigma}{\sqrt{n}} \right)$  (32)

$$E = (1.96) \left( \frac{0.5}{\sqrt{35}} \right) \approx 2.753953$$

5) calculate the interval

lower bound:  $x - E = 92 - 2.753953 \approx 39.8$  calls/month

upper bound:  $x + E = 92 + 2.753953 \approx 99.2$  calls/month

$$(39.8, 99.2)$$



the true population mean ~~is~~ is in this interval.

\* Find the sample size based on confidence Interv-  
also

$$\frac{E = Z_{\frac{\alpha}{2}} \left( \frac{\sigma}{\sqrt{n}} \right)}{Z_{\frac{\alpha}{2}}} = \frac{E}{Z_{\frac{\alpha}{2}}}$$

$$\sqrt{n} \cdot \frac{E}{Z_{\frac{\alpha}{2}}} = \frac{\sigma}{\sqrt{n}} \cdot \sqrt{n}$$

$$(\sqrt{n})^2 = \left( \frac{\sigma Z_{\frac{\alpha}{2}}}{E} \right)^2 \Rightarrow n = \left( \frac{\sigma Z_{\frac{\alpha}{2}}}{E} \right)^2$$

(33)

ExPs Determine the minimum sample size needed if we wish to be 90% confident that the sample mean is within  $\frac{\sigma}{2}$  units of the population mean. An estimate for the population std. of 8.9 is available from a previous study.

$$\sigma = 8.9$$

$$C = 0.90$$

$$E = 2$$

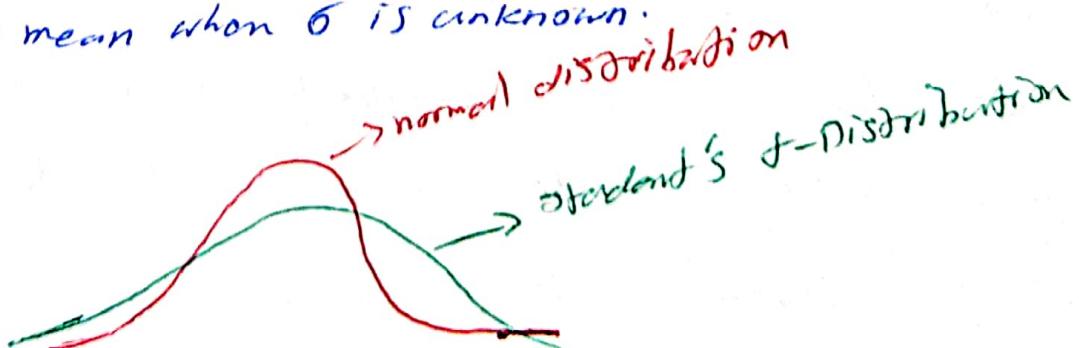
$$I_{\frac{\alpha}{2}} = I_{0.10} = 1.645$$

$$n = \left( \frac{I_{\frac{\alpha}{2}} \cdot \sigma}{E} \right)^2 = \left( \frac{1.645 \times 8.9}{2} \right)^2 = 47.73 \approx 48$$

the least size of sample.

### "Student's t-Distribution"

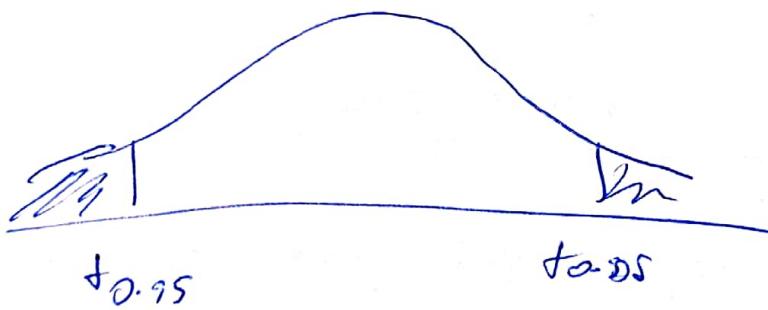
- when the population is normally distributed or the sample size is sufficiently large. in this case Student's t-distribution is a good option to calculate the margin of error for a population mean when  $\sigma$  is unknown.



- whence
- standard normal distribution has two parameters 1)  $\mu = 0$
  - 2)  $\sigma = 1$ , the t-distribution has only one parameter
  - the number of degrees of freedom  $\nu$  (df)
  - A t-distribution with fewer degrees of freedom has more area under the tails of the curve.

Exps: find the value of  $t$  for a t-distribution with 17 degrees of freedom such that the area under the curve to the left of  $t$  is 0.95. (when we denote  $t$ , the subscript denotes the area to the right, so we would be finding  $t_{0.95}$ )

$$\text{in excel} = T.INV(0.05, 17)$$



- find the value of  $t$  for a t-distribution with 17 df such that the area under the curve to the right of  $t$  is 0.10 - this would be indicated  $t_{0.9}$ . However, one would calculate the area to the left,  $1 - 0.1 = 0.9$

$$T.INV(0.9, 17)$$

in excel

$$(\bar{x} - E, \bar{x} + E) \text{ where } E = \left( t_{\frac{\alpha}{2}} \right) \left( \frac{s}{\sqrt{n}} \right)$$

$$\left( \bar{x} - \left( t_{\frac{\alpha}{2}} \right) \left( \frac{s}{\sqrt{n}} \right), \bar{x} + \left( t_{\frac{\alpha}{2}} \right) \left( \frac{s}{\sqrt{n}} \right) \right)$$

Ex: A marketing company wants to know the mean price of new vehicles sold in an up-and-coming area of town. The simple random sample of 756 cars has a mean of \$27.800 with a std. of \$1300. Construct a 95% confidence interval for the mean price of new cars sold in this area.

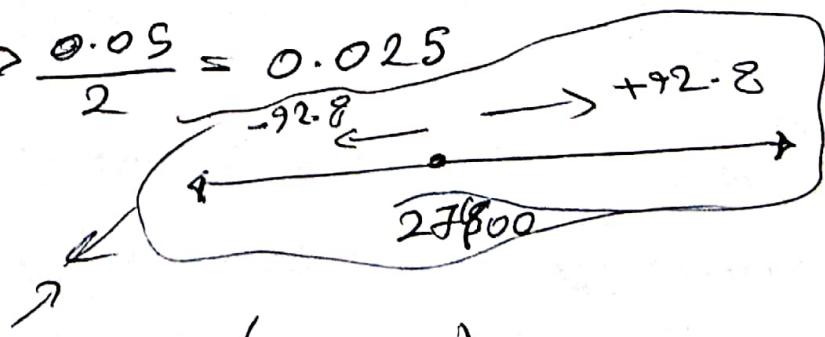
Step 1) find the point estimate: The point estimate for the population mean is the sample mean, which is \$27800.

