

# BIO 606: BIO STATISTICS

→ MTWT Lectures  
→ Quiz every week Thursday

by Prof. N.G. Prasad

KUNAL VERMA

MS18148

## Lecture - 1

(06-01-2022)

### Data in Biostatistics:

Data variables are properties of a sample population w.r.t. which they can be differentiated. Some of these variables can be "measured" and assigned a number. (continuous  $\in \mathbb{R}$  or discrete  $\in \mathbb{Z}$ )

- ① Nominal variable: variables which can't be measured and must be expressed qualitatively, i.e. we have some categories but no order to the categories. Example - Eye colors, variable, gender,

variable	<u>Eye color</u>	<u>Data</u>	"Categorical Variable"
Categories	→ Black	40	
	→ Blue	0	
	→ Brown	8	
	→ Green	2	

Questions like mean eye color are ill-defined.

- ② Ordinal variable: the order b/w the categories have some particular meaning i.e. the categories are ordered but the distance b/w the categories is unknown. Example - Ranks.

	<u>Rank</u>	<u>Data</u>
diff. categories	1	→ 100 marks
ordered	2	→ 99 marks
	3	
	⋮	
	n	

Based on scores of rank 1 & rank 2, one can't say anything about the score of rank 3. Hence, the distance b/w data is unknown.

③ Interval scale: The categories are ordered & the difference/distance b/w the categories is constant.

### Internal variable (Temp °C)

- 0 --- the zero of this scale is arbitrary.

1               $0^\circ\text{C} \not\Rightarrow$  no temperature.

2

3

4

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

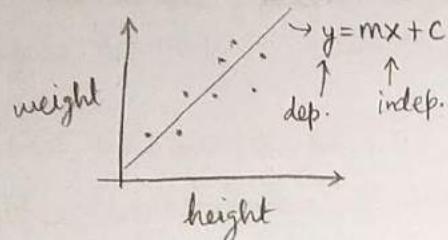
:

:

:

## Independent v/s dependent variables.

For example, we can ask if there's some correlation b/w heights & weights of individuals (statistically).



If we find a relation, then we say that the weight is linearly "dependent" on the height variable.

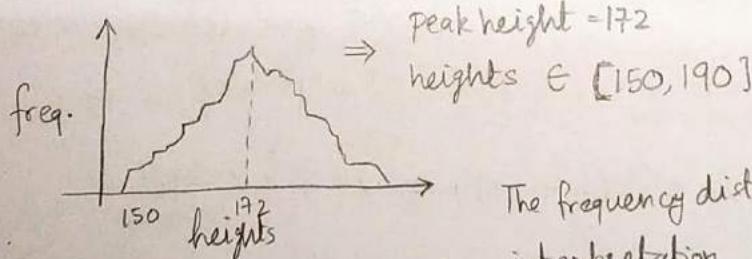
(not necessarily biologically, just according to the statistical model which explains the data)

Do not confuse statistical independence & dependence with biological causation.

## Frequentist view of statistics.

Let's say we measure the heights of a bunch of students ( $x_1, \dots, x_{1800}$ )

To make sense of these 1800 numbers isn't easy. We need to do something with that data to summarize the knowledge. One of the ways is creating graphs.



The frequency distribution of population has an interpretation

frequentist view of statistics → i.e. freq. of population ≡ probability distribution

Probability of a randomly collected person's height to be b/w 150 & 151 cm

= Area under the freq. distrib. graph from  $h=150\text{cm} \rightarrow h=151\text{cm}$ .

## Lecture-2

(10-01-2022)

Heights of IISERM students  $\{x_1, \dots, x_{120}\}$

However, as remarked earlier, raw data in itself is confusing and we need to do something to it to extract out useful info.

One of the ways is to use other numbers which represent the raw data. Those numbers should summarize the core knowledge of the raw data.

One of the ways is to use - Percentiles, Quartiles, Deciles.

(P)

Percentile is a number associated to a particular percentile rank.

$$\text{A percentile rank} = \frac{P}{100} \left( \frac{N+1}{N+1} \right) \quad \begin{matrix} \nearrow \text{no. of observations} \\ \text{"ordered"} \end{matrix}$$

Ex- 157, 159, 160, 161, 164, 165, 166, 167, 168, 169, 175, 181  
 1 2 3 4 5 6 7 8 9 10 11 12 13 → ranks (ordered)

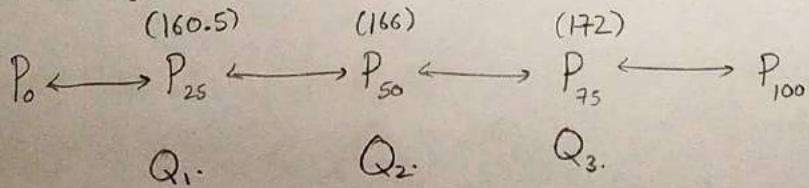
$$30^{\text{th}} \text{ percentile?} \quad P=30 \Rightarrow \text{Percentile Rank (PR)} = \frac{30}{100} (13+1)$$

$$= 0.3 \times 14 = \underline{4.2 \text{ Rank}}$$

$$\text{So, } 30^{\text{th}} \text{ percentile} = \text{Rank 4 data} + 0.2 \times (\text{Rank 5} - \text{Rank 4 data})$$

$$\Rightarrow 30^{\text{th}} \text{ percentile} = 161 + 0.2 (164 - 161) = 161 + 0.6 = \underline{\underline{161.6}}$$

We are usually interested in calculating  $P_{25}, P_{50}, P_{75}$ .



$$\text{For this data, } P_{50} = \text{Rank 7 data} = \underline{\underline{166}}, \quad P_{75} = \text{Rank 10.5 data} = 169 + 3 = \underline{\underline{172}}$$

$$P_{25} = \text{Rank 3.5 data} = 160 + 0.5 \times 1 = \underline{\underline{160.5}}$$

$$\text{So, } P_n = \text{rank} \lfloor PR \rfloor \text{ data} + \text{frac}(PR) \cdot [\text{rank}(\lfloor PR \rfloor + 1) \text{ data} - \text{rank}(\lfloor PR \rfloor) \text{ data}]$$

where  $PR \equiv \frac{n}{100} (N+1)$   $\rightarrow$  no. of observations.

Why is this  $N+1$ ? What is  $P_0$  &  $P_{100}$ ?

Another representative of data are Central Tendency measures and Dispersion measures.

The central tendency & dispersion measures give an idea of where the distribution is centred and how much is it spread around its central tendency.

### CENTRAL TENDENCY.

A central tendency representing the raw data must be a number which is closest to all the numbers simultaneously in the given data.

Say we have a set of nos.  $\{x_1, x_2, \dots, x_n\}$ . Say  $y$  is a good representative for central tendency. One of the ways of saying this is to do a sum over absolute differences.

$$|x_1 - y| = \Delta x_1$$

$$|x_2 - y| = \Delta x_2$$

$\vdots$

$$|x_n - y| = \Delta x_n$$

$$\text{Summed diff.} = X_\Delta$$

We want  $y$  to be such that  $X_\Delta$  is minimised.

It turns out, after doing the math,  $y = P_{50}$  (or  $Q_{12}$ ) which is also known as the median.

The median gives the minimum sum of abs. differences!

Another way to get a central tendency measure is to find a number  $\mu$  which minimizes the sum of squared differences. (instead of abs diff.)

$$(x_1 - \mu)^2 = \Delta x_1^2$$

$$(x_2 - \mu)^2 = \Delta x_2^2$$

:

$$(x_n - \mu)^2 = \Delta x_n^2$$

$$\underline{\text{Summed sq. diff}} = \underline{\Delta x^2}$$

So, now, we want a number  $\mu$  which minimizes the sum of squared differences  $\Delta x^2$ .

(arithmetic)

Consequently, again doing the math, one finds out that  $\mu = \text{mean/avg. of the distribution.} = \sum_{i=1}^N \frac{x_i}{N}$

In summary,

Mean  $\equiv$  minimizes the sum of squared differences.

Median  $\equiv$  minimizes the sum of absolute differences.

Another common measure of the central tendency can be obtained by observing which number appears the most frequently (might even be a qualitative category like eye color), then this is called the mode.

Mode  $\equiv$  the most common observation.

and median

Can't use mean for categorical variables, or for ordinal variables.

Can use mean & median for interval & ratio scale variables.

Mode is generally used for categorical & ordinal variables.

## Lecture - 3

(11-01-2022)

Meaning of  $k^{\text{th}}$  percentile? Three different def's -

- (1) smallest value that is larger than  $k\%$  of values.
- (2) smallest value that is equal to or larger than  $k\%$  of values.
- (3) value associated with  $k^{\text{th}}$  percentile rank.  $\rightarrow$  we use this def". Also valid for non-integer ranks.

Using def<sup>n</sup>(1),  $P_{100}$  is ill defined because  $\nexists$  any value  $>$  all the values.

$\therefore$  (1) can give slightly diff. values of percentile compared to (2) or (3).

Whatever the case,  $P_{50} \stackrel{!}{=} \text{Median}$ . (value for which sum of abs. differences is minimized)

And turns out we get the correct median when we define percentile with  $N+1$ .

Given  $\{x_1, x_2, \dots, x_n\}$

$$\text{Mean} = \frac{x_1 + x_2 + x_3 + \dots + x_n}{n}$$

Median = value at the centre of ordered list =  $P_{50}$   
position

Mean is more susceptible to extreme values (outliers) than median.

1, 2, 3, 4, 5

mean = 3

median = 3

1, 2, 3, 4, 50

mean = 12  $\rightarrow$  most values of distribution

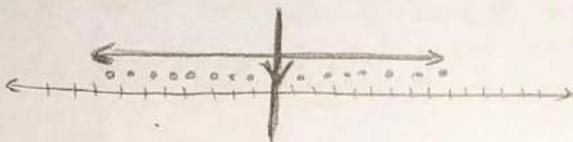
median = 3 aren't even close to 12.

3 is much closer to most of the distribution.

$\Rightarrow$  Mean is a bad representative when distribution has outliers!

Most of the times, we work with a small sample instead of the entire population. On an avg., we see that the sample mean is a better estimate of Population mean  $\mu$  compared to sample median acting as an estimate of Population median. (Efficiency of an estimate)

Say we know where the central tendency of a distribution is anchored. How do we estimate its spread about the centre?



150, 151, 153, 157, 160, 165, 167, 180

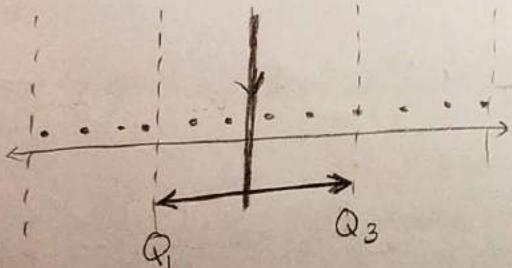
- The simplest estimate of dispersion is the Range.

$$\text{Range} \equiv \max(x_i) - \min(x_i) = 180 - 150 = 30.$$

However, the range is extremely susceptible to outliers & hence not a good representation for dispersion.

- A better alternative is the Inter-Quartile Range (IQR).

$$\text{IQR} \equiv Q_3 - Q_1 = P_{75} - P_{25}.$$



The outliers don't affect  $Q_3 - Q_1$  much since they dependent on intermediate values.

- Median absolute deviation (MAD)

$$150, 151, 153, 157, \mid 160, 165, 167, 180 \\ 158.5 \\ \text{Median} = M$$

To calculate median absolute deviation, we calculate  $|x_i - M| \forall x_i$ , and then find the median of this ordered list.

$$\Rightarrow \text{Median}(\{1.5, 1.5, 5.5, 6.5, \mid 7.5, 8.5, 8.5, 11.5\}) \\ = \underline{\underline{7}}$$

$$|150 - 158.5| = 8.5 = |D_1| \\ |151 - 158.5| = 7.5 = |D_2| \\ \vdots \\ 5.5 = |D_3| \\ 1.5 = |D_4| \\ 1.5 = |D_5| \\ 6.5 = |D_6| \\ 8.5 = |D_7| \\ 11.5 = |D_8|$$

$$\therefore \text{Median absolute dev.} = \text{Median of differences } \{|x_i - M| : \forall x_i\}$$

- Mean squared deviation.

$$150, 151, 153, 157, 160, 165, 167, 180$$

$$\mu \approx 160.4$$

(minimizes the  $\sum$  (squared deviation))

To calculate mean sq. deviation, calculate  $(x_i - \mu)^2 \forall x_i$  and then calculate the mean of this list.

$$(150 - 160.4)^2 = D_1^2 \\ \vdots \\ (180 - 160.4)^2 = D_8^2$$

$$\text{So, mean squared dev.} = \sum_{i=1}^n \frac{(x_i - \mu)^2}{n} = \frac{\sum_{i=1}^n D_i^2}{n} = \text{"Variance"} (\sigma^2) \\ (\text{MSD})$$

However, MSD is in squared units of the data, so we generally take a  $\sqrt$  to interpret the dispersion.

$$\sqrt{\text{Variance}} = \sigma = \text{Standard deviation} = \sqrt{\frac{\sum_{i=1}^n D_i^2}{n}}$$

∴ A distribution can be described by the mean & the std-dev. (or variance) or via median and M.A.D. (median absolute dev.) as the central tendency and dispersion measures respectively.

Insp. One important diff. b/w variance & std-dev. is that variances can simply be added whereas std-dev. can't be. (in what respect though?)

### Population & sample variance.

$$\text{Population variance} = \frac{\sum_{i=1}^n (x_i - \mu)^2}{n} \quad x_i \in \text{Population data} \\ \mu = \text{population mean}$$

However since we don't really have access to population data, we estimate the variance using a small sample.

$$\text{Sample variance} = \frac{\sum_{j=1}^m (x_j - \bar{x})^2}{m-1} \quad \bar{x} = \text{sample mean} \\ x_j \in \text{sample data}$$

$\bar{x}$  is an estimate of  $\mu$  which we obtain from a sample  $\subset$  population.

The reason why we have a  $(m-1)$  in the denominator instead of  $m$  is that sample variance, if defined with just  $m$ , will be an under estimate of population variance. (when  $m \ll n$ )

This is because by def",  $\bar{x}$  minimizes the sum  $\sum_{i=1}^m (x_i - \bar{x})^2$ , and in general,  $\bar{x} \neq \mu \Rightarrow \frac{\sum_{i=1}^m (x_i - \bar{x})^2}{m} \leq \frac{\sum_{i=1}^m (x_i - \mu)^2}{m}$ .  $\Rightarrow$  sample variance is always an underestimate so we use the  $m-1$  correction factor.

CHECK  
THIS

## Lecture-4

(12-01-2022)

Idea: A distribution can be summarized by specifying the central tendency and a measure of dispersion.

Given some population  $\{x_1, x_2, \dots, x_n\}$ , we can calculate  $\mu^2 \sigma^2$  (or  $\text{Var}[x]$ )

$$\sigma^2 = \text{Var}[x] = \sum_{i=1}^n \frac{(x_i - \mu)^2}{n}$$

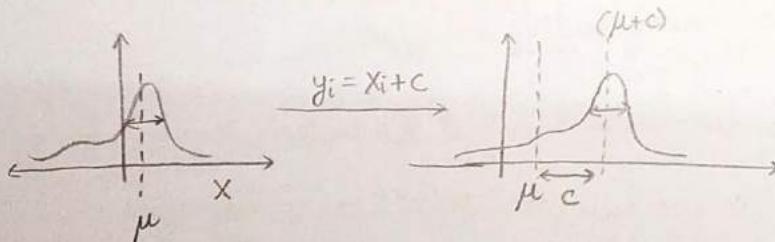
Let's define a new variable  $y_i \equiv x_i + c$ . Calculating  $\text{Var}[y]$

$$\text{Var}[y] = \text{Var}[x+c] = \sum_{i=1}^n \frac{(x_i + c - \tilde{\mu})^2}{n} \quad \text{where } \tilde{\mu} = \sum_{i=1}^n \frac{(x_i + c)}{n} = c + \mu$$

$$\Rightarrow (x_i + c - \tilde{\mu})^2 = (x_i + c - \mu - c)^2 = (x_i - \mu)^2$$

$$\therefore \text{Var}[x+c] = \sum_{i=1}^n \frac{(x_i + c - \tilde{\mu})^2}{n} = \sum_{i=1}^n \frac{(x_i - \mu)^2}{n} = \text{Var}[x]$$

$$\Rightarrow \underline{\text{Var}[x+c] = \text{Var}[x]}, \text{ but } \underline{\mu[x+c] = \mu[x] + c}$$



If we multiply by a scalar, i.e.  $y_i \equiv c x_i$

$$\tilde{\mu} = \mu[cx] = c \sum_{i=1}^n \frac{x_i}{n} = c\mu \Rightarrow \text{Var}[cx] = \sum_{i=1}^n \frac{(cx_i - c\mu)^2}{n} = c^2 \text{Var}[x]$$

$$\Rightarrow \underline{\text{Var}[cx] = c^2 \text{Var}[x]}, \text{ and } \underline{\mu[cx] = c\mu[x]}$$

- So,
- $\text{Var}[c+x] = \text{Var}[x]$ ,  $\mu[x+c] = c + \mu[x]$
  - $\text{Var}[cx] = c^2 \text{Var}[x]$ ,  $\mu[cx] = c\mu[x]$

$$\Rightarrow \boxed{\begin{aligned} \text{Var}[cx+b] &= \text{Var}[cx] = c^2 \text{Var}[x] \\ \mu(cx+b) &= b + \mu(cx) = b + c\mu[x] \end{aligned}}$$

Now, let's see what happens if we add variances of 2 variables.

$$\begin{aligned} \text{Var}[x+y] &= \sum_{i=1}^n \frac{(x_i + y_i - \mu_x - \mu_y)^2}{n} = \sum_{i=1}^n \frac{((x_i - \mu_x) + (y_i - \mu_y))^2}{n} \\ &= \sum_{i=1}^n \left( \frac{(x_i - \mu_x)^2}{n} \right) + \left( \frac{(y_i - \mu_y)^2}{n} \right) + 2 \left( \frac{(x_i - \mu_x)(y_i - \mu_y)}{n} \right) \\ &= \underline{\text{Var}[x] + \text{Var}[y]} \pm 2 \underline{\text{Cov}[x,y]} \end{aligned}$$

where  $\text{Cov}[x,y] = \sum_{i=1}^n \frac{(x_i - \mu_x)(y_i - \mu_y)}{n}$

If  $x$  &  $y$  are independent random variables,  $\text{Cov}[x,y] = 0$ .

and  $\text{Var}[x+y] = \text{Var}[x] + \text{Var}[y]$

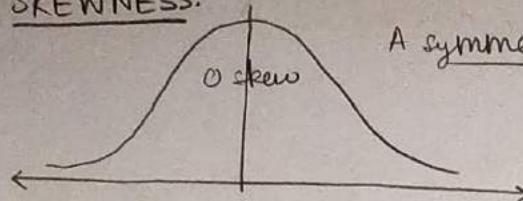
---

x ————— x ————— x ————— x ————— x

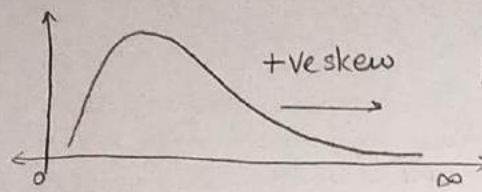
Quite often, describing the distribution via just the central tendency & the dispersion isn't enough. We use two more quantities - { SKEWNESS and KURTOSIS }

Describe the tails of the distribution.

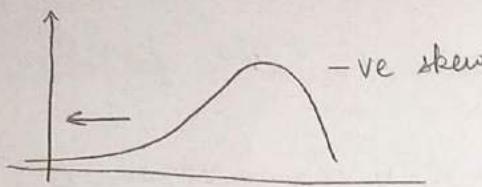
## SKEWNESS.



A symmetric distribution (about mean)  
has zero skewness.

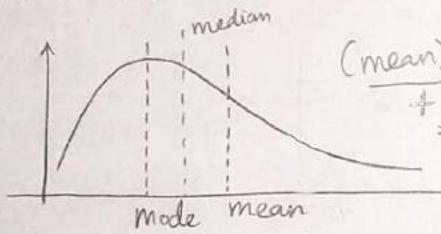


However, if one of the tails extends longer than the other, we have skewness. The graph on the left has +ve skewness.



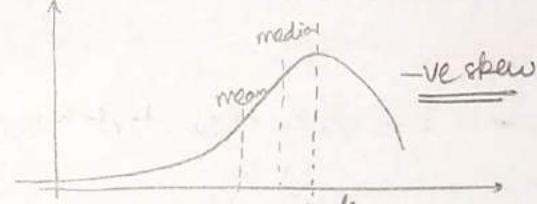
For a symmetric distribution,  $\text{mean} = \text{median} = \text{mode}$ .

For a skewed distribution,



(mean > median)  
+ve skew

(large values pull mean towards it)



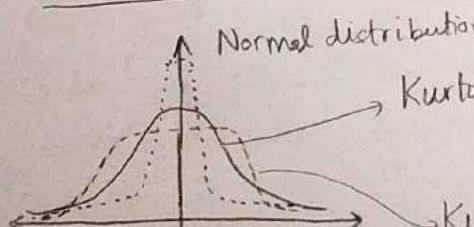
(median > mean)

-ve skew

why here?

(read wiki page on skewness)

## KURTOSIS. (measure of tail thickness)



Kurtosis (normal dist.) = 3

Kurtosis (↙)  $< 3$

Kurtosis (↗)  $> 3$

shorter/thinner tails.

longer/thicker tails.

## Excess kurtosis.

0

+ve

-ve

$K=3 \Rightarrow$  mesokurtic

$K>3 \Rightarrow$  leptokurtic

$K<3 \Rightarrow$  platykurtic

Therefore, we use the following to describe a distribution -

- (1) mean  
(2) variance  
(3) skewness  
(4) kurtosis
- } moments of a distribution. (Read about moments.)

The higher moments are important because there might exist data sets which have the exact same mean & variance, however they still might be drastically different distributions.

(Look up - Anscombe's quartet)

Given a population, the population mean  $\mu$ , the population variance  $\sigma^2$  are called the parameters of the population. (and constant)

If we instead collect a sample, and calculate sample mean  $\bar{X}$  and sample variance  $s^2$  are called a statistic. They are variable in the sense that they may depend upon the sample set we choose.

## Lecture - 5

(13-01-2022)

### BOX - WHISKER PLOT

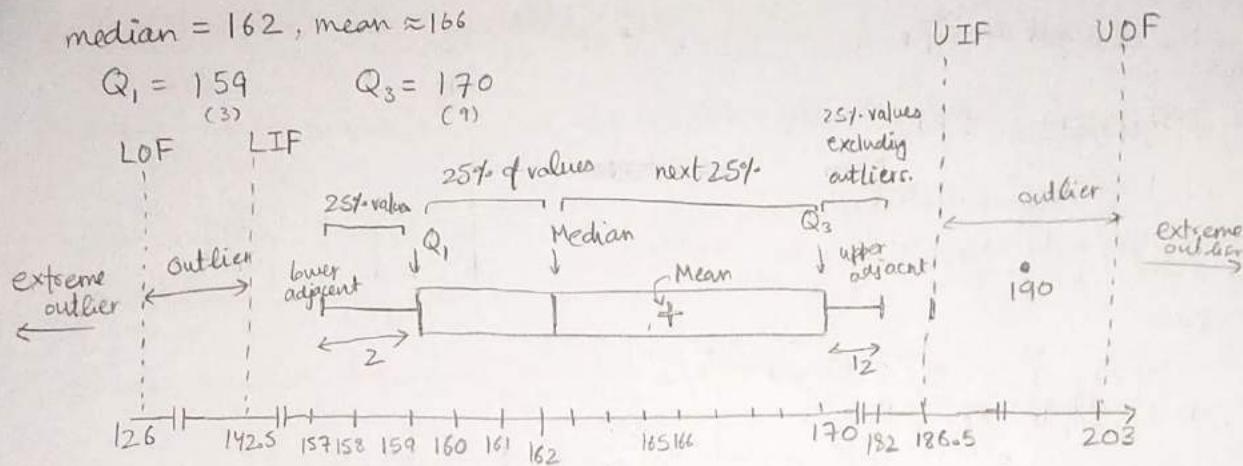
157, 158, 159, 159, 159, 162, 165, 167, 170, 182, 190

To make this plot, we need the median,  $Q_1$ , &  $Q_3$

median = 162, mean  $\approx$  166

$$Q_1 = 159 \quad (3)$$

$$Q_3 = 170 \quad (9)$$



$Q_1 \rightarrow$  lower hinge

$Q_3 \rightarrow$  upper hinge

$Q_3 - Q_1 \rightarrow H_{\text{spread}}$  (I.Q.R.)

Step  $\equiv H_{\text{spread}} \times 1.5$

Upper inner fence  $\equiv$  Upper hinge + Step

Upper outer fence  $\equiv$  Upper hinge + 2 · step.

Lower inner fence  $\equiv$  Lower hinge - step

Lower outer fence  $\equiv$  Lower hinge - 2 · step.

$$H_{\text{spread}} = 11$$

$$\text{step} = 11 \times 1.5 = 16.5$$

$$\text{UIF} = 186.5 = 170 + 16.5$$

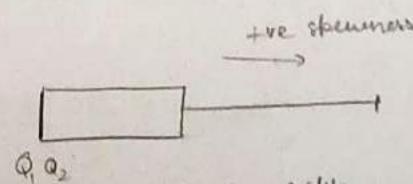
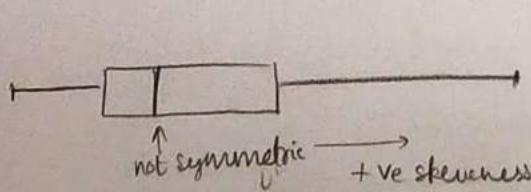
$$\text{UOF} = 203 = 170 + 33$$

$$\text{LIF} = 142.5 = 159 - 16.5$$

$$\text{LOF} = 126 = 159 - 33$$

Upper adjacent = value in the list  $<$  UIF but closest

Lower adjacent = value in the list  $>$  LIF but closest.



Lots of repetitive values in beginning

Interpreting the plot:

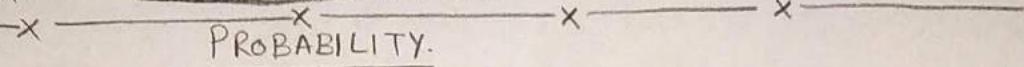
- There are two measures of skewness.

- if  $Q_1$  is closer to  $Q_2$  than  $Q_3$  is, then distribution is positively skewed.

if  $Q_1$  is further to  $Q_2$  than  $Q_3$ , then negatively skewed

- If mean > median  $\Rightarrow$  +ve skew. else -ve.

THESE TWO MEASURES NEED NOT BE CONSISTENT AND ARE DIFFERENTLY DEFINED.

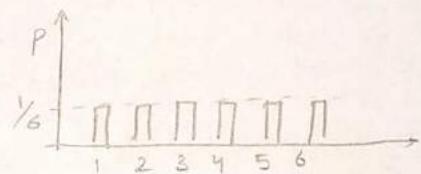
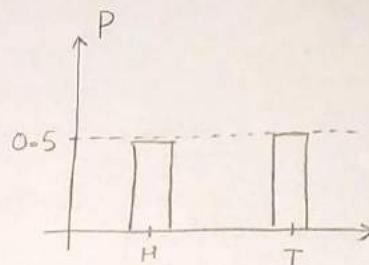


For the classical example of coin toss, the possible outcomes  $\{H, T\}$

$$P(H) = \frac{1}{2} = P(T)$$

For a die,  $\{1, 2, 3, 4, 5, 6\}$  is the sample space

$$P(i) = \frac{1}{6} \quad \forall i \in \text{sample space}$$



Let's now discuss the probabilities

of occurrence of more than one events.

We'll start with the simpler case of "independent" events.

Two independent events -

Probability of getting H on coin AND 6 on die = ?

Since H AND 6 are independent events  
(coin) (die)

$$\Rightarrow P(H \text{ AND } 6) = P(H) \cdot P(6) = \frac{1}{2} \cdot \frac{1}{6} = \frac{1}{12}$$

$$\text{Similarly } P(1 \text{ AND } 1) = \frac{1}{6} \cdot \frac{1}{6} = \frac{1}{36}$$

$$P(\text{even no. on die 1 AND odd no. on die 2}) = \frac{3}{6} \cdot \frac{3}{6} = \frac{1}{4}$$

## Lecture - 6

(17-01-2022)

$$P(\text{event 1 AND event 2}) = P(\text{event 1}) \cdot P(\text{event 2}) \quad \text{if 1 & 2 are independent.}$$

AND → multiplicative.

$$P(\text{event 1 OR event 2}) = P(\text{event 1}) + P(\text{event 2}) - P(\text{both})$$

$$\text{example- } P(H \text{ or even}) = P(H) + P(\text{even}) - P(\text{both}) = \frac{1}{2} + \frac{1}{2} - \frac{1}{2} \cdot \frac{1}{2} = \frac{3}{4}$$

$$P(\text{both}) = P(H) \cdot P(\text{even}) \quad \text{since H & even are independent events.}$$

Let's now consider dependent events.

