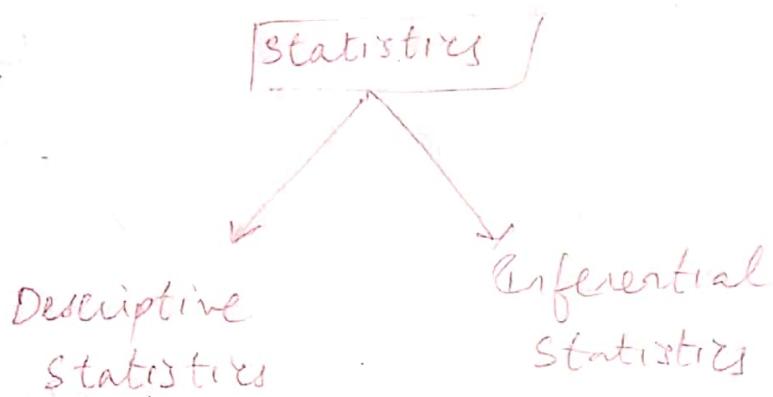


Statistics

Statistics is the science of collecting, organizing & analyzing data

* Data → facts or pieces of information

Ex: Age of students in your class.
30, 25, 60, 90, ... etc



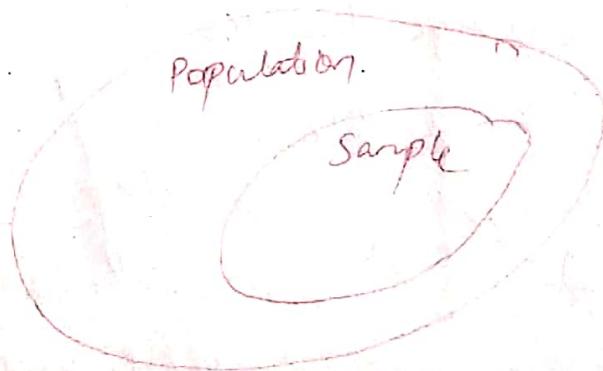
Descriptive: consist of organizing & summarizing data

Inferential: consist of using data you've measured to form conclusions.

Population & Sample

(1) Population: The group you're interested in studying

(2) Sample: a subset of population



Sampling methods

* Simple Random Sampling : Every member of population (N) has an equal chance of being selected for your sample.

* Stratified Sampling: The population (N) is split into non overlapping groups (strata), then simple random sampling is done on each group to form a sample (n).

* Systematic Sampling; Every n^{th} member of the population (N) is selected as a sample (A)

* Convenience Sampling: Easily obtained individuals from population (N) are placed in the sample
(h) also called (voluntary response sampling)

Types of Variable

Types of Variable (or many)
Variable: A property which takes any value
Ex: Age, Gender



Quantitative variables

Discrete

(specific or)
limited or infinite
values it can
take)

Continuous
(continuous many or infinite values)

Independent variable: It is any variable that is being manipulated (or) is changing.

Dependent Variable: Is any variable that is being measured.

Variable Measurement Scales

	Less Information
Nominal	
Ordinal	
Interval	
Ratios	More Info

Nominal → If data which is split into
Data categorical types is called
 Nominal data

Ordinal Data → Bi data i.e. in which order matter but distance between values does not (a) can be measured

External \rightarrow Order matters, i.e. Distance between
is equal & meaningful.
Natural zero is not present.
Ex: temperature.

Ratio! Order matters, distance between is equal & meaningful, And a natural zero does exist

Frequency distribution: lists each measured category & number of occurrences of occurrence for each category

Cumulative distribution: Adding of occurred frequencies is called cumulative frequency

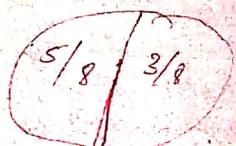
Representation of Qualitative Data (mostly DISCRETE)

- * Bar graphs: X-axis → character, Y-axis → Occurrence
- * Pie Charts: no contact b/w bars

Pie circle divided into sectors, each sector represent a category of data in proportion to the total amount of data collected

Ex: Pizza slices

Friend took $\frac{3}{8}$
You took $\frac{5}{8}$



Histogram Representation of Quantitative

- * Histogram → Bars are in contact
- * Stem & Leaf Plot used for discrete & continuous

For continuous histogram we converts into class & set up a Interval.

Range	freq
1-2	1
3-4	2

1 → lower limit
2 → upper limit

Class width = upper limit - lower limit

To Decide classes Ex: We have data of 30 students marks ranging from 1 to 50

$$30 / (\text{desired classes}) = 3 \rightarrow \text{class width.}$$

Arithmetic Mean = $\frac{\text{Sum of observation}}{\text{Number of observations}}$

Population Mean

(\bar{x})

$$\bar{x} = \frac{\sum x}{n} \quad (\text{number of sample})$$

$$\bar{x} = \frac{\sum x}{N}$$

(size of population)

Median: 1, 1, 2, 2, 2, 3, 3, 4, 4, 5, 5

$$\frac{2+3}{2} = \frac{5}{2} = 2.5$$

Mode: Most Repeated: 2

Bi-modal: Two modes; Multimodal: More than 2 modes

Mean is effected by Outliers

Median should be used for central tendency. (a) Mode depending upon Question (b) we can remove outliers.

Measures of Dispersion

Dispersion refers to how spread out a dataset is about the mean

Variance (σ^2)

Standard Deviation (S.D) (σ)

Definitional Formula
Population Variance

$$\sigma^2 = \frac{\sum (x - \mu)^2}{N}$$

Definitional Formula
Population SD (σ)

$$\sigma = \sqrt{\frac{\sum (x - \mu)^2}{N}}$$

(σ^2)
Computational Formula
Population Variance

$$\sigma^2 = \frac{\sum x^2 - (\sum x)^2}{N}$$

Computational Formula
Population S.D (σ)

$$\sigma = \sqrt{\frac{\sum x^2 - (\sum x)^2}{N}}$$

Definitional
Formula (σ^2)
Sample Variance

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

Definitional Formula
Sample S.D (σ)

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

Computational
Formula Sample
Variance (σ^2)

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

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

Percentile & Quartiles

Percentages = Meeting character
of Interest $\neq 100$

Total no. of observations

Dataset = 1, 2, 3, 4, 5

$$\frac{2}{5} \times 100 = 40\%$$

A Percentile is a value below which a certain percentage of observations lie

Ex: 2, 2, 3, 4, 5, 5, 5, 7, 8, 8, 8, 10, 11, 11, 12

Percentile rank $= \frac{\text{values below } x}{n} \times 100$

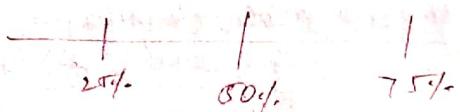
$$\text{Percentile Rank of } 10 \leftarrow \frac{16}{20} \times 100 = 80\%$$

What value exists at the percentile ranking of 25%.

$$\begin{aligned}\text{Value} &= \frac{\text{Percentile}}{100} \times (n+1) \\ &= \frac{25}{100} (20+1) = 5.25\end{aligned}$$

5.25, so we take average of 5th & 6th

Quartiles:



First Quartile - 25% $\rightarrow Q_1$

Second Quartile - 50% $\rightarrow Q_2$

Third Quartile - 75% $\rightarrow Q_3$

Five Number Summary: It is a method of summarizing a distribution of data

- * minimum
- * First Quartile (Q_1) (25%)
- * Median (50%)
- * Third Quartile (Q_3) (75%)
- * maximum

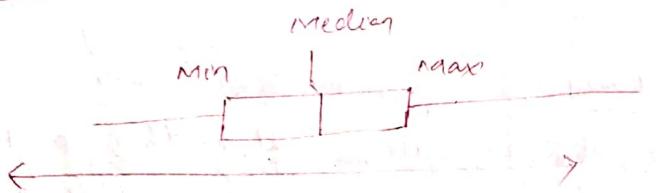
$$\text{Lower Fence} = Q_1 - 1.5(\text{IQR})$$

$$\text{Upper Fence} = Q_3 + 1.5(\text{IQR})$$

$$\text{IQR} = Q_3 - Q_1$$

Outliers can be removed if they don't fall between Upper Fence & Lower Fence.

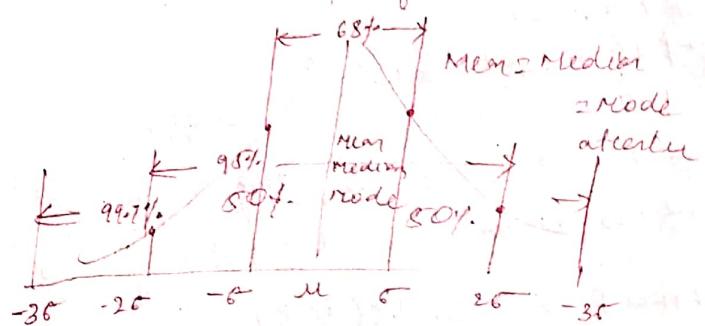
Boxplot:



Skewness: How a distribution of data leans

left skew; right skew, no skewness
Skewness is important to recognise because it has implications in hypothesis testing.