Step 2) find the critical value: since  $n=756$ , is sufficiently large, we use the student's t-distribution to calculate the critical value.

$$\alpha = 1 - 0.95 = 0.05 \rightarrow \frac{0.05}{2} = 0.025$$

$$df = n - 1 = 755$$

$$\text{so, } t_{0.025} \approx 1.963$$



Step 3)  $E = \left( t_{\frac{\alpha}{2}} \right) \left( \frac{s}{\sqrt{n}} \right) = (1.963) \left( \frac{1300}{\sqrt{756}} \right) \approx 92.811706$

Step 4)

$$\text{lower end point} = \bar{x} - E = 27902 - 92.8 \approx 27809 \quad (36)$$

$$\text{upper end point} = \bar{x} + E = 27993$$

Step 5)  $(27809, 27993)$

### Estimating population proportions

- since the  $P$  (population proportion) is often difficult to find, we use  $\hat{P}$  (sample proportion) as our point estimate  $\hat{P} = \frac{x}{n}$ . we will calculate a margin of error using the normal model and then add it on subtracting it from a point estimate.

$$(\hat{P} \geq 10, n(1-\hat{P}) \geq 10)$$

$$(\hat{P} - E, \hat{P} + E) \text{ where } E = Z_{\frac{\alpha}{2}} \sqrt{\frac{\hat{P}(1-\hat{P})}{n}}$$

$$\left( \hat{P} - Z_{\frac{\alpha}{2}} \sqrt{\frac{\hat{P}(1-\hat{P})}{n}}, \hat{P} + Z_{\frac{\alpha}{2}} \sqrt{\frac{\hat{P}(1-\hat{P})}{n}} \right)$$

EXP: A survey of 345 randomly selected students at one university found that 301 students think that there is not enough parking on campus. Find the 90% confidence interval for the proportion of all students at this uni who think that there is not enough parking on campus.

$$\hat{P} = \frac{x}{n} = \frac{307}{395} \approx 0.77$$

(37)

$$n = 395$$

$$c = 0.9$$

$$\alpha = 1 - 0.9 = 0.1 \Rightarrow \frac{\alpha}{2} = 0.05$$

$$\left( \hat{P} - Z_{\frac{\alpha}{2}} \sqrt{\frac{\hat{P} \cdot \hat{q}}{n}}, \hat{P} + Z_{\frac{\alpha}{2}} \sqrt{\frac{\hat{P} \cdot \hat{q}}{n}} \right)$$

$$\left( 0.87 + 1.96 \sqrt{\frac{0.87 \cdot 0.13}{395}}, 0.87 + 1.96 \sqrt{\frac{0.87 \cdot 0.13}{395}} \right)$$

$$\text{upper bound} = 0.87 + 0.0295 \approx 0.893$$

$$\text{lower bound} = 0.87 - 0.0295 \approx 0.842$$

$$0.843, 0.893$$

$$(0.843, 0.893)$$

Interpret the interval: with 90% confidence, we estimate that between 81.3% and 90.2% of all students at this uni think there is not enough parking on campus.

# "calculations with population proportions"

$$E = \frac{Z_{\alpha/2}}{2} \sqrt{\frac{\hat{p}(1-\hat{p})}{n}}$$

$$\Rightarrow \left( \frac{E}{Z_{\alpha/2}/2} \right)^2 = \frac{\hat{p}(1-\hat{p})}{n}$$

$$\Rightarrow n \cdot \left( \frac{E}{Z_{\alpha/2}/2} \right)^2 = \frac{\hat{p}(1-\hat{p})}{\underline{\underline{n}}}$$

$$\Rightarrow n = \hat{p} \hat{q} \left( \frac{Z_{\alpha/2}}{E} \right)^2$$

Expt: what is the minimum sample size needed for a 99% confidence interval for the population proportion if previous studies indicated  $P \approx 0.59$  and we desire no more than a 2% margin of error?

$$\hat{p} \approx 0.59$$

$$\alpha = 0.01 \quad \left. \right\} Z_{\alpha/2} = ?$$

$$\frac{\alpha}{2} = 0.005$$

$$E = 0.02$$

$$n = \hat{p} \hat{q} \left( \frac{Z_{\alpha/2}}{E} \right)^2$$

$$n = 0.59 \times 0.41 \left( \frac{2.575}{0.02} \right)^2$$

$$n \approx 4777.078725$$

$$n \approx 4778$$

## Comparing Population means ( $\sigma$ Known)

39

Comparing two samples:

- A) The samples are random and independent.
- B) both population std. are known.
- C) either  $n_1, n_2 \geq 30$  or both are normally distributed.

- Creating the confidence interval

1) find the point estimate:  $\bar{X}_1 - \bar{X}_2$

2) find the margin of error:  $E = Z_{\frac{\alpha}{2}} \sqrt{\frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}}$

3) the confidence interval,

$$\left( (\bar{X}_1 - \bar{X}_2) - Z_{\frac{\alpha}{2}} \sqrt{\frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}}, (\bar{X}_1 - \bar{X}_2) + Z_{\frac{\alpha}{2}} \sqrt{\frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}} \right)$$

a) interpret the confidence interval.

- If the CI contains 0, we must conclude that there is a possibility that the two population means are the same.  
- if the CI lies completely on the positive side of the line, the first mean > second mean  $\overline{\bar{X}_1} > \overline{\bar{X}_2}$

- if the CI is completely on the negative side of the line, the first mean < second mean

EXP: A researcher is looking at the study habits of college students. A random sample of 92 freshmen reported a mean study time of 15 hr/week. A random sample of 39 seniors reported a mean study time of 23 hr/week. Construct a 95% CI to estimate the true difference between the amount of time per week that both groups study. The Std. for freshmen = 4.7 and for seniors =

$$\bar{x}_1 - \bar{x}_2 = 15 - 23 = -8 \text{ hr/week}$$

$$\alpha = 1 - 0.95 = 0.05 \quad S = 2.$$

$$\frac{\alpha}{2} = 0.025$$

$$E = t_{\frac{\alpha}{2}} \sqrt{\frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}} \Rightarrow t = \frac{1.96}{0.025} \sqrt{\frac{(4.7)^2}{92} + \frac{(5.2)^2}{39}} \approx 2.169$$

$$(\bar{x}_1 - \bar{x}_2) - E, (\bar{x}_1 - \bar{x}_2) + E$$

$$(-10.769, -5.836)$$

Interpretations with 95% Confidence: freshmen study between 5.8 - 10.2 hr/week less than seniors.

### Comparing population means ( $\sigma$ unknown)

- The population variances are assumed to be unequal.

$$1) (\bar{x}_1 - \bar{x}_2) = ?$$

$$3) (\bar{x}_1 - \bar{x}_2) - E, (\bar{x}_1 - \bar{x}_2) + E$$

$$2) E = t_{\frac{\alpha}{2}} \sqrt{\frac{s_1^2}{n_1} + \frac{s_2^2}{n_2}}$$

Exps misty believes that the difficulty she is having in the class is the result of an inexperienced teacher. She believes that students in her class are receiving lower scores on their exams than they in other another class with a more experienced teacher.

$$\left\{ \begin{array}{l} n_1 = 11 \\ \bar{x}_1 = 75 \\ s_1 = 8 \end{array} \right. \quad \boxed{C = 90\%} \quad \left\{ \begin{array}{l} n_2 = 9 \\ \bar{x}_2 = 82 \\ s_2 = 5 \end{array} \right. \quad \text{both normally distributed}$$

$$(\bar{x}_1 - \bar{x}_2) = 75 - 82 = -7$$

$$t_{\frac{\alpha}{2}} = \frac{1 - 0.9}{2} = \frac{0.1}{2} = \boxed{0.05}$$

\*  $df =$  the number of degrees of freedom is the smaller of the values  $n_1 - 1 = 10$  and  $n_2 - 1 = 8$

$$\text{so, } df = 8$$

$$= \boxed{1.860}$$

$$+ t_{0.05}$$

$$E = \frac{t_{\alpha}}{2} \sqrt{\frac{s_1^2}{n_1} + \frac{s_2^2}{n_2}} = 1.860 \sqrt{\frac{8^2}{11} + \frac{5^2}{9}} \approx \boxed{5.953}$$

$$(\bar{x}_1 - \bar{x}_2) - E, (\bar{x}_1 - \bar{x}_2) + E$$

$$(-12.5, -7.5)$$

(42)

Interpretations with 90% confidence, the mean exam score for students in another class is better than ~~the~~ her class.