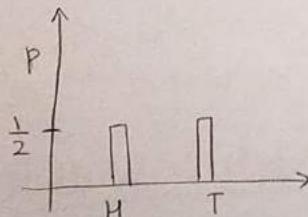
$$\text{We draw an A of spades. } P(A_{\text{spades}}) = \frac{1}{52}$$

$$\text{Now if we draw an A of hearts without putting A spade back, } P(A_2) = \frac{1}{51}$$

The probability of event 2 is dependent on event 1. This is known as conditional probability.  $P(B|A) = \text{prob. of B given A has happened}$

### Uniform distributions

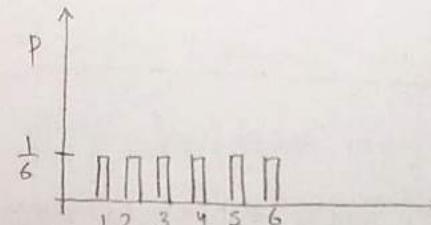
Toss of a fair coin



$$\mu = \frac{a+b}{2}$$

a, b are the endpts.

Throws of a die



$$\sigma^2 = \frac{n^2-1}{12} \quad \text{for the throw of die}$$

$$\mu = \frac{n+1}{2}$$

Technically,  $\neq$  a mean for something like a coin toss, but we can assign #'s to coin heads & tails. for ex- assign 0 to H & 1 to T.

$$\mu = \frac{0+1}{2} = 0.5 \quad \sigma^2 = \frac{0.5^2 + 0.5^2}{2} = 0.25$$

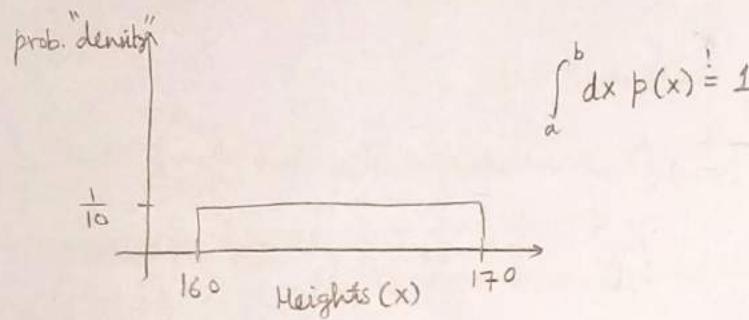
### Lecture - 7

(18-01-2022)

For a uniform <sup>discrete</sup> pdf,  $p(x = \text{any outcome}) = \frac{1}{n}$  where  $n = \# \text{ of possible outcomes}$

$$\mu = \frac{a+b}{2} \quad \text{and} \quad \sigma^2 = \frac{n^2-1}{12}$$

Similarly, we can also have a uniform continuous probability density f's.



$$\mu = \langle x \rangle = \int_a^b dx \times p(x) = \left( \frac{b+a}{2} \right)$$

$$\sigma^2 = \langle x^2 \rangle - \langle x \rangle^2 = \int_a^b dx \times x^2 p(x) - \left[ \int_a^b dx \times p(x) \right]^2 = \frac{(b-a)^2}{12}$$

### Bernoulli Trials / Binomial distribution.

Given a fixed no. of trials with a success & failure outcome which has a fixed prob of success / failure ( $p/q$ )

$$p(x = r \text{ successes}) = {}^n C_r p^r q^{n-r}$$

$${}^n C_r \equiv \binom{n}{r} = \frac{n!}{r!(n-r)!}$$

$n = \text{no. of trials}$

$r = \text{no. of success}$

$p = \text{prob. of success}$

$q = \text{prob. of failure} = 1-p$

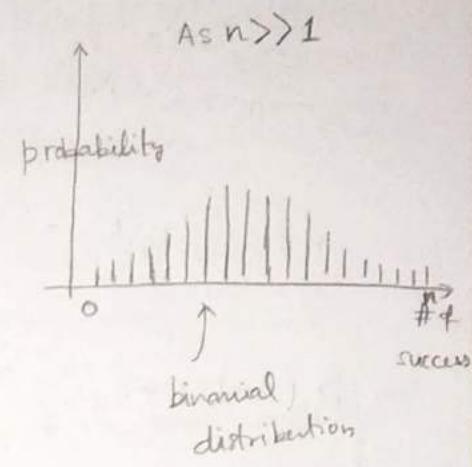
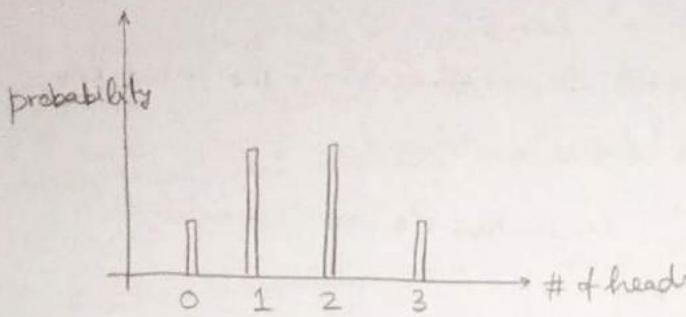
For  $n=3$  trials,

$$P(X=0 \text{ heads}) = \binom{3}{0} 0.5^0 0.5^2$$

$$P(X=1 \text{ head}) = \binom{3}{1} 0.5^1 0.5^2$$

$$P(X=2 \text{ heads}) = \binom{3}{2} 0.5^2 0.5^1$$

$$P(X=3 \text{ heads}) = \binom{3}{3} 0.5^3 0.5^0$$



$$\mu = np, \quad \sigma^2 = npq$$

for Binomial distrib<sup>n</sup>

X ————— X ————— X —————

## Lecture - 8

(19-01-2022)

If  $p \neq q$ , the binomial distribution isn't symmetric about the mean anymore. There is a non-zero skewness for  $p \neq q$ .

In general, binomial distribution is just a special case of multinomial distribution.

Multinomial distribution -

$n$  trials (independent)

Each trial has  $k$  different outcomes

$k=2$  case is called binomial distribution.

$p_i \rightarrow$  probability of each outcome  $i=1, 2, 3, \dots, k$ .

Multinomial distribution gives the probability of a particular combination of outcomes of the  $n$  trials.

Example.- Say for a certain cricket match of IND vs SA, we have 4 outcomes

- Win  $\rightarrow p_1 = 0.3$
- Loss  $\rightarrow p_2 = 0.1$
- Draw  $\rightarrow p_3 = 0.2$
- Tie  $\rightarrow p_4 = 0.4$

Now, if we want to answer what's the probability that IND will win 2 matches, lose 1 match, draw 0 & tie 2 matches? We do this via multinomial dist.

$$\text{So, } P(r_1=1, r_2=2, \dots, r_k=k) = \frac{n!}{r_1! r_2! \dots r_k!} p_1^{r_1} p_2^{r_2} \dots p_k^{r_k}$$

where  $r_1, r_2, \dots, r_k$  are the # of times outcomes 1, 2, ...,  $k$  occur respectively.

$$\sum_{i=1}^k r_i = n$$

$p_1, p_2, \dots, p_k$  are the probabilities of each outcome 1, 2, ...,  $k$  respectively

$$\sum_{i=1}^k p_i = 1$$

In the context of above problem,  $r_1=2, r_2=1, r_3=0, r_4=2$

$$r_1+r_2+r_3+r_4=5=n \Rightarrow P = \frac{5!}{2!1!0!2!} (0.3)^2 (0.1)^1 (0.2)^0 (0.4)^2$$

## Poisson Distribution.

Possibility of the # of events when they are rare.

Say we have  $k$  number of events happening in an interval (space/time), given that they occur at a fixed constant rate on average.

i.e. the average rate of an event is a fixed probability for a given interval, & scales with the size of the interval. These events are independent. And more than one event can't happen at the same instance.

⇒ Poisson distribution gives the probability of such  $X$  number of events happening in a given interval-

$$p(X) = \frac{e^{-\mu} \mu^X}{X!}$$

$\mu$  = mean rate of the event happening  
 $X$  = # of events of interest.

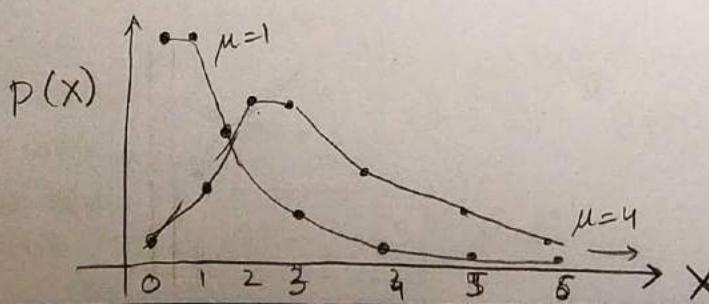
For example, if we encounter 2 delivery guys at the gate every 10 mins, then what's the prob. that we see 4 zomato guys in that 10 min interval?

$$\mu = 2, X = 4 \Rightarrow p(X=4) = \frac{e^{-2} 2^4}{4!}$$

↑                      ↑  
 avg rate              # of guys of interest  
 2 guys/10 mins

For a Poisson distribution,  $\boxed{\mu = 0.2}$

The skewness ↓ as  $\mu \uparrow$ .



## MEAN OF A BINOMIAL DISTRIBUTION.

Say, the total no. of trials is  $N$ . The probability of event  $E_1$ , which we label as success is  $p_1$  & probability of  $E_2$  (failure) =  $p_2 = 1 - p_1$ .

The # of successes in a series of trials is given the notation  $n_1$  & # of failures ( $E_2$ ) is  $n_2$ .

$$\sum_{i=1}^2 p_i = 1 \quad \text{and} \quad \sum_{i=1}^2 n_i = N$$

Now here we have 2 random variables,  $n_1$  &  $n_2$ , which represent the # of successes ( $E_1$ ) & # of failures ( $E_2$ ) respectively.

$$\Rightarrow \text{The average # of successes} = \langle n_1 \rangle = \langle n_1 \rangle = \sum_{n_1=0}^N n_1 \cdot P(X_1=n_1, X_2=n_2)$$

$$\langle n_1 \rangle = \sum_{n_1=0}^N n_1 \cdot {}^N C_{n_1} \cdot p_1^{n_1} \cdot p_2^{N-n_1} = \sum_{n_1=0}^N n_1 \cdot \frac{N!}{n_1! (N-n_1)!} \cdot p_1^{n_1} (1-p_1)^{N-n_1}$$

Relabelling  $n_1 = n$  &  $p_1 = p$  for now,

$$\langle n \rangle = \sum_{n=1}^N n \cdot \frac{N!}{n! (N-n)!} \cdot p^n (1-p)^{N-n} = \sum_{n=1}^N \frac{N!}{(n-1)! (N-n)!} \cdot p^n (1-p)^{N-n}$$

(the  $n=0$  term = 0)

$$= \sum_{n=1}^N n p \cdot \frac{(N-1)!}{(n-1)! (N-n)!} \cdot p^{n-1} (1-p)^{N-n} \quad \left| \begin{array}{l} n-1 \equiv \tilde{n} \\ N-1 \equiv \tilde{N} \end{array} \right.$$

$$= Np \sum_{\tilde{n}=0}^{\tilde{N}} \frac{\tilde{N}!}{\tilde{n}! (\tilde{N}-\tilde{n})!} \cdot p^{\tilde{n}} (1-p)^{\tilde{N}-\tilde{n}} = Np$$

$$\text{mean no. of successes } (E_1) \downarrow \quad (p + (1-p))^{\tilde{N}} = 1^{\tilde{N}} = 1$$

$$\therefore \boxed{\langle n_1 \rangle = N p_1}$$

mean # of failures ( $E_2$ )  
though we often don't talk about this.

Similarly,  $\boxed{\langle n_2 \rangle = N p_2}$

For a multinomial distribution, i.e.  $k$  outcomes (instead of 2 as in binomial) with  $N$  trials, the mean number of times outcome  $E_i$  occurs ( $i=1, 2, 3, \dots, k$ ) is given by -

$$\langle n_i \rangle = N p_i$$

where  $n_i$  is the # of times  $E_i$  occurs  
and  $p_i$  is the probability of each outcome

$$\text{Var}(n_i) = n_i p_i (1-p_i)$$

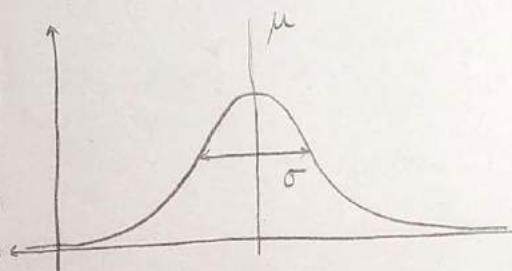

---

### Lecture - 9

(20-04-2022)

Gaussian / Normal distribution / Bell curve.

For  $np \approx nq \geq 5$ , Binomial  $\rightarrow$  Normal.



- Described by  $\mu, \sigma$ .
- CONTINUOUS probability distribution  
 $X \in (-\infty, \infty)$
- Symmetric about  $\mu \Rightarrow$ 
  - Skewness = 0
  - Excess Kurtosis = 0
- Centre is heavy, tails are light.

→ 68-95-99.7 rule ( $\mu \pm \sigma \approx 68\%$  of area,  $\mu \pm 2\sigma \approx 95\%$  of area,  $\mu \pm 3\sigma \approx 99.7\%$  of area)  
3σ rule

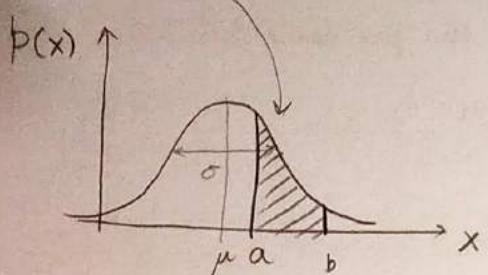
$$p(x) = \frac{1}{\sigma \sqrt{2\pi}} e^{-\frac{1}{2}(\frac{x-\mu}{\sigma})^2}$$

$$\int_{-\infty}^{\infty} dx p(x) = 1$$

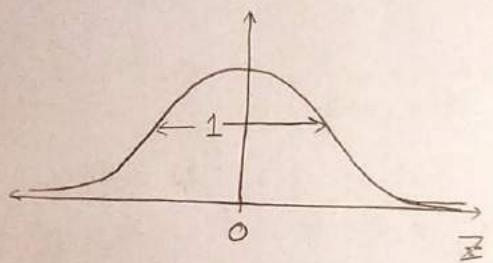
$$\langle x \rangle = \int_{-\infty}^{\infty} dx x p(x)$$

Probability that a normally distributed random variable (say, height) has values b/w  $a \& b$ -

$$P(X=a \rightarrow X=b) = \int_a^b dx \frac{e^{-(x-\mu)^2/2\sigma^2}}{\sigma\sqrt{2\pi}}$$



The standard normal distribution is defined as  $N(0, 1)$  i.e.  $\mu=0$ .  
 $\sigma=1$ .  
 The random variable is called a z-variate.



$$P_{\text{standard}}(z) = \frac{1}{\sqrt{2\pi}} e^{-z^2/2}$$

Any normal distribution can essentially be converted into a standard normal dist.

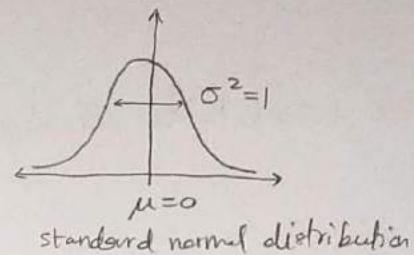
$$X \sim N(\mu, \sigma^2) \longrightarrow Z \sim N(0, 1)$$

This is achieved by  $\left(\frac{x_i - \mu}{\sigma}\right) = z_i \quad \forall i$

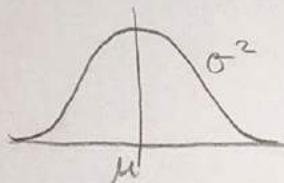
## Lecture-10

(31-1-22)

Standard Normal distribution  $\equiv N(\mu=0, \sigma^2=1)$



Given any normal distribution with  $\mu \neq 0$  &  $\sigma^2 \neq 1$ , we can convert it into a standard normal distribution



Say  $X \sim N(\mu, \sigma^2)$

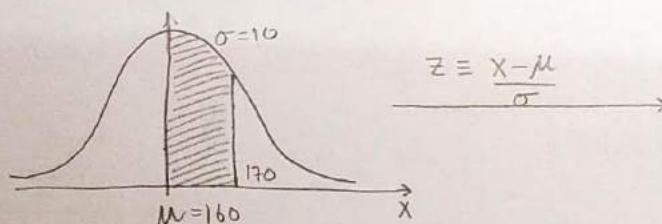
Then  $\frac{X-\mu}{\sigma}$  gives a new variable  $Z$

$$Z \equiv \frac{X-\mu}{\sigma} \text{ s.t. now } Z \sim N(0, 1).$$

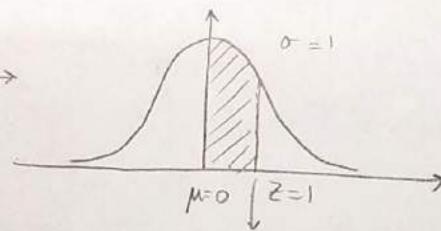
This can be useful for calculating probabilities given  $Z$ -variate tables.

Say heights are distributed as  $N(160, 100)$  i.e.  $\mu = 160$  &  $\sigma = 10$

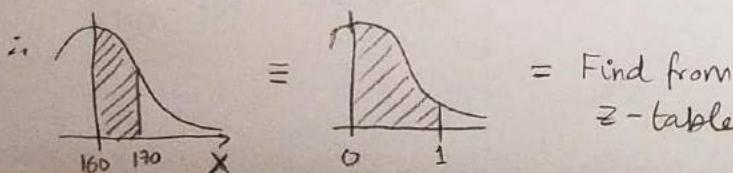
Then what's the probability of finding a height b/w 160 & 170cm.



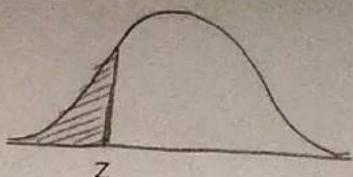
$$Z \equiv \frac{X-\mu}{\sigma}$$



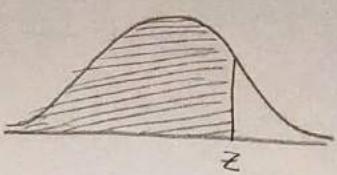
$$Z = \frac{170-\mu}{\sigma} = \frac{170-160}{10} = \frac{10}{10} = 1$$



Negative z-score table calculates

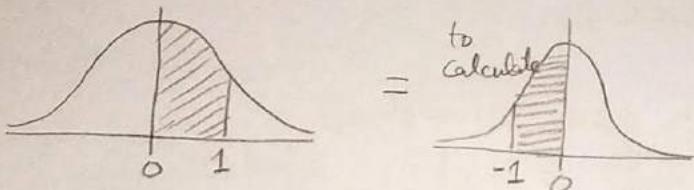


Positive z-score table calculates



But since gaussian is symmetric,

to calculate



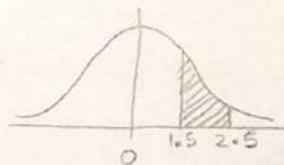
Given a negative z-score table,  $P(160\text{cm} - 170\text{cm}) = Z(0) - Z(-1)$

Given a positive z-score table,  $P(160\text{cm} - 170\text{cm}) = Z(1) - Z(0)$

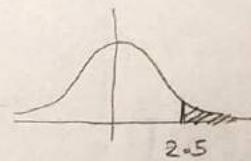
Similarly, let's now find  $P(175\text{cm} - 185\text{cm}) = ?$

$$\mu=160 \Rightarrow \frac{175-160}{10} = 1.5 \quad \frac{185-160}{10} = 2.5$$

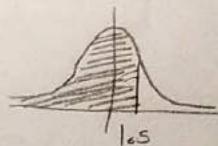
$$\therefore P(175\text{cm} - 185\text{cm}) = \int_{1.5}^{2.5} dz N_z(0,1) =$$



$$P(\text{atleast } 185) = \int_{185}^{\infty} dx N_x(160, 10^2) = \int_{2.5}^{\infty} dz N_z(0,1) =$$



$$P(\text{atmost } 175) = \int_{-\infty}^{175} dx N_x(160, 10^2) = \int_{-\infty}^{1.5} dz N_z(0,1) =$$



## Inferential statistics

Distribution of heights in IISERM students? Possible to take data of all the students, but samples may give some estimates.

We'll have to estimate a measure of

Central tendency (sample mean)  
Dispersion (sample variance)

$\bar{x}$  would be an estimate of  $\mu$ .

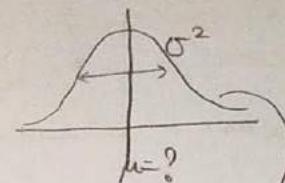
$s^2$  would be an estimate of  $\sigma^2$ .

However, they might or might not be close to  $\mu$  &  $\sigma^2$ .  
(depend on the sample)

Say we have a population whose data  $x$  is distributed normally

$$x \sim N(\mu, \sigma^2)$$

$\mu \rightarrow$  unknown,  $\sigma^2 \rightarrow$  known



take a sample of size  $n$ .

From a sample of population (size  $n$ ), we calculate the sample mean  $\bar{x}$

$$\bar{x} = \sum_{i=1}^n \frac{x_i}{n} \quad (\text{an estimate of popn mean})$$

The quantity which we want to estimate (here  $\mu$ ) is referred to as the estimand.

The rule for calculating the estimand (here  $\bar{x} = \sum_{i=1}^n \frac{x_i}{n}$ ) is referred to as the estimator.

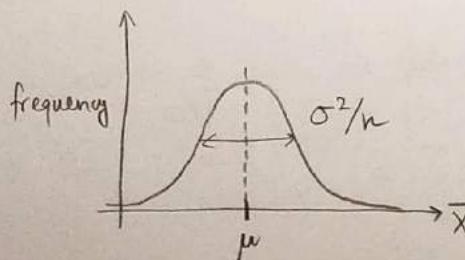
The numerical value of the  $\bar{x}$  is known as the estimate.

Properties of the estimator -

- should be unbiased. (the mean of all  $\bar{x}$  estimates  $\rightarrow \mu$ ) and accurate.
- should be precise (there shouldn't be too much variation in  $\bar{x}$  estimates).
- as  $n \rightarrow$  popn size, then  $\bar{x} \rightarrow \mu$ .

But how are the estimates of  $\bar{x}$  distributed?

If we calculate  $\bar{x}$  infinitely many times & plot it, then  $\bar{x} \sim N(\mu, \frac{\sigma^2}{n})$



where  $n =$  sample size.

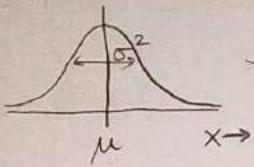
## Lecture-12

(02-02-2022)

Sampling distribution of mean.

$\mu \rightarrow$  unknown

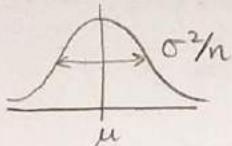
$\sigma^2 \rightarrow$  known



$n \rightarrow$  take samples of size  $n$  & estimate  $\bar{x}$

$\bar{x}_1, \bar{x}_2, \dots, \bar{x}_m$

The distribution of  $\bar{x}$ 's looks like



$$x \sim N(\mu, \sigma^2) \Rightarrow \bar{x} \sim N(\mu, \sigma^2/n)$$

$$\frac{\sigma}{\sqrt{n}} \equiv \text{standard error of the mean} = S_{\bar{x}}$$

We can also normalize the original normal distribution of the data

$x \sim N(\mu, \sigma^2)$  to make it a standard normal distribution

$$z = \frac{x - \mu}{\sigma} \Rightarrow z \sim N(0, 1)$$

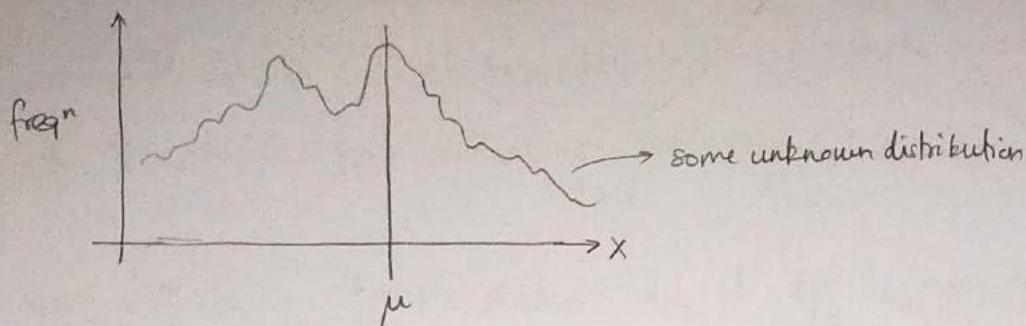
Similarly, we can also renormalize the sample distribution of means & convert it into a standard normal distribution, but this will require the change of variables

$$z = \frac{\bar{x} - \mu}{\sigma/\sqrt{n}} \Rightarrow z \sim N(0, 1)$$

A standard normal distribution  $\int_{-\infty}^{\infty} dz N_z(0, 1) = \int_{-\infty}^{\infty} dz \frac{e^{-x^2/2}}{\sqrt{2\pi}} = 1$

Let's now take the case when Heights ~ unknown distribution.

(x)



Now take samples of size  $n$  ( $n \sim \text{Large}$ ) -  $\bar{x}_1, \bar{x}_2, \bar{x}_3, \dots, \bar{x}_n$  ---

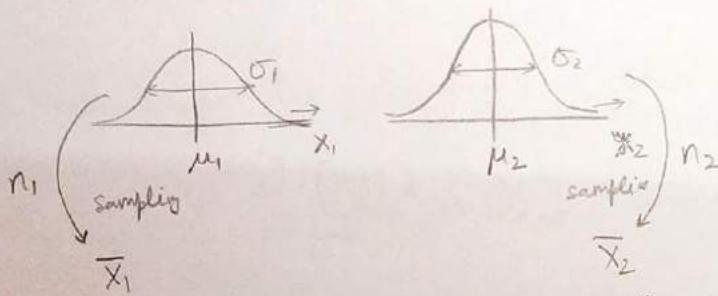
What about the sample distribution of  $\bar{x}$ 's?

As long as your sampling size is large i.e.  $n \gg 1$ , then the  $\bar{x}$  distribution will approximate a normal distribution  $N(\mu, \sigma/\sqrt{n})$  where  $\mu$  &  $\sigma$  are the mean & std. dev. of our unknown distribution.

The above is a consequence of the Central Limit Theorem.

Sampling Distribution of "difference b/w means"

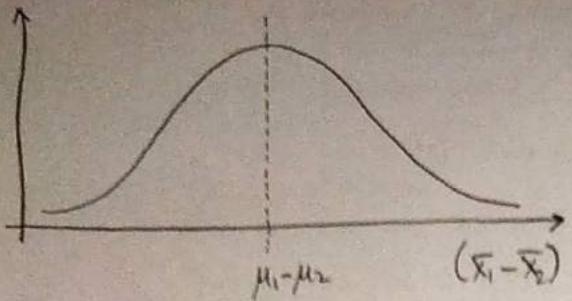
Say  $X_1 \sim N(\mu_1, \sigma_1)$  and  $X_2 \sim N(\mu_2, \sigma_2)$  are two distributions.



diff b/w sample means  $\equiv (\bar{x}_1 - \bar{x}_2)$

Doing this a number of times, we get the distribution

$(\bar{x}_1 - \bar{x}_2)_1, (\bar{x}_1 - \bar{x}_2)_2, (\bar{x}_1 - \bar{x}_2)_3, \dots, (\bar{x}_1 - \bar{x}_2)_{\infty}$



$$\text{mean of } (\bar{x}_1 - \bar{x}_2) = \mu_1 - \mu_2$$

$$\text{Var}(\bar{x}_1 - \bar{x}_2) = \text{Var}(\bar{x}_1) + \text{Var}(\bar{x}_2)$$

Now  $\text{Var}(\bar{x}_1) = \frac{\sigma_1^2}{n_1}$  and  $\text{Var}(\bar{x}_2) = \frac{\sigma_2^2}{n_2}$

$$\text{So, } \text{Var}(\bar{x}_1 - \bar{x}_2) = \left( \frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2} \right) \Rightarrow \text{std-dev of } (\bar{x}_1 - \bar{x}_2) = \sqrt{\frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}}$$

$$\text{So, } \Rightarrow (\bar{x}_1 - \bar{x}_2) \sim N\left(\mu_1 - \mu_2, \sqrt{\frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}}\right)$$

To convert this into a normal distribution, we define

$$z = \frac{(\bar{x}_1 - \bar{x}_2) - (\mu_1 - \mu_2)}{\sqrt{\frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}}} \Rightarrow z \sim N(0, 1)$$

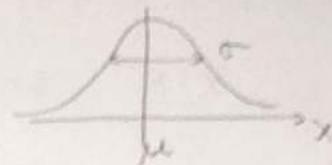
## Lecture - 13

(03-02-2022)

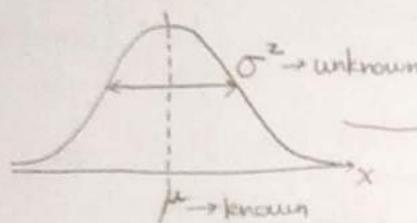
Till now, we've mainly explored how distributions of sample means approach  $\bar{X} \sim N(\mu, \sigma^2/n)$  where  $n$  = sample size.

Let's move onto the next sample statistic and ask, let's say we have a population with  $X \sim N(\mu, \sigma^2)$ .

Say  $\mu \rightarrow$  known, but  $\sigma^2 \rightarrow$  unknown.



Now we need an estimator function to estimate the variance.