Distributions of most continuous random variables will follow the shape of normal curve



Graph changes direction at inflection points

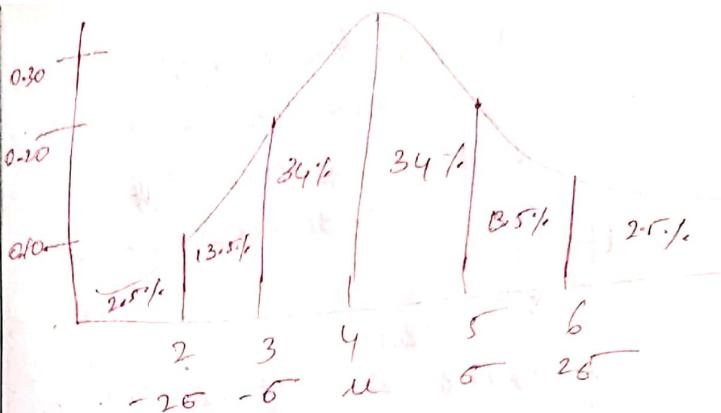
Empirical Rule:

- 6 to 6 → 68% of data
- 20 to 20 → 95% of data
- 30 to 30 → 99.7% of data

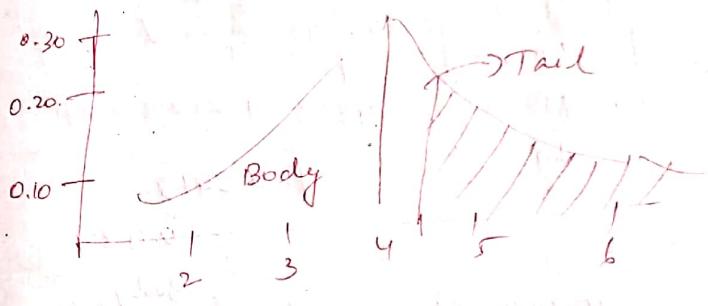
Z-scores

Are standardised values that can be used to compare scores in different distributions.

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



Let's say we have random variable 'x' distributed as such



What percentage of scores fall above

$$4.25\% \quad Z = \frac{x - \mu}{\sigma} = \frac{4.25 - 4.00}{1} = 0.25$$

From Mean \uparrow distance

$$\text{Area} = 0.5987$$

Area is Body

$$1 - 0.5987 = 0.4013$$

In United States Avg IQ is 100 with SD of 15
What percentage of the population would you expect to have an IQ

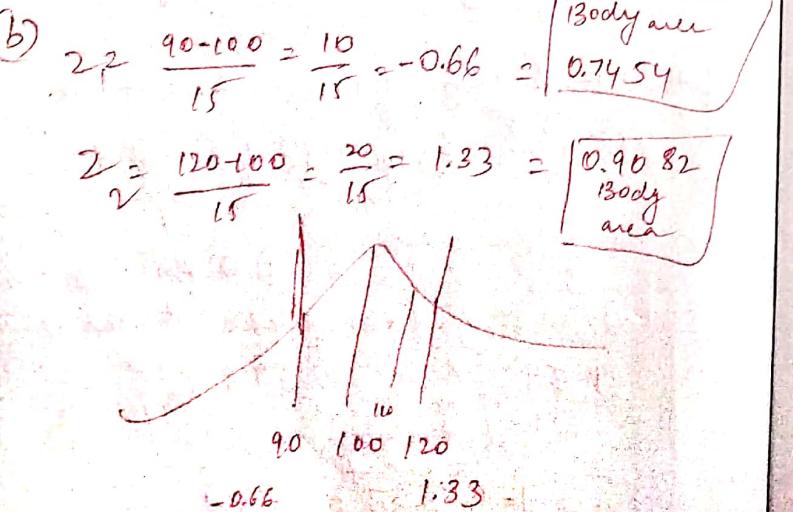
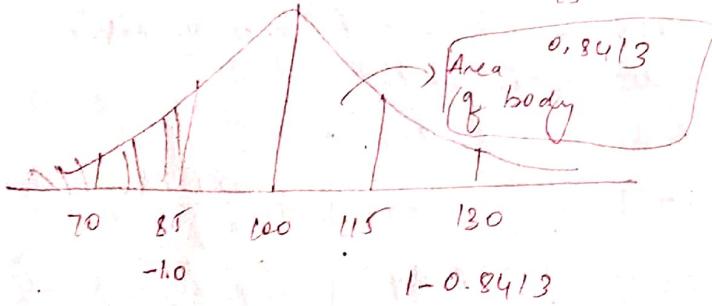
- (a) Lower than 85
(B) Between 90 & 120

$$\mu = 100$$

$$\sigma = 15$$

(a)

$$Z = \frac{x - \mu}{\sigma} = \frac{85 - 100}{15} = -1.0$$



1, Area under tail 25.46%

2, Area under tail 9.18%

Between area is 65.36%

$$\begin{array}{r} 1 \\ 25.46 \\ 9.18 \\ \hline 34.64 \end{array}$$

$$\begin{array}{r} 1.00 \\ 100.00 \\ 34.64 \\ \hline 65.36 \end{array}$$

Probability: Is a measure of the likelihood of the event.

Prob = $\frac{\text{no of ways an event can occur}}{\text{no of possible outcomes}}$

(s)
Sample Space: Collection of all possible outcomes.

Probabilities will always be between 0 & 1
An event with a probability 0 is impossible
An event with a probability 1 is certain.

Addition Rule

Two events are "mutually exclusive" if they cannot occur at same time.

Ex: Rolling a dice and getting 6, & 1 at a time.

For mutually exclusive events (Addition Rule)

$$P(A \text{ or } B) = P(A) + P(B)$$

Probability of choosing a card that is heart
(or) a king.

$$P(\text{Heart}) = \frac{13}{52}; P(\text{King}) = \frac{4}{52}$$

$$P(\text{King and Heart}) = \frac{1}{52}$$

For Non Mutually exclusive events (Addition Rule)

$$P(A) + P(B) - P(A \cap B) = P(A \text{ or } B)$$

$\cap \rightarrow \text{and}$; $\cup \rightarrow \text{or}$

$$\boxed{P(\text{Heart and King}) = 0.308}$$

Multiplication Rule

Two events are independent if they do not effect one another.

Ex: Rolling a '5' & then Rolling a '6'

Dependent: Two events are dependent if they do affect one another.

Ex: Drawing a king & then drawing a queen from deck of cards without putting the king back.

What is probability of rolling a 5 & then a 3 with a normal six sided die?

$$P(A \cap B) = P(A) * P(B)$$

$$P(3 \text{ and } 5) = \cancel{\frac{1}{6} \times \cancel{\frac{1}{6}}} = \frac{1}{36}$$

What is probability of drawing a king & then drawing a queen, from a deck of cards

$$P(B/A) = \frac{P(A \cap B)}{P(A)}$$

$$P(\text{King and Queen}) = \frac{4}{52} \times \frac{4}{51} = \frac{16}{2652}$$

Multiplication rule for Dependent events is conditional probability.

Permutation:

Ex: You visit a zoo with 6 animals & ask you to record the first

3 animals you see

Tiger, Lion, Monkey, Zebra, Walrus, Snake

$$6 \times 5 \times 4 = 120$$

With permutations order matters.

$$\text{Total no of objects} \leftarrow nPr = \frac{n!}{(n-r)!} = 6P_3 = 120$$

no of objects you are picking.

Combinations: With combination order doesn't matter

$$\text{Total no of objects} \leftarrow nCr = \frac{n!}{(n-r)!r!} = 6C_3 = 20$$

no of objects you are picking

Random Variable: Is a variable which has its value determined by a probability experiment.

Discrete Random Variable: Is a random variable which has a finite no of values.

Continuous Random Variable: Is a random variable which has an infinite no of values.

Probability Distribution: Displays probabilities associated with all possible outcomes

Discrete: Probability Distribution for one roll of a six sided die?

Number	P(Number)
1	$\frac{1}{6} = 0.167$
2	$\frac{1}{6} = 0.167$
3	$\frac{1}{6} = 0.167$
4	$\frac{1}{6} = 0.167$
5	$\frac{1}{6} = 0.167$
6	$\frac{1}{6} = 0.167$

Probability Histogram: Is a histogram with possible values on the x axis, & probabilities on y axis.

x	P(x)
1	0.10
2	0.15
3	0.35
4	0.40



Mean & Expected Value of Discrete Random Variable :-

Discrete Random Variable

x = Number of strokes to complete course

Below is probability distribution for a golfer on a par 3 hole

P(x)	x	$P(x)$
0.10	1	0.10
0.60	2	0.30
1.38	3	0.45
0.60	4	0.15

$$\text{mean formula: } \mu_x = \sum [x * P(x)]$$

$$\boxed{\mu_x = 2.65}$$

The mean we just calculated of 2.65 is an expected value.

If we were to take a large enough sample of this golfer's performance on par 3 holes we expected his mean to approach

$$2.65$$

This is an Example of law of large numbers.

Variance & S.D of Discrete Random Variable

Variance formula

$$\sigma^2 = \sum [x^2 * P(x)] - \mu_x^2$$

x	$P(x)$	x^2	$x^2 * P(x)$
1	0.10	1*1=1	0.10
2	0.30	4	1.20
3	0.45	9	4.05
4	0.15	16	2.40
			7.75

$$\mu_x = 2.65$$

$$\sigma^2 = 7.75 - (2.65)^2$$

$$\text{Variance } \boxed{\sigma^2 = 0.73}$$

$$\text{S.D. } \boxed{\sigma = 0.05}$$

Law of Large Numbers

As a probability experiment is performed many times, the observed value (mean) will arrive at the expected value (mean).