Var

Comparing two pop. means (6 unknown, Equal)

- The population variances are assumed to be equal (use pooled variance)

1) Point estimate :  $\bar{x}_1 - \bar{x}_2$

2) margin of error where  $df = n_1 + n_2 - 2$

$$E = \frac{t_{\alpha/2}}{\sqrt{\frac{(n_1-1)s_1^2 + (n_2-1)s_2^2}{n_1+n_2-2} \left( \frac{1}{n_1} + \frac{1}{n_2} \right)}}$$

- 3) Interval confidence:

$$\left( (\bar{x}_1 - \bar{x}_2) - \frac{t_{\alpha/2}}{\sqrt{\frac{(n_1-1)s_1^2 + (n_2-1)s_2^2}{n_1+n_2-2} \left( \frac{1}{n_1} + \frac{1}{n_2} \right)}}, (\bar{x}_1 - \bar{x}_2) + E \right)$$

$\downarrow$   
 $E$

Exps: numbers of hours per week spent studying for a particular mathematics class were gathered from randomly selected students. One of the questions asked if the respondent usually chooses a seat at the front of the classroom or at the back. Assume that both population distributions are approx. normal. Complete a 90% confidence interval

Front of classroom preferred:

12, 10, 17, 9, 7, 11, 12, 10, 8, 9, 5

Back of classroom preferred:

9, 11, 5, 8, 9, 5, 10, 7, 8, 7, 5, 2, 9, 6, 8

$$\begin{cases} n_1 = 11 \\ \bar{x}_1 = \frac{\sum_{i=1}^n x_i}{n} \approx 9.9 \\ s_1 = \sqrt{\frac{\sum (x_i - \bar{x})^2}{n-1}} = 3.86 \end{cases}$$

$$(\bar{x}_1 - \bar{x}_2) \approx 3.64$$

$$df = n_1 + n_2 - 2 = 29$$

$$E = c = 90\% = 0.9$$

$$\alpha = 1 - 0.9 = 0.1$$

$$\frac{\alpha}{2} = 0.05$$

$$\frac{\alpha}{2} = 0.05$$

$$\begin{cases} n_2 = 15 \\ \bar{x}_2 \approx 6.27 \\ s_2 \approx 2.46 \end{cases}$$

$$\frac{t_{0.05}}{\sqrt{29}} \leq 1.71$$

$$E = 5.45$$

$$(\bar{x}_1 - \bar{x}_2) - E, (\bar{x}_1 - \bar{x}_2) + E$$

$$(-2.8, 9.0)$$

# Comparing two population means ( $\sigma$ unknown, Paired)

- The samples are random and dependent.  
our data are dependent.
- The paired difference for any pair of data values is given by  $d = x_2 - x_1$
- The mean of the paired differences (point estimate) is given by  $\bar{d} = \frac{\sum d_i}{n}$
- The margin of error  $E = \frac{t_{\alpha/2}}{f} \sqrt{\frac{\sum (d_i - \bar{d})^2}{n-1}}$   
*(for the differences)*
- The std. is  $\frac{\sum (d_i - \bar{d})^2}{n-1}$  . we devide it by  $\sqrt{n}$  further to get the standard error.
- subtract and add the  $E$  for the point estimate  $(\bar{d} - E, \bar{d} + E)$

Ex: The amounts of home utility bills are given for two consecutive months for 5 different homes in a neighborhood. Find a 90% confidence interval for the mean difference in utility bills from March to April.

(4B)

March	April	diff.	Sample size is too small but it is just a practice!
	$x_2$		$x_2 - x_1$
119.75	127.06	7.31	$\rightarrow x_2 - x_1$
62.43	79.09	16.67	
202.39	189.25	-12.89	
97.88	99.64	1.76	
66.07	68.52	2.51	

$$\text{mean of Diff.} = 0.67$$

$$\text{Std of Diff.} = 8.46$$

$$n = 5$$

$$\alpha \text{-Level} = 0.9$$

$$\text{critical value } (t_{\frac{\alpha}{2}}) = 2.731$$

$$\text{Standard Error} = 3.78$$

$$E = 8.07359$$

$$\text{lower bound} = -7.39$$

$$\text{upper bound} = 8.735$$

(16)

"Comparing population proportions"  
(two pop.s)

the criteria

- 1) equal chance for all possible samples to be chosen.
- 2) The samples are independent.
- 3) The conditions for a Binomial distribution are met (~~Part 2~~)
- 4) The sample sizes are large enough to ensure that each distribution has at least 10  $\hat{P}$  and 10  $\hat{1-P}$ .

$$n_1 \hat{P}_1 \geq 10, n_1 (1-\hat{P}_1) \geq 10$$

$$n_2 \hat{P}_2 \geq 10, n_2 (1-\hat{P}_2) \geq 10$$

$$(1-P)$$

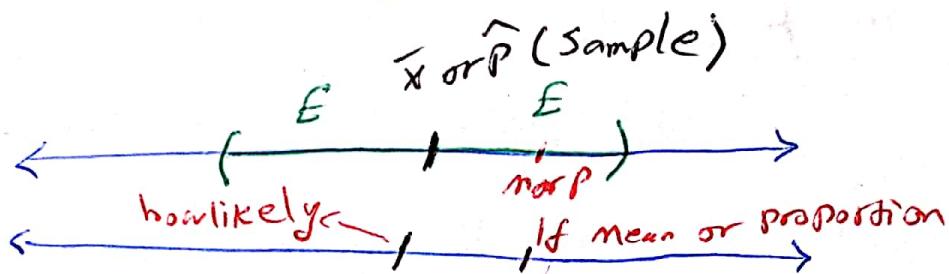
a) The point estimate =  $\hat{P}_1 - \hat{P}_2$

b) margin of Error  $E = Z_{\alpha/2} \sqrt{\frac{\hat{P}_1(1-\hat{P}_1)}{n_1} + \frac{\hat{P}_2(1-\hat{P}_2)}{n_2}}$

c)  $((\hat{P}_1 - \hat{P}_2) - E, (\hat{P}_1 - \hat{P}_2) + E)$

# "Introduction to Hypothesis Testing"

- we use confidence intervals to give us a range in which we feel confident that the true population parameter will fall - many times, however, we want to test a claim that is being made about a parameter.



- In science, a hypothesis is usually a claim about how the world works. In statistics, a hypothesis involves a claim about the numerical value of a population parameter, such as the population mean or proportion

Eg: Based on historical data, the board of Education for one school district believes that the percentage of high school sophomores considering dropping out of school is 10%. A high school counselor in the district claims that this percentage is too high. ~~He~~ collects a random sample of 50 sophomores and finds that 8.5% of them have considered dropping out of school. Perform a hypothesis test with  $\alpha = 0.02$ .