$\xrightarrow{n}$  take a sample of size  $n$   
 $(x_1, x_2, \dots, x_n)$

$$\downarrow \text{calculate } s^2 = \frac{\sum_{i=1}^n (x_i - \bar{x})^2}{n-1}$$

Let's say we repeat this sampling process multiple times.

$$s_1^2, s_2^2, s_3^2, \dots, s_m^2, \dots$$

Here, the sampling distribution of  $s^2$  doesn't follow any distribution, but

$$\frac{s^2 \cdot (n-1)}{\sigma^2} \sim \chi_{n-1}^2$$

But what if our initial distribution is some unknown distribution?

Unfortunately, it turns out that for the variance, there is no equivalent central limit theorem. If the original population is NOT normally distributed

then  $\frac{s^2(n-1)}{\sigma^2}$  does NOT follow a  $\chi^2$  distribution.

Let's now take the case  $\sigma$  is unknown,  $\mu$  is known but we'd like to figure out what's the distribution of a certain estimand.

$$\text{So, } X \sim N(\mu, \sigma^2) \Rightarrow \bar{X} \sim N(\mu, \sigma^2/n)$$

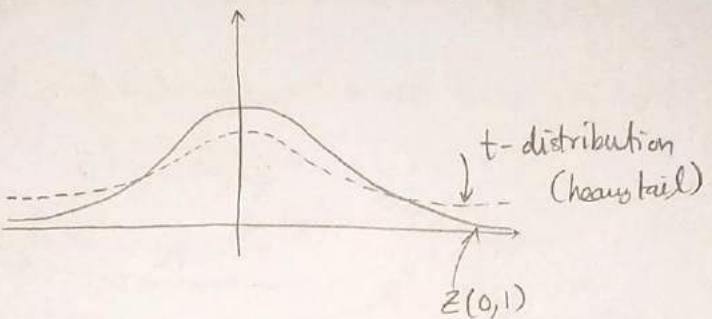
$$\Rightarrow Z = \frac{\bar{X} - \mu}{\sigma/\sqrt{n}} \sim N(0, 1)$$

But since  $\sigma$  is unknown, we replace it by the sample std dev.  $s$

i.e.  $\frac{\bar{X} - \mu}{s/\sqrt{n}}$  which ends up following a student t-distribution

with  $(n-1)$  d.o.f.s

$$\frac{\bar{X} - \mu}{s/\sqrt{n}} \sim t_{n-1}$$



Low d.o.f. ( $n \sim \text{small}$ )  $\Rightarrow$  Heavy tails of t-distribution.

High d.o.f. ( $n \sim \text{large}$ )  $\Rightarrow$  tends to a Z-distribution

### Lecture - 14

(14-02-2022)

#### Student's t-distribution.

Say we start with a normal distribution (standard)  $Z$ , and we divide that by  $\frac{X^2}{\text{d.o.f.}}$ , then

$$\frac{Z}{\sqrt{X^2/\text{d.o.f.}}} \sim \text{t-distribution (t}_{\text{d.o.f.}}\text{)}$$

We know that the distribution  $\frac{\bar{X} - \mu}{\sigma/\sqrt{n}}$  with  $\bar{X}$  as the random variable follows a Z-distribution.

Similarly, we know  $\frac{s^2(n-1)}{\sigma^2}$  follows a  $\chi^2$  distribution with s as the random variable (with  $n-1$  d.o.f.)

$$\frac{\bar{X} - \mu}{\sigma/\sqrt{n}} \sim Z \quad \text{and} \quad \frac{s^2(n-1)}{\sigma^2} \sim \chi_{n-1}^2$$

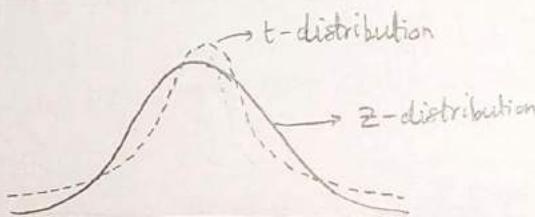
$$\text{So, } \frac{Z}{\sqrt{\left(\frac{\chi^2_{n-1}}{n-1}\right)}} = \frac{\bar{X} - \mu}{\frac{\sigma/\sqrt{n}}{\sqrt{s^2/\sigma^2}}} = \frac{(\bar{X} - \mu)}{\frac{s/\sqrt{n}}{\sigma/\sqrt{n}}} \sim t_{n-1} \rightarrow$$

A t-distrib<sup>n</sup> is obtained  
ONLY if our original  
data is normally  
distributed.

$\Rightarrow$  The t-distribution with  $(n-1)$  d.o.f.s is given by

$$t_{n-1} = \frac{\bar{X} - \mu}{S/\sqrt{n}}$$

$\bar{X}$  = sample mean  
 $S$  = sample std. dev.  
 $n$  = sample size



The t-distribution depends on  
the # of d.o.f.s =  $n-1$

The  $t_{n-1}$  is our new random variable now. The probability density function is a bit complicated but not needed much

$$p(t) = \frac{\Gamma\left(\frac{n}{2}\right)}{\sqrt{(n-1)\pi} \Gamma\left(\frac{n-1}{2}\right)} \cdot \left(1 + \frac{t^2}{n-1}\right)^{-\frac{n-1}{2}}$$

Why does the t-table find values of t- corresponding to some area instead of the area corresponding to some t-value & d.o.f.s?

Now earlier we looked at the random variable "difference of means" and how that followed a Z-distribution when renormalized properly.

Similarly, if we look at the distribution of "difference of means" when  $\sigma$  is unknown, we can see another t-distribution emerging.

Say  $X_1 \sim N(\mu_1, \sigma_1^2)$  &  $X_2 \sim N(\mu_2, \sigma_2^2)$ , and eventually

With  $\sigma_1$  &  $\sigma_2$  unknown,  $t_1 = \frac{\bar{X}_1 - \mu_1}{S_1 / \sqrt{n_1}}$  and  $t_2 = \frac{\bar{X}_2 - \mu_2}{S_2 / \sqrt{n_2}}$

follow t-distribution

Difference b/w means =  $(\bar{X}_1 - \bar{X}_2)$

$$\text{Now } \mu_{\bar{X}_1 - \bar{X}_2} = E[\bar{X}_1 - \bar{X}_2] = E[\bar{X}_1] - E[\bar{X}_2] = \mu_1 - \mu_2$$

$$S_{\bar{X}_1 - \bar{X}_2}^2 = \text{Var}[\bar{X}_1 - \bar{X}_2] = \text{Var}[\bar{X}_1] + \text{Var}[\bar{X}_2]$$

$$S^2 = \frac{S_1^2}{n_1} + \frac{S_2^2}{n_2} \Rightarrow S_{\bar{X}_1 - \bar{X}_2} = \sqrt{\frac{S_1^2}{n_1} + \frac{S_2^2}{n_2}}$$

$$\text{So, } t = \frac{(\bar{X}_1 - \bar{X}_2) - (\mu_1 - \mu_2)}{\sqrt{\frac{S_1^2}{n_1} + \frac{S_2^2}{n_2}}} \sim t_{n_1+n_2-2}$$

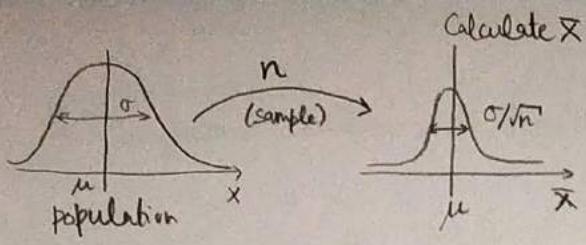
We can also write this as

$$t = \frac{(\bar{X}_1 - \bar{X}_2) - (\mu_1 - \mu_2)}{\sqrt{S_p^2 \left( \frac{1}{n_1} + \frac{1}{n_2} \right)}}$$

$$S_p^2 = \frac{S_1^2(n_1-1) + S_2^2(n_2-1)}{n_1+n_2-2}$$

↗ pooled variance

↑ t-distribution with  $(n_1+n_2-2)$  d.o.f.



However, estimating  $\mu$  from  $\bar{x}$  generally has a lot of error.  $\bar{x}$  is known as a point estimate of  $\mu$ .

### We will now discuss INTERVAL ESTIMATES.

Instead of giving a simple point estimate i.e. a single answer, we can specify an interval in which  $\mu$  will lie (with some probability). This idea is also known as a "Confidence Interval".

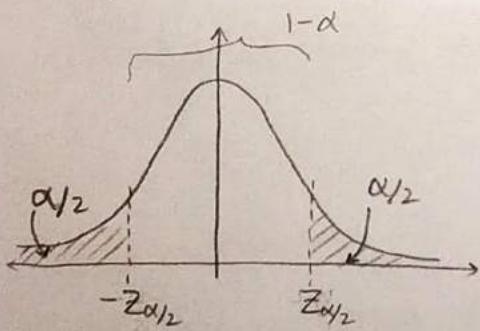
Confidence interval takes the following form -  $(\bar{x} \pm \text{margin of error})$

We can demand the probability with which we want to have  $\mu$  in some interval. If we have an error  $\alpha$ , then the prob. of  $\mu$  lying in the confidence interval =  $1-\alpha$ .

so, a confidence interval  $\equiv (1-\alpha) \cdot 100\% \text{ C.I.}$

### Confidence intervals for Z-distribution:

Say our error is  $\alpha$ . Compute  $\alpha/2$



Calculate Z value ( $Z_{\alpha/2}$ ) corresponding to area of  $\alpha/2$  b/w  $Z_{\alpha/2}$  and  $+\infty$  on  $Z(0,1)$ .

Then

$$P(Z \in [-Z_{\alpha}, Z_{\alpha}]) = 1-\alpha$$

Since the  $\left(\frac{\bar{X}-\mu}{\sigma/\sqrt{n}}\right)$  distribution is also z-distributed, the probability that this number  $\in [-Z_{\alpha/2}, Z_{\alpha/2}]$  is

$$P\left(-Z_{\alpha/2} \leq \frac{\bar{X}-\mu}{\sigma/\sqrt{n}} \leq Z_{\alpha/2}\right) = 1-\alpha$$

$$\Rightarrow P\left(-\frac{\sigma}{\sqrt{n}} Z_{\alpha/2} + \bar{X} \leq \mu \leq \frac{\sigma}{\sqrt{n}} Z_{\alpha/2} - \bar{X}\right)$$

$$= P\left(\bar{X} - \frac{\sigma}{\sqrt{n}} Z_{\alpha/2} \leq \mu \leq \bar{X} + \frac{\sigma}{\sqrt{n}} Z_{\alpha/2}\right) = 1-\alpha$$

So, with a probability of  $(1-\alpha)$ , the mean  $\mu$  lies in the confidence interval of  $\left[\bar{X} - \frac{\sigma}{\sqrt{n}} Z_{\alpha/2}, \bar{X} + \frac{\sigma}{\sqrt{n}} Z_{\alpha/2}\right] !!!$

$$\boxed{\text{Margin of Error} \equiv \frac{\sigma}{\sqrt{n}} Z_{\alpha/2}}$$

The width of the confidence interval remains the same  $\frac{2\sigma}{\sqrt{n}} Z_{\alpha/2}$  independent of the sample chosen. The ranges of confidence intervals only change because of the sample mean  $\bar{X}$ , so the confidence interval may or may not capture  $\mu$ . The probability of  $\mu$ -capture =  $(1-\alpha)$

## Lecture-15

(15-02-2022)

The confidence interval is always calculated for a population parameter.

For the pop<sup>n</sup> mean  $\mu$ , we found out that the confidence interval was of the form -

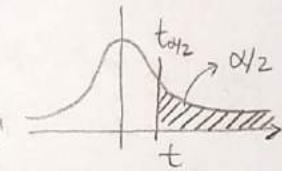
$$\bar{X} \pm Z_{\alpha/2} \frac{\sigma}{\sqrt{n}}$$

The length of C.I. depends upon  $\sigma \rightarrow$  std-dev. of pop<sup>n</sup>  
 $n \rightarrow$  sample size  
 $\alpha \rightarrow$  level of error  
 $\Rightarrow 1-\alpha \rightarrow$  level of confidence.

Higher the  $(1-\alpha)$ , the longer is the length of C.I.

Now let's say one also doesn't know what the pop<sup>n</sup> std dev  $\sigma$  is.  
Then we estimate  $\sigma \approx s$  (sample std dev.), and we replace the z-distribution with the t-distribution.

$$\Rightarrow \text{C.I. if } \sigma \text{ is unknown} = \bar{X} \pm t_{\frac{\alpha}{2}, n-1} \frac{s}{\sqrt{n}}$$



∴ We can also see that the length of the C.I. changes depending on the value of sample std-dev.  $s$ !

Let us now discuss Confidence intervals for the parameter of difference in means b/w 2 distributions  $X_1$  &  $X_2$ .

Assuming  $\sigma$  is known, then the  $(1-\alpha) \cdot 100\%$  C.I. of difference b/w 2 means  $\mu_1 - \mu_2$  will be-

$$\text{C.I. of } (\mu_1 - \mu_2) = (\bar{x}_1 - \bar{x}_2) \pm Z_{\frac{\alpha}{2}} \sqrt{\frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}}$$

(σ-known)

But in the case when  $\sigma$  is unknown, we write the C.I. as

$$\text{C.I. of } (\mu_1 - \mu_2) = (\bar{x}_1 - \bar{x}_2) \pm t_{\frac{\alpha}{2}, n_1+n_2-2} \sqrt{s_p^2 \left( \frac{1}{n_1} + \frac{1}{n_2} \right)}$$

(σ-unknown)

$$\text{where } s_p^2 \equiv \frac{s_1^2(n_1-1) + s_2^2(n_2-1)}{n_1+n_2-2}$$

Since we calculate C.I.s for population parameters, we can do a similar thing for pop<sup>n</sup> variance  $\sigma^2$ .

Say  $X \sim N(\mu, \sigma)$  where  $\sigma$  or  $\sigma^2$  is unknown.

Just like  $\bar{X}$  was a point estimate for  $\mu$ , the sample variance  $s^2$  is a point estimate for  $\sigma^2$ .

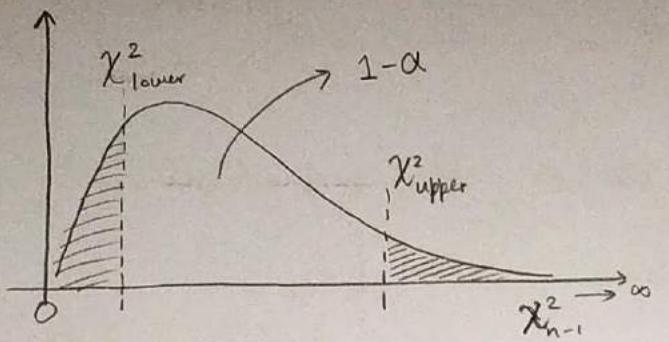
$$s^2 \approx \sigma^2$$

and our parameter of interest was  $\bar{X}$

For the distribution where  $\sigma$  was known, we knew that  $\frac{(\bar{X}-\mu)}{\sigma/\sqrt{n}}$  followed a Z-distribution & we used that to get a C.I. for  $\mu$ .

For a distribution where  $\sigma$  was unknown, we know that

$$\frac{(n-1)s^2}{\sigma^2} \sim \chi^2_{n-1}$$



$X^2_{n-1}$  is a normalized distribution!

We now try to find  $X^2_{\text{lower}}$  s.t. the area before it =  $\alpha/2$ .

similarly, we find a  $X^2_{\text{upper}}$  s.t. the area beyond it =  $\alpha/2$

∴ We can now say that if our error is  $\alpha$ , then the probability of lying in the confidence of  $(1-\alpha) \cdot 100\%$  is

$$P(X^2_{\text{lower}} < X^2_{n-1} < X^2_{\text{upper}}) = 1-\alpha$$

$$\text{so, } P\left(X^2_{\text{lower}} < \frac{(n-1)s^2}{\sigma^2} < X^2_{\text{upper}}\right) = 1-\alpha$$

$$\Rightarrow P\left(\frac{s^2(n-1)}{X^2_{\text{upper}}} < \sigma^2 < \frac{s^2(n-1)}{X^2_{\text{lower}}}\right) = 1-\alpha$$

∴ The value of  $\sigma^2$  lies between  $\left[\frac{s^2(n-1)}{X^2_{\text{upper}}(n-1)}, \frac{s^2(n-1)}{X^2_{\text{lower}}(n-1)}\right]$  with

a probability of  $(1-\alpha)$ .

↳ level of confidence.

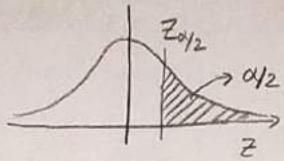
### Example 1

$$X \sim N(\mu, \sigma=20)$$

Sample size  $n = 10$ . Say we take a sample & get  $\bar{X} = 25$ .

Find the 95% C.I. for  $\mu$ .

95% C.I.  $\Rightarrow \alpha = 0.05$  error.



Also,  $\sigma$  is known but  $\mu$  isn't, so the sample will follow a normal distribution.

$$P\left(-Z_{\alpha/2} \leq \frac{\bar{X}-\mu}{\sigma/\sqrt{n}} \leq Z_{\alpha/2}\right) = 1-\alpha$$

$$\begin{aligned} \Rightarrow \text{C.I.} &= \bar{X} \pm Z_{\alpha/2} \frac{\sigma}{\sqrt{n}} = 25 \pm Z_{0.025} \cdot \frac{20}{\sqrt{10}} \\ &= 25 \pm 2 Z_{0.025} \cdot \frac{\sigma}{\sqrt{n}} = 25 \pm 2(1.96) \cdot 3.14 \\ &= 25 \pm 12.31 \end{aligned}$$

### Example 2 Let's say even $\sigma$ is unknown now.

$$X \sim N(\mu, \sigma^2)$$

$$n=10, \bar{X}=25, S=5.$$

Let's now calculate the 95% C.I.

Now we'll use the t-distribution to figure this out.

$$\text{C.I.} = \bar{X} \pm t_{\alpha/2}^{(n-1)} \frac{S}{\sqrt{n}} = 25 \pm t_{0.025}^{d.o.f=9} \cdot \sqrt{\frac{25}{10}}$$

$$\text{Now } t_{0.025} \approx 2.26216 \quad \text{and } \sqrt{\frac{5}{2}} \approx 1.58$$

$$\text{So, C.I. (95\%)} \text{ of } \mu = 25 \pm 3.574$$

Example 3.

Say  $X \sim N(\mu, \sigma=20)$  like before.

$\bar{X}=25$ , and the Margin of Error for 95% C.I. is 5.

Find  $n=?$

$$\text{So, M.o.E.} = Z_{\alpha/2} \frac{\sigma}{\sqrt{n}} \Rightarrow n = \left\lceil \left( \frac{Z_{\alpha/2} \cdot \sigma}{\text{M.o.E.}} \right)^2 \right\rceil = \left\lceil \left( \frac{1.96 \times 20}{5} \right)^2 \right\rceil = \left\lceil 61.46 \right\rceil = 62$$

## Lecture - 16.

(17-02-2022)

### HYPOTHESIS TESTING.

① Let's say we have a data set of heights  $X$ . We take a sample & calculate the sample mean  $\bar{X}$ .

We now "hypothesize" that  $\mu_{\text{height}} = \mu_H$

Assuming this pop<sup>n</sup> is normally distributed, is there any relation of hypothesized mean  $\mu_H$  with  $\bar{X}$  i.e. does your data support your hypothesis?

Model  $\equiv$  Hypothesized distribution (or distribution parameters)

So, we want to calculate the probability that given some data, our model is some XYZ

$$P(\text{Model XYZ} | \text{data})$$

However, most of the times we know  $P(\text{data} | \text{model})$

② Taking another example where we have 2 means  $\bar{X}_1$  &  $\bar{X}_2$ , we calculate the difference b/w means  $\bar{X}_1$  &  $\bar{X}_2$ . If there is a significant difference, we say  $\bar{X}_1$  &  $\bar{X}_2$  don't come from the same distribution. If the difference is not significant, then we say they do come from the same one.

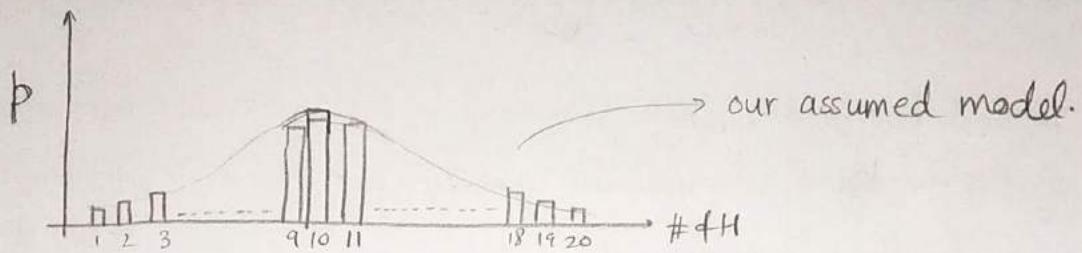
Again, we want to know  $P(\text{Model} | \text{Data})$

↓  
same distrib<sup>n</sup>?      ↓  
 $\bar{X}_1, \bar{X}_2$

Let's now take a concrete example we can work with-

We have a coin. Is the coin fair or not? We only toss it 20 times.

We begin with the assumption that the coin is fair. Then the distrib<sup>n</sup> for the # of heads looks like-



⇒ we carry forward the assumption that  $p=q=0.5 \Rightarrow$  Binomial distrib<sup>n</sup> with  $p=q=0.5$

∴ If our assumption is true, then our most likely event is  $\# \text{ of } H = 10$ .

As our data starts deviating from the expected # of Heads (=10), the  $P(\text{data} | \text{model})$  decreases further, & our model becomes more and more unlikely.

Now say we get 19 Heads out of 20 flips, an extremely rare event if our model is correct (the coin is fair).

We can argue that the rarity of this event  $\Rightarrow$  our assumption was wrong  
& our model was incorrect.

But we could also have disregarded it as a genuine rare event since  $P \neq 0$ .

Where do we draw the line?

We decide that boundary.

So in this experiment, we assume the model of a fair coin (20 flips), take the data & calculate  $\bar{X}$  from the data, and now we ask,

what's the probability that we get such a data outcome given our model?

So we end up calculating  $P(\text{data} | \text{model})$  instead of  $\cancel{P(\text{model} | \text{data})}$

To take the final decision, we have to define 'What is rare?' before rejecting the hypothesis. However, finally, we take a SUBJECTIVE decision to decide the bounds of rarity.

Frequentist statistics  $\longrightarrow$  calculate  $P(\text{data} | \text{model})$

Bayesian statistics  $\longrightarrow$  calculate  $P(\text{model} | \text{data})$

NOTES

$(\bar{X})$

$P(\bar{X})$

- If the sample obtained has a probability of occurrence less than the pre-specified threshold  $P(\bar{X}_0)$  i.e.  $P(\bar{X}) < P(\bar{X}_0)$ , given  $H_0$  is true, then difference between sample & null hypothesis is statistically significant.

We then reject  $H_0$  & accept  $H_A$ .

- Hypothesis Tests based on statistical significance are another way of expressing Confidence Intervals!

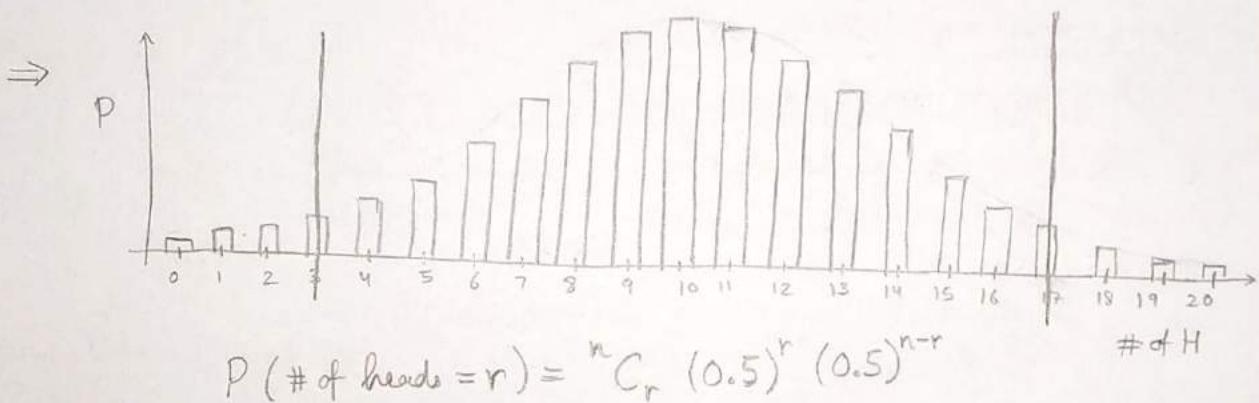
Hypothesis Testing (contd.)

Ideally, we want to calculate  $p(\text{Model} | \text{data})$  where  $\text{model} \sim \text{distrib}^n$   
 i.e. we want to find certain parameters of a distrib<sup>n</sup> given the data.

However, actually, we end up calculating  $p(\text{data} | \text{model})$  in frequentist statistics. We assume that the model is known & try to see if the data fits the model.

Example - Is the coin fair? (we can only flip  $n=20$  times before making a decision.)

We say, assume that the coin is fair  $p=q=0.5$ . This is our assumed / hypothesized model.



We put bars at  $\# \text{ of } H = 3$  and  $17$ , and say that if  $\# \text{ of } H > 17$  or  $< 3$ , then the events are so rare that upto some error, we can say that the coin is not fair.

However,  $\# \text{ of } H = 0, 1, 2, 18, 19, 20$  are still valid possibilities, although extremely unlikely. So, we go forward with the hypothesis with an error of

$$\epsilon = P(0) + P(1) + P(2) + P(18) + P(19) + P(20)$$

If our data says  $3 \leq \# \text{ of Heads} \leq 17$ , then we cannot reject the Hypothesis that the coin is fair. However, we reject our hypothesis if  $\# \text{ of Heads} > 17$  or  $< 3$ , but our decision can still be wrong with an error

$$\varepsilon = \sum_{i \in \{0, 1, 2, 18, 19, 20\}} P(i) \rightarrow \begin{array}{l} \text{extremely low possibility.} \\ \text{Acceptable error.} \end{array}$$

So, three main steps in hypothesis testing are -

- ① Assume a model.
- ② Make the boundaries of acceptability of a Hypothesis. These are decided by choosing what level of error is okay with us.
- ③ Reject or not reject the hypothesis depending on your error & data.

### Z-Test for a single mean.

Say heights  $\sim N(\mu, \sigma=20)$   $\mu$  is unknown.

Make a hypothesis that  $\mu = 170$  cm.

Take a sample of  $n=16$  students, and we get  $\bar{X}=175$  cm.

Now "ideally" we'd have wanted to find what's the probability that  $\mu=170$  cm

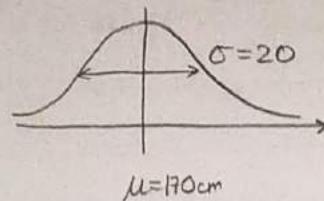
IF  $\bar{X}=175$  cm i.e.  $P(\mu=170 \text{ cm} | \bar{X}=175 \text{ cm}, n=16) = ?$

However, we end up calculating the probability that  $\bar{X}$  <sup>can be</sup> as extreme as 175 cm when  $\mu=170$  cm, i.e.  $P(\bar{X} \text{ is as extreme as } 175 \text{ cm} | \mu=170 \text{ cm})$

We start by taking a null hypothesis  $H_0$ .

$H_0: \mu = 170 \text{ cm}$  is what we hypothesize. (with a normal distribution)

So, hypothesized distribution

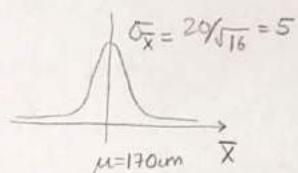
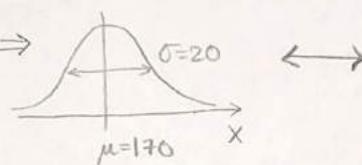


We also have an alternative hypothesis  $H_A$ .

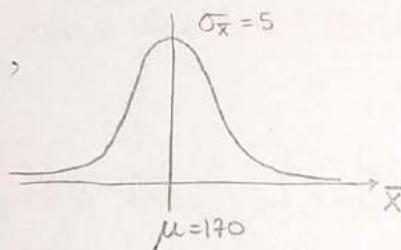
$H_A: \mu \neq 170 \text{ cm}$  is our alternate hypothesis.

We accept the type I error  $\alpha = 0.05$ . (5%)

Assuming that  $\mu = 170 \text{ cm}$  is true,  $\Rightarrow$



Now given,

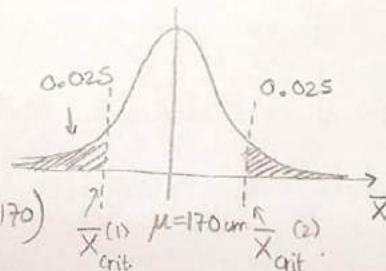


how do we decide the boundaries on this where we say a certain value of  $\bar{x}$  obtained is within our acceptable range & we can't reject  $H_0$ ?