Ex: Suppose we toss a coin once & find probability of getting head is

$$P(H) = \frac{1}{2} = 0.5 \rightarrow \text{Expected mean value.}$$

If a coin is tossed as timer then the overall value of the probability of head will eventually reach 0.5 (as it is expected)

Binomial Distribution: An experiment is a binomial distribution if

- * It is repeated fixed number of times
- * Trials are independent
- * Trials have mutually exclusive outcomes either success or failure
- * Probability of success is same for all trials.

$$P(x) = nC_x p^x (1-p)^{n-x}$$

In a recent survey we found 85% of houses have high speed Internet. If you take sample of 18 households, what is the probability that exactly 15 will have high speed Internet?

$$P(x) = nC_x p^x (1-p)^{n-x}$$

$$n=18; x=15; p=85\% = 0.85$$

$$P(15) = 18C_{15} (0.85)^{15} (1-0.85)^{18-15}$$
$$= 0.239$$

At least 15

$$P(x \geq 15) = P(15) + P(16) + P(17) + P(18)$$

$$P(x \geq 15) = 0.718$$

Mean & S.D of binomial Distribution

$$\mu_x = np ; \sigma_x = \sqrt{np(1-p)}$$
$$= (18)(0.85) ; \sigma_x = \sqrt{(18)(0.85)(1-0.85)}$$
$$\boxed{\mu_x = 15.3} \quad \boxed{\sigma_x = 1.515}$$

Poisson Distribution

Is used when computing probability of a certain number of success within a specified interval.

An Experiment follows Poisson if:-

- * Probability of two successes in small enough interval is 0%.
- * Probability of a success is the same for any two intervals which share the same length.
- * Success are independent of successes in other intervals.

At a theme park, there is a roller coaster that sends an average of 3 cars through its circuit every minute between 6pm & 7pm. A random variable, X , represents the number of roller coaster cars to pass through the circuit between 6pm & 6:10pm.

What is probability that 35 cars will pass through the circuit between 6pm & 6:10pm?

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

↑ success
time interval
e ↑ rate constant

$\lambda = 3$, $x = 35$, $t = 10 \text{ minutes}$

$$e = \text{constant} = 2.718$$

$$P(35) = \frac{(\lambda)(10)^{35}}{35!} \times e^{-(\lambda)(10)}$$

$$P(35) = 0.045$$

Mean & S.D of Poisson Random Variables

$$\text{Mean} \\ \mu_x = \lambda t$$

$$\text{S.D} \\ \sigma_x = \sqrt{\mu_x}$$

$$x = \text{number of success} = 35 \\ t = \text{a length of time} = 10 \text{ minutes} \\ \lambda = \text{avg no. of success in} \\ \text{interval} = 3$$

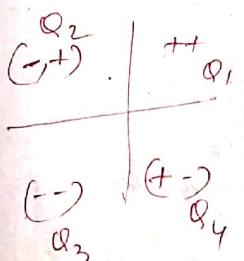
$$e = 2.718$$

$$\mu_x = (\lambda)(t)$$

$$\boxed{\mu_x = 30}$$

$$\sigma_x = \sqrt{30}$$

$$\boxed{\sigma_x = 5.477}$$



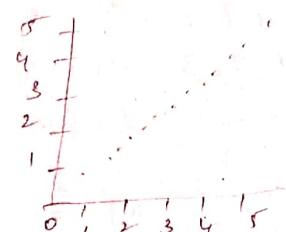
Scatter Plot: Are a method of graphically displaying bi-variate data.



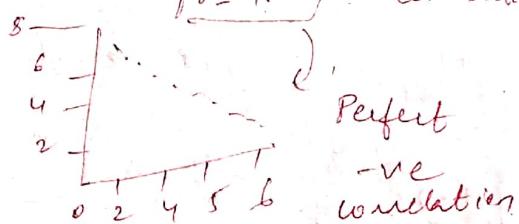
Pearson's Correlation (r)

It measures strength of linear relationship between two variables

Pearson's (r) is always between -1 & 1

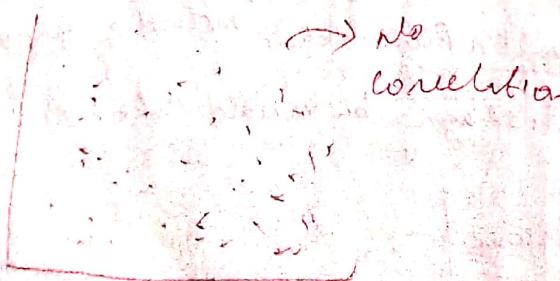


$r = 1.0 \rightarrow$ Perfect +ve correlation



Perfect -ve correlation

r=0



\rightarrow No correlation

Pearson's r subject	(X) Age	(Y) Income	(XY)
1	18	15,000	270,000
2	25	24,000	725,000
3	57	68,000	3,876,000
4	45	52,000	2,340,000
5	26	32,000	832,000
6	64	80,000	5,120,000
7	37	41,000	1,517,000
8	40	45,000	1,800,000
9	24	26,000	624,000
10	33	33,000	1,089,000
	369	421,000	18,193,000

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

$$r = 18193000 - \frac{(369)(421000)}{10}$$

$$\sqrt{15629 - \left(\frac{369}{10}\right)^2} \cdot \sqrt{212890000000 - \left(\frac{421000}{10}\right)^2}$$

$$= 0.99 \quad \text{Strong correlation}$$

Pearson's 'r' Hypothesis testing (r)
 Researchers collect 10 individual data
 to test if age & Income are related
 $\alpha = 0.05$

SL#	Age	Income
1	18	15k
2	25	29k
3	57	68k
4	45	52k
5	26	32k
6	64	80k
7	37	41k
8	40	41k
9	24	26k
10	33	33k

$H_0 : \rho = 0$ $\alpha = 0.05$

$H_1 : \rho \neq 0$ $df = n - 2$

$$10 - 2 = 8$$

using r-table

using df, $r \rightarrow$ column

$$\text{Row } r_{0.05} = 0.632$$

Using 'r' formula

and find

'r' = 0.19 So we reject
null hypothesis

Spearman's Correlation (γ_s)

So far Pearson's 'r' which measures the relationship between two continuous (interval or ratio) variables

Spearman's correlation is used when

* Measuring the relationship between two ordinal variables

* Measuring relationship between two variables that are related but not linearly.

* To calculate spearman's correlation, we must first rank scores

X	Y
2	21
5	17
8	14
11	10
15	5
16	3

X	Y
1	6
2	5
3	4
4	3
5	2
6	1

So lower score in X has lowest rank & similarly in Y lowest has lowest rank.

Now we find ' γ_s '

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

$$r_s = \frac{56 - \frac{(21)(21)}{6}}{\sqrt{\frac{91 - (21)^2}{6}} \cdot \sqrt{\frac{91 - (21)^2}{6}}} = -1.00$$

$$r_s = -1.00$$

Perfect -ve relationship

The reason we find correlation is to find linear regression.

Linear Regression

If we know that two variables are strongly correlated we can use one variable to predict other.

linear
Regression
function

$\hat{y} = \beta_1 x + \beta_0$	↑	y-intercept & Regression line
$\beta_1 = r \frac{s_y}{s_x}$	↑	$s_y \rightarrow S.D(s_y)$
$\beta_0 = \bar{y} - \beta_1 \bar{x}$	↑	$s_x \rightarrow S.D(s_x)$

Stronger your correlation (that is, the closer r is to $-1(0)$), the more accuracy your prediction will be.

s.no	x	y	sof from previous kai! Pearson we get $r = 0.99$
1	18	15k	
2	25	29k	
3	31	6.8k	
4	45	6.8k	
5	48	52k	
6	26	32k	
7	37	80k	
8	40	49k	
9	24	49k	
10	33	26k	

$$s_x = \sqrt{\frac{\sum x_i^2 - \frac{(\sum x_i)^2}{n}}{n-1}} = \sqrt{\frac{15629 - \frac{136161}{10}}{10-1}} = 14.96$$

$$s_y = \sqrt{\frac{\sum y_i^2 - \frac{(\sum y_i)^2}{n}}{n-1}} = \sqrt{\frac{19902.26 - \frac{19902.26}{10}}{10-1}} = 19.902.26$$

$$\beta_1 = (0.99) \frac{19902.26}{14.96} = 1317.06$$

$\rightarrow y$ Mean

$$B_0 = \bar{y} - B_1 \bar{x}$$

$\rightarrow x$ Mean

$$B_0 = (4200) - (1317.06)(36.9)$$

$$\boxed{B_0 = 6499.51}$$

$$\boxed{\hat{y} = 1317.06x + 6499.51}$$

For Ex:

Predict yearly income of someone who is 33 yr old

$$\hat{y} = (1317.06)(33) + 6499.51$$

$$= 36963.47$$

Correlation vs Causation

Causation means that one variable causes something to happen to another variable.

something to happen is another variable.

To say that two things are correlated is to say that they share some kind of relationship

In order to imply causation, a true experiment must be performed where subjects are randomly assigned to different conditions.

Parameter, statistic & sampling error

Parameter: A characteristic that describes a population
Because it's not possible to measure whole population, parameters are most often estimated.

Ex Population mean, μ , Population variance, σ^2

For instance we have to measure mean of height of people on earth that could be a parameter but that's impossible so why we measure statistic.

A characteristic that describes a sample is called a statistic.

Statistics are most often used to estimate the value of unknown parameters.

Ex Sample Mean: \bar{x} Sample Variance: s^2

Sampling error: Is any difference that exists between statistic & its corresponding parameter.

Ex: We use statistic to find an estimate population mean i.e. 6.8. However when do manually we find it to 7.0 so this dispersion is called sampling error. $7.0 - 6.8 = 0.2 \rightarrow$ Sampling error is 0.2

Distribution of Sample Mean (statistic)

Is a probability distribution for all possible values of a sample mean computed from a sample of size n .