The null hypothesis ( $H_0$ ):  $H_0: \hat{P} = 0.70$

The alternative hypothesis ( $H_A$ ):  $H_A: \hat{P} \stackrel{(a)}{\neq} 0.70$

$H_A$  = in general means  $\hat{P} \neq P_H$  [or < or >]

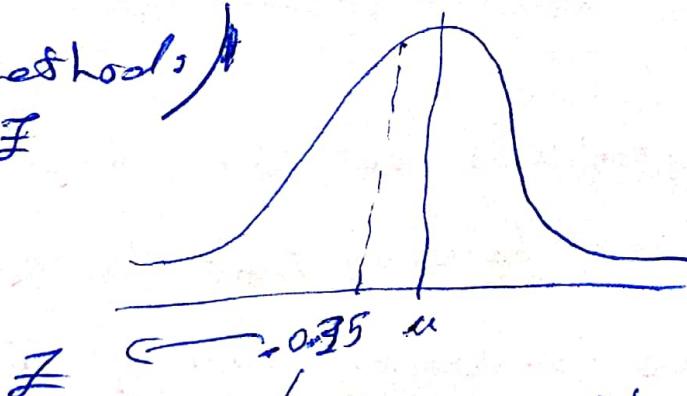
In this case  $\hat{P} < 0.70$

from sample

first steps calculate Z-test based on  $\bar{x}$  or  $\hat{P}$

$$Z = \frac{\hat{P} - P}{\sqrt{\frac{P(1-P)}{n}}} \Rightarrow Z = \frac{0.085 - 0.7}{\sqrt{\frac{0.7(0.9)}{50}}} = -0.3536$$

(The first method)  
using  $Z$

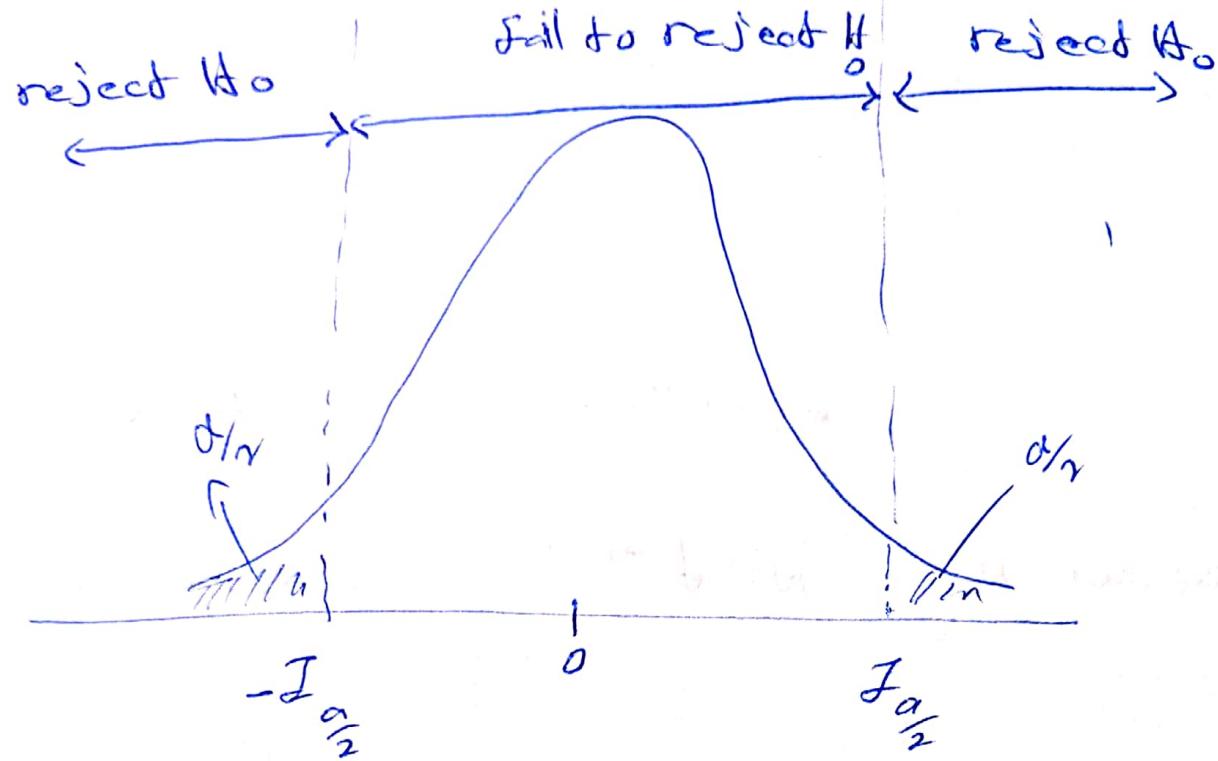


(level of significance) to the left of mean.

$\alpha = 0.02$  based on the level of significant with the null hypothesis, then we reject the null hypothesis

if we use excel: norm.S.Inv(0.02) = [-2.051]

so because our  $Z$ -score is  $-0.3536$ , so we fail to reject the null hypothesis.



The second method (using  $\hat{P}$ )  
p-value

- if  $p < \alpha$ , we reject the null hypothesis

- if  $p > \alpha$ , we fail to reject the  $H_0$

using  $I = -0.3536$  and the fact that we are running a left-tailed test

$$P(I < -0.3536) = 0.3678$$

since  $p = 0.3678$  and  $\alpha = 0.02$  and  $P > \alpha$ , we fail to reject the null hypothesis -

## Errors in hypothesis Testing

		$H_0$ is true	$H_0$ is false
Decision	Reject $H_0$	Type I error	Correct decision
	Fail to reject $H_0$	Correct decision	Type II error

a question, Technically, we can choose the level of significance " $\alpha$ " to be any value . so why wouldn't we choose the level of confidence to be <sup>as</sup> ~~a~~ as possible ?

The probability of making a Type II error is represented by the  $\beta$ . The smaller the probability of making a Type II error, the larger the probability of making a type I error will be.

so you must choose an " $\alpha$ " small enough to control the risk of a Type I error while not making it too small, which would increase the chance of making a Type II error.

## "Hypothesis testing for means ( $\sigma$ known) (right-tailed)"

Ex: The state education department is considering introducing new initiatives to boost the reading levels of fourth graders. The mean reading level of fourth graders in the state over the last 5 years was a Lexile reader measure of 850 L. The developers of a new program claim that their technique will raise the mean reading level of fourth graders. To access the impact of their initiative, the developers were given permission ~~of~~ to implement their ideas in the classrooms. At the end of the pilot study, a simple random sample of  $n = 2000$  fourth graders had a mean reading level of 856 L. It is assumed that the population standard deviation is  $\sigma = 98 L$ . Using a  $\alpha = 0.05$  level of significance, should the findings of the study convince the education department of the validity of the developers' claim?

solution:

step 1) check conditions

- random sample ✓
- $\sigma$  known ✓ ( $\sigma = 98$ )
- $n > 30$  or pop. dist. approx. nor. mat.
- ( $n = 2000$ )

step 2) state the hypotheses

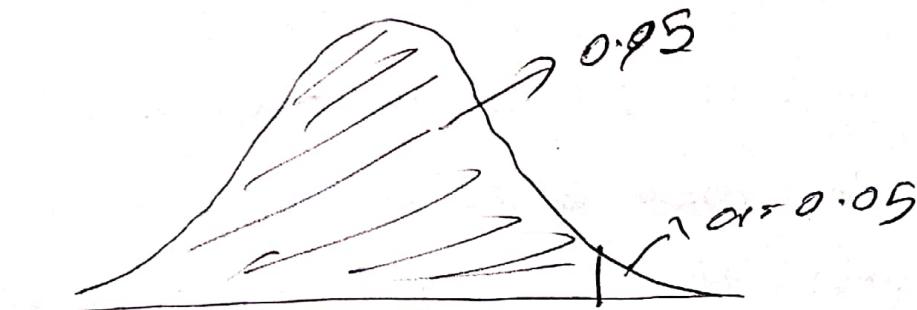
$$\begin{cases} H_0 = \mu = 850 \\ H_A = \mu > 850 \end{cases}$$

step 3) calculate the necessary sample statistics.

$$Z = \frac{\bar{x} - \mu}{\sigma/\sqrt{n}} = \frac{256 - 250}{98/\sqrt{1000}} = 1.936$$

- critical value:  $1-\alpha = 1-0.05 = 0.95$

use<sup>ing</sup> excel;  $\text{NORM.S.INV}(0.95) = 1.645$



$$p\text{-value} = P(Z > 1.936) = 0.0269$$

in excel:  $\text{T-NORM.S.DIST}(1.936)$

step 4) conclusion: with  $p=0.0269 < \alpha=0.05$ , we reject the null hypothesis. with 95% confidence, there is evidence to suggest that the mean Lexile reader measure will increase.