We divide  $\alpha = 0.05$   $\rightarrow 0.025$   
 $\rightarrow 0.025$

(since dist<sup>n</sup> is symmetric)

(since  $\mu \neq 170$  is our  $H_A$ , so it can be  $< 170$  or  $> 170$ )

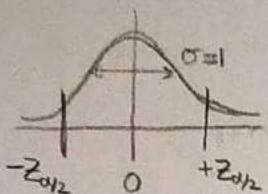


- If our calculated  $\bar{x}$  from the sample  $\in [\bar{x}_{\text{crit.}}^{(1)}, \bar{x}_{\text{crit.}}^{(2)}]$ , then we do not reject the hypothesis that  $\mu = 170 \text{ cm}$ !
- However, if  $\bar{x} > \bar{x}_{\text{crit.}}^{(2)}$  or  $\bar{x} < \bar{x}_{\text{crit.}}^{(1)}$ , then we reject  $H_0$ .

Now we make the switch to Z-distr<sup>b</sup> to figure out  $\bar{X}_{\text{critical}}^{(1,2)}$ .

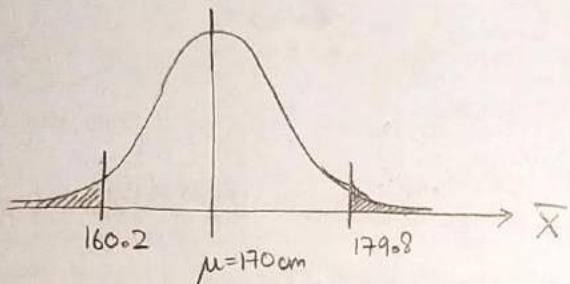
$$\frac{\bar{X} - \mu}{\sigma/\sqrt{n}} = Z \sim N(0, 1) \Rightarrow$$

$$\begin{aligned}\bar{X}_{\text{critical}}^{(1)} &= \mu - \frac{\sigma}{\sqrt{n}} Z_{\alpha/2} \\ \bar{X}_{\text{critical}}^{(2)} &= \mu + \frac{\sigma}{\sqrt{n}} Z_{\alpha/2}\end{aligned}$$



Usually,  $\alpha = 0.05$  so  $Z_{\alpha/2} = Z_{0.025} = 1.96$

For  $\mu = 170 \text{ cm}$  and  $\frac{\sigma}{\sqrt{n}} = 5$ ,  
(hypothesized)  $\bar{X}_{\text{critical}}^{(1)} = 160.2 \text{ cm}$   
 $\bar{X}_{\text{critical}}^{(2)} = 179.8 \text{ cm}$



Now if our  $\bar{X} \in [160.2, 179.8] \Rightarrow H_0$  can't be rejected & is rejected if it is outside the critical extremes.

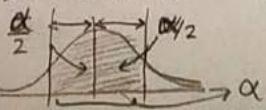
Another way to check if our data, and hence  $\bar{X}_{\text{obs}}$  rejects or doesn't reject  $H_0$  is to calculate  $Z_{\text{obs}}$  corresponding to  $\bar{X}_{\text{obs}}$  and see if  $Z_{\text{obs}} \in [-1.96, +1.96]$   
 $([-Z_{\alpha/2}, +Z_{\alpha/2}])$

Isn't "error" a slightly misleading word for  $\alpha$  here?

Shouldn't the allowed error be this?

i.e. how much error do we allow

that  $\bar{X}$  is some error away from  $\mu$ ?



## Lecture-18

(22-02-2022)

We label " $\alpha$ " as an error because  $100\cdot\alpha\%$  of the times, we might make an error in our decision EVEN IF our Hypothesis is TRUE.

In the whole game of Hypothesis testing, we're comparing  $\bar{X}_{\text{calc}}$  with  $\bar{X}_{\text{critical}}^{(1,2)}$  (or  $Z_{\text{calc}}$  with  $Z_{\text{critical}}^{(1,2)}$ ) and seeing if it's less or more extreme.

In the beginning we stated that in frequentist stats, we calculate

$$P(\text{data being atleast as extreme as critical pts.} \mid \text{model } H_0) = "p"\text{-value}$$

$\hookrightarrow$  when  $H_0$  is true

However, we haven't yet calculated it in our heights  $X \sim N(\mu, \sigma=20)$  example where we took  $H_0: \mu=170$  as the null hypothesis with  $\alpha=0.05$ .

There are two ways of making decisions about  $H_0$  -

- ① Check if  $Z_{\text{calc}}$  is as extreme as  $Z_{\text{critical}}^{(1,2)}$ . (which we did)
- ② Compare  $\alpha$  vs  $p$ .

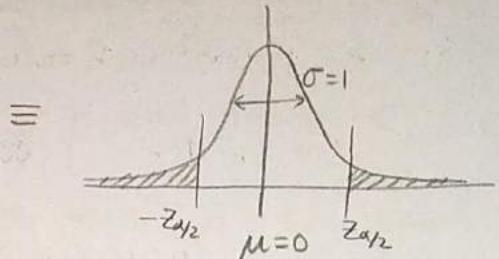
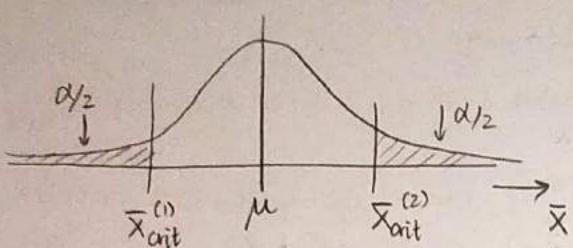
If  $p < \alpha \Rightarrow$  reject  $H_0$ .

If  $p > \alpha \Rightarrow$  do not reject  $H_0$ .

Let us calculate the  $p$ -value for the same example with  $X \sim N(\mu, \sigma=20)$  and  $H_0: \mu=170$  with  $\alpha=0.05$

If we take a sample and  $\bar{X} = 175$ . ( $n=16$ )

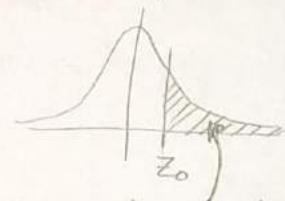
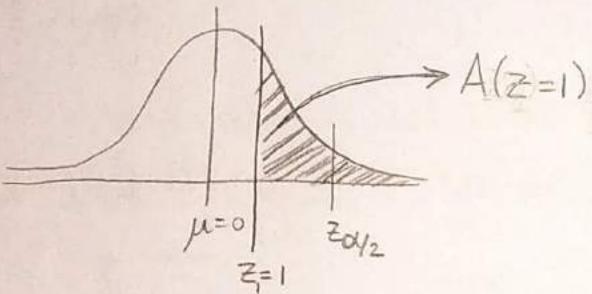
$$Z_{\text{calc}} = \frac{\bar{X} - \mu}{\sigma/\sqrt{n}} = \frac{175 - 170}{20/\sqrt{16}} = \frac{5}{5} = 1$$



$z$ -distribution

$$z \equiv (\bar{x} - \mu) / (\sigma/\sqrt{n})$$

Now,  $Z_{\text{calc}} = +1$ . So if we calculate the area b/w  $Z_{\text{calc}}$  and  $+\infty$



area-value correspondence  
to  $Z_0$ .

$$A(z=+1) = 0.1587$$

p-value is the area associated with  $(Z_{\text{calc}} \text{ to } \infty) \times 2$  when we put our error  $\alpha$  on both the tails. (TWO-TAILED TEST)

$$\text{so, } p = 0.1587 \times 2 = \underline{0.3174}$$

$$\text{Now } p = 0.3174 > \alpha = 0.05 \Rightarrow p > \alpha.$$

The reason why  $p > \alpha$  is because  $Z_{\text{calc}} < Z_{\text{critical}} \Rightarrow A(z_{\text{calc}}) > A(z_{\text{crit}})$

Since  $p > \alpha \Rightarrow$  we can't reject  $H_0$ .

However, suppose that  $\bar{x} = 185 \Rightarrow Z_{\text{calc}} = 3$

Now  $A(Z=3) = 0.0013$

$$\Rightarrow p = 0.0013 \times 2 = 0.0026$$

$$\text{But } p = 0.0026 < \alpha = 0.0050 \Rightarrow p < \alpha$$

Since  $p < \alpha \Rightarrow H_0$  is rejected.

The reason why  $p < \alpha$  is b/c  $Z_{\text{calc}} > Z_{\text{crit}} \Rightarrow A(Z_{\text{calc}}) < A(Z_{\text{crit.}})$

$$\begin{matrix} " & " \\ p/2 & \alpha/2 \end{matrix}$$

If  $Z_{\text{calc}} < 0$ , then take  $|Z_{\text{calc}}|$  and calculate area from  $|Z_{\text{calc}}|$  to  $+\infty$  and compare that to  $\alpha/2$ . (or calculate area b/w  $Z_{\text{calc}}$  &  $-\infty$ )

Ques: why did we distribute  $\alpha/2$  &  $\alpha/2$  in both the tails?

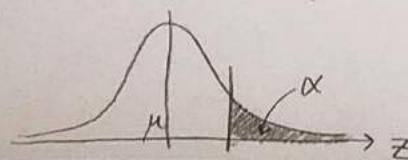
Ans: Not because  $N(\mu, \sigma)$  is symmetric, but because our alternate hypothesis  $H_A$  was defined as  $H_A: \mu \neq 170$  (not  $<$  or  $>$  170)

So our alternate hypothesis works even if  $\mu > 170$  or  $\mu < 170$ .

since we are okay with  $\mu$  being  $< 170$  or  $> 170$ , we put the error in both the tails.

So,  $H_A: \mu \neq 170 \Rightarrow$  Two tailed test!

However, if we had  $H_A: \mu > 170$  cm. In that case, we put all the error in a single tail s.t. area beyond  $Z_{\text{critical}}$  is  $\alpha$ .



$$z = (\bar{x} - \mu) / (\sigma / \sqrt{n})$$

In such a case, if we see  $Z_{\text{calc}} > Z_\alpha$ , then we say  $H_0$  is rejected.

If  $Z_{\text{calc}} < Z_\alpha$ , then  $H_0$  is NOT rejected.

∴ The directionality of alternative hypothesis  $H_A$  carries forward to rejecting or not rejecting  $H_0$ , where we put all the error in one tail.

same procedure if  $H_A: \mu < 170 \text{ cm}$ , with all  $\alpha$  in left-tail.

When all the error  $\alpha$  is sent into a single tail  $\Rightarrow$  ONE-TAIL TEST!

★ For a one-tail test, we won't multiply the area by 2 to find  $p$ .

## Lecture - 19

(23-02-2022)

### Discussion on $\alpha$ -error:

Say heights  $\sim N(\mu, \sigma=20)$  where  $\mu$  is unknown. The size of the sample is known. ( $n=16$ ). And we calculate the  $\bar{X}$ .

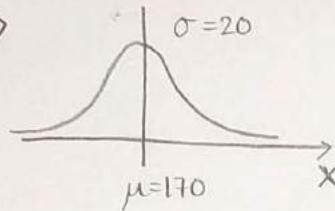
We are interested in testing the hypothesis that  $\mu=170$ . So we start with-

Let  $H_0$  be true where  $H_0: \mu=170$ .

assumption

$$H_0: \mu=170 \Rightarrow$$

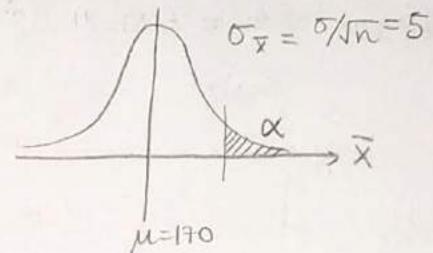
$$H_A: \mu > 170.$$



$$\Rightarrow n=16$$

Distribution if  $H_0$  is true

We are analyzing the one-tail example here



① IF  $H_0$  is really TRUE-

- then we don't reject  $H_0$  if  $\bar{X}$  lies less than the  $\bar{X}_{\text{critical}}$ .
- we reject it if it lies beyond  $\bar{X}_{\text{critical}}$ .

However, these values of  $\bar{X}$  are still a valid outcome of the distribution.

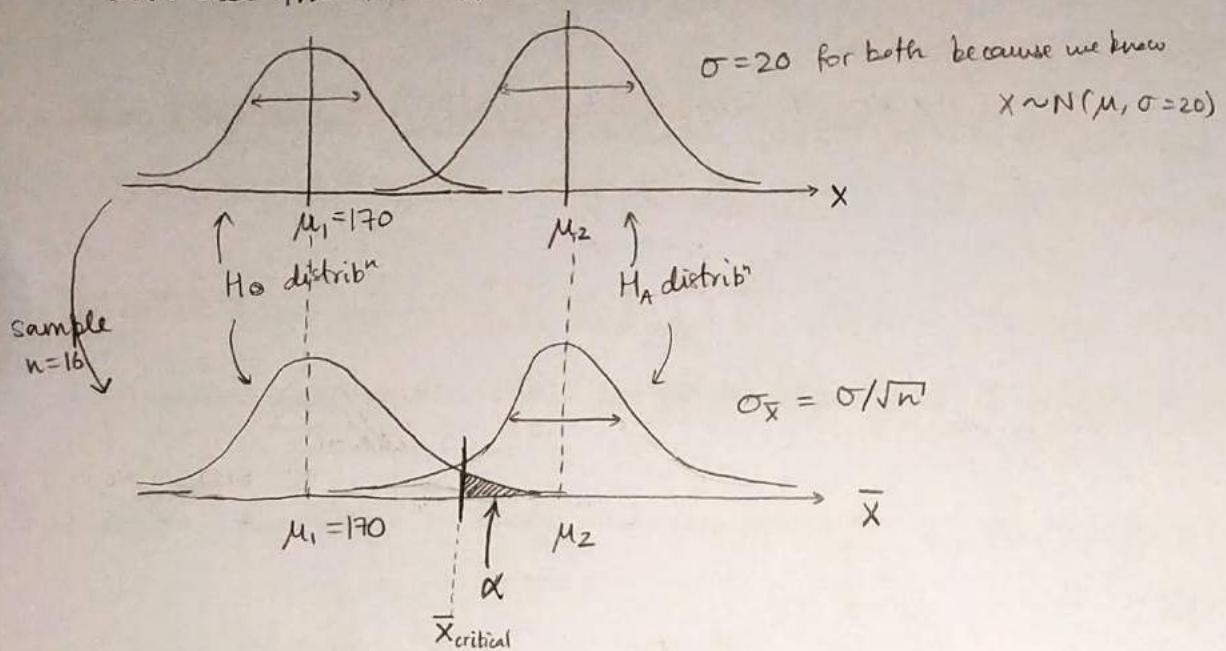
So, with probability  $\alpha$ , we make the error of rejecting  $H_0$  EVEN WHEN IT IS REALLY TRUE. This is known as type I error.

Let us now consider the scenario that in reality,  $H_0$  is false (or that  $\mu \neq 170 \text{ cm}$ ), but we as statisticians do not have that intel-

What do we do in this case?

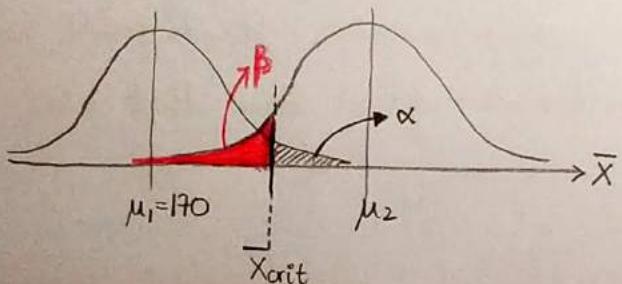
Say the "actual" value of  $\mu$  was  $\mu_2$  (which we are not aware of).

Let's call this the  $H_A$  distrib'.



② Now, when  $H_0$  is really FALSE -

- if  $\bar{X} > \bar{X}_{\text{critical}}$  i.e. lies in the shaded band with area  $\alpha$ , then we reject null hypothesis  $H_0$  because it's too extreme.  
So, we get  $\bar{X} > \bar{X}_{\text{critical}} \Rightarrow H_0$  is rejected, which is good because in reality  $H_0$  was false.
- However, if  $\bar{X} < \bar{X}_{\text{critical}}$  i.e. if it lies outside the  $\alpha$ -area band, then we do not reject  $H_0$ . This is a problem (or an error) however, since in reality  $H_0$  is false, but we still didn't reject it



Since our actual distribution from which the samples are drawn is the

$H_A$  distribution, then the probability that we DO NOT REJECT  $H_0$  DESPITE  $H_0$  being false in reality is -

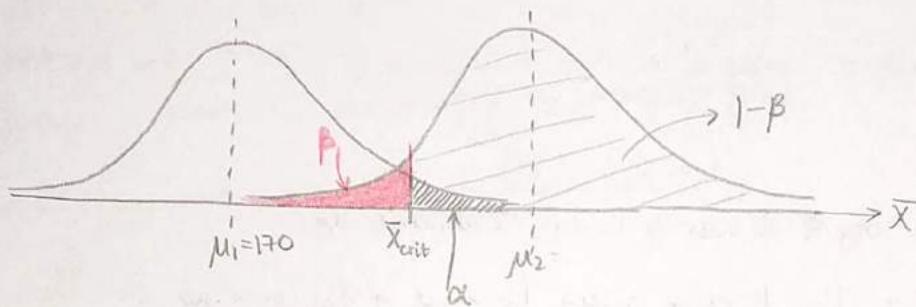
$$\int_{-\infty}^{\bar{x}_{\text{crit}}} d\bar{x} N_{\bar{x}}(\mu_2, \sigma_{\bar{x}}^2 = 5) \equiv \beta$$

And we define it as the  $\beta$ -error.

So,  $H_0: \mu = 170$

$H_A: \mu > 170$

	$H_0$ TRUE	$H_0$ FALSE
Reject $H_0$	$\alpha \cdot 100\%$ times Type I error.	✓ $1 - \beta$
Do not reject $H_0$	✓ $1 - \alpha$	$\beta \cdot 100\%$ times Type-II error



$\beta$  = probability that you don't reject  $H_0$  when  $H_0$  is actually false.

$1 - \beta$  = probability that you reject  $H_0$  when  $H_0$  is actually false.

$$\beta = P(\text{not reject } H_0 \mid H_0 \text{ is actually FALSE})$$

$$\Rightarrow 1 - \beta = P(\text{reject } H_0 \mid H_0 \text{ is actually FALSE})$$

POWER OF A TEST

$$\alpha = P(\text{reject } H_0 \mid H_0 \text{ is actually TRUE})$$

$$\Rightarrow 1 - \alpha = P(\text{not reject } H_0 \mid H_0 \text{ is actually TRUE})$$

Now we come back to the question which always haunted me.

why can't I make  $\alpha$  as small as possible if it's an error?

Because reducing  $\alpha \Rightarrow \bar{X}_{\text{critical}}$  becomes more and more extreme.

The more extreme  $\bar{X}_{\text{critical}}$  becomes, then in cases where  $H_0$  is actually FALSE, the probability that  $H_0$  is NOT REJECTED even when  $\bar{X}$  was extreme as hell will increase.

Geometrically, as  $\alpha \downarrow$ ,  $\bar{X}_{\text{crit}}$  becomes more extreme, and the value of  $\beta \uparrow \Rightarrow (1-\beta) \downarrow$  which is a measure of the "power of the test".

$\therefore$  Decreasing  $\alpha$  comes at a cost of increasing  $\beta$ . (when  $n$  is fixed.)

Ques Is there a way of decreasing  $\alpha$  &  $\beta$  simultaneously?

Ans: Yes! since the  $\beta$  error arises because of the overlap b/w the hypothesized ( $H_0$ ) and the actual distribution ( $H_A$ ), if we increase  $n$ ,

$\sigma_{\bar{X}} = \sigma / \sqrt{n}$  decreases, and so does the overlap. Hence we can decrease both  $\alpha$  &  $\beta$  errors simultaneously  $\Rightarrow$  larger  $(1-\beta) \Rightarrow$  larger power of test.

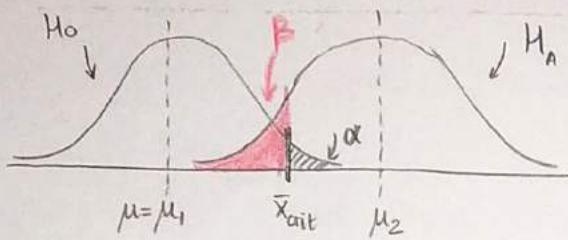
$\therefore$  "Power of the test" is directly dependent on sample size  $n$ .

## Lecture - 20

(24-02-2022)

As discussed earlier, for a one-tailed test with  $H_0: \mu = 170$  & alternative hypothesis as  $\mu > 170$ , we get all the  $\alpha$  in the right tail.

In the case when our  $H_0$  is actually false (although we don't know it), there might exist a different  $\mu = \mu_2$  of the actual distribution.



Corresponding to a given  $n$ , decreasing  $\alpha \Rightarrow$  increasing  $\beta$ .

However, changing  $n$  can give us a handle on both  $\alpha$  &  $\beta$ .

Converting to a z-distribution

$$Z_\alpha, Z_\beta > 0$$

$$Z_\alpha = \frac{\bar{x}_{\text{crit}} - \mu_1}{\sigma/\sqrt{n}}$$

$$Z_\beta = - \frac{(\bar{x}_{\text{crit}} - \mu_2)}{\sigma/\sqrt{n}}$$

$$\text{Now, } Z_\alpha + Z_\beta = \underbrace{\left( \frac{\mu_2 - \mu_1}{\sigma} \right) \sqrt{n}}_{\text{effect size / Cohen's D}} = \text{Non-centrality parameter } \delta.$$

(standardized difference b/w 2 means)

$$\Rightarrow n = \frac{(Z_\alpha + Z_\beta)^2}{[(\mu_2 - \mu_1)/\sigma]^2} \quad \text{For one-tailed test.}$$

HW- do this for a two-tailed test. ( $\alpha \rightarrow \alpha/2$ )

But how do I know  $\mu_2$  beforehand? That itself is the quantity of interest.

The quantity  $(\mu_2 - \mu_1)$  has to be addressed subjectively.

We have to ask, "what to me is a meaningful difference?"

So,  $(\mu_2 - \mu_1)$  is not the exact difference b/w hypothesized & actual mean. It's a meaningful difference, a difference which is significant enough for us to test the hypothesis for.

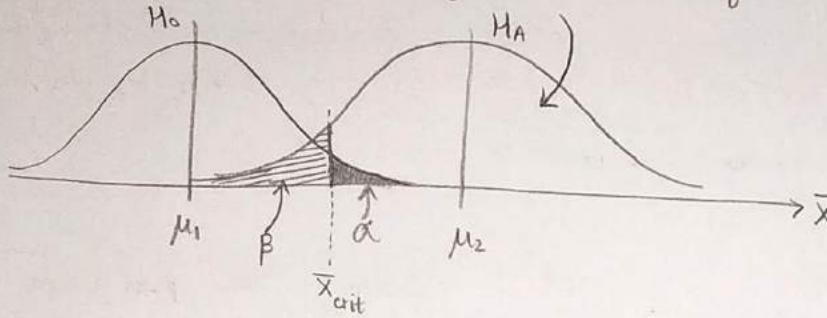
In general, we want to go for high power (high  $1-\beta$ ) tests.

## Lecture - 21

(28-02-2022)

The errors  $\alpha, \beta$ , "p-value", all of these are inter-related by the parameter  $(1-\beta)$  or the power of the test.

$(1-\beta)$  = probability that we reject  $H_0$  when it's false.



We want to be in the  $(1-\beta)$  area when our null hypothesis is false.

This area depends on the amount of overlap. The overlap here is in our control b/c  $\sigma_{\bar{x}} = \sigma/\sqrt{n}$ , so if  $n \uparrow$ ,  $\sigma_{\bar{x}} \downarrow$ , and the overlap goes  $\downarrow$   
 $\Rightarrow$  power of test  $(1-\beta) \uparrow$ .

So we can decide the optimum "n" for our desired  $\alpha, \beta$  values.

Say we want •  $(1-\beta) = 0.9$  or  $\beta = 0.1$

•  $\alpha = 0.01$

Now  $n = \left[ \frac{Z_{\alpha} + Z_{\beta}}{\frac{(\mu_2 - \mu_1)}{\sigma}} \right]^2$  for a one-tailed test.

$(\mu_2 - \mu_1)$  is NOT the difference b/w actual & hypothesized mean. We wouldn't be doing hypothesis testing if we knew  $\mu_2$ . So,  $\mu_2 - \mu_1$  is a minimum "meaningful difference".

Let's say we keep our minimum meaningful difference at  $\mu_2 - \mu_1 = 5$  and find out a sample size "n" value corresponding to it.

If in reality  $\mu_2 - \mu_1 = 7$  now, then  $Z_{\alpha} + Z_{\beta} \propto n \cdot (\mu_2 - \mu_1)^2$

$\Rightarrow Z_{\beta} \propto (\mu_2 - \mu_1)^2$ . So if  $(\mu_2 - \mu_1) \uparrow$  for fixed n,  $Z_{\beta} \uparrow \Rightarrow \beta \downarrow$  and the power of test  $(1 - \beta) \uparrow$ .

So, once our  $(\mu_2 - \mu_1)$  is fixed to a min. difference, having  $\mu_2 - \mu_1 > \text{diff.}_{\min}$  only improves the power of the test!

However, by choosing a large enough n according to above prescription, we can make any difference statistically significant, which apparently is bad.

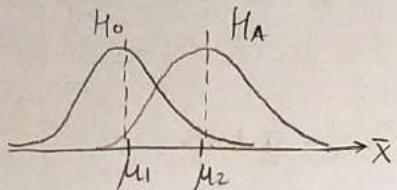
There is a difference b/w statistical significance vs Biological significance!

## Misusing statistical significance

### • Case - 1

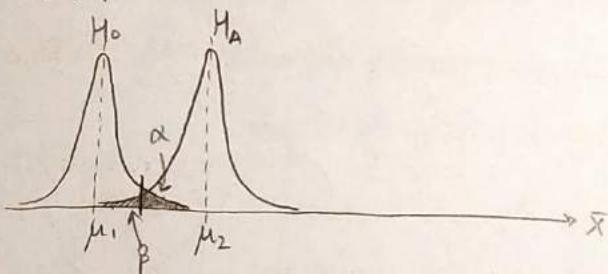
Let's say the hypothesized mean is  $\mu = \mu_1$  & the actual mean (unknown to us) is  $\mu = \mu_2$ . Say  $|\mu_2 - \mu_1| \ll 1$  or maybe it is extremely small.

In that case, our distribution would look like this for intermediate "n"-



since  $\mu_1$  &  $\mu_2$  are not very different, a null hypothesis  $H_0: \mu = \mu_1$  won't be rejected since there is a large overlap b/w the values of both the distributions. So, the result  $\bar{x}$  drawn from  $H_0$  distrib<sup>n</sup> ( $H_0: \mu > \mu_1$ ) won't have a statistically significant difference.

However, by ramping up the value of n, we can make the distributions sharply peaked & reduce overlap.



$$\text{Since } n = \left[ \frac{Z_{\alpha} + Z_{\beta}}{(\mu_2 - \mu_1)/\sigma} \right]^2$$

If  $n \uparrow \Rightarrow Z_{\beta} \uparrow \Rightarrow \beta \downarrow \Rightarrow (1-\beta)$  power  $\uparrow$

so even though the difference b/w  $\mu_2$  &  $\mu_1$  was NOT significant, by choosing a high enough n & a high powered  $(1-\beta)$  test, one can make it look like the  $\bar{x}$  generated from the sample to have a statistically significant difference from  $\mu_1$ . (i.e. we consider it as a rare event assuming  $H_0$  is true)  
 $\bar{x} > \bar{x}_{\text{crit}}$  or  $p < \alpha \Rightarrow$  rejecting  $H_0$ .

∴ Mathematically, a high-power test is always better because higher the power, the higher the accuracy of discerning small differences.

But it can be used to create a biased impression depending on how it's presented.

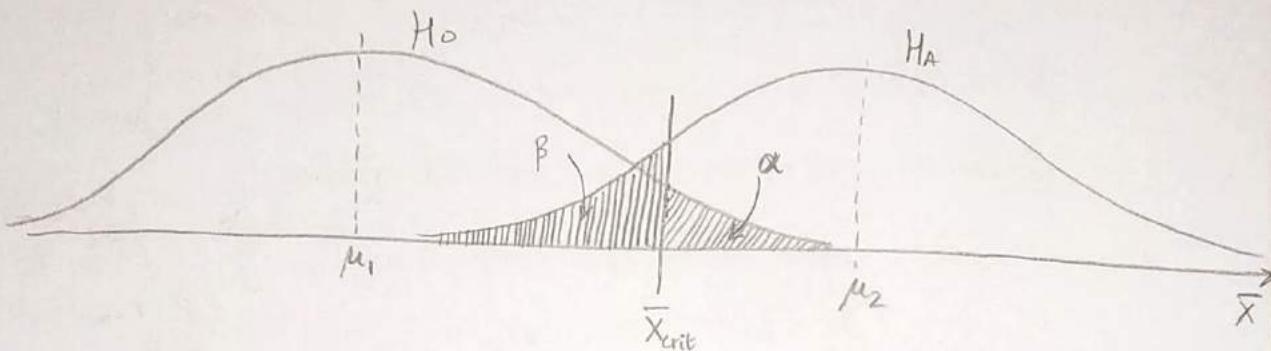
• Case-2

Let's say  $|\mu_2 - \mu_1| \gg 1$  or a significant difference exists b/w hypothesized & actual mean.