Ex: A statistics class has 6 students, ages displayed below. Construct a sampling distribution of the mean of age for samples ($n=2$)

Ages: 18, 18, 19, 20, 20, 21

Sample Sample Mean

18, 18 18

18, 19 18.5

18, 20 19

18, 20 19

18, 21 20

18, 19 18.5

18, 20 19

18, 21 19

19, 20 19.5

19, 20 19.5

19, 20 19.5

19, 21 20

20, 20 20

20, 21 20.5

20, 21 20.5

All Possible
Combination
& samples of
size $n=2$

Sample mean	Frequency	Probability
18	1	1/15
18.5	2	2/15
19	4	4/15
19.5	3	3/15
20	3	3/15
20.5	2	2/15
	15	

$$\text{Mean} = \frac{18+18+19+20+20+21}{6}$$

$$M = 19.33$$

So it's accurate.

From above most sample Mean method we did we will get population mean close to 19 & least population mean is 18 so it's less mean.

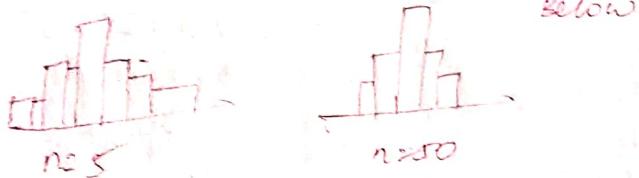
Remember larger our sample size is
the closer our sample mean should
be to the population mean. This is
law of large numbers.

(large sample better
population estimate.)

Standard Error of Mean

There are sampling distribution of

Mean
Below



As sample size increases, S.D decreases
(σ_x)

Larger sample size means less dispersion
Standard deviation of sampling distribution
is also known as the standard error of
the mean

$$\sigma_x = \frac{\sigma}{\sqrt{n}} \rightarrow \text{S.D. of population}$$

Standard error
of mean
↓
sample
size

Central Limit Theory (CLT)

It states that regardless of the shape of
the population distribution, the distribution
of sample means will be approximately
normal.

The distribution of sample mean will become
more normal as its sample size increases

Good rule of thumb: Sample distributions will
be approximately normal if their sample
size is $n=30$ or larger

i.e. Sample Distribution always will be
normal distribution. \Rightarrow (CLT)

Sample Proportion

Let say we want to measure how many
left handed exist in a population obtained
by it is impossible so we take a sample of
500 people & create a proportion.

In our sample, we found 75 are left handed

$$\hat{p} = \hat{P} = \frac{x}{n} = \frac{75}{500} = 0.15 = 15\% \text{ of population is left handed}$$

Distribution of Sample Proportion

We can find out the distribution of the sample proportion if the sample size is less than 5% of the total population size.

If $np(1-p) \geq 10$, the distribution of sample proportion is approximately normal.

S.D. of Sampling Distribution

$$\sigma_{\hat{p}} = \sqrt{\frac{p(1-p)}{n}}$$

Ex: In a sample of 500 individuals, 75 are left handed. Describe distribution of the sample proportion.

Sample size

Yes 500 is less than 5% of entire world population.

$$500(0.15)(1-0.15) \geq 10$$

$63.75 \geq 10$ so normal

$$\text{S.D.} \quad \sigma_{\hat{p}} = \sqrt{\frac{(0.15)(1-0.15)}{500}} = 0.016$$

$$\boxed{\mu_{\hat{p}} = 0.15} \quad \boxed{\sigma_{\hat{p}} = 0.016}$$

Confidence Intervals about the Mean, Population, S.D. known:-

Point estimate

Value of any statistic that estimates the value of a parameter is called a point estimate.

Ex:

$$\bar{x} \rightarrow \mu$$

$$\bar{x} = 2.95 \quad \mu = 3.00$$

We rarely know if our point estimate is correct because it is merely an estimation of the actual value.

Confidence Intervals

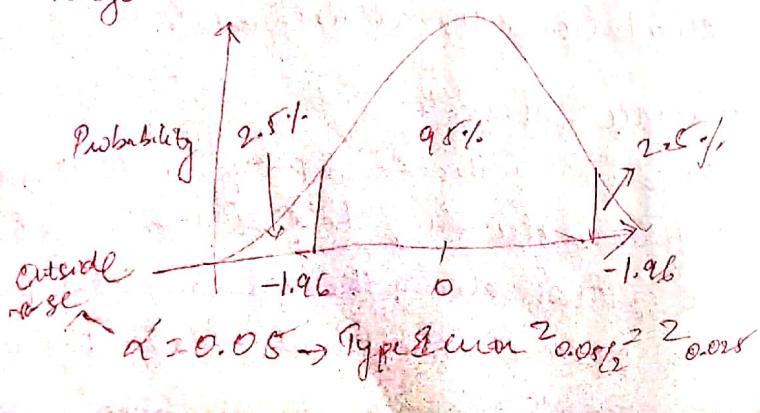
Because of this discrepancy, we construct confidence intervals to help estimate what the actual value of the unknown population mean is $\rightarrow \text{Point Estimate} \pm \text{Margin of Error}$

- first of all how to
- * How confident we want to be with our assessment
 - * Population Standard deviation
 - * How large our sample is.

Confidence Intervals about the mean

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

Let's say we have a distribution like this & we want 95% to be confidence & a we take a chance of 5% of it to be outside the range



Ex: On the verbal section of the SAT, S.D is known to be 100. A sample of 25 test-takers has a mean of 520. Construct a 95% confidence interval about the mean.

$$\bar{x} \pm 2\alpha/2 \frac{\sigma}{\sqrt{n}} \quad 2\alpha/2 = 2.025 \\ \bar{x} = 520 \quad \sigma = 100$$

$$520 \pm (2.025)(100) = 480.8, 559.2$$

↓ ↓
Lower Bound Upper Bound

i.e. I'm 95% confident that mean SAT score is between 480.8 & 559.2

Calculating Required sample size to Estimate Population Mean :-

We can calculate what sample size we will need in order to have a certain margin of error.

$$n = \left(\frac{2\alpha/2 \cdot \sigma}{E} \right)^2$$

Ex: On the verbal section of the SAT S.D is known to be 100. What size sample would we need to construct a 95% confidence interval with a margin of error of 20?

$$n = \left(\frac{1.96}{20} \right)^2 = 96.04$$

A sample size of 97 is needed to create a 95% confidence interval with a margin of error of 20.

Student's t-Distribution

When performing any type of test (or) analysis using Z-score it is required that σ (Population S.D.) is already known.

In real life, this is hardly ever the case. It is almost always impossible for us to know the S.D. of the population from which our sample is drawn.

So how we perform an analysis when we don't know S.D. of population. So we use: Student's t-Distribution

$$Z = \frac{\bar{x} - \mu}{\frac{\sigma}{\sqrt{n}}}; t = \frac{\bar{x} - \mu}{\frac{s}{\sqrt{n}}}$$

Degrees of freedom

As you can see formula is same major difference comes from Degrees of Freedom

Student's t-Distribution has $n-1$ degrees of freedom

Remember we are no longer given population standard deviation. Instead we must estimate it with Sample S.D.

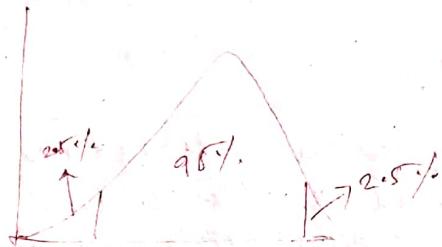
Sample S.D. is itself a random variable. By calculating Sample S.D., it is given a fixed value & thus one less value is free to vary.

D.F change has the probability distribution looks. Probability distribution of t has more dispersion than normal distribution so we use Table.

Confidence Intervals about Mean, Population S.D unknown.

On the verbal section of the SAT, a sample of 25 test takers has a mean of 520 with S.D. of 80. Construct a 95% Confidence interval about the mean.

$$\bar{x} \pm t_{\alpha/2} \left(\frac{s}{\sqrt{n}} \right)$$



$$\alpha = 0.05$$

$$t_{\alpha/2} = t_{0.025} = t_{0.025}$$

because it's a t-test

$$df = n - 1 = 25 - 1 = 24$$

Column

$$t = 2.0639$$

$$520 \pm 2.0639 \left(\frac{80}{\sqrt{25}} \right) = (486.978, 553.022)$$

80% confidence interval
mean \pm 2.0639 $\left(\frac{80}{\sqrt{25}} \right)$
Lower bound: 486.978
Upper bound: 553.022

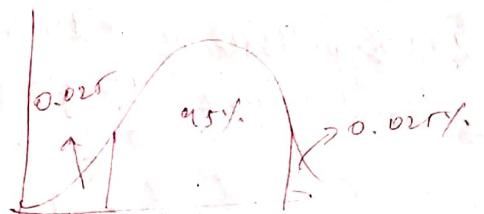
Confidence Interval for Population proportion

In a recent poll of 200 households, it was found that 152 households had at least one computer. Estimate the proportion of households in the population that have at least one computer.

$$\hat{p}_p = \hat{p} = \frac{152}{200} = 0.76$$

$$\hat{\sigma}_p = \sqrt{np(1-p)} = \sqrt{(200)(0.76)(1-0.76)} \geq 10 \\ 36.48 \geq 10$$

$$\hat{p} \pm 2t_{\alpha/2} \sqrt{\frac{\hat{p}(1-\hat{p})}{n}} \\ 2t_{0.025} = 1.96$$



$$0.76 \pm 1.96 \sqrt{\frac{(0.76)(1-0.76)}{200}}$$

$$95\% \text{ confidence: } 0.701 \text{ to } 0.819$$

Lower Bound: 0.701
Upper Bound: 0.819

Type I & Type II :-

Type I: We reject Null hypothesis when in reality it is true.

Type II: we retain Null Hypothesis when in reality it is False.