## Hypothesis testing for means ( $\sigma$ -known) 2-tailed

(53)

- A recent study showed that the mean number of children for women in europe is 1.5. A global watch group claims that German women have a mean fertility rate that is different from the mean for all of europe. To test its claim, the group surveyed a simple random sample of 128 German women and found that they had a mean fertility rate of 1.9 children. The population std. is assumed to be 0.8. Is there sufficient evidence to support the claim made by the global watch group at the 90% level of confidence?

$$\sigma = 0.8$$

$$n = 128$$

$$\mu = 1.5 \rightarrow \text{pop. mean}$$

$$\text{Confidence} = 90\% = 0.90$$

$$\bar{x} = 1.9 \rightarrow \text{observed mean (sample mean)}$$

$$\begin{cases} H_0: \mu = 1.5 \\ H_1: \mu \neq 1.5 \end{cases}$$

$$\alpha = 1 - 0.9 = 0.1$$

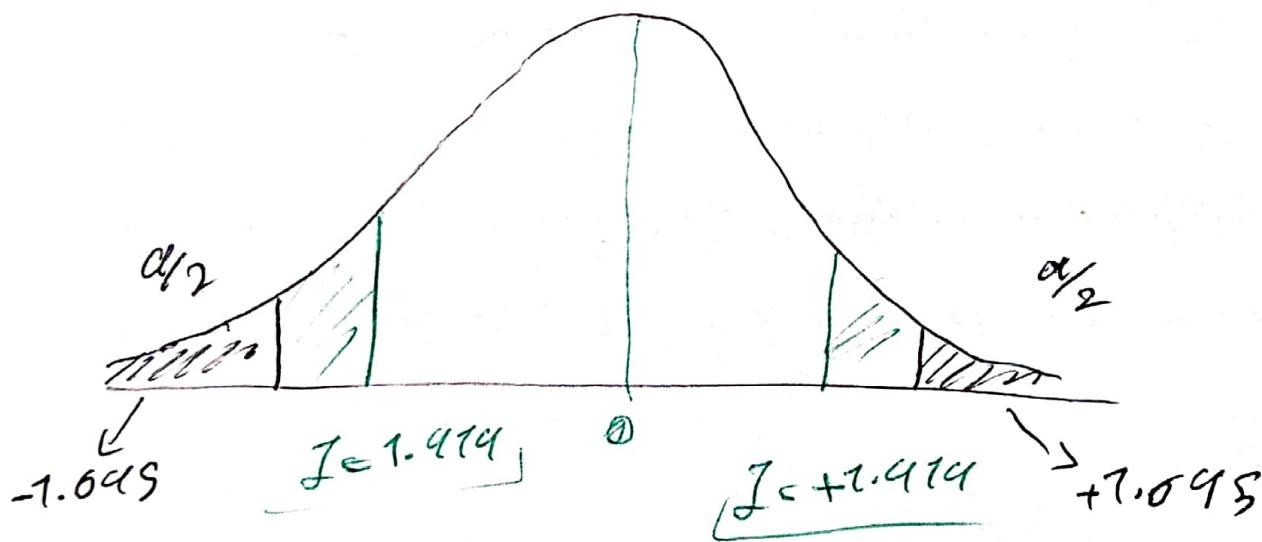
$$Z = \frac{\bar{X} - \mu}{\frac{\sigma}{\sqrt{n}}} = \frac{1.4 - 1.5}{\frac{0.8}{\sqrt{128}}} = \frac{-0.1}{\frac{0.8}{\sqrt{128}}} = \boxed{-1.479}$$

- critical value for rejection region:

$$\text{NORM.S.INV}(0.95) = \pm 1.959$$

$$\text{P-value: } 2 \times P(Z < -1.479) = 0.1573$$

2-tailed



$P = 0.1573 > \alpha = 0.1$  so we fail to reject the null hypothesis.

# Hypothesis testing for population means ( $\sigma$ unknown)

## t-tailed

Ex: nurses in a large teaching hospital have complained for many years that they are overworked and understaffed. The consensus among the nursing staff is that the mean number of patients per nurse each shift is 8.0. The hospital administrators claim that the mean is lower than 8.0. To prove their point to the nursing staff, the administrators gather information from a simple random sample of 19 nurses' shifts. The sample mean is 7.5 patients per nurse with a standard deviation of 1.1 patients per nurse. Test the claim using  $\alpha=0.025$  and assume that the number of patients per nurse has a normal distribution.

step1) check conditions

$\mu = 8.0$ $n = 19$ $s = 1.1$	$H_0: \mu = 8$ $H_a: \mu < 8$	Random sample ✓ $\sigma$ -unknown ✓ $n < 30$ but pop. dist. normal ✓
--------------------------------------	----------------------------------	----------------------------------------------------------------------------

$$\begin{array}{c} \text{step2)} \\ \begin{cases} H_0: \mu = 8 \\ H_a: \mu < 8 \end{cases} \end{array}$$

$$\alpha = 0.025$$

Step 3)

$$t = \frac{\bar{x} - \mu}{\frac{s}{\sqrt{n}}} = \frac{7.5 - 8}{\frac{1.1}{\sqrt{19}}} = -1.981$$

$$df = 19 - 1 = 18$$

critical value for the rejection

~~excel:~~  $T.INV(0.025, 18) = -2.701$

P-value is  $P(t < -1.981) = 0.0375$

$\rightarrow$  excel:  $T.DIST(-1.981, 18)$

Step 4)  $P = 0.0375 > \alpha = 0.025$

Interval  $(\bar{x} \pm t_x \left( \frac{s}{\sqrt{n}} \right))$

we fail to reject the null hypothesis.

with 97.5% confidence.

### "Hypothesis testing for population proportions"

(1-Tailed)

EXP: The local school board has been advertising that 65% of voters favor a tax increase to pay for a new school. A local politician believes that less than 55% of his constituents favor this tax increase. To test his belief, his staff asked a simple random sample of 50

(57)

of his constituents whether they favor the tax increase and 27 said that they would vote in favor of the tax increase. if the politician wishes to be 95% confident in his conclusion, does this information support his belief?

1)  $\begin{cases} \text{- random sample } \checkmark \\ \text{- Binomial distribution } \checkmark \\ n=50, 50(0.65)=32.5 \geq 10 \quad \& \quad n(1-p)=50(0.35)=17.5 \geq 10 \end{cases}$

2)  $\begin{cases} H_0: p=0.65 \\ H_a: p < 0.65 \end{cases}$

$p = 0.65$   
 $n = 50$   
 $\bar{x} = 27$   
 $\hat{p} = \frac{27}{50} = 0.54$   
 $\alpha = 1 - c = 1 - 0.95 = 0.05$