If we choose a very low  $n$ , it is easy to see that since  $n = \left\lceil \frac{(Z_\alpha + Z_\beta)}{(\mu_2 - \mu_1)/\sigma} \right\rceil$

$$n \downarrow \Rightarrow Z_\beta \downarrow \Rightarrow \beta \uparrow \Rightarrow (1-\beta) \text{ power of test} \downarrow$$

As the power  $(1-\beta) \downarrow$ , the overlap b/w distrib's increases



$\beta \equiv P(\text{not rejecting } H_0 \text{ even when it's wrong})$  is very high.

If  $\beta$  is high, it's very likely that we find a sample  $\bar{X}$  (sampled from  $H_A$ ) and mistake it to be sampled from  $H_0$  & not rejecting  $H_0$  hypothesis since  $\bar{X} < \bar{X}_{\text{crit}}$  or  $p > \alpha$  is a likely possibility.

This would make  $\bar{X}$  look like a statistically INSIGNIFICANT difference.

In all of this, we use the following def<sup>n</sup> of statistical significance

Def<sup>n</sup>: A result is said to be statistically significant if the chances of that particular event happening, assuming  $H_0$  is true, is extremely unlikely (or in simple terms,  $p < \alpha$ ).

Therefore -

- ① If  $|\mu_2 - \mu_1| \ll 1$ , we can make  $\bar{X}$  look like a statistically significant result (assuming  $H_0$  is true) i.e. an extremely unlikely event in  $H_0$  distrib<sup>n</sup>, by increasing the sample size n.
- ② If  $|\mu_2 - \mu_1| \gg 1$ , we can make  $\bar{X}$  look like a statistically INsignificant result (assuming  $H_0$  is true) i.e. as if it came from  $H_0$  distrib<sup>n</sup>, by decreasing the sample size n.

⇒ A high power test and a statistically significant difference } can be used  
A low power test and a statistically INsignificant difference } to fudge the  
inferences of the data.

Coming up with a measure of "biologically" meaningful difference.

$$\text{Effect size} = \frac{\mu_2 - \mu_1}{\sigma} = \xrightarrow{\text{actual}} \xrightarrow{\text{hypothesized}} \text{Cohen's D}/8$$

However, when we don't have access to actual mean, we find  $\bar{x}$  and replace  $\mu_2$  with it (and write  $\mu_1 = \mu_H$ )

$$\text{Effect size} \equiv \boxed{\frac{\bar{x} - \mu_H}{\sigma} = \text{Cohen's d}} \quad (\text{doesn't depend on } n)$$

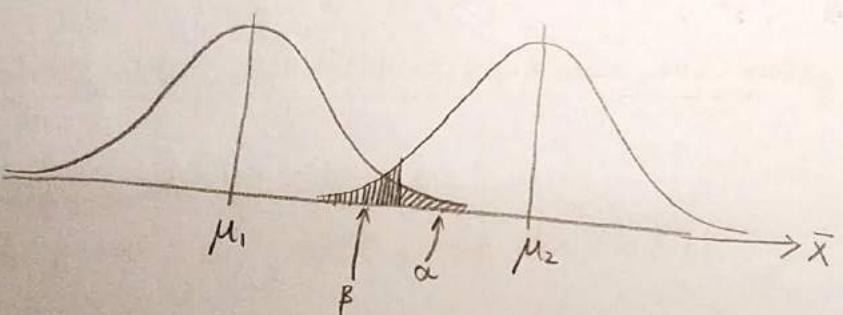
Thumb rules  $d = 0.2 \rightarrow \text{low diff.}$

$d = 0.2 \text{ to } 0.5 \rightarrow \text{intermediate diff.}$

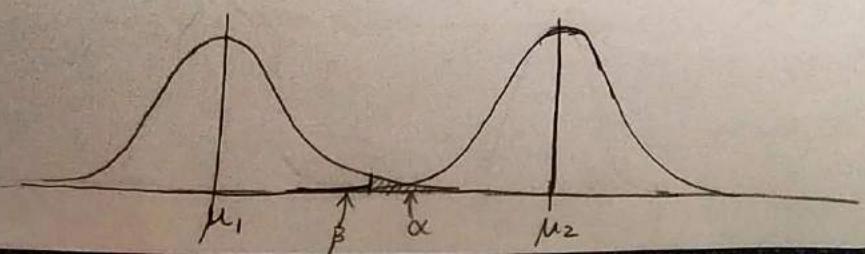
$d > 0.6 \rightarrow \text{large diff.}$

Still very subjective. Not valid everywhere.

Let's take an example where  $H_0$  is false.



To increase power ( $1-\beta$ ), we increase the sample size  $n$ .



Now if we do high powered test (given  $H_0$  is false), what happens to the p-value v/s  $\alpha$ -value comparison?

Ques Are we more likely to see  $p < \alpha$  ( $\bar{x} > \bar{x}_{\text{crit.}}$ ) i.e.  $H_0$  being rejected, when we have a high powered test compared to a low powered test?

Ans: Yes. We are more likely to see  $p < \alpha$  because as  $n \uparrow$ , the  $\beta \downarrow$  and  $(1-\beta) \uparrow$  hence decreasing the overlaps b/w the  $H_0$  &  $H_A$  distributions. If the overlap decreases, then we are less likely to make the mistake of not rejecting  $H_0$  even when it's false.

So, we prefer high powered tests.

### Reproducibility Crisis.

Reproducing results of an experiment is likely when the tests we do are high powered. If  $(1-\beta)$  is low however, then it's very likely that a hypothesis was not rejected purely by chance.

Need to repeat analysis for TWO-TAILED TEST.

Typically, we need a higher  $n$  (for the same power) in a TWO-TAILED compared to a ONE TAILED TEST.

## Lecture - 22

(02-03-2022)

All of our prior discussion involved one sample hypothesis tests, where  $\sigma$  is a known quantity. We used this to test our hypotheses concerning the means of the distribution.

But what if  $\sigma$  is unknown?

Say Heights ( $x$ )  $\sim N(\mu, \sigma)$ . We want to check the hypothesis that

$$H_0: \mu = 175 \text{ cm.}$$

The only thing we know about  $x$  is that it's normally distributed.

Take a sample of size  $n = 5$ .

If we take  $H_A: \mu \neq 175 \text{ cm.}$  (two-tailed test.)

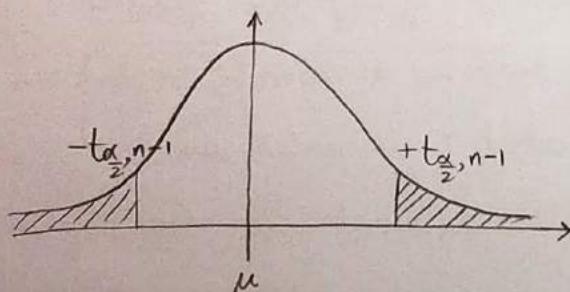
Fix  $\alpha = 0.05$ .  $\begin{cases} \alpha/2 = 0.025 \\ \alpha/2 = 0.025 \end{cases}$  to be put in both the tails.

$$\sigma \text{ known} \Rightarrow \frac{\bar{x} - \mu}{\sigma/\sqrt{n}} = z \sim N(0, 1)$$

for  $n \geq 30$ ,  $t \approx z$ .

$$\sigma \text{ unknown} \Rightarrow \frac{\bar{x} - \mu}{s/\sqrt{n}} = t \sim t_{n-1}$$

$$\pm t_{0.025, 4} = \pm 2.78$$



If calculated  $t$  is more extreme than  $\pm t_{\frac{\alpha}{2}, n-1} = \pm 2.78$ , then we say that there is a significant difference & we can reject  $H_0$ .

Say the sample gives  $\bar{X}=168$  and  $s=5$ .

$$\Rightarrow t = \frac{168 - 175}{5/\sqrt{5}} = -3.13$$

Now  $-3.13$  is more extreme than  $\pm t_{\text{crit}}$ , which is a significantly different value.  $\therefore$  We reject  $H_0$  (which means that at  $\alpha=0.05$ , there is a significant difference b/w observed & hypothesized mean).

But here we were blindly given  $n$  (the sample size). Usually, we decide on a  $\alpha, \beta$  which gives a sample size estimate.

For the case when  $\sigma$  was known,  $n = \left[ \frac{Z_\alpha + Z_\beta}{(\mu_2 - \mu_1) / \sigma} \right]^2$

For the case when  $\sigma$  is unknown (two-tailed)

$\sigma \rightarrow s, Z \rightarrow t$

$$n = \left[ \frac{t_{\alpha/2, n-1} + t_{\beta, n-1}}{(\mu_2 - \mu_1) / s} \right]^2$$

For a two-tailed test though,  $n = \left[ \frac{t_{\alpha/2, n-1} + t_{\beta, n-1}}{(\mu_2 - \mu_1) / s} \right]^2$

$\underline{\alpha \rightarrow \alpha/2}$

However, this is a problem since calculating  $t_{\frac{\alpha}{2}, n-1}$  &  $t_{\beta, n-1}$  requires knowledge of  $n$ , which we are supposed to estimate. So this is a recursive eqn. We also can't know what  $s$  is unless we know the sample size  $n$ .

so, we usually need some prior estimate of  $s$  based on some previous data to estimate  $n$ . Once we have that, we start with any random value for  $n$ , calculate  $t_{\alpha/2, n-1}$  and  $t_{\beta, n-1}$  and use it to find  $\left( \frac{t_{\alpha/2, n-1} + t_{\beta, n-1}}{(\mu_2 - \mu_1)/s} \right)^2 = n$

Now our input & output would generically be different. Use this new calculated value of  $n$  again in the formula, find the new  $n$  & compare with the old one. Repeat till you converge to a  $n$  s.t. input  $n \approx$  output  $n$ .

So, usually, we report a # of such parameters for our statistical tests.

$$t_{\alpha, df} = ?$$

$$\alpha = ?$$

$$d_{\text{realised}} = \frac{\bar{X} - \mu_H}{s} = ?$$

$$\beta = ?$$

$$\beta = ?$$

$\bar{X}$  = sample mean

$s$  = sample std dev

However, all of this includes a big assumption, that the  $X$  is normally distributed  $X \sim N(\mu, \sigma)$ . This is because  $t$ -variable is defined as follows-

$$t = \frac{\bar{X} - \mu}{\sqrt{\frac{\chi^2}{n-1}}} \xrightarrow{\text{z-dist.}} \text{z-dist.}$$

since  $\frac{\bar{X} - \mu}{\sigma/\sqrt{n}}$  is a z-distrib<sup>n</sup>

and  $\frac{\sigma^2(n-1)}{\sigma^2} = \chi^2$  is a  $\chi^2$ -distrib<sup>n</sup>

for  $X \sim N(\mu, \sigma)$  CRUCIAL!

$$\Rightarrow \frac{\frac{\bar{X} - \mu}{\sigma/\sqrt{n}}}{\sqrt{\frac{s^2(n-1)}{s^2(n-1)}}} \sim t_{n-1}$$

so,  $\frac{s^2(n-1)}{\sigma^2} \sim \chi^2_{n-1}$  only for  $X \sim N(\mu, \sigma) \Rightarrow \frac{\bar{X} - \mu}{s/\sqrt{n}} \sim t_{n-1}$  only for  $X \sim N(\mu, \sigma)$

## Lecture - 23

(03-03-2022)

Assumptions in One sample mean Hypothesis testing with the t-test-

- samples are independent.
- The original population should be normally distributed  $X \sim N(\mu, \sigma^2)$

t-test is "robust" in the sense that tiny perturbations from normal distribution don't affect the results of t-test too much. In fact, more precisely, as long as the distrib<sup>n</sup> remains symmetric under perturbations, t-test is fairly accurate.

Severe skew is bad for a t-test. But as long as distrib<sup>n</sup> is symmetric, like even a uniform distrib<sup>n</sup>, t-test is robust.

One sample test = determine, using your sample mean, if that sample came from a hypothesized distrib<sup>n</sup> with hypothesized mean  $\mu_H$ .

However, we don't usually do that in experiments, and most of the times we compare two different samples. (two sample test)

Two independent sample hypothesis test-

say height of males  $\sim N(\mu, \sigma^2 = 10)$   $\sigma_1 = \sigma_2 = 10$ .

height of females  $\sim N(\mu, \sigma^2 = 10)$

We want to check if mean height of males & females is same or different.  
i.e. if  $\mu_1 \stackrel{?}{=} \mu_2$

Say we take samples of size  $n_1$  &  $n_2$  from pop<sup>n</sup> of males & females respectively. From these, we calculate  $\bar{X}_1$  &  $\bar{X}_2$ .

Our null hypothesis is  $H_0: \mu_1 - \mu_2 = 0$  (diff. of means = 0)

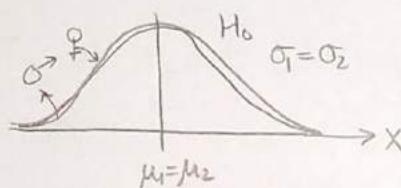
Alternate hypothesis  $H_A: \mu_1 - \mu_2 \neq 0$  → other options are  $\mu_1 - \mu_2 > 0$   
(two-tailed)  
or  $\mu_1 - \mu_2 < 0$

$$\alpha = 0.05$$

(one-tailed tests)

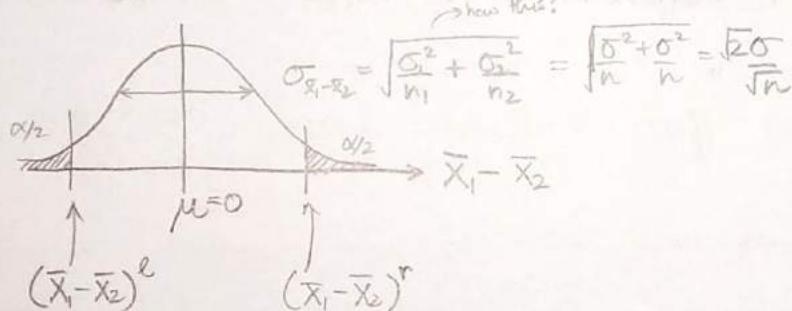
Choosing a value of  $\beta$ , we can calculate  $n_1$  &  $n_2$  & use that as sample size.  
(think about how to do this?)

If  $H_0$  is true, then the distrib<sup>n</sup> for ♂  
(M) and ♀  
(F) is the same.



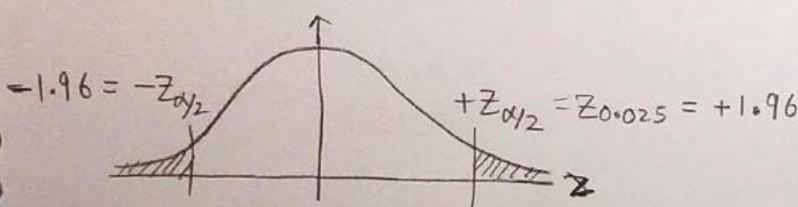
⇒ distrib<sup>n</sup> of  $(\bar{X}_1 - \bar{X}_2)$  should have  $\mu = 0$  if  $H_0$  is true.

$$\sigma_1 = \sigma_2 \\ n_1 = n_2$$



renormalize into z-distrib<sup>n</sup>

$$z = \frac{(\bar{X}_1 - \bar{X}_2) - (\mu_1 - \mu_2)}{\sqrt{2} \sigma / \sqrt{n}}$$



Say  $\bar{X}_1 = 175 \text{ cm}$   $\bar{X}_2 = 172 \text{ cm}$

$$n_1 = n_2 = 9$$
$$\sigma_1 = \sigma_2 = 10$$

$$\Rightarrow Z_{\text{calc}} = \frac{(175 - 172) - (0)}{\sqrt{2} \cdot 10 / \sqrt{9}}$$
$$= \frac{(3 - 0)}{\sqrt{2} \cdot 10 / 3} = \frac{9}{\sqrt{2} \cdot 10} \approx 0.64$$

Now  $0.64 \in [-1.96, +1.96]$  i.e. it is not as extreme as  
Z critical values.

$\Rightarrow$  The difference b/w observed ( $\bar{X}_1 - \bar{X}_2$ ) and hypothesized ( $\mu_1 - \mu_2$ ) is  
not significant enough.  $\therefore$  We do not reject  $H_0$ .

Thus, given the data, it is a valid possibility that  $\mu_1 = \mu_2$ . But we're not sure.

The samples must be independent for an independent sample test to be valid!

This assumption affects the calculation of  $\sigma_{\bar{X}_1 - \bar{X}_2} = \sqrt{\frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}}$

---

Now say we have a case where we don't know what  $\mu_1$  &  $\mu_2$  are  
AND we don't know what  $\sigma_1$  &  $\sigma_2$  are (but  $\sigma_1 = \sigma_2$ ).

Now how do we setup the test to check if the two popn's are identical?

So,  $H_0: \mu_1 - \mu_2 = 0$

$$\begin{aligned} X_1 &\sim N(\mu_1, \sigma_1) \\ X_2 &\sim N(\mu_2, \sigma_2) \end{aligned} \quad \sigma_1 = \sigma_2 = ?$$

So, estimate  $\sigma_1 \approx s_1$ , and  $\sigma_2 \approx s_2$

$$\text{Males} \xrightarrow{n_1} \bar{x}_1, s_1 \quad \text{Females} \xrightarrow{n_2} \bar{x}_2, s_2$$

so given  $\sigma_1 = \sigma_2$ , what should we use as  $s_1 = s_2 = s$ ?

so, we use an average

$$s_p \equiv \frac{s_1^2(n_1-1) + s_2^2(n_2-1)}{n_1+n_2-2}$$

↑  
pooled  
variance

For two sample test,

$$d_{\text{effect size}} = \frac{(\bar{x}_1 - \bar{x}_2) - (\mu_1 - \mu_2)}{s_p}$$

### Two independent sample testing.

Say  $X_1$  &  $X_2$  are random variables representing the two samples.

$$(X_1) \quad X_1 \sim N(\mu_1, \sigma^2)$$

samples are assumed  
to be independent.

$$n_1 = n_2$$

$$(X_2) \quad X_2 \sim N(\mu_2, \sigma^2)$$

$$\sigma_1 = \sigma_2 = \sigma \text{ but unknown}$$

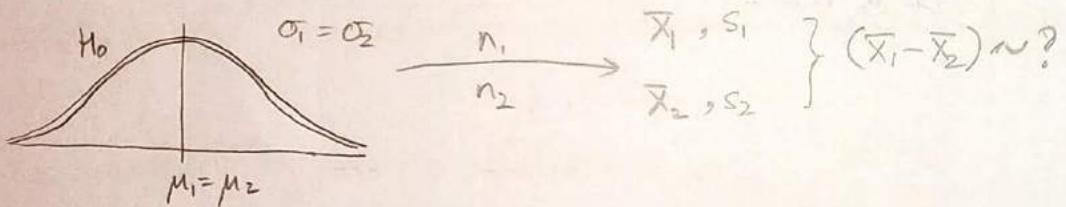
We want to test whether the means  $\mu_1$  &  $\mu_2$  of these distrib's are same or not.

$$H_0: \mu_1 - \mu_2 = 0 \quad (\text{no difference b/w } X_1 \text{ & } X_2 \text{ distrib's})$$

$$H_A: \mu_1 - \mu_2 \neq 0 \quad (\text{different distrib's}) \rightarrow \text{two-tailed.}$$

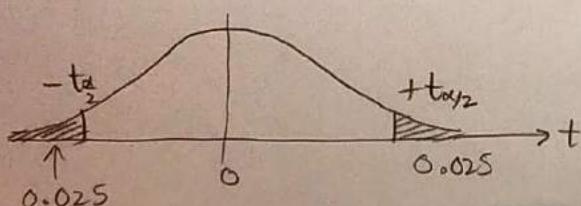
Choose  $\alpha = 0.05$

According to the  $H_0$  distribution,  $X_1$  &  $X_2$  distrib's overlap



$$\Rightarrow \bar{X}_1 - \bar{X}_2 \text{ has a mean of } \mu_1 - \mu_2 = 0 \text{ and } S = \sqrt{S_p^2 \left( \frac{1}{n_1} + \frac{1}{n_2} \right)}$$

$$\therefore \frac{(\bar{X}_1 - \bar{X}_2) - (\mu_1 - \mu_2)}{\sqrt{S_p^2 \left( \frac{1}{n_1} + \frac{1}{n_2} \right)}} \sim t_{n_1+n_2-2} \quad (\text{t-distrib'n})$$



$$\text{If } n_1 = n_2 = 5, \text{ then } \frac{(\bar{X}_1 - \bar{X}_2) - 0}{\sqrt{s_p^2 \left( \frac{1}{5} + \frac{1}{5} \right)}} \sim t_8 \quad d.o.f. = 8$$

$$t_{8, \frac{\alpha}{2}=0.025} = 2.306$$

$\therefore t_{\text{critical}}$  is  $\pm 2.306$

Knowing this, if we pull out samples & get

$$\begin{aligned}\bar{X}_1 &= 175, s_1 = 10 \\ \bar{X}_2 &= 168, s_2 = 8\end{aligned}$$

$$t_{\text{calc}} = \frac{(175 - 168)}{\sqrt{82 \cdot \frac{2}{5}}} = \underline{\underline{1.22}}$$

$$\begin{aligned}s_p^2 &= \frac{s_1^2(n_1-1) + s_2^2(n_2-1)}{n_1+n_2-2} \\ &= 4 \cdot \frac{164}{82}\end{aligned}$$

$\therefore t_{\text{calc}} \in [-t_{\frac{\alpha}{2}}, +t_{\frac{\alpha}{2}}] \Rightarrow$  We do not reject  $H_0$  (or we say that the difference of means isn't statistically significant!)

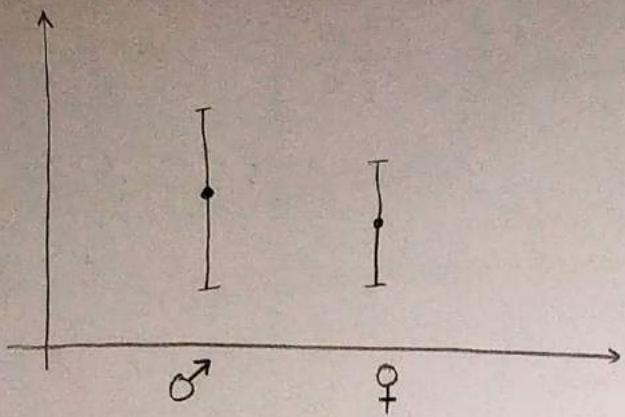
$$\boxed{\text{Cohen's } d = \frac{(\bar{X}_1 - \bar{X}_2) - \mu_H}{s_p}} = \frac{(\bar{X}_1 - \bar{X}_2)}{s_p} = \frac{7}{\sqrt{82}} \approx 0.883$$

$$\left( s_p^2 = \frac{\sum_i (\bar{X}_1 - X_{1,i})^2 + \sum_j (\bar{X}_2 - X_{2,j})^2}{n_1+n_2-2} \rightarrow \text{new type of variance + the entire sample set of } 18^2 \right)$$

Good way to plot the results of this data analysis?

Mean, 95% C.I.

Contd. on next page.

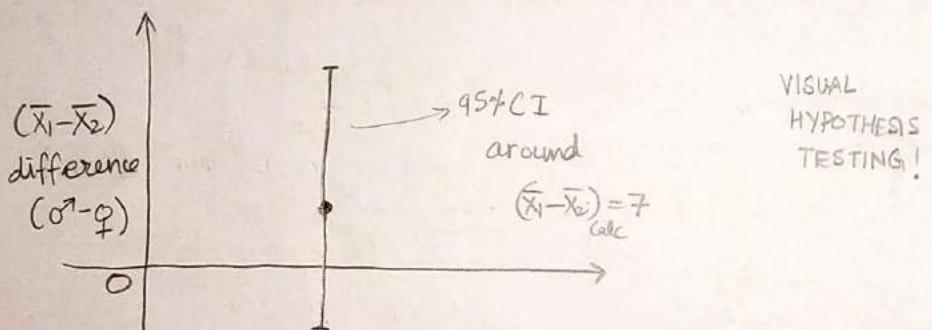


→ But this plot doesn't give any info. on the statistical test to determine differences in means.

$$\text{For } \text{♂, 95\% C.I. is } \Rightarrow \bar{x}_1 \pm t_{\frac{\alpha}{2}, n_1-1} \cdot \frac{s_1}{\sqrt{n_1}}$$

$$\text{For } \text{♀, 95\% C.I. is } \Rightarrow \bar{x}_2 \pm t_{\frac{\alpha}{2}, n_2-1} \cdot \frac{s_2}{\sqrt{n_2}}$$

A much better way is to plot the difference ( $\text{♂} - \text{♀}$ )



The 95% CI around  $(\bar{x}_1 - \bar{x}_2)$  will be given by

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

$$= 7 \pm (2.306) \cdot (5.78)$$

$$= 7 \pm 13.21$$

SINCE  $\bar{x}_1 - \bar{x}_2 = 0$  IS CONTAINED IN THE 95% C.I.

$\Rightarrow H_0$  CANNOT BE REJECTED!

## Lecture - 25.

(08-03-22)

Once the samples become correlated and aren't independent anymore, our assumptions are violated as well as the variance  $s \neq \sqrt{s_p^2 \left(\frac{1}{n_1} + \frac{1}{n_2}\right)}$  anymore since we'll also have covariance terms now.

In such cases, we do paired sample t-tests.

Example - a Before-After experiment.

Use the same set of individuals for the sample. Take data from them before taking a drug and after taking a drug.

<u>Before</u>	<u>After</u>		<u>Difference (B-A)</u>
$x_{11}$	$x_{21}$	$\xrightarrow{\text{take a difference}}$	$\Delta_1$
$x_{12}$	$x_{22}$		$\Delta_2$
$x_{13}$	$x_{23}$		$\Delta_3$
$\vdots$	$\vdots$		$\vdots$

mean =  $\bar{x}_\Delta$   
 std. dev. =  $s_\Delta$

To do the Hypothesis, take  $H_0: \mu_{\text{before}} - \mu_{\text{after}} = 0$  (meds have no effect)

$\therefore$  restating it,  $H_0: \text{mean}(\Delta) = \bar{x}_\Delta = 0$  (or  $\mu_B - \mu_A = 0$ )

$H_A: \text{mean}(\Delta) \neq 0$  (or  $\mu_B - \mu_A \neq 0$ )

$\alpha = 0.05$

$\rightarrow$  expected diff. b/w  $\mu_A$  &  $\mu_B$  assuming  $H_0$  is true.

$$\text{Now, } \frac{\bar{x}_\Delta - 0}{s_\Delta / \sqrt{n}} \sim t_{n-1}$$

Find  $t_{\text{calc}} = \frac{\bar{x}_\Delta - 0}{s_\Delta / \sqrt{n}}$  from the sample and compare with the critical value  $t_{\frac{\alpha}{2}, n-1}$ .

Three samples.

4x3  
2x1

$$A \sim N(\mu_1, \sigma_1)$$

$$\sigma_1 = \sigma_2 = \sigma_3 = \sigma \text{ (unknown)}$$

$$B \sim N(\mu_2, \sigma_2)$$

$$C \sim N(\mu_3, \sigma_3)$$

Take samples of sizes  $n_1 = n_2 = n_3 = n$

$$n_1 \rightarrow \bar{X}_1, S_1$$

$$n_2 \rightarrow \bar{X}_2, S_2$$

$$n_3 \rightarrow \bar{X}_3, S_3$$

Is at least one of these means different from the other?

One possibility is to do pair-wise t-tests

$$A \text{ vs. } B \quad H_0^{(1)}: \mu_1 - \mu_2 = 0 \quad (\alpha)$$

$$B \text{ vs. } C \quad H_0^{(2)}: \mu_2 - \mu_3 = 0 \quad (\alpha)$$

$$C \text{ vs. } A \quad H_0^{(3)}: \mu_3 - \mu_1 = 0 \quad (\alpha)$$

3 different pair-wise tests.

Total error in this pairwise test =  $3\alpha$  ???

$$\text{Bon Ferroni correction} = \frac{\alpha}{\# \text{ of pairwise tests}} = \alpha'$$

so, to keep the total error =  $\alpha$ , we put  $\alpha$  value for all  $H_0^{(i)}$  to be  $\alpha' = \frac{\alpha}{3}$  in our example.

But decrease in  $\alpha \Rightarrow$  increase in  $\beta \Rightarrow$  Power ( $1-\beta$ ) decreases.

So, we'll need to increase sample sizes "n" to compensate.

But this method isn't very realistic as pairwise test for 4 samples would demand  $\alpha' = \alpha/6$  & 6 separate tests. Better method - ANOVA

## Lecture-26.

(21-03-2022)

<u>A</u>	<u>B</u>	<u>C</u>
5	5	7
5	7	8
6	7	8
7	6	9
6	5	9
<hr/>		
$\bar{X} = 5.8$	6	8.2

One way to compare the measure & compare the yields from fertilizers A, B, and C is to do pairwise tests

$A \text{ vs. } B$     }  
 $B \text{ vs. } C$     }  
 $C \text{ vs. } A$     } Is at least one of them different?

A better method is ANOVA (Analysis of variance)

ANOVA allows us to partition the variances in the problem, but it addresses questions about the means.

Assumptions-

- ① Populations are normally distributed,  $A \sim N(\mu_1, \sigma)$
- ② Variances are equal.  $B \sim N(\mu_2, \sigma)$
- ③ Samples are independent.  $C \sim N(\mu_3, \sigma)$

Based on the # of factors (explanatory variables), we can have -

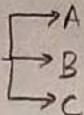
- One-Factor (One-way) ANOVA
- Two-Factor (Two-way) ANOVA
- ⋮
- $n$ -Factor ( $n$ -way) ANOVA

Each factor can have multiple levels.

Say our yield was only dependent on one factor - fertilizers.