Type II

Effect size:

A measure of the strength of an effect

After doing statistical analysis if we reject null hypothesis we calculate effect size to determine strength of the effect

$$\text{Cohen's } d = \frac{\text{Mean difference}}{\text{S.D.}}$$

Different sample take in each analysis

$$\text{Cohen's } d' = \frac{(120 - 100)}{30} = 0.66$$

If we have 2 cases

Ex: Medicine A → mean 120
medicine B mean 100

of 'd'
when
0.2 → small effect
0.5 → medium effect
0.8 → large effect

Power:

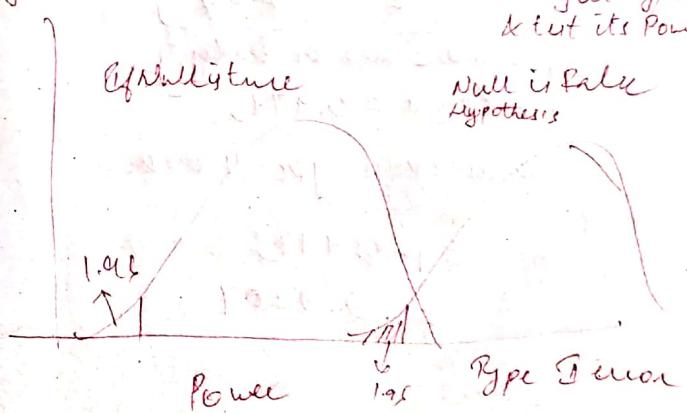
is the probability of correctly rejecting Null hypothesis

Beta: (Type II)

is the probability of incorrectly retaining the null hypothesis

Power Ex:

A new medication that claims to improve typing ability is currently being tested. The average person types at 30 words per minute with a standard deviation of 16. The medication is expected to increase average wpm to 46. A sample of 16 individuals is taken to determine if the medication improves typing ability at $\alpha = 0.05$. Here we assume we reject hypothesis & test its Power



we find standard error

$$\sigma_M = \frac{\sigma}{\sqrt{n}} = \frac{15}{\sqrt{30}} = 4$$

$$z_{\alpha} = 0.05 \text{.} \\ = z \text{ score (com)} \\ (1.96) (\sigma_M)$$

$$= 1.96 (4) \\ = 7.84$$

soil mean is at a distance

of 7.84 from mean

$$30 + 7.84 = (37.84) \rightarrow N$$

with mean

$$z = \frac{\bar{x} - \mu}{\sigma_M} = \frac{37.84 - 46}{4}$$

$$z = -2.04$$

from 2 score area in body

$$\text{Power} = 0.9793$$

Remaining is Type I error

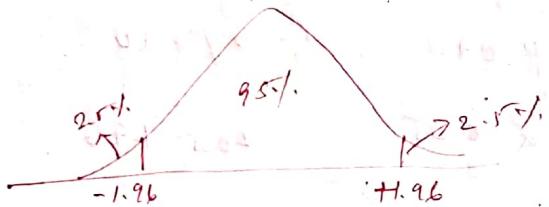
$$= 1 - 0.9793$$

$$= 0.0207$$

One Sample z-test

In population, Avg 218 at 100 & with S.D of 15. A team of scientists wants to test a new medication to see if it has either a positive or negative effect on intelligence or no effect at all. A sample of 30 participants who have taken the medication has a mean of 140. Did the medication affect intelligence?

$$H_0: \mu = 100 ; H_1: \mu \neq 100 ; \alpha = 0.05$$



$$z_{0.05} = 1.96$$

$$z = \frac{\bar{x} - \mu}{\sigma_M} = \frac{140 - 100}{\frac{15}{\sqrt{30}}} = 14.60$$

$$z = 14.60$$

so we reject null H_0

$$z = 14.60 ; p < 0.05$$

One Sample z-test for proportions

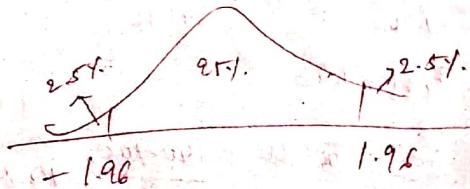
A survey claims that 9 out of 10 doctors recommend aspirin for their patients with head aches. To test this claim, a random sample of 100 doctors is to be obtained. Of these 100 doctors, 12 indicate that they recommend aspirin. $\alpha = 0.05$

9 out of 100 so proportion $\hat{p} = 0.82$

$$H_0: p = 0.90 ; H_1: p \neq 0.90$$

$$\alpha = 0.05$$

$$z_{0.05} = 1.96$$



$$Z_0 = \frac{\hat{p} - p_0}{\sqrt{\frac{p_0(1-p_0)}{n}}} = \frac{0.82 - 0.90}{\sqrt{\frac{0.90(1-0.90)}{100}}} = -2.667$$

So we reject null hypothesis

$$Z = -2.667, P < 0.05$$

F-table

Note: For t-test if its unpaired we use F-test i.e
 $F = \frac{\text{highest } S^2}{\text{smallest } S^2}$ & F.D.F ; If P-value < Critical value
 (reject)

One Sample t-Test (when S.D is unknown)

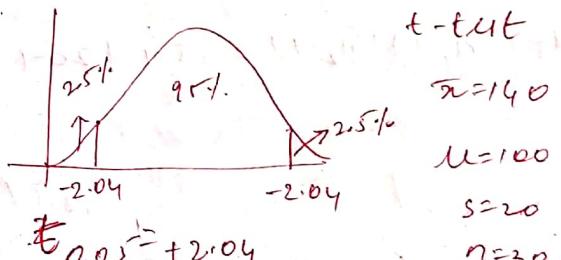
In the population, Avg $\bar{x} = 100$, a group scientist check it has effect. A sample of 30 participants who have taken medication has a mean of 140, S.D of 20, $\alpha = 0.05$

$$H_0: \mu = 100 ; H_1: \mu \neq 100$$

$$\alpha = 0.05$$

$$df = n - 1 = 30 - 1 = 29 \rightarrow \text{row}$$

column



t-tail

$$\bar{x} = 140$$

$$\mu = 100$$

$$S = 20$$

$$n = 30$$

$$t = 0.05 = 2.04$$

$$t = \frac{\bar{x} - \mu}{\frac{s}{\sqrt{n}}} = \frac{140 - 100}{\frac{20}{\sqrt{30}}} = \frac{40}{3.65} = 10.96$$

So we reject Null hypothesis

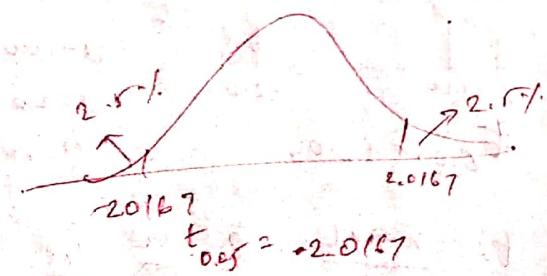
$$t > 10.96, P < 0.05$$

Independent Sample t-test

A statistic teacher wants to compare his two classes to see if they performed any differently on the tests he gave that semester. Class A had 25 students with an average score of 70, S.D of 15, class B had 20 students with an average score of 74 standard deviation 25. Using $\alpha=0.05$ did two classes perform differently or

$$H_0: \mu_A = \mu_B ; H_1: \mu_A \neq \mu_B \quad \alpha=0.05$$

$$df = (n_1 - 1) + (n_2 - 1) = 25 - 1 + 20 - 1 = 43$$



$$t = \frac{\bar{x}_1 - \bar{x}_2}{\sqrt{\frac{(sp)^2}{n_1} + \frac{(sp)^2}{n_2}}} \quad sp^2 = \frac{SS_1 + SS_2}{df_1 + df_2}$$

$$df_1 = n_1 - 1 = 25 - 1 = 24$$

$$df_2 = n_2 - 1 = 20 - 1 = 19$$

$$SS_1 = s_1^2 (df_1) = (15)^2 (24) = 5400$$

$$SS_2 = s_2^2 (df_2) = (25)^2 (19) = 11875$$

$$sp^2 = \frac{5400 + 11875}{24 + 19} = 401.74$$

$$t = \frac{(70 - 74)}{\sqrt{\frac{401.74}{25} + \frac{401.74}{20}}} = \frac{-4}{\sqrt{36.16}} = -0.67$$

so we do not reject null hypothesis

$$t = -0.67, p > 0.05$$

Confidence Intervals for Independent Samples

t-test:

After t-test, confidence can then be constructed to estimate how large the mean difference is we do this only when we reject null hypothesis

Ex: We use

$$\begin{aligned} \bar{x}_1 &= 28.00 & s_1 &= 2.00 & n_1 &= 30 \\ \bar{x}_2 &= 20.00 & s_2 &= 3.00 & n_2 &= 30 \end{aligned}$$

Construct 95% confidence interval for the difference of these two means

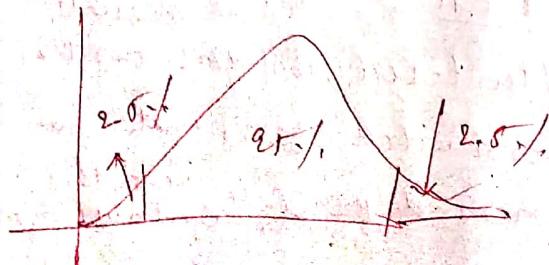
Confidence 2 point estimate & Margin of error

$$\text{Lower bound: } (\bar{x}_1 - \bar{x}_2) - t_{\alpha/2} * \sqrt{\frac{s_1^2}{n_1} + \frac{s_2^2}{n_2}}$$

$$\text{Upper bound: } (\bar{x}_1 - \bar{x}_2) + t_{\alpha/2} * \sqrt{\frac{s_1^2}{n_1} + \frac{s_2^2}{n_2}}$$

$$\text{U.B.: } (28 - 20) - t_{\alpha/2} * \sqrt{\frac{2^2}{30} + \frac{3^2}{30}} = 6.65$$

$$\text{L.B.: } (28 - 20) + t_{\alpha/2} * \sqrt{\frac{2^2}{30} + \frac{3^2}{30}} = 9.35$$



$$df_1 = 30 - 1 = 29 \quad t_{0.025} = 2.0412$$

$$df_2 = 30 - 1 = 29$$

We are 95% confident mean difference between sample 1 & sample 2 is between 6.65 & 9.35