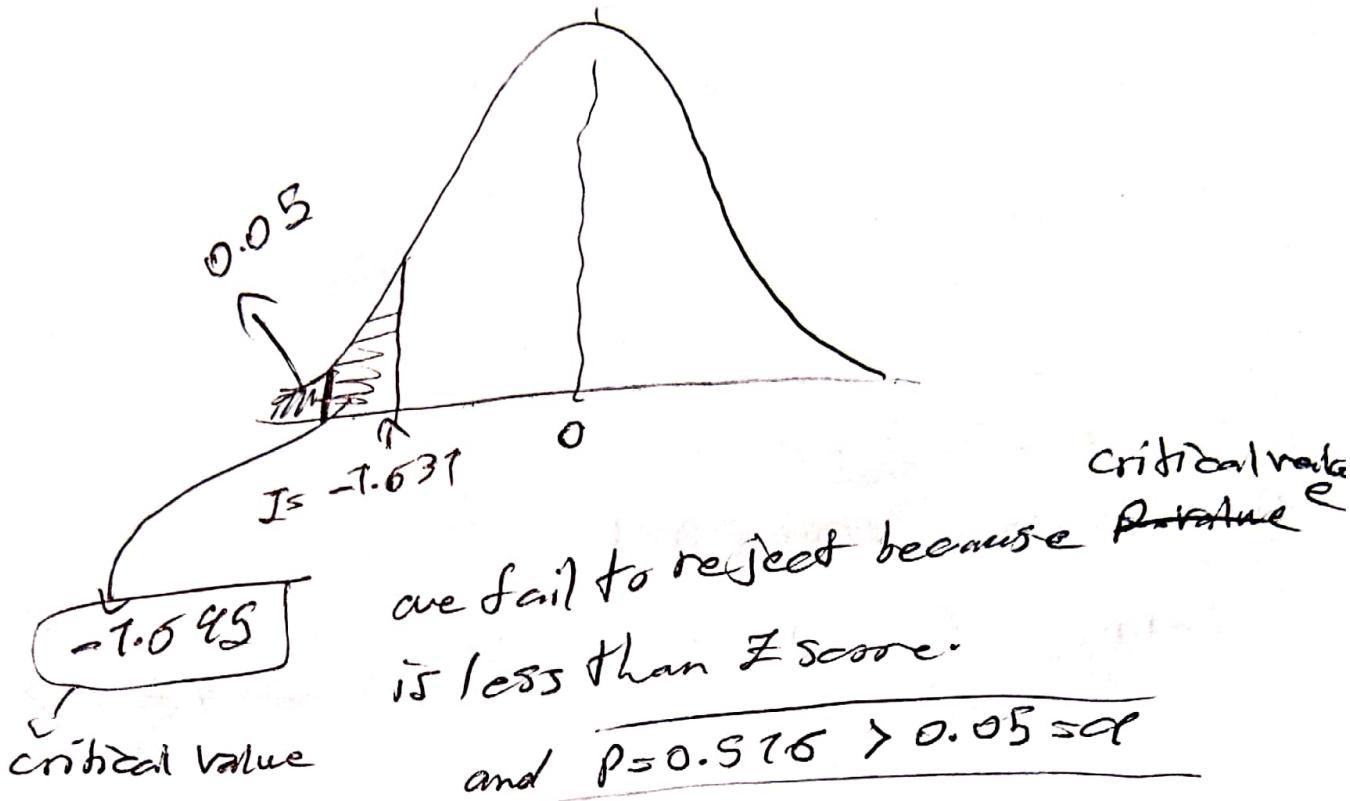
3)  $Z = \frac{\hat{p} - p}{\sqrt{\frac{p(1-p)}{n}}} = \frac{0.54 - 0.65}{\sqrt{\frac{(0.65)(0.35)}{50}}} = -1.645$

critical value for rejection region

excl  $\Rightarrow$  norm. S.JNR(0.05) = -1.645

$P(Z < -1.645) = 0.0512 \rightarrow$  norm. S.DJST (-1.631)

4)



## 5) finding Confidence Intervals:

$$\hat{P} \pm E, \sqrt{\frac{\hat{P}\hat{q}}{n}} \Rightarrow 0.59 \pm 1.645 \left( \sqrt{\frac{0.59(0.41)}{50}} \right)$$

$$\Rightarrow 0.59 \pm 0.116$$

$$\Rightarrow (0.474, 0.656)$$

with 95% confidence, we can claim that the true proportion of voters favor a tax increase falls between 42.4% and 55.6%. since the hypothesized value of 65% falls in our interval, it fails to support failing to reject the null hypothesis.

Hypothesis testing for 2 sample means ( $\sigma$  known)  
 (1-tailed) (59)

$$\text{Test Statistic: } Z = \frac{(\bar{x}_1 - \bar{x}_2) - (\mu_1 - \mu_2)}{\sqrt{\frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}}}$$

$(\mu_1 - \mu_2)$  = hypothesized difference, it's often "0"

$(\bar{x}_1 - \bar{x}_2)$  = 1st sample mean - 2<sup>nd</sup> sample mean.

Expt: A drug manufacturer claims that its new cholesterol drug, when used together with a healthy diet and exercise plan, lowers a patient's total cholesterol level by over 20 points more than simply changing a patient's diet ~~and~~ exercise regimen. To test the claim, a sample of 55 patients with high cholesterol is chosen at random from a cardiologist's list of patients to take the drug in addition to changing their diet and exercise plans. Over the course of three months, this group lowers its total cholesterol level by a mean of 19.7 points. The population std. for this group is 6.8 points. Another 55 patients with high cholesterol lowered their level by a mean of 23.1 points. The population std. for this group is 5.3 points. Test the drug manufacturer's claim.

(50)

claim using a 0.01 level of significance.

- 1)  $\begin{cases} \text{- Both groups are random samples - both are randomly selected.} \\ \text{- Both known } \sigma_1 = 5.8, \sigma_2 = 5.3 \\ \text{- Both } n_1, n_2 \geq 30 \quad n_1 = 55, n_2 = 55 \end{cases}$

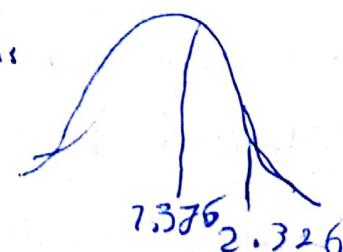
2) 
$$\begin{cases} H_0: \mu_1 - \mu_2 = 20 \\ H_1: \mu_1 - \mu_2 > 20 \end{cases}$$
  
 $\bar{x}_1 = 99.7, \bar{x}_2 = 23.1 \quad \underline{\alpha = 0.01}$   
 $\sigma_1 = 5.8 \quad \sigma_2 = 5.3$   
 $n_1 = 55 \quad n_2 = 55$

3) 
$$Z = \frac{(\bar{x}_1 - \bar{x}_2) - (\mu_1 - \mu_2)}{\sqrt{\frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}}} = \frac{(99.7 - 23.1) - (20)}{\sqrt{\frac{(5.8)^2}{55} + \frac{(5.3)^2}{55}}}$$

$$\Rightarrow Z = 7.376$$

critical value for rejection region:

CV: norm. S. Inv(0.99) = 2.326



P-value:  $P(Z > 7.376) = 0.0849$

$$\boxed{P > \alpha} \Rightarrow (0.08 > 0.01)$$

Sail to reject  $H_0$

# Determining statistical significance for the Pearson Correlation Coefficient.

- Is our linear relationship statistically significant?