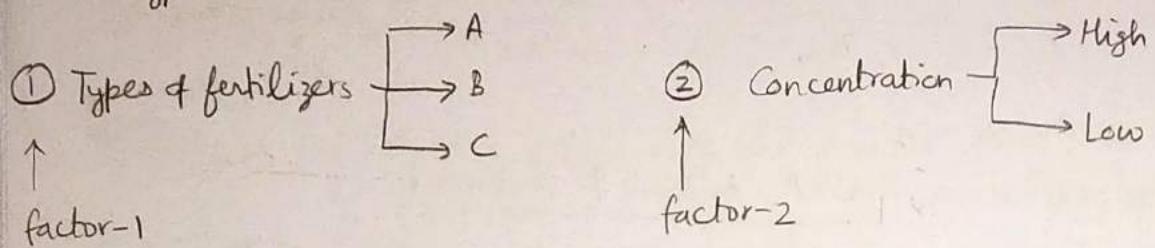
(one-way ANOVA)

And we have 3 different fertilizers



So under fertilizers, we have 3 levels.

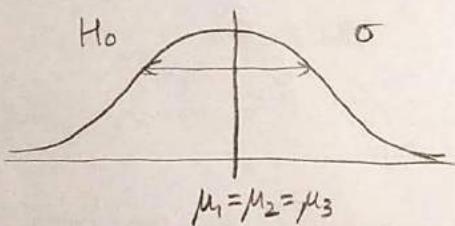
Let's now add another factor, say concentration of fertilizers along with the type.



This would be a two-factor ANOVA.

For the one-factor example (fertilizers A, B, C), the problem is setup as follows:

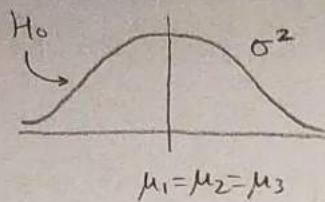
$$H_0: \mu_1 = \mu_2 = \mu_3 = \mu. \quad (\sigma \text{ is same, but unknown})$$



$$H_A: \mu_i \neq \mu_j \quad \text{for at least one } i \neq j$$

$$\alpha = 0.05$$

We'll come back to this setup in a while.



How can we estimate value of  $\sigma^2$ ?

$\sigma^2$  can be estimated in two ways -

① First method estimates  $\sigma^2$

② Second method estimates  $\sigma^2$  only if  $H_0$  is TRUE.

Won't give  $\sigma^2$  if  $H_0$  is FALSE. In fact, the estimate  $> \sigma^2$ .

#### • Method 1 : Variance Within

From our data on first page, we have the data from a sample for fertilizers A, B & C. We pick up sample & calculate  $\bar{X}$  and  $S^2$ .

$$n_1 \rightarrow \bar{X}_1, S_1^2$$

$$n_2 \rightarrow \bar{X}_2, S_2^2$$

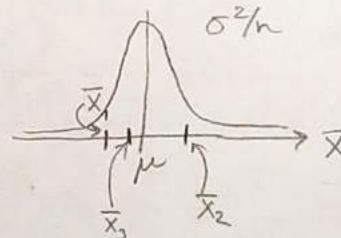
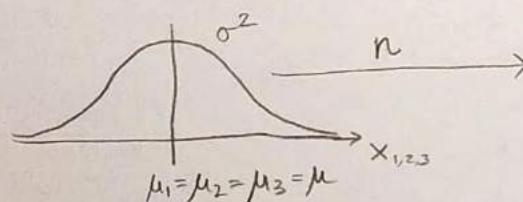
$$n_3 \rightarrow \bar{X}_3, S_3^2$$

Now all three should estimate  $\sigma^2$ .  
So which one do we choose?

⇒ We use  $S_{\text{pooled}}^2$  using  $S_1^2, S_2^2, S_3^2$  to estimate  $\sigma^2$ .

#### • Method 2 : Between Variance

Let's say  $H_0$  is true. Then taking samples of size  $n_1 = n_2 = n_3 = n$  for all three



Now  $\frac{\sigma^2}{n}$  = between variance

$\therefore (\text{b/w variance}) \cdot n \approx \sigma^2$

$$\text{b/w variance} = \text{variance of sample means}$$

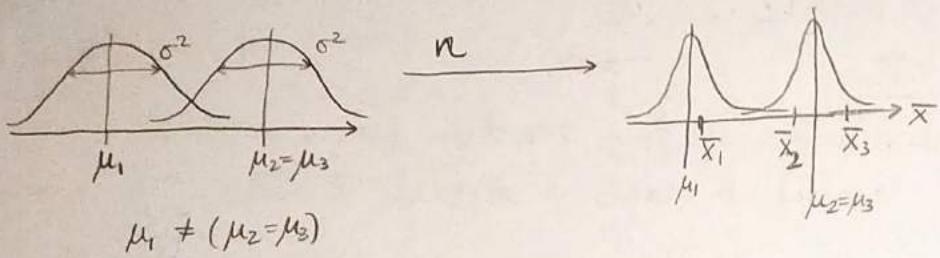
$$S = \text{Var}(\bar{X}_1, \bar{X}_2, \bar{X}_3)$$

$$= S_{\bar{X}}$$

$\bar{X}_{1,2,3}$  = sample mean for fertilizer A, B, C.

$$\Rightarrow \underbrace{n S_{\bar{X}}^2 \approx \sigma^2}_{\text{variance between}}$$

However if  $H_0$  is FALSE,



The Variance between  $\equiv \text{Var}(\bar{X}_1, \bar{X}_2, \bar{X}_3) = S_{\bar{X}}^2 > \frac{\sigma^2}{n}$  always in this case because the  $\bar{X}_1$  will be further away from  $\bar{X}_2$  &  $\bar{X}_3$  than if it would have come from the same distribution.

$\therefore$  If  $H_0$  is true, then  $n S_{\bar{X}}^2$  is a good estimate for  $\sigma^2$ .  
 If  $H_0$  is false, then  $n S_{\bar{X}}^2$  overestimates  $\sigma^2$ .

so, in the end, we'll have  $S_{\text{pooled}}^2 (S_1^2, S_2^2, S_3^2) = V_B$

$$\text{and } n \cdot S_{\bar{X}}^2 \equiv V_B$$

	$H_0$ is TRUE	$H_0$ is FALSE
and $\frac{V_B}{V_W}$	$\frac{\sigma^2}{\sigma^2} = 1$	$\frac{\delta + \sigma^2}{\sigma^2} > 1$

b/c  $V_B$  is an overestimate of  $\sigma^2$  when  $H_0$  is false.

But when is  $\frac{V_B}{V_W} > 1$  due to  $H_0$  being FALSE or just by chance?

Sol<sup>n</sup>: ratio of variances  $\sim F$ -distrib<sup>n</sup>

The F-test here will be one-tailed test because one of the tails close to  $\frac{V_B}{V_W} = F = 0$  will always support  $H_0$  but we begin to doubt

$H_0$  as  $\frac{V_B}{V_W} = F > 1$ .

## Lecture - 27

(22-03-2022)

Variance within - an average of the variance within each group i.e. for fertilizers A, B, C we get  $s_A^2, s_B^2, s_C^2$  and  $s_{\text{pooled}}^2 = \frac{s_{\text{pooled}}^2}{3} (s_A^2, s_B^2, s_C^2)$

Variance between - a measure of variance across/between the group means.

This gives us an accurate measure of variance only when  $H_0$  is true.

( $V_B \geq \sigma^2$  if  $H_0$  isn't true)

ANOVA Table

A	B	C
5	5	7
5	7	8
6	7	8
7	6	9
6	5	9
$\bar{X}_A = 5.8, \bar{X}_B = 6, \bar{X}_C = 8.2$		

$\Rightarrow$

Source	S.S.	Dof.	MS	F	P
Factor (between)					
Error (within) variance					
Total					

$$\bar{X}_{\text{grand}} = 6.67$$

S.S (sum of squares) for within variance

$$\begin{aligned} &= (5-5.8)^2 + (5-5.8)^2 + \dots + (6-5.8)^2 + \\ &+ (5-6)^2 + (7-6)^2 + \dots + (5-6)^2 + = 9.6 \\ &+ (7-8.2)^2 + (8-8.2)^2 + \dots + (9-8.2)^2 \end{aligned}$$

$$SS_{\text{within}} = SS_{\text{error}} = \sum_{i \in \{A, B, C\}} \sum_{j=1}^5 (x_{ij} - \bar{x}_i)^2$$

Now,  $SS_{\text{between}} = \text{sum of squares for b/w variance} = n \cdot \sum_{i \in \{A, B, C\}} (\bar{x}_i - \bar{x}_{\text{Grand}})^2$

$$= 5 \cdot [(5.8 - 6.67)^2 + (6 - 6.67)^2 + (8.2 - 6.67)^2] = 17.7$$

Here  $n_1 = n_2 = n_3$

If that wasn't the case,

$$SS_{\text{between}} = \sum_{i \in \{A, B, C\}} n_i (\bar{x}_i - \bar{x}_{\text{Grand}})^2$$

Also,  $SS_{\text{total}} = \sum_{i \in \{A, B, C\}} \sum_{j=1}^5 (x_{ij} - \bar{x}_{\text{Grand}})^2$

$$\begin{aligned} SS_{\text{total}} &= (5 - 6.67)^2 + (5 - 6.67)^2 + (6.67 - 6)^2 + \dots \\ &+ \dots + (9 - 6.67)^2 + (9 - 6.67)^2 \end{aligned}$$

It turns out that

$$SS_{\text{total}} = SS_{\text{between}} / \text{factor} + SS_{\text{error}}$$

$$\therefore SS_{\text{total}} = 27.3$$

The next thing in the table is d.f. (degrees of freedom)

Each group has  $n-1$  d.f.

∴ Total d.o.f for error is a sum of  $n-1 + n-1 + n-1$  since we are summing over all of them

$$df_{\text{error}} = n-1 + n-1 + n-1 = 3(n-1) = 3 \cdot (5-1) = 12 \text{ here}$$

In general,

$$df_{\text{error}} = k \cdot (n-1)$$

$k = \# \text{ of groups or levels in the factor}$   
 $n = \# \text{ of entries in a group in the factor}$

degree of freedom for / between the "factors"

$$df_{\text{factor}} = df_{\text{between}} = k-1$$

$$df_{\text{factor}} = 2 \text{ here.}$$

$$df_{\text{total}} = k \cdot n - 1$$

$$\text{and } df_{\text{total}} = df_{\text{error}} + df_{\text{factor}}$$

Roughly speaking,  $\frac{\sum_i (x_i - \bar{x})^2}{n-1} \rightarrow \text{d.f.} \approx \text{variance measure}$   $\rightarrow \text{s.s.}$

$$\therefore \text{Variance} \approx MS \text{ (Mean square)} \equiv \frac{S.S.}{d.f.}$$

$$MS_{\text{error}} = \frac{SS_{\text{error}}}{df_{\text{error}}}$$

$$MS_{\text{total}} = \frac{S.S._{\text{total}}}{d.f._{\text{total}}}$$

$$MS_{\text{factor}} = \frac{SS_{\text{factor}}}{df_{\text{factor}}}$$

$$\text{but } MS_{\text{total}} \neq MS_{\text{error}} + MS_{\text{factor}}$$

Observe that  $M S_{\text{error}} = S_{\text{pooled}}^2$

$$M S_{\text{error}} = \frac{\sum_{i=1}^k \sum_{j=1}^n (x_{ij} - \bar{x}_i)^2}{k \cdot (n-1)}$$

and this is  $\equiv S_{\text{pooled}}^2$

Now  $MS \approx \sigma^2$

$$\text{So, } M S_{\text{error}} \approx \sigma^2 \approx S_{\text{error}}^2$$

↑  
also written as  $S_{\bar{x}}^2$  or  $S_{\text{within}}^2$

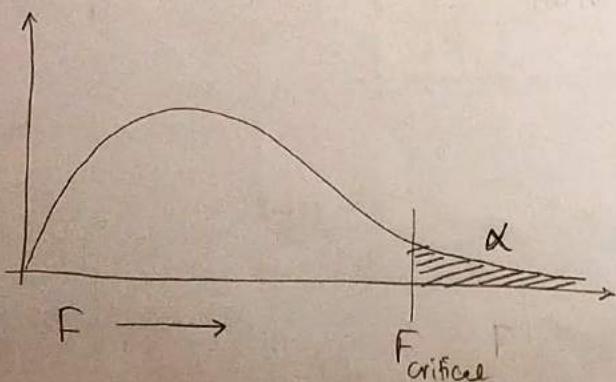
$$\text{and } M S_{\text{factor}} \approx \sigma^2 \approx n S_{\text{factor}}^2$$

↑  
also written as S.E. or  $S_{\text{b/w}}^2$

So, the F-value is calculated as

$$\frac{M S_{\text{factor}}}{M S_{\text{error}}} = F_{\text{num df, den df}} = F_{\text{f.d.f., e.d.f.}}$$

In our case      d.f. factor = 2       $M S_{\text{factor}} = 8.85$        $\Rightarrow F_{2,12} = \frac{M S_{\text{factor}}}{M S_{\text{error}}} \approx 11$   
                      d.f. error = 12       $M S_{\text{error}} = 0.8$



Now ideally  $F=1$  for  $H_0$  to not be rejected. But it occasionally might be  $>1$ . But how much  $>1$  is acceptable before we consider A, B, C to be different distrib's?

We accept an error of  $\alpha$  & say that with probability  $\alpha$ , we have rejected  $H_0$  even when it's true, so we generally keep it small.

For  $df_{num, den} = 2, 12 \Rightarrow F_{critical} = 3.89$  for  $\alpha = 0.05$

Now  $F_{calc} > F_{crit}$  ( $11 > 3.89$ )  $\Rightarrow p < \alpha$

$\therefore$  We reject our null hypothesis  $H_0$  i.e. atleast one of the data groups are different from the others.

We put all our  $\alpha$  in the right tail only because  $F < 1$  or  $\sigma_{\text{factor}}^2 < \sigma_{\text{error}}^2$

$\Rightarrow$  The variance b/w group means  $\bar{X}_A, \bar{X}_B, \bar{X}_C$  is less than the variance within the groups themselves  $\Rightarrow$  the distrib's must be fairly close to each other. In fact, it is a much stronger evidence for our  $H_0$ .

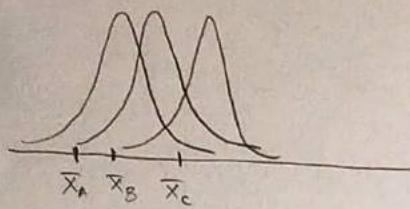
The  $MS_{\text{error}}$  &  $MS_{\text{factor}}$  measures effectively partition the variance of the system. This can be used to calculate the effect size

$$\text{Effect size} = \eta^2 = \frac{SS_{\text{factor}}}{SS_{\text{total}}}$$

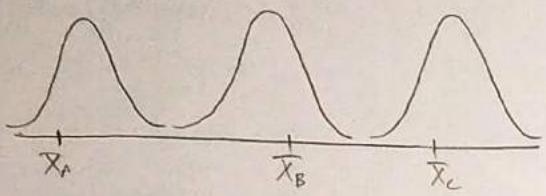
The higher the  $\eta^2$ , the higher are the chances of rejecting  $H_0$ .

$SS_{\text{factor}}$  = ss value calculated between the factor means.

$SS_{\text{error}}$  = ss value considered as an error present within a factor's group.



Low between variance  
(factor)  
⇒ low  $\eta^2$



High between variance  
(factor)  
⇒ high  $\eta^2$

In summary,

$$m = \# \text{ of groups}$$

$$n_i = \# \text{ of entries in group } i$$

Source	S. S.	d.f.	M. S.
Factor (between)	$\sum_{i=1}^m n_i (\bar{x}_i - \bar{X}_{\text{grand}})^2$	$m-1$	$MS_{\text{factor}} = \frac{SS_{\text{factor}}}{df_{\text{factor}}}$
Error (within)	$\sum_{i=1}^m \sum_{j=1}^n (x_{ij} - \bar{x}_i)^2$	$m(n-1)$	$MS_{\text{error}} = \frac{SS_{\text{error}}}{df_{\text{error}}}$
Total	$\sum_{i=1}^m \sum_{j=1}^n (x_{ij} - \bar{X}_{\text{grand}})^2$	$m \cdot n - 1$	$MS_{\text{total}} = \frac{SS_{\text{total}}}{df_{\text{total}}}$

$$SS_{\text{total}} = SS_{\text{factor}} + SS_{\text{error}}$$

$$df_{\text{total}} = df_{\text{factor}} + df_{\text{error}}$$

$$F_{\text{calc}} = \frac{MS_{\text{factor}}}{MS_{\text{error}}} . \text{ compare with } F_{df_{\text{factor}}, df_{\text{error}}}^{\text{crit}(\alpha)}$$

## Lecture - 28

(23-03-2022)

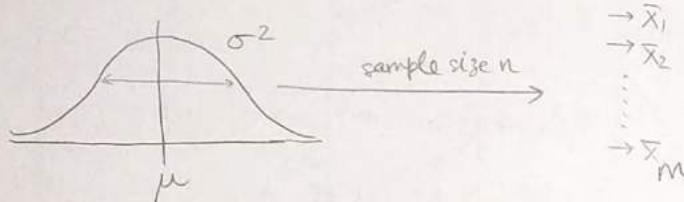
In ANOVA test, if we reject  $H_0$ , we do Post-Hoc tests.

Post hoc tests -

- Tukey's honest significant difference test. (HSD test)

The central idea is to figure which of the pairs (out of AB, BC, AC) are significantly different.

Tukey's HSD test depends on Tukey's studentized range. ( $Q$ -distrib")



$\rightarrow \bar{x}_1$   
 $\rightarrow \bar{x}_2$   
 $\vdots$   
 $\rightarrow \bar{x}_m$

$m = \#$  of groups in data table

$$\bar{x}_{\max} - \bar{x}_{\min} \rightarrow \text{range } (\bar{x})$$

Then

$$\frac{\bar{x}_{\max} - \bar{x}_{\min}}{\sqrt{\frac{M_{\text{Swithin}}}{n}}} \sim Q_{m, df_w}$$

A	B	C	
5	5	7	
5	7	8	
6	7	8	
7	6	9	
6	5	9	$\bar{x}_a$
$\bar{x}$	5.8	6	8.2
			6.67

$$\text{here, } df_w = 3.(5-1) = 12$$

$$k/m = 3 \\ n=5$$

At  $\alpha=0.05$ , we find the critical value  $Q_{\text{crit}}$  at (3, 12)

$$Q_{3,12}^{\text{crit}(\alpha=0.05)} = 3.773$$

We define HSD =  $Q_{k, df_w}^{\text{crit}(\alpha)} \cdot \sqrt{\frac{MS_w}{n}}$

$$Q_{3,12}^{\text{crit}(\alpha=0.05)} = 3.733 \quad \text{and} \quad MS_w = 0.8 \quad \text{and} \quad n=5$$

$$\Rightarrow HSD = 3.733 \sqrt{\frac{0.8}{5}} = 1.5$$

Now, we do pairwise differences & see if it's  $>$  or  $<$  HSD.

$$\Delta_{AB} = |5.8 - 6| = 0.2 < HSD \quad \left. \begin{array}{l} \\ \end{array} \right\} A-B \text{ difference isn't significant.}$$

$$\Delta_{BC} = |6 - 8.2| = 2.2 > HSD \quad \left. \begin{array}{l} \\ \end{array} \right\} B-C \text{ and } A-C \text{ differences}$$

$$\Delta_{AC} = |5.8 - 8.2| = 2.4 > HSD \quad \left. \begin{array}{l} \\ \end{array} \right\} \text{are significant.}$$

ANOVA Model (ways of interpreting the observed data)

$$\bullet \text{Factor effects} \rightarrow X_{ij} = \mu + \underbrace{T_i}_{\substack{\text{factor term} \\ \uparrow \\ \bar{X}_{\text{grand}}}} + \underbrace{\varepsilon}_{\substack{\text{error} \\ \uparrow}} \quad i \in \{A, B, C\}$$

$$\varepsilon \sim N(0, \sigma)$$

$$\text{In our case } M \approx \bar{X}_{\text{grand}} = 6.67$$

$$T_i = \bar{X}_i - \bar{X}_{\text{grand}}$$

$$\text{and } \varepsilon \sim N(0, \sigma) \text{ where } \sigma = V_i = S_i^2$$

Due to this way of defining things, we can write null hypothesis as

$$H_0: T_i = 0 \quad \forall i \quad \text{and} \quad H_A: T_i \neq 0 \text{ for any } i$$

## INTRO. TO TWO-FACTOR ANOVA

We have two factors in our example now -

Factors	# of levels
Fertilizers	3 → A → B → C
Concentration	2 → High → Low

Fertilizers / Concentration	A	B	C	$\bar{X}_{..j..}$
Low Conc.	5	5	7	$X_{ijk}$
	5	7	8	$i \in \{A, B, C\}$
	6	7	8	$j \in \{low, high\}$
	7	6	9	$k \in \{1, 2, 3, 4, 5\}$
	6	5	9	
$\bar{X}_{ij..}$	5.8	6	8.2	6.67
	7	7	11	$\bar{X}_{i..} = \text{mean over all } i = \hat{i}$
	7	9	12	$\bar{X}_{.j..} = \text{mean overall } j = \hat{j}$
	8	9	12	
	9	8	13	
	8	7	13	$\bar{X}_{ij..} = \text{mean over all values with } i = \hat{i} \text{ & } j = \hat{j}$
$\bar{X}_{i...}$	7.8	8	12.2	9.33
	6.8	7	10.2	8
	Mean of all A's	Mean of all B's	Mean of all C's	$\bar{X}_{grand} = \bar{X}_{...}$

$$\therefore \bar{X}_{ijk} = \mu + \tau_i + \beta_j + \underbrace{\tau_i \beta_j}_{\text{interaction effect.}} + \varepsilon$$

$$\mu \approx \bar{X}_{\text{grand}} = \bar{X} \dots$$

$$\tau_i \approx (\bar{X}_{i\dots} - \bar{X} \dots)$$

$$\beta_j \approx (\bar{X}_{\cdot j\dots} - \bar{X} \dots)$$

### Lecture-29

(24-03-2022)

#### ANOVA

Two-factor ANOVA. (Refer to the data on prev. page)

ANOVA table in this case looks like.

Source	S.S.	df.	MS	F	P
Factor 1 (fertilizer)	72.8	2	36.4		
Factor 2 (concentration)	53.06	1	53.06		
Interaction $F_1 \times F_2$	6.67	2	3.345		
Error	19.2	24	0.8		
Total	151.7	29	5.23		

For  $n$ -factor ANOVA with  $n \geq 2$ , we don't have a single null hypothesis

We have more than one  $H_0$ , each one for  $F_1, F_2$  and  $F_1 \times F_2$

① Let's start with the error row.

$$\begin{aligned} i &\in \{A, B, C\} \\ j &\in \{\text{low, high}\} \\ k &\in \{1, 2, 3, 4, 5\} \end{aligned}$$

$$SS_{\text{error}} = \sum_{i \in \{A, B, C\}} \sum_{j \in \{l, h\}} \sum_{k=1}^5 (x_{ijk} - \bar{x}_{ij.})^2$$

$$\begin{aligned} SS_{\text{error}} &= (5-5.8)^2 + (5-5.8)^2 + \dots + (5-6)^2 + (7-6)^2 + \dots + (7-8.2)^2 + (8-8.2)^2 \\ &\quad + \dots + (7-7.8)^2 + (7-7.8)^2 + \dots + (7-8)^2 + (9-8)^2 + \dots + (11-12.2)^2 \\ &\quad + (12-12.2)^2 + \dots + (13-12.2)^2 = 19.2 \end{aligned}$$

$SS_{\text{error}}$  ≡ sum of squared errors for all groups.

② For Factor 1, the S.S. can be calculated as

$$SS_{F_1} = \sum_{i \in \{A, B, C\}} (\bar{x}_{i..} - \bar{x}_{...})^2 \cdot n_i \cdot (\# \text{ of levels in } F_2)$$

$\uparrow$   
sample size = 5

$$= \frac{5}{10} \cdot [ (6.8 - 8)^2 + (7-8)^2 + (10.2 - 8)^2 ] = 72.8$$

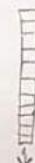
③ For Factor 2,

$$SS_{F_2} = \sum_{j \in \{l, h\}} (\bar{x}_{.j.} - \bar{x}_{...})^2 \cdot n_j \cdot (\# \text{ of levels in } F_1) = \begin{cases} 5.3 [(6.67 - 8)^2 \\ \quad + (9.33 - 8)^2] \\ = 53.06 \end{cases}$$

$\uparrow$   
sample size = 5

$n_j \cdot (\# \text{ of levels in } F_1)$  = # of entries over which the average  $\bar{x}_{.j.}$  is performed

$n_i \cdot (\# \text{ of levels in } F_2)$  = # of entries over which the average  $\bar{x}_{i..}$  is done



④ What about the interaction  $F_1 \times F_2$ ?

$$\overline{X}_{ij\cdot} = \underbrace{\overline{X}_{\dots}}_{\substack{\uparrow \\ \text{central mean}}} + \underbrace{(\overline{X}_{i\dots} - \overline{X}_{\dots})}_{\substack{\uparrow \\ \text{deviation due to} \\ \text{fertilizer } i}} + \underbrace{(\overline{X}_{\cdot j\dots} - \overline{X}_{\dots})}_{\substack{\uparrow \\ \text{deviation due to} \\ \text{concentration } j}} + \text{interaction effects}$$

$$(T_i) \qquad \qquad \qquad (\beta_j)$$

$$\therefore \text{Interaction effect} = \overline{X}_{ij\cdot} - [\mu + T_i + \beta_j] = \overline{X}_{ij\cdot} - [\overline{X}_{\dots} + \overline{X}_{i\dots} + \overline{X}_{\cdot j\dots} - \overline{X}_{\dots}]$$

and a measure of interaction variance would be

$$SS_{F_1 \times F_2} = n \sum_{ij} (obs_{ij} - exp_{ij})^2$$

$$= n \sum_{i \in \{A, B, C\}} \sum_{j \in \{L, H\}} \left[ \overline{X}_{ij\cdot} - (\overline{X}_{i\dots} + \overline{X}_{\cdot j\dots} - \overline{X}_{\dots}) \right]^2$$

For  $i=A, j=L$

$$\overline{X}_{\dots} = 8$$

$$\text{so, } \overline{X}_{A, L, \cdot}^{(\text{expected})} = 8 + T_A + \beta_L = 8 + (6.8 - 8) + (6.67 - 8)$$

$$= 5.47$$

$$\overline{X}_{A, L, \cdot}^{(\text{observed})} = 5.8$$

$$\Rightarrow \Delta \overline{X}_{A, L} = \text{interaction}_{A, L} = (\overline{X}_{A, L}^{obs} - \overline{X}_{A, L}^{exp})$$

$$\Rightarrow SS_{F_1 \times F_2} (\text{int}) = n \sum_{i \in \{A, B, C\}} \sum_{j \in \{L, H\}} \Delta \overline{X}_{ij}^2 = 6.67$$

⑤ Total measure

$$SS_{\text{total}} = \sum_{i, j, k} (X_{ijk} - \overline{X}_{\dots})^2 = 151.7$$

Let's move onto # of d.o.f.s now

①  $df_{\text{error}} = \underbrace{(5-1) + (5-1) + \dots + (5-1)}_{6 \text{ times}} = 6 \cdot (5-1) = 24$

$$df_{\text{error}} = (n-1) \cdot l \cdot m$$

$n$  = sample size

$l$  = # of levels for  $F_2$  (conc)

$m$  = # of levels for  $F_1$  (fertilizer)

②  $df_{F_1} = \underbrace{m-1}_{\substack{\text{# of levels of } F_1}} = 2$

③  $df_{F_2} = \underbrace{l-1}_{\substack{\text{# of levels of } F_2}} = 1$

④  $df_{F_1 \times F_2} = df_{F_1} \times df_{F_2} = 1 \cdot 2 = 2$

⑤  $df_{\text{total}} = df_{\text{error}} + df_{F_1} + df_{F_2} + df_{F_1 \times F_2} = 29$

Mean squared values  $MS = \frac{SS}{df}$

$$MS_{F_1} = 36.4$$

$$MS_{F_2} = 53.06$$

$$MS_{F_1 \times F_2} = 3.345$$

$$MS_{\text{error}} = 0.8$$

$$MS_{\text{total}} = 5.23$$

Now we calculate F values

$$F_{F_1} = \frac{MS_{F_1}}{MS_{\text{error}}}$$

$$F_{F_2} = \frac{MS_{F_2}}{MS_{\text{error}}}$$

There are two different types of factor

→ fixed factors

→ random factors

### Lecture - 30

(28-03-2022)

For a  $n \geq 2$  factor ANOVA,  $\exists$  a series of hypotheses.

$H_0^{(1)}$ : Factor 1 has no effect.

$H_0^{(2)}$ : Factor 2 has no effect.

$H_0^{(1 \times 2)}$ : No interaction exists b/w the two factors.

	F	df	Critical F ( $F_{\alpha}^{\alpha}$ )	p	
Factor 1	$\frac{MS_{F_1}}{MS_e} = 45.5$	2	$F_{2, 24}^{0.05} = 3.4$	$p < \alpha$	Reject
Factor 2	$\frac{MS_{F_2}}{MS_e} = 66.3$	1	$F_{1, 24}^{0.05} = 4.26$	$p < \alpha$	Reject
$F_1 \times F_2$	$\frac{MS_{F_1 \times F_2}}{MS_e} = 4.16$	2	$F_{2, 24}^{0.05} = 3.4$	$p < \alpha$	Reject
Error		24			

Factors → Fixed

→ Random

description  
of the levels

This gives rise to different types  
of ANOVA tests.