Effect size: (After rejecting H_0)

$$\bar{x}_1 = 28.00, s_1 = 2.00, n_1 = 30$$

$$\bar{x}_2 = 20.00, s_2 = 3.00, n_2 = 30$$

$$s_p^2 = 6.5 \quad df = 58 \quad t = 12.15$$

$$\text{Cohen's } d = \frac{\bar{x}_1 - \bar{x}_2}{\sqrt{s_p^2}} = \frac{28 - 20}{\sqrt{6.5}} = 3.14$$

$d=0.2 \rightarrow \text{small effect}$

$d=0.5 \rightarrow \text{medium effect}$

$d=0.8 \rightarrow \text{large effect}$

3.14 indicate a very large

Our means are likely very different.

$$\text{Other way: } \eta^2 = \frac{t^2}{t^2 + df} = \frac{(12.15)^2}{(12.15)^2 + 58} = .718$$

$\eta^2 \leq 0.01 \rightarrow$ small effect

$\eta^2 \geq 0.09 \rightarrow$ medium effect

$\eta^2 \geq 0.25 \rightarrow$ large effect

0.718 indicates a very large effect

Our means are likely very different

Required is the variance shared between two conditions but it is

similar to Cohen's d' [to measure effect size]

Dependent Samples t-Test:

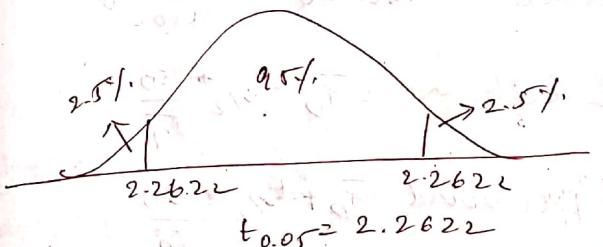
Researchers want to test a new anti-hunger weight loss pill. They have 10 people rate their hunger both before & after taking the pill Does the pill do anything? $\alpha = 0.05$

Before	After	Difference
9	7	-2
10	6	-4
7	5	-2
5	4	-1
1	4	3
9	6	-3
6	7	1
8	5	-3
7	7	0

$$H_0: \mu_{\text{before}} = \mu_{\text{after}} \quad \alpha = 0.05$$

$$H_1: \mu_{\text{before}} \neq \mu_{\text{after}} \quad n = 10$$

$$df = N - 1 = 10 - 1 = 9$$



$$t = \frac{\bar{x}_D}{S_D / \sqrt{n}} = \frac{\bar{x}_D}{S_D} \cdot \sqrt{n}$$

$$\bar{x}_D = \frac{2+4+2+1+3+(-1)+(-2)+1+3+0}{10} = 1.7$$

$$S_D = \sqrt{\frac{\sum x^2 - (\sum x)^2}{n-1}} = \sqrt{\frac{49 - (17)^2}{9}} = 1.49$$

$$t = \frac{1.7}{1.49 / \sqrt{10}} = 3.61$$

$t = 3.61$ so we reject null hypothesis.

Some rejected null hypothesis & find confidence intervals for Dependent sample t-test

$$\bar{x}_D = 1.7 ; s_D = 1.49 \quad n = 10$$

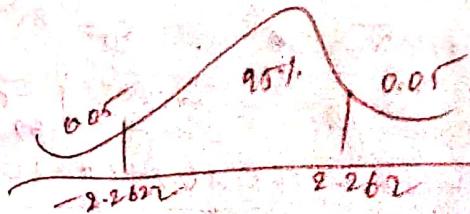
$$t = 3.61$$

95% confidence

Confidence Interval estimate = Margin of error

$$\text{Lower bound} = \bar{x}_D - t_{\alpha/2} * \frac{s_D}{\sqrt{n}}$$

$$\text{Upper bound} = \bar{x}_D + t_{\alpha/2} * \frac{s_D}{\sqrt{n}}$$



$$df = N - 1 = 10 - 1 = 9$$

$$LB = 1.7 - 2.262 * \frac{1.49}{\sqrt{10}} = 0.634$$

$$UB = 1.7 + 2.262 * \frac{1.49}{\sqrt{10}} = 2.76$$

Effect size: allows us to measure the magnitude of real differences & we do this only after rejecting H_0

$$\bar{x}_D = 1.7, s_D = 1.49, n = 10, t = 3.61$$

$$\text{Cohen's } d = \frac{\bar{x}_D}{s_D} = \frac{1.7}{1.49} = 1.14$$

$d = 0.2$ small effect

$d = 0.5$ medium effect

$d = 0.8$ large effect. i.e. means are very different

$$R^2 \text{ squared } R^2 = \frac{t^2}{t^2 + df} = \frac{(3.61)^2}{(3.61)^2 + 9} = 0.59.$$

Amount of variability shared between the groups

$R^2 = 0.01$ small effect

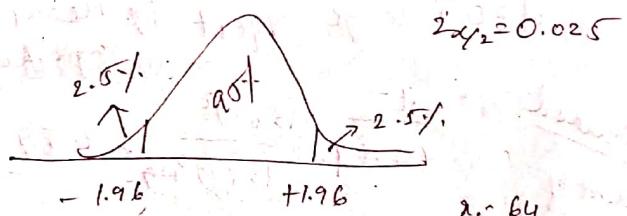
$R^2 = 0.09$ medium effect

$R^2 = 0.25$ large effect

Z-test for Proportions, Two Samples

Researchers want to test the effectiveness of a new anti-anxiety medication. In clinical testing, 64 out of 200 people taking the medication report ~~symptoms of anxiety~~ symptoms of anxiety. Of the people receiving a placebo, 92 out of 200 report symptoms of anxiety. Is the medication any differently than the placebo? Test this claim using alpha $\alpha = 0.05$

$$H_0: p_1 = p_2 ; H_1: p_1 \neq p_2 \quad \alpha = 0.05$$



$$n_1 = 200, \hat{p}_1 = \frac{64}{200} = 0.32, n_2 = 200$$

$$\hat{p}_2 = \frac{92}{200} = 0.46$$

$$Z = \frac{(\hat{p}_1 - \hat{p}_2)}{\sqrt{\hat{p}(1-\hat{p})} \cdot \sqrt{\frac{1}{n_1} + \frac{1}{n_2}}}$$

$$\hat{p} = \frac{x_1 + x_2}{n_1 + n_2}$$

$$\hat{p} = \frac{64 + 92}{200 + 200} = 0.39$$

$$Z = \frac{0.32 - 0.46}{\sqrt{(0.39)(1-0.39)} \cdot \sqrt{\frac{1}{200} + \frac{1}{200}}} = -2.869$$

$$Z = -2.869$$

so we reject null hypothesis.

$$Z = -2.869, p < 0.05$$

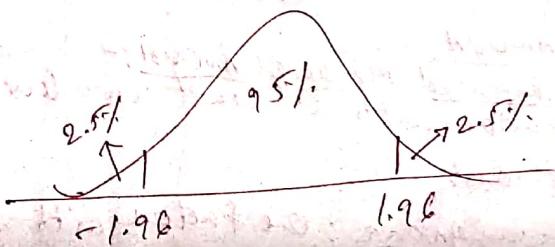
Confidence Interval \rightarrow 95% confidence

$$\alpha = 5\%, n_1 = 200, n_2 = 200$$

$$\hat{p}_1 = \frac{64}{200} = 0.32; \hat{p}_2 = \frac{92}{200} = 0.46$$

$$L.B = (\hat{p}_1 - \hat{p}_2) - 2\sqrt{\frac{\hat{p}_1(1-\hat{p}_1)}{n_1} + \frac{\hat{p}_2(1-\hat{p}_2)}{n_2}}$$

$$U.B = (\hat{p}_1 - \hat{p}_2) + 2\sqrt{\frac{\hat{p}_1(1-\hat{p}_1)}{n_1} + \frac{\hat{p}_2(1-\hat{p}_2)}{n_2}}$$



$$L.B = (0.46 - 0.32) - 1.96 \sqrt{\frac{(0.46)(1-0.46)}{200} + \frac{(0.32)(1-0.32)}{200}}$$

$$U.B = (0.46 - 0.32) + 1.96 \sqrt{\frac{(0.46)(1-0.46)}{200} + \frac{(0.32)(1-0.32)}{200}}$$

$$L.B = 0.045 ; U.B = 0.235$$

We are 95% confident mean difference between two proportion is between 0.045 & 0.235.

T ANOVA

Is a statistical method used to compare the means of two or more groups

* Factor (Variable) Ex: Gender

* Levels (Categorical) Ex: ~~Male~~ Male & Female

Example: Factor - "Dosage"

Levels - "0mg, 5mg, 10mg"

Types of ANOVA

Repeated Measures ANOVA:-
* One factor with at least two levels

dependent

* One way Anova: One factor with at least two levels one independent

One way Anova Example

0mg	50 mg	100 mg
9	3	4
6	2	3
:	:	2
1	1	1

Repeated Measures ANOVA Example

Day 1	Day 2	Day 3
9	7	4
8	6	3
:	:	2
1	1	1

Factorial ANOVA:- Two or more factors (each with at least two levels), levels can be either independent, dependent or both mixed

	Day 1	Day 2	Day 3
Men	9	7	4
	8	6	3
	7	6	2
Women	1	1	1

Assumptions in ANOVA

1. Normality of Sampling Distribution of Means

→ Distribution of sample means is normally distributed

2. Independence of Errors

Errors between cases are independent of one another.