Step 1) write your hypothesis:

$$\begin{cases} H_0: \rho = 0: \text{there is no significant linear relationship.} \\ H_1: \rho \neq 0: \text{there is a significant linear relationship.} \end{cases}$$

Step 2) Test statistics  $t_s = \frac{r}{\sqrt{\frac{1-r^2}{n-2}}}$

$$P\text{-value: } T_{\text{DIST}} \cdot 2T(t_s, df) \Rightarrow df = n-2$$

Step 3) Draw your conclusion:

- Reject if  $P < \alpha$ , so there is a significant linear relationship
- fail to reject if  $P > \alpha$ , there is not a significant linear relationship.

$$r_s = \frac{n \sum x_i y_i - (\sum x_i)(\sum y_i)}{\sqrt{n \sum x_i^2 - (\sum x_i)^2} \sqrt{n \sum y_i^2 - (\sum y_i)^2}}$$

→ Excel  
 $\text{correl}(x_i, y_i)$   
 or  
 $\text{pearson}(x_i, y_i)$

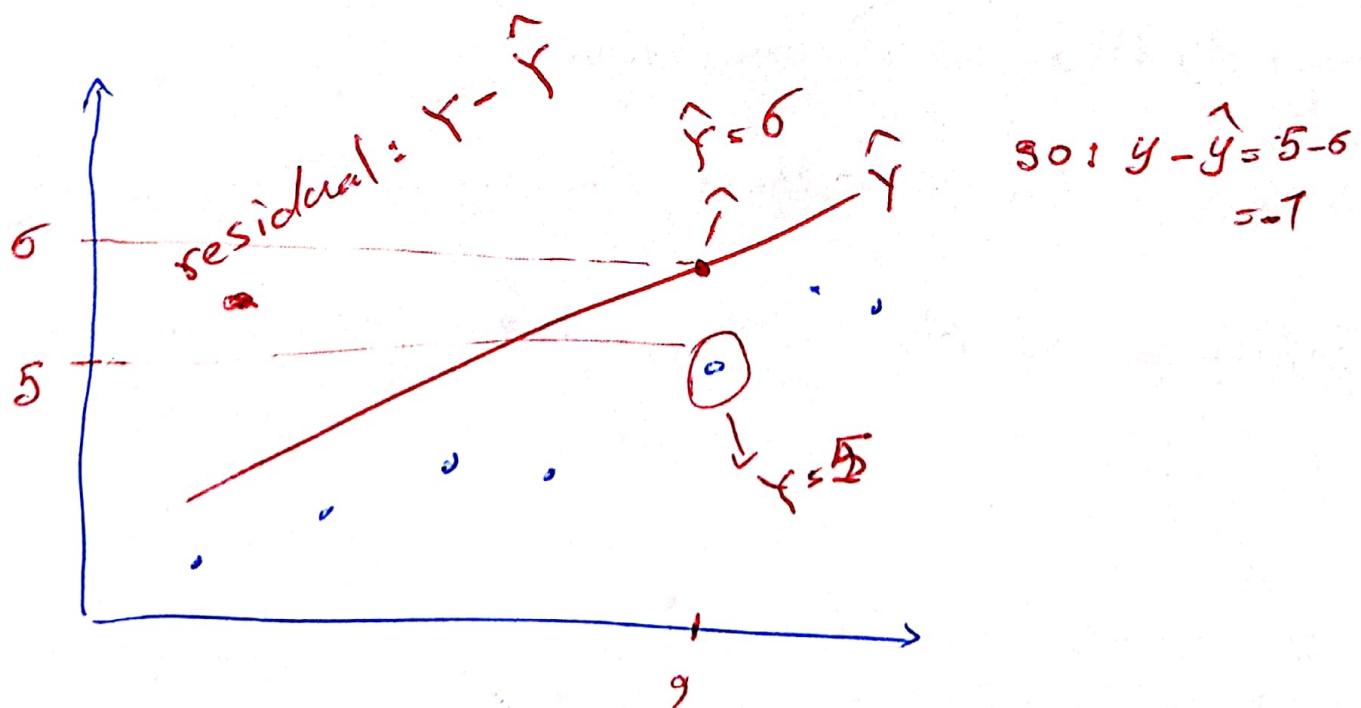
## "The Least Squares Regression Line" (LSRL)

- The LSRL is the line for which the distance between each data point and the predicted data point (located on the LSRL) is the smallest. For each value, we will round to three decimal places.

$$\hat{Y} = b_0 + b_1 x$$

slope:  $b_1 = \frac{n \sum x_i y_i - (\sum x_i)(\sum y_i)}{n \sum x_i^2 - (\sum x_i)^2}$

Intercept:  $b_0 = \frac{\sum y_i}{n} - b_1 \frac{\sum x_i}{n}$



a) with  $p = 0.773 > 0.05 = \alpha$  we fail to reject the null hypothesis.

point estimate ( $\bar{x}$ )

$$\text{confidence interval } (902 - 898) \pm 1.761 \sqrt{\frac{103^2}{18} + \frac{95^2}{18}}$$

$$\Rightarrow 94 \pm 57.23$$

$$(-17.23, 105.23)$$

with  $\alpha = 5\%$  we believe that the tax returns at Smith CPA are between \$17.23 less and \$105.23 more on average than those at Jones and Company CPA. since the hypothesis value of 0 falls in our interval, this supports failing to reject the null hypothesis.

### Hypothesis Testing for 2 Sample ( $\sigma$ unknown) - equal var

$$df = n_1 + n_2 - 2$$

$$t = \frac{(\bar{x}_1 - \bar{x}_2) - (\mu_1 - \mu_2)}{\sqrt{\frac{(n_1-1)s_1^2 + (n_2-1)s_2^2}{n_1 + n_2 - 2}} \cdot \sqrt{\frac{1}{n_1} + \frac{1}{n_2}}}$$

## 1 Hypothesis Testing for 2 sample means (Paired or dependent) (6 unknown)

- mean of pair differences:  $\bar{d} = \frac{\sum d_i}{n}$  where each  $d_i$   $d_i = y_i - x_i$

$$t = \frac{\bar{d} - \mu_d}{S_d / \sqrt{n}} \text{ where } S_d = \sqrt{\frac{\sum (d_i - \bar{d})^2}{n-1}}$$

## 2 Hypothesis testing for 2 sample proportions

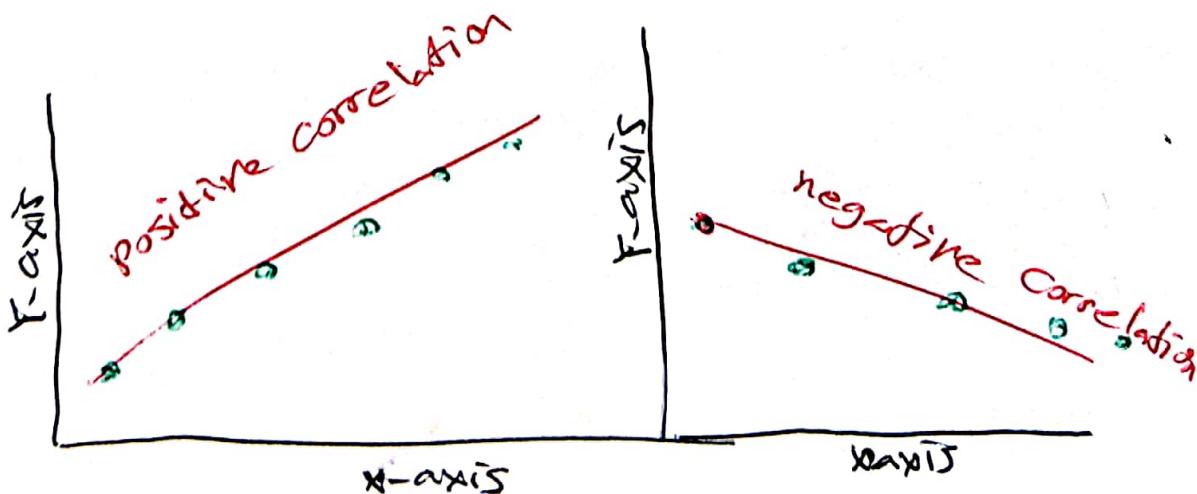
$$n_1 \hat{P}_1 \geq 10, n_1(1-\hat{P}_1) \geq 10, n_2 \hat{P}_2 \geq 10, n_2(1-\hat{P}_2) \geq 10$$