Expected MS table

	Model I (Fixed)	Model II (Random)	Model III (Mixed)
$F_1$	$\sigma_{F_1}^2 + \sigma_{\text{error}}^2$	$\sigma_e^2 + \sigma_{F_1}^2 + \sigma_{F_1 \times F_2}^2$	$\sigma_e^2 + \sigma_{F_1}^2 + \sigma_{F_1 \times F_2}^2$
$F_2$	$\sigma_{F_2}^2 + \sigma_{\text{error}}^2$	$\sigma_e^2 + \sigma_{F_2}^2 + \sigma_{F_1 \times F_2}^2$	$\sigma_e^2 + \sigma_{F_2}^2$
$F_1 \times F_2$	$\sigma_{F_1 \times F_2}^2 + \sigma_{\text{error}}^2$	$\sigma_e^2 + \sigma_{F_1 \times F_2}^2$	$\sigma_e^2 + \sigma_{F_1 \times F_2}^2$
Error	$\sigma_{\text{error}}^2$	$\sigma_e^2$	$\sigma_e^2$

• FIXED (MODEL-I)

If  $H_0$  is true, then  $\sigma_{F_i}^2 \rightarrow 0 \Rightarrow F_i = \frac{\sigma_{F_i}^2 + \sigma_e^2}{\sigma_e^2} \approx 1$  } For factor  $i$   
 If  $H_0$  is false, then  $\frac{\sigma_{F_i}^2 + \sigma_e^2}{\sigma_e^2} > 1 \Rightarrow F_i > 1$  }  $i \in \{1, 2\}$

$$F_{F_1 \times F_2} = \frac{\sigma_{F_1 \times F_2}^2 + \sigma_e^2}{\sigma_e^2}$$

so, for fixed (model-I) ANOVA, our F-values are calculated as

$$F_I = \frac{\sigma_{F_1/F_2/F_1 \times F_2}^2 + \sigma_e^2}{\sigma_e^2} = \frac{\text{M.S.}_{F_1/F_2/F_1 \times F_2}}{\text{M.S.}_{\text{Error}}}$$

• RANDOM (MODEL-II)

For model-II (random), the F-values are measured a little differently since  $MS_{F_1}$  &  $MS_{F_2}$  also have an interaction component now.

For Factors 1 & 2

$$F_{F_i} = \frac{\underline{MS_{F_i/F_2}}}{\underline{MS_{F_i \times F_2}}} = \frac{\sigma_e^2 + \sigma_{F_i \times F_2}^2 + \sigma_{F_i}^2}{\sigma_e^2 + \sigma_{F_i \times F_2}^2}$$

If  $H_0^{(i)}$  is true,  $\sigma_{F_i}^2 \rightarrow 0$  &  $F_{F_i} \approx 1$ .

If  $H_0^{(i)}$  is false,  $\sigma_{F_i}^2 \neq 0$  &  $F_{F_i} > 1$

For the interaction though,

$$F_{F_i \times F_2} = \frac{\underline{MS_{F_i \times F_2}}}{\underline{MS_{\text{error}}}} = \frac{\sigma_e^2 + \sigma_{F_i \times F_2}^2}{\sigma_e^2}$$

### MIXED (MODEL-III) ANOVA

Say  $F_1 \rightarrow \text{fixed}$  &  $F_2 \rightarrow \text{random}$

$$F_1 - \sigma_e^2 + \sigma_{F_1}^2 + \sigma_{F_1 \times F_2}^2$$

$$F_{F_1} = \frac{\underline{MS_{F_1}}}{\underline{MS_{F_1 \times F_2}}} = \frac{\sigma_e^2 + \sigma_{F_1}^2 + \sigma_{F_1 \times F_2}^2}{\sigma_e^2 + \sigma_{F_1}^2}$$

$$F_2 - \sigma_e^2 + \sigma_{F_2}^2$$

$$F_{F_2} = \frac{\underline{MS_{F_2}}}{\underline{MS_e}} = \frac{\sigma_e^2 + \sigma_{F_2}^2}{\sigma_e^2}$$

$$F_1 \times F_2 - \sigma_e^2 + \sigma_{F_1 \times F_2}^2$$

$$F_{F_1 \times F_2} = \frac{\underline{MS_{F_1 \times F_2}}}{\underline{MS_e}} = \frac{\sigma_e^2 + \sigma_{F_1 \times F_2}^2}{\sigma_e^2}$$

$$\text{error} - \sigma_e^2$$

Since the numerator & denominator in  $F$  depend upon the type of factors, the critical  $F$ 's dependence on #d.f.'s also changes!

## F-values

### Model-I (fixed)

$F_1$

$$\frac{MS_{F_1}}{MS_e}$$

### Model-II (random)

$$\frac{MS_{F_1}}{MS_{F_1 \times F_2}}$$

Model-III  
 $(F_1 \rightarrow \text{fixed})$   
 $(F_2 \rightarrow \text{random})$   
Mixed

$$\frac{MS_{F_1}}{MS_{F_1 \times F_2}}$$

$F_2$

$$\frac{MS_{F_2}}{MS_e}$$

$$\frac{MS_{F_2}}{MS_{F_1 \times F_2}}$$

$$\frac{MS_{F_2}}{MS_e}$$

$F_1 \times F_2$

$$\frac{MS_{F_1 \times F_2}}{MS_e}$$

$$\frac{MS_{F_1 \times F_2}}{MS_e}$$

$$\frac{MS_{F_1 \times F_2}}{MS_e}$$

## Critical-F values.

### Model-I (fixed)

$F_1$

$$F^{\alpha}_{df_{F_1}, df_e}$$

### Model-II (random)

$$F^{\alpha}_{df_{F_1}, df_{F_1 \times F_2}}$$

Model-III (mixed)  
 $(F_1 - \text{fixed}, F_2 - \text{random})$

$F_2$

$$F^{\alpha}_{df_{F_2}, df_e}$$

$$F^{\alpha}_{df_{F_2}, df_{F_1 \times F_2}}$$

$$F^{\alpha}_{df_{F_1}, df_{F_1 \times F_2}}$$

$F_1 \times F_2$

$$F^{\alpha}_{df_{F_1 \times F_2}, df_e}$$

$$F^{\alpha}_{df_{F_1 \times F_2}, df_e}$$

$$F^{\alpha}_{df_{F_1 \times F_2}, df_e}$$

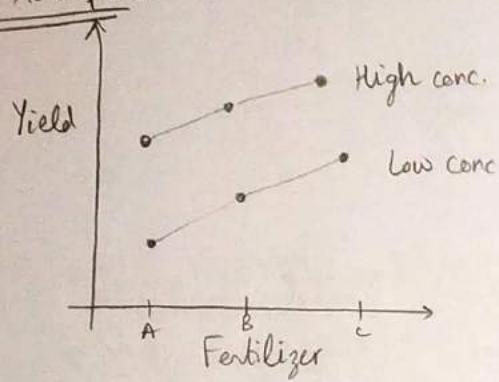
In our example, we rejected  $H_0^{(F_1 \times F_2)}$ ,  $H_0^{(F_1)}$  and  $H_0^{(F_2)}$

Rejection of  $H_0^{(F_1)}$   $\Rightarrow$  Atleast one of the three fertilizers gives a different mean yield.

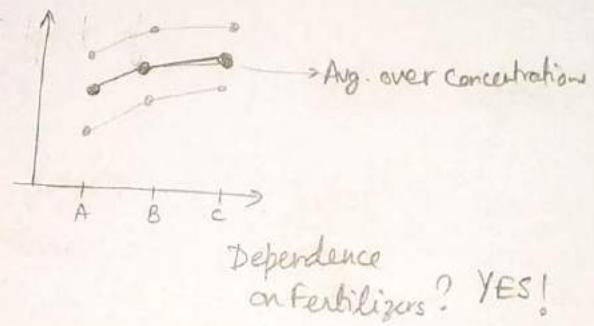
Rejection of  $H_0^{(F_2)}$   $\Rightarrow$  Atleast one of the concentration has a different mean yield.

Rejection of  $H_0^{(F_1 \times F_2)}$   $\Rightarrow$  ~~an~~ interaction b/w the factors.

### Example 1



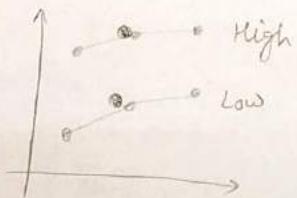
$\Leftrightarrow$  Effect of fertilizer? Average over concentrations.



Dependence on fertilizers? YES!

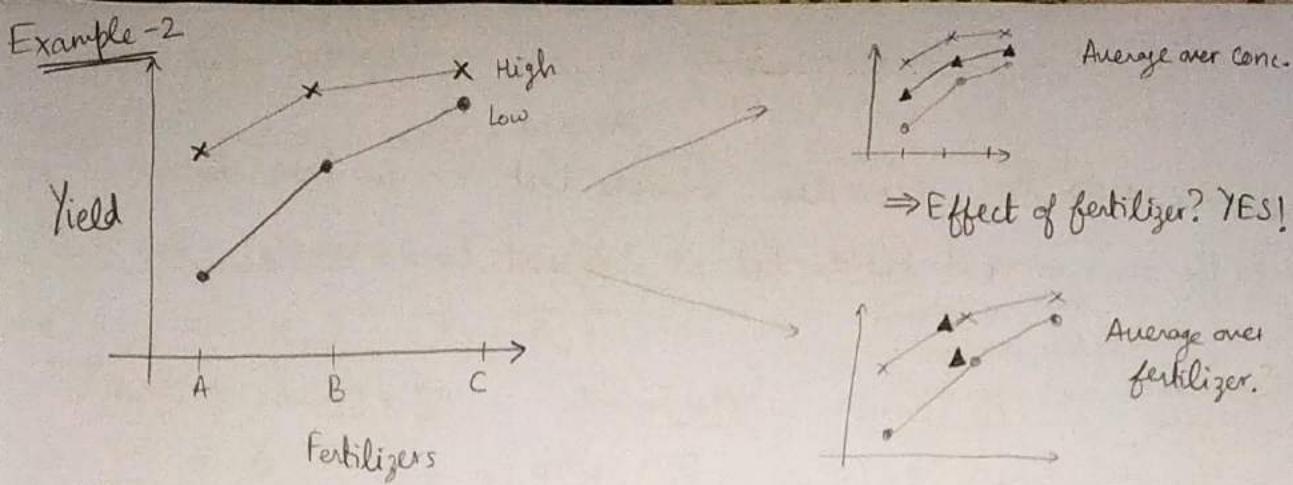
If the lines are parallel, then  $\nexists$  any interaction b/w  $F_1$  &  $F_2$ .

$\Leftrightarrow$  Effect of conc.? Avg. over fertilizers



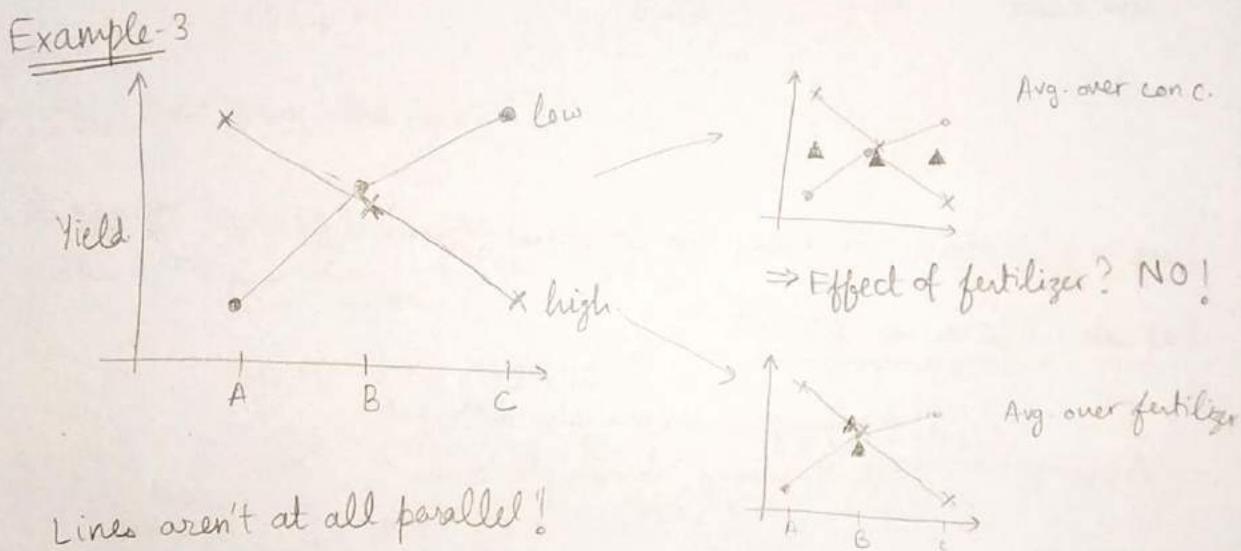
Dependence on concentration? YES!

Effect of  $F_1$ ? Average over  $F_2$ .



But lines aren't parallel.

⇒ ∃ an interaction between  
concentration and fertilizer.  
 $(F_2)$                                $(F_1)$



Lines aren't at all parallel!

⇒ ∃ a significant interaction.

So, by effect of fertilizer averaged over the levels of other factor (conc.) doesn't exist. Same for concentration's effect

The reason is that the strong interaction cancels the average effect of each other

Lecture - 31  
 (29/03/2022)

Let's consider this dependent sample test i.e. we take observations for the same group of people but at different time intervals.

<u>Subjects</u>	<u>Pre-exercise</u>	<u>After 3 months</u>	<u>After 6 months</u>	<u>subject Mean</u>
1	45	50	55	50
2	42	42	45	43
3	36	41	43	40
4	39	35	40	38
5	51	55	59	55
6	44	49	56	49.7
Time mean	42.8	45.3	49.7	

$$\text{Grand mean } \mu_g = 45.9$$

Due to dependent samples, we don't use standard ANOVA. We use Repeated Measures ANOVA.

- Assumption
- (1) samples are normally distributed
  - (2) The variance remains the same.

→ ANOVA table

<u>Source</u>	<u>S.S.</u>	<u>d.f.</u>	<u>M.C.</u>	<u>F</u>	<u>P</u>	
Time						Time - Factor
Subjects						
Error						
Total						

In one-factor ANOVA,  $SS_{\text{within}} = SS_{\text{error}}$

But in repeated measures ANOVA,  $SS_w = SS_{\text{error}} + SS_{\text{subject}}$

### S.S calculations.

- $SS_T = \frac{n_s}{\# \text{ of subjects}} \sum_t (\bar{X}_t - \bar{X}_G)^2$   
=  $6 \cdot [ (42.8 - 45.9)^2 + (45.3 - 45.9)^2 + (49.7 - 45.9)^2 ]$   
= 143.44
- $SS_{\text{subject}} \equiv SS_s = \frac{n_t}{\# \text{ of time points}} \sum_s (\bar{X}_s - \bar{X}_G)^2$   
=  $3 [ (50 - 45.9)^2 + (43 - 45.9)^2 + \dots + (49.7 - 45.9)^2 ]$   
= 658.3
- $SS_{\text{error}} = SS_w - SS_{\text{sub.}}$

Remember,

$$\begin{aligned} SS_w &= \sum_{j \in \text{subj.}} \sum_{i \in \text{time}} (x_{ij} - \bar{x}_{ij})^2 \\ &= (45 - 42.8)^2 + (42 - 42.8)^2 + \dots + (44 - 42.8)^2 \\ &\quad + (50 - 45.3)^2 + (42 - 45.3)^2 + \dots + (49 - 45.3)^2 \\ &\quad + (55 - 49.7)^2 + (45 - 49.7)^2 + \dots + (56 - 49.7)^2 \\ &= 715.5 \end{aligned}$$

$$\Rightarrow SS_{\text{error}} = 715.5 - 658.3 = \underline{57.2}$$

$$\Rightarrow SS_{\text{total}} = \frac{SS_T + SS_s}{SS_w} + SS_e = SS_T + SS_w = 143.44 + 715.5 = \underline{858.9}$$

Dummy variable t takes a sum over all time intervals (pre, 3 month, 6 month)

## degrees of freedom (d.f. calculation)

- $d.f.T = \# \text{ of time points} (\# \text{ of time columns}) - 1$   
 $= 3 - 1 = \underline{2}$
- $d.f. \text{ subj} = \# \text{ of subjects} - 1 = 6 - 1 = \underline{5}$
- $d.f. \text{ error} = \underbrace{(n-1)}_{df_T} \cdot \underbrace{(k-1)}_{df_{\text{subj}}} = 2 \cdot 5 = \underline{10}$

looks like an interaction  
from 2-factor ANOVA

$$\Rightarrow df_{\text{Total}} = 2 + 5 + 10 = 17$$

## M.S. calculations.

$$MS_i \equiv \frac{SS_i}{df_i}$$

- $MS_T = \frac{143.4}{2} = \underline{71.7}$
- $MS_{\text{error}} = \frac{57.2}{10} = \underline{5.72}$
- $MS_{\text{subj}} = \frac{658.3}{5} = \underline{131.6}$
- $MS_{\text{total}} = \frac{858.9}{17} = \underline{50.53}$

## F-ratio calculations.

$$F_T = \frac{MS_T}{MS_{\text{error}}} = 12.5$$

$$F_{\text{subj}} = \frac{MS_{\text{subj}}}{MS_{\text{error}}} = 23$$

and  $F_T^{\text{crit}} = F_{2,10}^{\alpha} = 4.1$   
 $(\alpha=0.05)$

$$F_s^{\text{crit}} = F_{5,10}^{\alpha} = 3.33$$
 $(\alpha=0.05)$

∴ In both cases, we reject the null hypotheses  $H_0$ .

⇒ We conclude that there is a significant effect of time as well as subjects!

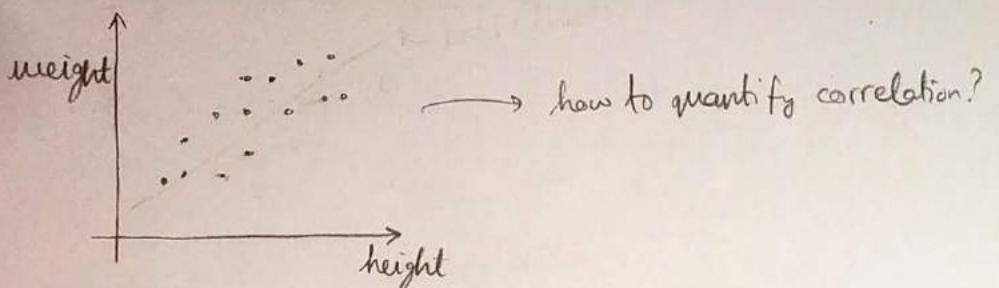
( data for at least one time is different from the others )  
( data for at least one of the subjects is diff. from the others )

## Lecture - 32

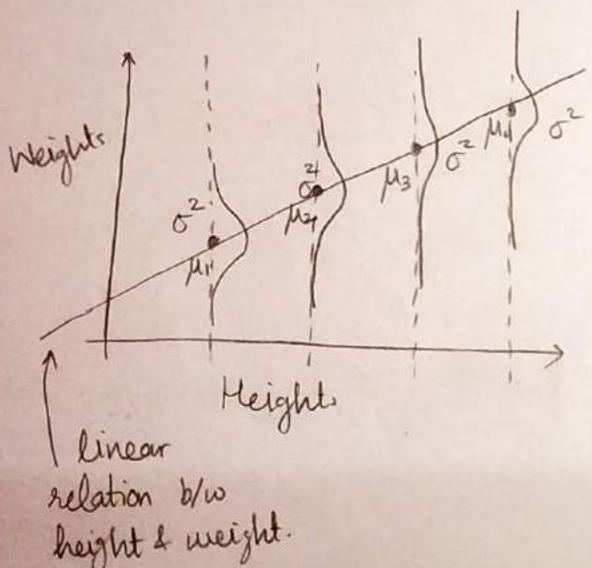
( 30 - 03 - 2022 )

### Correlation

We want to test if height & weight are correlated.



⇒ Pearson's product moment correlation coefficient.  
(Pearson's "r")



If we look at the entire pop,

at every value of height, we have a normal distrib'n of weights.

However, the  $\sigma^2$  of the weight distrib'n at each value of height is same.

Assumptions in correlation calc.:-

- ① relationship is linear ( $y = m \cdot x + c$ )
- ②  $\forall x \exists y$  s.t.  $y \sim N(\mu(x), \sigma^2)$
- ③  $\sigma^2$  is the same  $\forall x$ .
- ④ Samples are independent.

However, we don't have the access to means  $\mu(x)$  for every value of  $x$  (height). We only have access to samples from the distrib<sup>n</sup>  $y \sim N(\mu(x), \sigma^2)$ , and we have to estimate the best fit line using this. This is called regression.

Let's do the regression on the following x-y dataset.  
and correlation

$$\begin{array}{c}
 \begin{array}{cc}
 \bar{x} & \bar{y} \\
 2 & 8 \\
 3 & 10 \\
 4 & 10 \\
 5 & 12 \\
 6 & 13 \\
 7 & 14 \\
 \hline
 \bar{x} = 4.5 & \bar{y} = 11.16
 \end{array}
 &
 \begin{array}{l}
 \sum_i (x_i - \bar{x})^2 = 17.5 \\
 \sum_i (y_i - \bar{y})^2 = 24.86 \\
 \sum_i (x_i - \bar{x})(y_i - \bar{y}) = 20.5 \\
 \Rightarrow r = 0.98
 \end{array}
 \end{array}$$

Pearson's coeff

$$r = \frac{\sum_i (x_i - \bar{x})(y_i - \bar{y})}{\sqrt{\sum_i (x_i - \bar{x})^2 \sum_j (y_j - \bar{y})^2}}$$

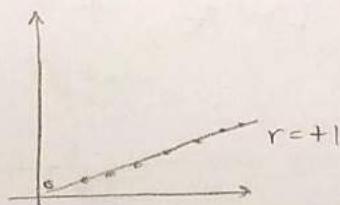
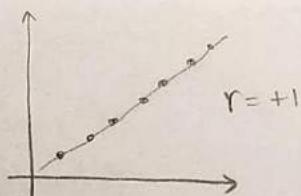
$$r \in [-1, 1]$$

$r = -1 \Rightarrow$  Perfect negative relation

$r = 0 \Rightarrow$  No relation.

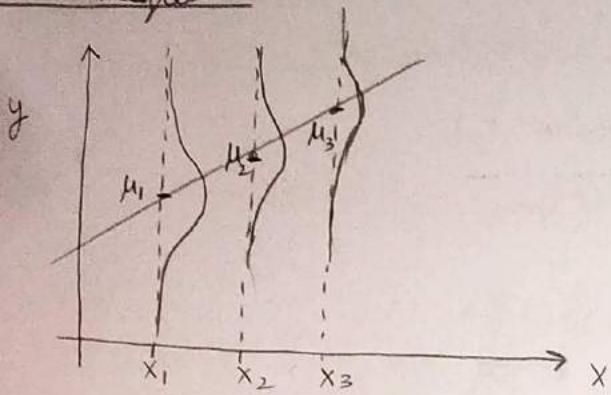
$r = +1 \Rightarrow$  Perfect positive relation

Disclaimer:  $r \neq$  slope.  $r$  represents the strength of correlation.  
(if datapts. sit on the line or not.)



Lecture - 33  
 (31-03-2022)

Linear Regression.



For every value  $x$ , the corresponding  $y \in N(\mu(x), \sigma^2)$ . This means that the data is not necessarily linear when we look at a sample & is generally scattered.

We wish to get the line joining means. It is given by

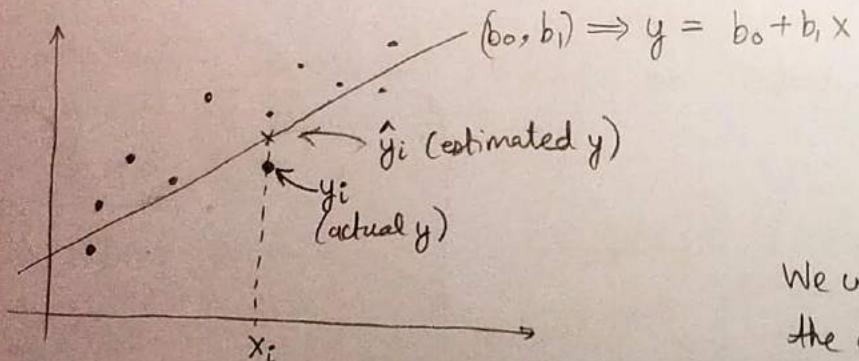
$$My = \beta_0 + \beta_1 x$$

↑                   ↑  
 y-intercept      slope

But we don't know what  $\mu_y$ 's are. So our goal is estimating  $\beta_0$  &  $\beta_1$  without knowing  $\mu_y$ 's.

Best-Fit line

Let's say our guess for  $(\beta_0, \beta_1)$  is  $(b_0, b_1)$



$$\Rightarrow \text{minimize } \sum_i (y_i - \hat{y}_i)^2$$

We want to reduce the difference for all  $y$ 's &  $\hat{y}$ 's

The best fit line is defined to give least squared deviation.

One can derive that the estimands  $b_0$  &  $b_1$  are -

$$b_0 = \bar{y} - b_1 \bar{x}$$

$$b_1 = \frac{\sum_i (x_i - \bar{x})(y_i - \bar{y})}{\sum_i (x_i - \bar{x})^2}$$

For yesterday's dataset ,  $\sum_i (x_i - \bar{x})(y_i - \bar{y}) = 20.5$

$$\sum_i (x_i - \bar{x})^2 = 17.5$$

$$\Rightarrow b_1 = \frac{20.5}{17.5} = 1.17$$

$$\Rightarrow b_0 = (11.16) - (1.17) \cdot (4.5) = 5.89$$

$$b_0 \approx \beta_0$$

$$b_1 \approx \beta_1$$

## Lecture - 34

(04-04-2022)

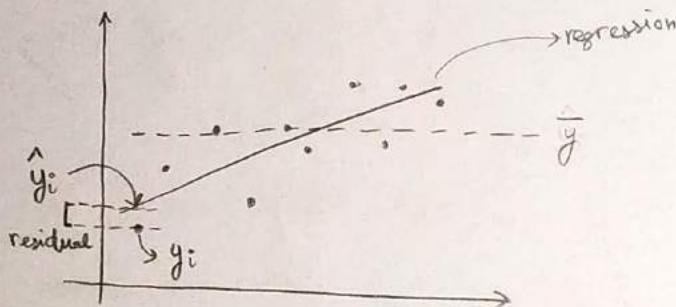
Ultimately, we are using linear regression to find the relationship b/w the  $x$  &  $y$  variables. We can pose this question in terms of a hypothesis test, whether  $\beta_1 = 0$  or  $\beta_1 \neq 0$

(independent) (dependent)

We are estimating  $\beta_1$  by  $b_1$ .  $\beta_1 \approx b_1$

$$\text{so, } H_0: \beta_1 = 0$$

$$H_A: \beta_1 \neq 0$$



F-table.

<u>Source</u>	<u>ss</u>	<u>df</u>	<u>MS</u>	<u>F</u>
Regression				
Residual (error)				
Total				

### SS calculation

$$\bar{y} = \langle y \rangle = \sum_{i=1}^n \frac{y_i}{n}$$

- $SS_{res} = \sum_i (y_i - \hat{y}_i)^2$
- $SS_{reg} = \sum_i (\hat{y}_i - \bar{y})^2$
- $SS_{total} = SS_{reg} + SS_{res}$

$\hat{y} \rightarrow$  predicted  $y$  (by regression)  
 $y \rightarrow$  observed  $y$  values.

### degrees of freedom

df for regression is always 1.

- $df_{res} = n - 2$
  - $df_{reg} = 1$
- $$\Rightarrow df_{tot} = n - 2 + 1 = n - 1$$

### MS calculation

$$MS_{reg} = \frac{SS_{reg}}{df_{reg}}$$

$$MS_{res} = \frac{SS_{res}}{df_{res}}$$

$$MS_{total} = \frac{SS_{total}}{df_{total}}$$

### F-values

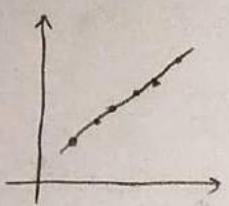
$$F_{\text{reg}} = \frac{MS_{reg}}{MS_{res}} \quad \text{and} \quad F^{\text{crit.}} = F_{1, n-2}^{\alpha}$$

In ANOVA, your independent variable is a categorical variable.

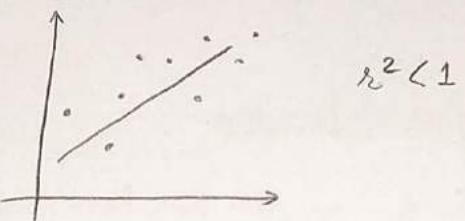
In regression, the independent variable is continuous.

We define something analogous to effect size for regression as well

$$r^2 = \frac{SS_{\text{reg}}}{SS_{\text{total}}}$$



$$\begin{aligned} SS_{\text{residue}} &= 0 \\ \Rightarrow r^2 &= 1 \end{aligned}$$



The notation is suggestive &  $r$  is literally the Pearson's coefficient!

$$r = \sqrt{r^2} = \sqrt{\frac{SS_{\text{reg}}}{SS_{\text{tot}}}}$$

(Pearson's coeff)

All our tests till now have been parametric tests which try to estimate a parameter of the population assuming some distribution.