3. Absence of Outliers

Outlying scores have been removed from the data set

4. Homogeneity of Variance

Population variances in different levels of each independent variable are equal

Hypothesis in ANOVA

ANOVA with one factor ("A", three levels)

$$H_0: \mu_{A_1} = \mu_{A_2} = \mu_{A_3} \quad \leftarrow \text{Main effect}$$

$$H_1: \text{not all means are equal}$$

ANOVA with two factors (A & B, each 3 levels)

$$H_0: \mu_{A_1} = \mu_{A_2} = \mu_{A_3}$$

$$H_1: \text{not all means are equal}$$

$$H_0: \mu_{B_1} = \mu_{B_2} = \mu_{B_3}$$

$$H_1: \text{not all means are equal}$$

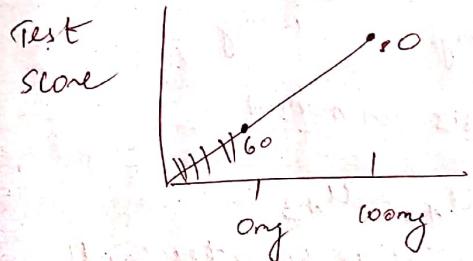
$$H_0: \text{Interaction absent}$$

$$H_1: \text{Interaction present}$$

← Interaction effect

Ex: ANOVA with One factor) → Main effect

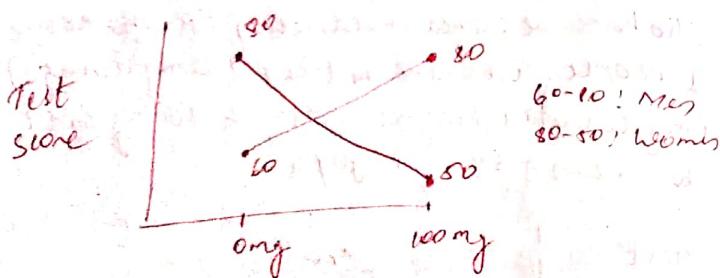
Pretend we are comparing test score of people who have received medication (100mg dosage) & people who have ~~not~~ taken (0mg dosage). 0mg condition has mean 60 & 100mg condition has mean of 80. In graph:



Interaction effect in ANOVA:-

$H_0: \mu_{A1} = \mu_{A2} = \mu_{A3}$	$H_0: \mu_{B1} = \mu_{B2} = \mu_{B3}$	$H_0: \text{Interaction absent}$
$H_1: \text{not all means are equal}$	$H_1: \text{not all means are equal}$	$H_1: \text{Interaction present}$

Here we have a factorial ANOVA with two factors: dosage (ong & 100ng) & gender (Men & Women). In the 0ng dosage condition, men have a mean of 80 while women have a mean of 50. This could be represented in a graph like this:



Pac-Hoc Analysis in ANOVA

$$H_0: \mu_{A1} = \mu_{A2} = \mu_{A3}$$

If we reject the Null hypothesis all we know that there is a difference somewhere among the group

Additional tests called Pac-Hoc test can be done to determine where differences lie

F distribution in ANOVA:-

When doing an ANOVA we calculate an 'F' statistic. It is similar to other statistics such as 'Z' & 't'

$$F = \frac{\text{Treatment Differences}}{\text{Random Differences}}$$

$$F = \frac{\sigma^2_{\text{High}}}{\sigma^2_{\text{Low}}} \quad F \rightarrow \text{is always +ve}$$

If there are no treatment differences (that is, there is no actual effect), we expect F to be 1. If there are treatment differences, we expect F to be greater than 1.

F statistic has its own one tailed distribution much like how the "Z" & "t" statistic have their own separate distributions.

One way ANOVA

One factor with at least two levels,
levels are independent

Ex: Researchers want to test a new anti-anxiety medication. They split participants into 3 conditions (0mg, 50mg, 100mg), then ask them to rate their anxiety level on a scale of 1-10. Are there any differences between 3 conditions using alpha = 0.05

0mg	50mg	100mg
9	7	4
8	6	3
7	6	2
8	7	3
8	8	4
9	7	3
8	6	2

$H_0: \mu_{0\text{mg}} = \mu_{50\text{mg}} = \mu_{100\text{mg}}$ $n=7$
 $H_1:$ not all μ_i 's are equal $N=21$
 $\alpha=0.05$

$$N=21 \quad , \quad n=7$$

$$df_{\text{Between}} = a-1 = 3-1 = 2$$

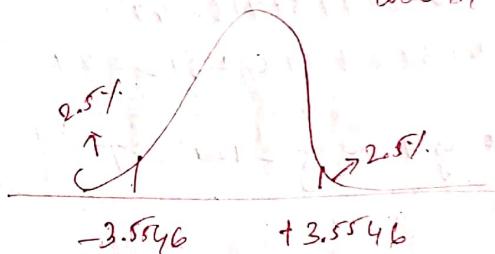
$$df_{\text{Within}} = N-a = 21-3 = 18$$

$$df_{\text{Total}} = N-1 = 21-1 = 20$$

To look up critical value we take

$$(df_{\text{Between}}, df_{\text{Within}}) = (2, 18)$$

↓ ↓
Row column Row



If F is greater than 3.5546 we reject

H_0

	SS	df	MS	F
Between	98.67	2	49.34	
Within	10.29	18	0.57	
Total	108.95	20		

$$SS_{\text{between}} \quad SS_{\text{within}} \quad SS_{\text{Total}}$$

$$SS_{\text{between}} = \frac{\sum (\bar{x}_i)^2}{n} - \frac{T^2}{N}$$

$$SS_{\text{between}} = \frac{(57)^2 + (47)^2 + (21)^2}{7} - \frac{T^2}{N}$$

$$0 \text{mg group: } 9+8+7+8+8+9+8 = 57$$

$$50 \text{mg group: } 7+6+6+7+8+7+6 = 47$$

$$100 \text{mg group: } 4+3+2+3+4+3+2 = 21$$

$$\boxed{0 \text{mg} + 50 \text{mg} + 100 \text{mg} = 125 = T}$$

$$SS_{\text{between}} = \frac{(57)^2 + (47)^2 + (21)^2}{7} - \frac{(125)^2}{21} = 98.67$$

$$SS_{\text{within}} = \sum Y^2 - \sum (\bar{x}_i)^2$$

$$\sum Y^2 = \frac{(57)^2 + (47)^2 + (21)^2}{7}$$

$$\sum Y^2 = 9^2 + 8^2 + 7^2 + 8^2 + 8^2 + 9^2 + 8^2 + 7^2 + 6^2 + 6^2 \\ + 7^2 + 8^2 + \dots + 2^2 = 853$$

$$SS_{\text{within}} = \frac{853 - (57)^2 + (47)^2 + (21)^2}{7}$$

$$= 10.29$$

$$SS_{\text{Total}} = \sum Y^2 - \frac{T^2}{N} = 853 - \frac{(125)^2}{21} \\ = 108.95$$

$$F = \frac{MS_{\text{between}}}{MS_{\text{within}}} \quad MS_{\text{between}} = \frac{98.67}{2} = 49.34$$

$$MS_{\text{within}} = \frac{10.29}{18} = 0.57$$

$$F = \frac{49.34}{0.57} = 86.56$$

As $3.5546 > 86.56$ so we reject

Null Hypothesis

Conclusion:

3 conditions differed significantly
on anxiety level $F(2, 18) = 86.56$,
 $p < 0.05$

Effect size for One way ANOVA

ANOVA Measures means are different from one another if it's large (or) small so we use effect size it helps us to measure

Continue to previous problem

The Most common measure of effect size for a One way ANOVA is

Eta Squared

$$\rightarrow \eta^2 = \frac{SS_{\text{between}}}{SS_{\text{Total}}} = \frac{98.67}{108.95} = 0.91$$

Eta Squared

91% of the total variance is accounted for by the treatment effect.

So we conclude there is a large meaningful difference between 3 groups.

Post Hoc Tests for One way ANOVA:-

So after doing One way ANOVA all we know that groups are different in some way

But Post Hoc Test tell us how they are different H_0 is different (or) H_0 is not different (or) H_{100} is different vice versa & where they are different

Two different tests

* Tukey HSD

* Scheffe

In Tukey test is more liberal than Scheffe test. In Tukey we have greater chance of rejecting H_0 & making Type I error.

In Scheffe we have less chance of rejecting Null hypothesis & less chance of making error.

In Tukey we might find what we want but make a mistake as the chance of

making mistake is high

In Scheffe we have less chance of finding what we want & also less chance of making mistake

Tukey HSD

With this test we are interested in examining mean differences

Omg.	50 mg.	100 mg.
9	7	4
8	6	3
7	6	2
8	7	3
8	8	4
9	7	3
8	6	2

$$\bar{x}_{\text{Omg}} = 8.14 ; \bar{x}_{\text{50mg}} = 6.71$$

$$\bar{x}_{\text{100mg}} = 3.00$$

$$HSD = q \sqrt{\frac{MS_{\text{within}}}{n}}$$
$$= q \sqrt{\frac{0.57}{7}}$$

$n=7$
 q -table

Here we use q table & only df within
k = 3 groups we are comparing &
 $\alpha = 0.05$

$$q = 3.61$$

$$HSD = 3.61 \sqrt{\frac{0.57}{7}} = 1.03$$

$$\bar{x}_{\text{Omg}} = 8.14 ; \bar{x}_{\text{50mg}} = 6.71 ; \bar{x}_{\text{100mg}} = 3.00$$

$$\bar{x}_{\text{Omg}} - \bar{x}_{\text{50mg}} = 8.14 - 6.71 = 1.43$$

$$\bar{x}_{\text{Omg}} - \bar{x}_{\text{100mg}} = 8.14 - 3.00 = 5.14$$

$$\bar{x}_{\text{50mg}} - \bar{x}_{\text{100mg}} = 6.71 - 3.00 = 3.71$$

Any mean above 1.03 are apart i.e.
are different.