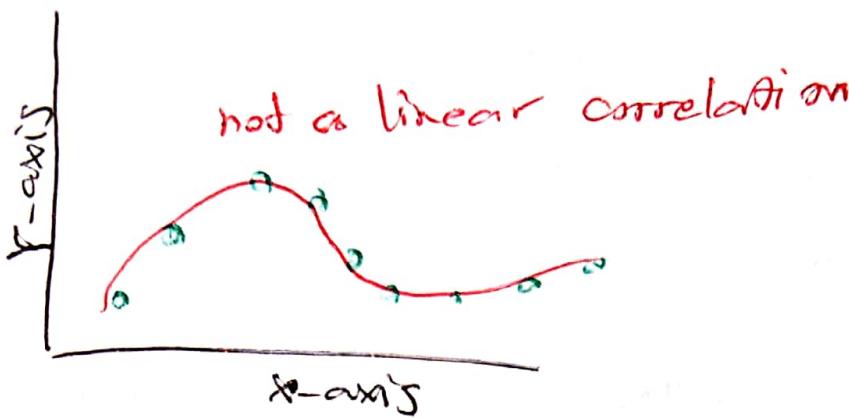
$$J = \frac{(\hat{P}_1 - \hat{P}_2) - (P_1 - P_2)}{\sqrt{\bar{P}(1-\bar{P}) \left( \frac{1}{n_1} + \frac{1}{n_2} \right)}} \text{ where } \bar{P} = \frac{x_1 + x_2}{n_1 + n_2}$$

## Scatter plots and Correlation

- A scatter plot is a graphical display of sets of quantitative values, typically in groups of 2. The two variables are related, as they are recorded from the same subject.

- an explanatory variable, or independent variable, is the variable whose change in value influences a change in the value of another variable.
- A response variable, or the dependent variable, is the variable that has its value change as a result of a change in the value of the explanatory variable.
- ~~one~~ In general, in a scatter plot, as the explanatory variable increases (as one moves to the right on the x-axis), what happens to the response variable (the values plotted on the y-axis)?
  - Does it increase? Then we have a positive correlation.
  - Does it decrease? There is a negative correlation.
  - Is a relationship difficult to determine? Then there is likely no linear correlation.





- The Pearson correlation coefficient,  $P$ , is the parameter that measures the strength of a linear relationship between two quantitative variables in a population. This coefficient for a sample is denoted by  $r$ , it always takes a value between -1 and 1, inclusive ( $-1 \leq r \leq 1$ )

$$r = \frac{n \sum x_i y_i - (\sum x_i)(\sum y_i)}{\sqrt{n \sum x_i^2 - (\sum x_i)^2} \sqrt{n \sum y_i^2 - (\sum y_i)^2}}$$

a) confidence interval

$$(22.7 - 23.1) \pm 2.326 \sqrt{\left( \frac{6.8^2}{55} + \frac{5.3^2}{55} \right)}$$

$$27.6 \pm 2.298$$

$$\boxed{(18.9, 29.3)}$$

At the 0.01 level of significance ( $\alpha$ ), we believe the true difference in post-test cholesterol levels to be between 18.9 and 29.3 points. Since the hypothesis value of 20 falls in our interval, this supports failing to reject  $H_0$ .

~~H<sub>0</sub>~~

~ Hypothesis Testing for 2 Sample means <sup>H</sup>

( $\sigma$  unknown) - unequal variances

$$\text{test statistic: } t = \frac{(\bar{x}_1 - \bar{x}_2) - (\mu_1 - \mu_2)}{\sqrt{\frac{s_1^2}{n_1} + \frac{s_2^2}{n_2}}}$$

Expt suppose that the Smith CPA firm in Chicago claims that its clients receive larger tax refunds, on average, than clients of its competitor, Jones and Company CPA, located on the other side of Chicago. To test the claim

15 clients from the Smith firm are randomly selected and found to have a mean dark refund of \$992 with a std. of \$103. At Jones and Company, a random sample of 18 clients are surveyed and found to have a mean refund of \$898 with a std. of \$95. Test the claim made by the Smith firm at the 0.05 level of significance. Assume that both population distributions are approx. normal.

- 1) both samples are randomly selected.
- 2) - both  $\sigma_1$  and  $\sigma_2$  are unknown.
- 3) both pop. dist. are normal
- 4) - equal variances.

$$2) \begin{cases} H_0: \mu_1 = \mu_2 \\ H_A: \mu_1 > \mu_2 \end{cases} \quad \left| \begin{array}{l} \bar{x}_1 = 992 \\ S_1 = 103 \\ n_1 = 15 \end{array} \right. \quad \left| \begin{array}{l} \bar{x}_2 = 898 \\ S_2 = 95 \\ n_2 = 18 \end{array} \right.$$

$$3) t = \frac{(992 - 898) - 0}{\sqrt{\frac{103^2}{15} + \frac{95^2}{18}}} = 1.260$$

$$\begin{aligned} df_{n_1} &= 15 - 1 = 14 \\ df_{n_2} &= 18 - 1 = 17 \end{aligned}$$

critical value:  $T.INV(0.95) \approx 1.757$

p-value:  $P(t > 1.260) = 0.113$

T.DIST.RT()

- The local school board wants to evaluate the relationship between class size and performance on the state achievement test. It decides to collect data from various schools in the district, and the data from a sample of eight classes are shown in the following table. Each pair of data represents the class size and corresponding average score on the achievement test for one class.

a) Determine if there is a significant linear relationship between class size and average test score at the 0.05 level of significance.

b) If there the relationship is significant, find the LSRL for these data.

a) ?

b) ?

class size	score
15	85.3
17	86.2
18	85
20	82.7
21	87.9
24	78.8
26	79.3
29	72.1

## "Predicting with and Interpreting Values of The LSRL" (69)

Do not Predict if:

- The data do not fall in a linear pattern when graphed on a scatter plot.
- The correlation coefficient is not statistically significant.
- You want to predict a value outside the range of sample data (extrapolation).
- The population is different than that from which the sample data were drawn.

## "Prediction Intervals for Linear regression"

- we would like to create an interval of data in which we feel confident the true value would lie between.

(E)  $\rightarrow$  margin of error

$$\hat{y} \pm \left( t_{\frac{\alpha}{2}} \times S_e \right) \sqrt{1 + \frac{1}{n} + \frac{n(x_0 - \bar{x})^2}{n(\sum x_i^2) - (\sum x_i)^2}}$$

$\hat{y}$   $\rightarrow$  predicted  $y$  (Plug given  $x$  into LSRL equation)

$t_{\frac{\alpha}{2}}$   $\rightarrow$  critical value ( $T.INV.2T(\alpha, df)$  where  $df = n - 2$ )

$s_e$  → standard error of  $\hat{y}$  → the root of the SSE divided by  
our sample the degrees of freedom,  $\sqrt{\text{SSE}/\text{df}}$

SSE → sum of squared (the sum of the squares of the re-  
errors residuals,  $\sum (y_i - \hat{y}_i)^2$ )

$y - \hat{y}$  → residual (subtract the predicted value for each  $x$   
(using LSRL) from the actual value for each  $x$ )

n → number of pairs data.

### "multiple Regression"

- in multiple regression models, we are interested in whether multiple explanatory variables influence the response variable. for instance, we know that a student's age has a significant effect on his or her reading level. Do other factors play a role, as well? we will let:

$x_1$ : represent the age of the student

$x_2$ : represent the years of experience of the child's teacher.

$x_3$ : represent the educational level of the child's Parent(s).

(56)

The final regression equation will look like this:

$$\hat{y} = b_0 + b_1 x_1 + b_2 x_2 + \dots + b_k x_k$$

and our hypothesis will incorporate all variables:

$$H_0: \beta_1 = \beta_2 = \dots = \beta_k = 0$$

$H_a$ : At least one coefficient does not equal 0.