Non-parametric tests - ??  
(distrib<sup>n</sup>-free tests)

## Lecture - 35

(05/04/2022)

A      B

122

123

135

150

153

160

173

175

180

185

190

say we don't have any reason to believe that  
the pop<sup>n</sup> distrib<sup>n</sup> ~ N.

⇒ we do distribution-free tests. (non-parametric  
tests)

Based upon ranking of data.

→ A & B are two independent samples & the pop<sup>n</sup> distrib<sup>n</sup> is  
unknown!

The first test which dealt with this was the Wilcoxon's rank sum test.  
It was improved by Mann & Whitney ⇒ Mann-Whitney U-test

This is a non-parametric equivalent of 2 independent-sample t-test.

<u>A</u>	<u>B</u>	<u>Rank A</u>	<u>Rank B</u>	
122		1		$H_0$ : Median are same. (median A = median B)
123		2		$H_A$ : Median are different.
135		3		
150		4		
	153		5	OR
	160		6	
	173		7	$H_0$ : distrib <sup>n</sup> s are same.
175		8		$H_A$ : distrib <sup>n</sup> s are different.
	180		9	
	185		10	$\alpha=0.05$
	190		11	
		<u><math>R_1 = 18</math></u>	<u><math>R_2 = 48</math></u>	$R_1, R_2$ = Rank Sums

Calculate  $U_1, U_2$

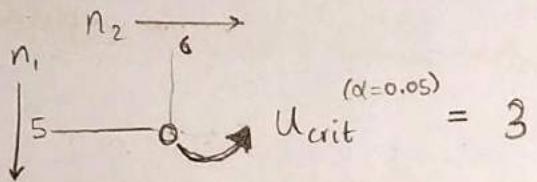
$$U_1 = n_1 n_2 + \frac{n_1(n_1+1)}{2} - R_1 = 27$$

$$U_2 = n_1 n_2 + \frac{n_2(n_2+1)}{2} - R_2 = 3$$

$n_1, n_2 \rightarrow$  sample sizes of A, B  
 $R_1, R_2 \rightarrow$  rank sums.

$$U_{\text{calc}} \equiv \min(U_1, U_2) = 3$$

$U_{\text{crit}}$  is calculated from a Mann-Whitney U-table. ( $\alpha = 0.05$ )



$U_{\text{calc}} > U_{\text{crit}}^{\alpha=0.05} \Rightarrow$  Do not reject  $H_0$ .

$U_{\text{calc}} \leq U_{\text{crit}}^{\alpha=0.05} \Rightarrow$  Reject  $H_0$ .

Here  $U_{\text{calc}}(3) \leq U_{\text{crit}}(3) \Rightarrow$  We reject  $H_0$ .  
 $H_0$

Also notice,  $U_1 + U_2 = n_1 \cdot n_2$

Let's take a simple example where  $n_A = n_1 = 3$  &  $n_B = n_2 = 3$ .

Both groups have the same # of ranks.

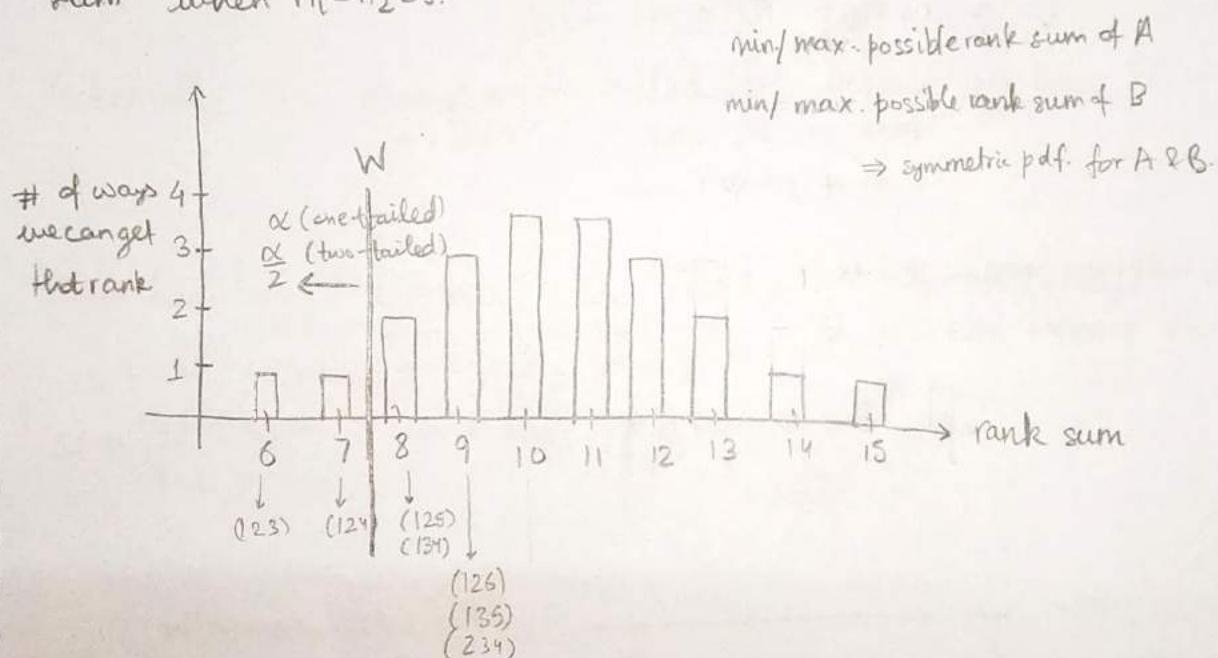
Depending upon the ranks in A or B, the total ranks  $R_1$  &  $R_2$  can vary.

Example-

A	B
1	4
2	5
3	6
<u>6</u>	<u>15</u>

A	B
1	3
2	5
4	6
<u>7</u>	<u>16</u>

Let's now find the # of ways we can get a particular rank sum when  $n_1 = n_2 = 3$ .



Probability distrib' of rank sums assuming  $H_0$  is true.

If we have  $n_1 = 3$  and  $n_2 = 4$  ( $n_1 + n_2$ )

<u>A</u>	<u>B</u>	min. rank sum of A = 6 (123)
1	4	min. rank sum of B = 10 (1234)
2	5	
3	6	max. rank sum of A = 18 (567)
	7	max. rank sum of B = 22 (4567)

for A, ranks  $\in [6, 18]$

for B, ranks  $\in [10, 22]$

In such a case, we can't have a W.

The U values

$$U_i = \underbrace{n_1 n_2 + \frac{n_i(n_i+1)}{2}}_{\text{max possible rank sum of group } i} - \underbrace{R_i}_{\text{the observed rank sum}}$$

minimum possible value of  $U_i = 0$

$$\begin{aligned} \text{max. possible value of } U_i &= \text{max. rank-sum} - \text{min. ranksum} \\ &\quad \text{of group } i \quad \text{of group } i \\ &= \text{range of rank-sums} = 18 - 6 \text{ (for A)} = 12 \\ &\quad = 22 - 10 \text{ (for B)} \end{aligned}$$

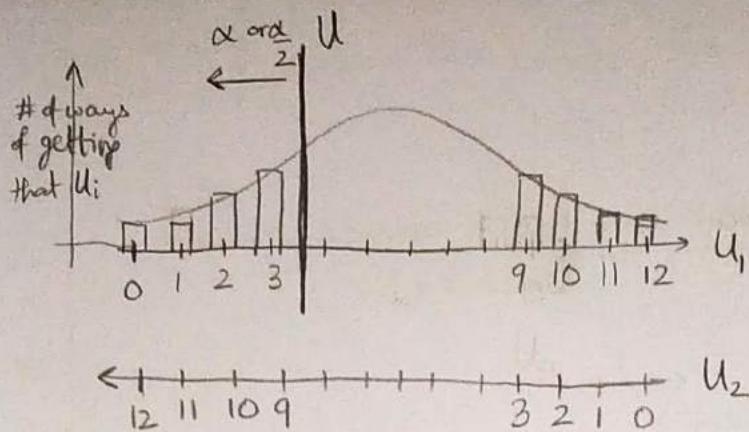
So,  $\xrightarrow[0]{12}$  The scales are now the same now

$U_i$        $i=A,B$

$$\text{If } U_1 = 0, \quad U_2 = n_1 n_2 - U_1 = 12 - 0 = 12.$$

$$\text{If } U_1 = 12, \quad U_2 = n_1 n_2 - U_1 = 12 - 12 = 0$$

$$U_1 + U_2 = n_1 n_2$$



The drawback of Wilcoxon was exactly that it needed  $n_1 = n_2$ .  
 Mann-Whitney U renormalizes both groups to the same scale  
 hence making hypothesis testing achievable.

Technically, U is always a one tailed test because we have  $U = \min(U_1, U_2)$

If we had defined  $U = \max(U_1, U_2)$  instead of min, the only change  
 is that the condition for rejecting  $H_0$  gets reversed!

## Lecture - 36

(06 - 04 - 2022)

Non-parametric test equivalent to a paired sample t-test.  
(dependent.)

- Wilcoxon's signed rank test

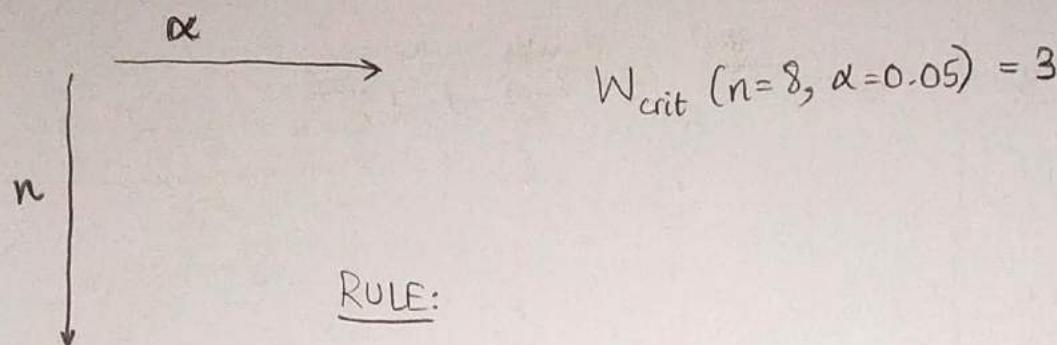
<u>Before</u>	<u>After</u>	<u><math> \Delta </math></u>	<u>Ranks</u>	<u>sign (B-A)</u>
85 — 75	10	3	+	
70 — 50	20	6	+	
40 — 50	10	3	-	
65 — 40	25	7	+	
80 — 20	60	8	+	$n=8$
75 — 65	10	3	+	$= \sum S.R.$
55 — 40	15	5	+	
20 — 25	5	1	-	

If we get a tie, we do an average of the ranks the numbers were supposed to take, so here we have 10 at rank 2,3 & 4, so we take  $\frac{2+3+4}{3} = \frac{9}{3} = 3 \leftarrow \text{avg. rank}$ . The next rank to be taken would be 5.

Signed rank sum  $\rightarrow - (3+1) = 4$   
 $\rightarrow + (3+6+7+8+3+5) = 32$

$$W_{\text{calc}} = \min (R_+, R_-) = \underline{\underline{4}}$$

$W_{\text{crit}}$   $\alpha=0.05$  gives us the critical value from a W-table



$W_{\text{calc}} \leq W_{\text{crit}} \Rightarrow \text{Reject } H_0.$

$W_{\text{calc}} > W_{\text{crit}} \Rightarrow \text{Do not reject } H_0.$

Here  $W_{\text{calc}} (4) > W_{\text{crit}} (3) \longrightarrow \text{we don't reject } H_0.$

The null hypothesis here is -

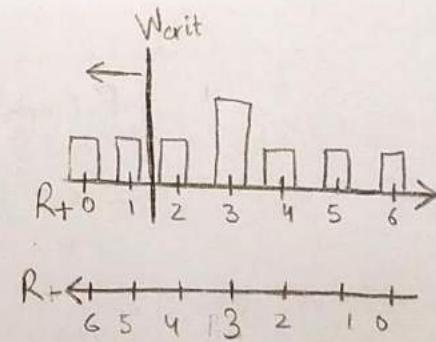
$H_0$ : The distrib's are the same.

$H_A$ : The distrib's are NOT the same

Why does this work?

Assume  $n=3$  & the  $\Delta$  values are distinct.

Ranks	sign possibilities ( $B-A$ )							
	+	-	+	-	+	-	+	-
1	+	-	+	-	+	-	+	-
2	+	+	-	-	+	+	-	-
3	+	+	+	+	-	-	-	-
$R_+$	6	5	4	3	3	2	1	0
$R_-$	0	1	2	3	3	4	5	6



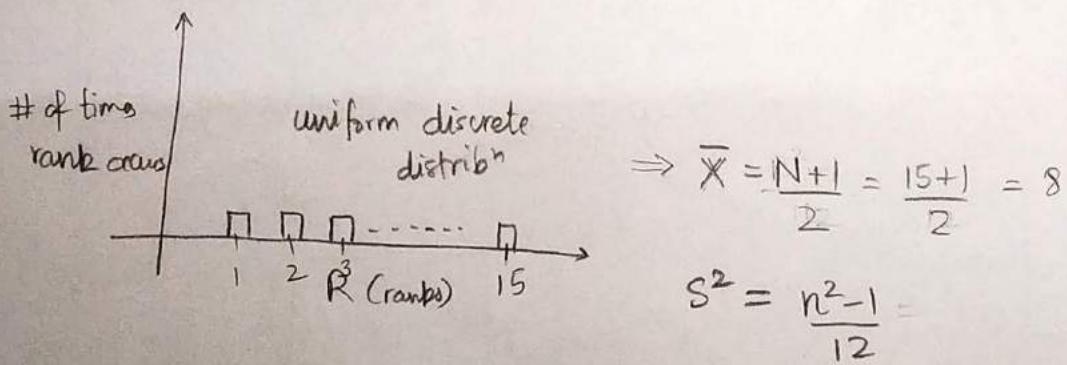
# Non parametric equivalent of 2-sample ANOVA

## Kruskal-Wallis ANOVA

			<u>Ranks</u>		
<u>A</u>	<u>B</u>	<u>C</u>	<u>A</u>	<u>B</u>	<u>C</u>
120	121		1	2	11
122			3	5	12
123	124		4	7	13
125	126		6	8	14
128	127		9	10	15
			23	32	65
			<u>Rank sums</u>		
			$\bar{R}_j$	4.6	6.4
			j=A,B,C		13
			$\bar{R}_{..} = 8$		

$$\bar{R}_{.j} = \frac{\sum_{ij} R_{ij}}{n_j} \rightarrow \text{mean ranks}$$

$$\bar{R}_{..} = \frac{\sum_{ij} R_{ij}}{n_i \times \# \text{ of gps.}} = \text{grand mean } \bar{R}_G$$



The ratio  $\frac{S^2}{\sigma^2} \sim \chi^2_{n-1}$  where  $S^2 = \text{sum of squares}$

$$\frac{\sum_j n_j (\bar{R}_{.j} - \bar{R}_{..})^2}{\frac{N^2-1}{12}} \sim \chi^2_{k-1}$$

$$\Rightarrow \frac{\sum_j n_j \left( \bar{R}_{.j} - \frac{N+1}{2} \right)^2}{\frac{N^2-1}{12}} \sim \chi^2_{k-1}$$

$k = \# \text{ of groups}$

$\frac{N^2-1}{12} = \text{variance of the discrete uniform distrib.}$

$$\boxed{\sum_j \frac{n_j (\bar{R}_{.j} - \frac{N+1}{2})^2}{\frac{N^2-1}{12}} \cdot \left( \frac{N-1}{N} \right) = H \sim \chi^2_{k-1}}$$

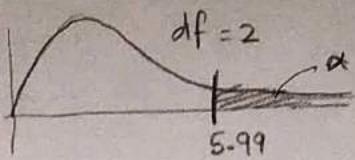
A simplified form looks like-

$$H = \frac{12}{N(N+1)} SS_{R \text{ b/w gp.}}$$

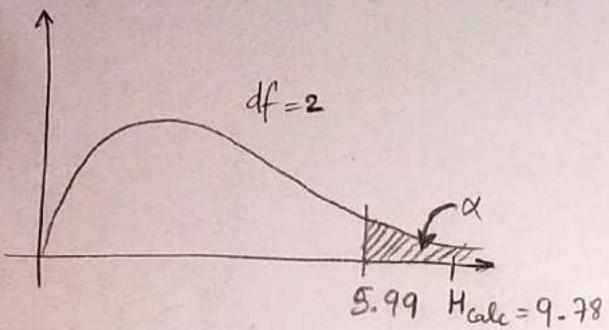
$$N = n \cdot k \\ = 5 \cdot 3 = 15$$

$$\begin{aligned} \text{Here, } SS_{R \text{ b/w group}} &= 5 \cdot (4.6 - 8)^2 + 5 \cdot (6.4 - 8)^2 + 5 \cdot (13 - 8)^2 \\ &= S [(4.6 - 8)^2 + (6.4 - 8)^2 + (13 - 8)^2] \\ &= 195.6 \end{aligned}$$

$$\Rightarrow H = \frac{12}{15 \cdot 16} \times 195.6 = 9.78$$



Now  $H \sim \chi^2_{k-1}$  where  $k=3 \Rightarrow \chi^2_{\text{crit}} = \chi^2_2 \text{ } (\alpha=0.05)$   
 $= 5.99.$



$$H_{\text{calc}} (9.78) > \chi^2_{\text{crit}} (5.99)$$

$\therefore$  We reject  $H_0$ .

where

$H_0$ : All three groups have the same ranks.

$H_A$ : The ranks of atleast one of the groups is different from others.

## Lecture -37

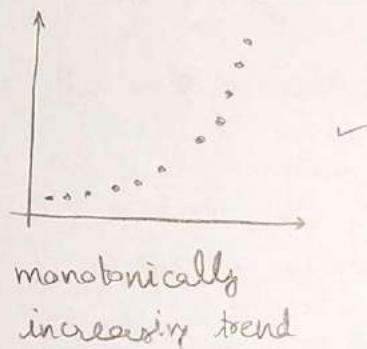
(07 - 04 - 2022)

Non-parametric equivalent of correlation & regression.

The reason for doing this test is because we don't know if the underlying probability distrib<sup>n</sup> ~ N.

- Spearman rank correlation test.

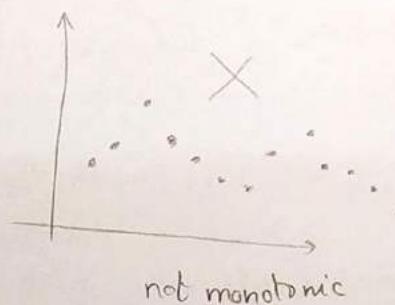
Here we assume that the two variables we're testing for are related monotonically



monotonically  
increasing trend



monotonically  
decreasing trend.



not monotonic

<u>Height</u>	<u>Weight</u>	<u>R<sub>H</sub></u>	<u>R<sub>W</sub></u>	<u>Diff. b/w ranks (D)</u>	<u>D<sup>2</sup></u>
150	45	1	1	0	0
151	55	2	3	-1	1
160	52	3	2	+1	1
165	57	4	5	-1	1
167	56	5	4	+1	1
172	60	6	6	0	0
175	65	7	7	0	0
180	70	8	8	0	0

$$\sum_i D_i^2 = 4$$

Spearman rank correlation  $\rho = 1 - \frac{6 \sum_i D_i^2}{N(N^2-1)}$

N = # of paired samples  
= 8

$$\rho = 1 - \frac{6 \cdot 4}{8 \cdot 63} = 0.95$$

Just like Pearson's correlation coeff,  $\rho \in [-1, 1]$

↑  
strong positive correlation  
strong negative correlation

Logic:

$$\min (\sum_i D_i^2) = 0 \Rightarrow \text{perfect correlation}$$

$$\begin{array}{cc} R_W & R_H \\ 1 & 1 \\ 2 & 2 \\ \vdots & \vdots \\ 8 & 8 \end{array}$$

max ( $\sum_i D_i^2$ ) is obtained when

$$= 2 \cdot [7^2 + 5^2 + 3^2 + 1^2]$$

$$= 168$$

$$\begin{array}{cc} R_H & R_W \\ 1 & 8 \\ 2 & 7 \\ 3 & 6 \\ \vdots & \vdots \\ 7 & 2 \\ 8 & 1 \end{array}$$

$\Rightarrow$  perfect negative correlation

For a given  $N$ ,

$$\max \left( \sum_i D_i^2 \right) = 2 \cdot \sum_{i=1}^{\lfloor N/2 \rfloor} (2n)^2 = \frac{N(N^2-1)}{3}$$

So, lowest  $\sum_i D_i^2 = 0 \Rightarrow$  perfect positive correlation.

Highest  $\sum_i D_i^2 = \frac{N(N^2-1)}{3} \Rightarrow$  perfect negative correlation

So we can define a measure of correlation

$$\tilde{\rho} = \frac{\sum D_i^2}{\left( \frac{N(N^2-1)}{3} \right)} \leftarrow \begin{array}{l} \text{calculated} \\ \leftarrow \max \end{array}$$

$\rightarrow = 0$  for perfect +ve correlation  
 $\rightarrow = 0.5$  for no correlation  
 $\rightarrow = 1$  for perfect -ve correlation

We map this into a scale similar to Pearson's coeff by following

$$\rho = 1 - 2\tilde{\rho} \leftarrow \begin{array}{l} = 1 \text{ for perfect +ve corr.} \\ = 0 \text{ for no corr.} \\ = -1 \text{ for perfect -ve corr.} \end{array}$$

Hence,

$$\boxed{\rho = 1 - \frac{6 \sum_i D_i^2}{N(N^2-1)}}$$

## Statistical Tests for Nominal data. (named)

### Chi-squared tests-

- ① Goodness of fit.
  - ② Independence
  - ③ Homogeneity.
- $\approx \chi^2_{df}$        $(n \geq 5)$   
                                ↑  
                                  sample size.

#### • Goodness of fit

used to check whether an observed variable fits a pre determined proportion.

	<u>AA</u>	<u>Aa</u>	<u>aa</u>
expected proportions	0.25	0.5	0.25
observed values	30	50	24

→ are these deviations significant?

In the sense of significance of these deviations, we do a goodness of fit.

#### • Independence

take a single popn & measure 2 variables - x & y.

Are these independent?

#### • Homogeneity

measure the same variable in two populations.

Are the pop's homogeneous w.r.t. that variable?

## Lecture-38

(11-04-22)

### $\chi^2$ tests for categorical / nominal data.

$n_i \geq 5$

- Goodness of fit
- Independence
- Homogeneity

} all of them  
use the same test statistic

$$\sum_i \frac{(O_i - E_i)^2}{E_i} \sim \chi^2_{df}$$

#### ① Goodness of fit.

To check that the observed data fits well with the predicted values or not.

Typical monohybrid cross

<u>F<sub>0</sub></u>	AA	X	aa
		↓	
<u>F<sub>1</sub></u>	Aa		
		↓	
<u>F<sub>2</sub></u>	AA	Aa	aa
		1 : 2 : 1	← expected ratio
	↓	↓	↓
But we get observed data as	30	65	20 ← observed distrib'

The null hypothesis is described as: ( $\alpha = 0.05$ )

H<sub>0</sub>: The observed data fits well with the predicted model.

H<sub>A</sub>: The goodness of fit is bad!

To get the expected distrib', find total sample size =  $30 + 65 + 20 = 115$

<u>Observed</u>	AA	Aa	aa	→ total = 115
	30	65	20	
<u>Expected</u>	$\frac{1}{4} \times 115$	$\frac{2}{4} \times 115$	$\frac{1}{4} \times 115$	
	$\approx 28$	$\approx 59$	$\approx 28$	→ to keep the sum

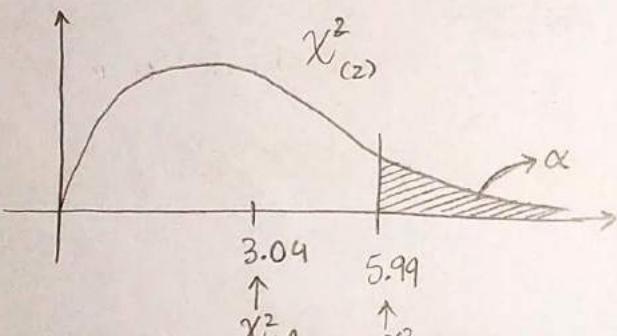
(each group has  $n_i \geq 5$  so we can use  $\chi^2$  test)

$$\chi^2_{\text{calc}} = \frac{(30-28)^2}{28} + \frac{(65-59)^2}{59} + \frac{(20-28)^2}{28} \approx 3.04$$

$$df = k-1 = 3-1 = 2$$

↑  
# of grps.

$$\chi^2_{(2)} (\alpha=0.05) = 5.99$$



Since  $\chi^2_{\text{calc}} (3.04) < \chi^2_{\text{crit}} (5.99) \Rightarrow \text{We can't reject } H_0$ .

∴ The level of fit is good enough.

- $\chi^2$  test for Independence

We want to see if 'smoking' and 'cancer' are linked or independent of each other.

- [  
 Ho: Smoking & Cancer are independent of each other.  
 Ha: They are dependent.

	No Cancer.	Cancer.	total
Non-smokers	50	10	60
Previous smokers	50	25	75
Current smokers	50	28	78
	150	63	213

This is the observed data.

What are the expected values though?

start by assuming the two factors are independent. Then we should get a 1:1 ratio b/w both factors.

$$\text{Individuals without any cancer} = \frac{150}{213}$$

in the popn

$$\text{Individuals who are non-smokers} = \frac{60}{213}$$

in the popn

$$\text{Ideally, the proportion of non-smokers without cancer (assuming they are independent)} = \frac{150}{213} \times \frac{60}{213}$$

$$\Rightarrow \text{Expected # of non-smokers without cancer} = \frac{150 \times 60}{213} \approx 42$$

Similarly, we can do this for each of them

$$\text{Expected \# of non-smokers with cancer} = \frac{63}{213} \times \frac{60}{213} \times 213 \approx \underline{\underline{18}}$$

prob. of  
cancer      prob. of  
non-smoker

$$\text{Expected \# of prev-smokers without cancer} = \frac{150}{213} \times \frac{75}{213} \times 213 \approx \underline{\underline{53}}$$

prob. of  
no cancer      prob. of  
prev-smoker

$$\text{Expected \# of prev-smokers with cancer} = \frac{63}{213} \times \frac{75}{213} \times 213 \approx \underline{\underline{22}}$$

prob. of  
cancer      prob. of  
prev-smoker

$$\text{Expected \# of smokers without cancer} = \frac{150}{213} \times \frac{78}{213} \times 213 \approx \underline{\underline{55}}$$

prob. of  
no cancer      prob. of  
smoker

$$\text{Expected \# of smokers with cancer} = \frac{63}{213} \times \frac{78}{213} \times 213 \approx \underline{\underline{23}}$$

prob. of  
cancer      prob. of  
smoker

	No cancer		Cancer	
	Obs.	Exp.	Obs.	Exp.
Non-smoker	50	42	10	18
Previous smoker	50	53	25	22
Current smoker	50	55	28	23
Total	150		63	213

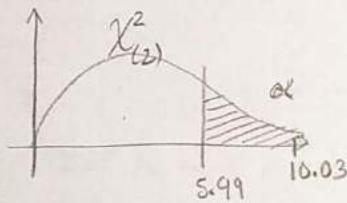
$$\begin{aligned}\chi^2_{\text{calc}} &= \sum_i \frac{(O_i - E_i)^2}{E_i} \\ &= \frac{(50-42)^2}{42} + \frac{(50-53)^2}{53} + \frac{(50-55)^2}{55} + \frac{(10-18)^2}{18} + \frac{(25-22)^2}{22} + \frac{(28-23)^2}{23} \\ &= \underline{\underline{10.03}}.\end{aligned}$$

degrees of freedom =  $(r-1) \cdot (c-1)$  =  $(3-1)(2-1) = \underline{\underline{2}}$

$\uparrow$  # of rows       $\downarrow$  # of columns

$$\Rightarrow \chi^2_{\text{crit}} (\alpha=0.05) = \chi^2_{(2)} (\alpha=0.05) = 5.99$$

since  $\chi^2_{\text{calc}} (10.03) > \chi^2_{\text{crit}} (5.99) \Rightarrow \text{Reject } H_0$ .



$\therefore$  Smoking & Cancer are related  
and not independent variables.

- $\chi^2$  test for Homogeneity.

Ex- Risk of ankle injury in two pop's

→ those who play sports  
→ those who don't play sports.

$H_0$ : The pop's are homogeneous w.r.t. no. of injuries

# of injuries	<u>Sportsperson</u>		<u>Not sportsperson (normies)</u>	<u>total</u>
	0-1	2-3	> 3	
0-1	50	50	10	60
2-3			25	75
> 3	50		28	78
<u>total</u>	150		63	213

We have the same dataset as before, and because the test runs in exactly the same way, we again reject the  $H_0$ .

No cancer, Cancer  $\longrightarrow$  sportsperson, non sportsperson

smoker categories  $\longrightarrow$  # of injury categories

The only difference b/w test of homogeneity & independence lies in-

### Independence

measure 2 variables (smoking, cancer)  
in same popn.

### Homogeneity

measuring same variable (injuries)  
in 2 popns (sportsperson, not sportsperson)