Scheffe

$$\text{Comparison} = \frac{MS_{\text{between}}}{MS_{\text{within}}}$$

MS within is same as in One way ANOVA

MS_{between} is different, now Scheffé is a pair wise test so we have to calculate

$$3 \quad SS_{\text{between's}} \quad (0, 50) \quad (0, 100) \quad (50, 100)$$

(3)

$$SS_{\text{between}} = \frac{\sum (\bar{x}_{ij})^2}{n} - \frac{T^2}{N}$$

For (0, 50)

$$SS_{\text{between}} = \frac{57^2 + 47^2}{7} - \frac{104^2}{14} = 7.14$$

$$MS_{\text{between}} = \frac{SS_{\text{between}}}{df_{\text{between}}} = \frac{7.14}{2} = 3.57$$

$$F_{\text{comparison}} = \frac{3.57}{0.57} = 6.26$$

We are going to use same F value we used in One way ANOVA

(0, 50) $F_{\text{comparison}} > F_{\text{value}}$ so we reject.

Null Hypothesis: to say 0mg & 50mg groups different from one another.

Similarly for (0, 100) & (50, 100)

Repeated Measures ANOVA

One Factor with at least two levels, levels are dependent.

By saying that levels are dependent, it means that they share variability in some way. It's almost same except one additional calculation we must perform to account for this shared variability.

Same example of One way ANOVA but Medication is given once a week & twice another time. Once ~~twice~~ weeks.

$\alpha = 0.05$

Before	week 1	week 2
9	7	4
8	6	3
7	6	2
8	7	3
8	8	4
9	7	3
8	6	2

$$H_0: \mu_{\text{before}} = \mu_{\text{week 1}} = \mu_{\text{week 2}}$$

$$H_1: \text{Not all } \mu's \text{ are equal}$$

$\alpha = 0.05$

$$N=21$$

$$S=7$$

$$df_{\text{Between}} = a-1 = 3-1 = 2$$

$$df_{\text{within}} = N-a = 21-3 = 18$$

$$df_{\text{subject}} = S-1 = 7-1 = 6$$

$$df_{\text{error}} = df_{\text{within}} - df_{\text{subject}} = 18 - 6 = 12$$

$$df_{\text{Total}} = N-1 = 21-1 = 20$$

To look up critical value

$$df_{\text{between}} = a-1 = 3-1 = 2$$

$$df_{\text{error}} = df_{\text{within}} - df_{\text{subject}}$$

$$= 18 - 6 = 12$$

$$(df_{\text{between}}, df_{\text{error}}) = (2, 12)$$

~~row~~ column
(column, row)

$$F_{\text{value}} = 3.8853$$

If F is greater than 3.8853 we reject null hypothesis

	SS	d.f.	MS	F
Between	98.67	2	49.34	224.27
within	108.96	18		
subject	7.62	6		
Error	2.67	12	0.22	
Total	108.96	20		

$$SS_{\text{between}} = \frac{\sum (\sum a_{ij})^2}{S} - \frac{T^2}{N}$$

$$\text{Before group} = 9+8+7+8+18+19+8 = 57$$

$$\text{week 1 group} = 7+6+6+7+8+7+6 = 47$$

$$\text{week 2 group} = 4+3+2+3+4+3+2 = 21$$

$$\begin{aligned} SS_{\text{between}} &= \frac{(57)^2 + (47)^2 + (21)^2}{21-3} - \frac{T^2}{N} \\ &= \frac{(57)^2 + (47)^2 + (21)^2 - (125)^2}{7} - \frac{125^2}{21} \\ &= 98.67 \end{aligned}$$

$$SS_{\text{within}} = \sum Y^2 - \frac{\sum (\sum a_{ij})^2}{S}$$

$$\sum Y^2 = 9^2 + 8^2 + \dots + 2^2 = 853$$

$$SS_{\text{within}} = 853 - \frac{(57)^2 + (47)^2 + (21)^2}{7}$$

$$= 10.29$$

$$SS_{\text{subject}} = \frac{\sum (\sum S_i)^2}{a} - \frac{T^2}{N}$$

$$\text{Subject one: } 9+7+4=20$$

$$\text{Subject Two: } 8+6+3=17$$

$$\text{Subject seven: } 8+6+2=16$$

$$SS_{\text{subject}} = \frac{(20)^2 + (17)^2 + (15)^2 + (18)^2 + (20)^2 + (19)^2}{12}$$

$$= \frac{(125)^2}{24} = 7.62$$

$$\text{Error} = \text{Within-Subjects} = 10.29 - 7.62 \\ = 2.67$$

$$\text{Total} = \text{Between} + \text{Within} \\ = 98.67 + 10.29 = 108.96$$

$$F = \frac{MS_{\text{between}}}{MS_{\text{error}}}$$

$$MS_{\text{between}} = \frac{98.67}{2} = 49.34$$

$$MS_{\text{error}} = \frac{2.67}{12} = 0.22$$

$$F = \frac{49.34}{0.22} = 224.27$$

$$F\text{ value} = 3.8853$$

$$F = 224.27$$

so we reject Null Hypothesis

Conclusion:

3 conditions differed significantly on anxiety level, $F(2, 12) = 224.27$
 $P < 0.05$

Factorial ANOVA, Two independent factors

Two factors with at least two levels each, levels are independent.

Factorial ANOVA (with independent factors) is kind of like the One way ANOVA, except now you're dealing with more than one independent variable.

Researchers want to test a new anti anxiety medication. They measure anxiety of 36 participants on different dosages of the medication: 0mg, 50mg, 100mg. Participants are also divided based on what school they are attending. Which researcher hypothesis will also effect anxiety levels. Anxiety is rated on a scale of 1-10, with 10 being "high anxiety" & 1 being "low anxiety". Use alpha = 0.05 to conduct your analysis.

	0mg	50mg	100mg
High school students	3 4 5 3 4 3	5 4 3 5 5 5	7 8 7 8 7 7
College students	1 2 1 1 2	5 4 3 5 4	9 8 9 8 7 9

$$H_0: \mu_{high\ school} = \mu_{college}$$

$$H_1: \mu_{high\ school} \neq \mu_{college}$$

$$H_0: \mu_{0mg} = \mu_{50mg} = \mu_{100mg}$$

$$H_1: \text{not all dosage means}$$

H_0 : An interaction is absent

H_1 : An interaction is present

$$\alpha = 0.05$$

$$df_{school(A)} = a-1 = 2-1 = 1$$

$$df_{dosage(B)} = b-1 = 3-1 = 2$$

$$df_{school \times dosage(A \times B)} = (a-1)(b-1) = (2-1)(3-1) = 2$$

$$df_{error} = N - ab = (36) - (2)(3) = 30$$

$$df_{Total} = N - 1 = 36 - 1 = 35$$

in F table
 School ($df_{school}(A)$, df_{evar}): $(1, 50)$ c.v = 4.17

Dosage ($df_{dosage}(B)$, df_{evar}): $(2, 50)$ c.v = 3.32

School x Dosage ($df_{school \times dosage}(AB)$, df_{evar})
 $= (2, 50)$
 c.v = 3.32

3' hypothesis, so we have 3 decision rules

[School] \rightarrow F is greater than 4.17, reject Null hypothesis

[Dosage] \rightarrow F is greater than 3.32, reject Null hypothesis

[Interaction] \rightarrow F is greater than 3.32, reject null hypothesis.

	SS	df	MS	F
School	2.25	1	2.25	4.166
Dosage	175.17	2	87.59	162.20
Interaction	17.16	2	8.58	15.89
Evar		30	0.54	
Total	210.75	35		

SS school SS dosage SS school x dosage SS evar

Total

$$SS_{school} = \frac{\sum (\sum a_{ij})^2 - T^2}{(b)(n)} \quad \frac{N}{N}$$

$$\text{High school sum} = 3+4+5+3+4+3+5+4+3 \\ + 5+5+5+7+8+7+8+7+7 = 93$$

$$\text{College sum} = 1+2+1+4+1+2+1+4+3+5+4 \\ + 9+8+9+8+7+9 = 84$$

$$SS_{school} = \frac{(93)^2 + (84)^2}{(3)(6)} - \frac{(177)^2}{36} = 2.25$$

$$SS_{dosage} = \frac{\sum (\sum b_{ij})^2 - T^2}{(a)(n)} \quad \frac{N}{N} = \frac{(30)^2 + (33)^2 + (64)^2}{(2)(6)} - \frac{(177)^2}{36} = 175.17$$

$$\text{Omg sum} = 3+4+5+3+4+3+1+2+1+1+1+2 = 30$$

$$\text{Song sum} = 5+4+3+5+5+5+4+3+5+5+4 = 53$$

$$\text{100mg sum} = 7+8+7+8+7+7+9+8+9+8+7+9 = 94$$

Test statistic

$SS_{\text{school} \times \text{dosage}}$ (Interaction)

$$= \frac{\sum (\sum a_i b_i)^2}{n} - \frac{\sum (\sum a_i)^2}{(b)(n)} - \frac{\sum (\sum b_i)^2}{(a)(n)} + \frac{T^2}{N}$$

$$= \frac{(22)^2 + (8)^2 + (27)^2 + (6)^2 + (4)^2 + (8)^2}{6} - \frac{(3)^2 + (6)^2}{(3)(6)} - \frac{(30)^2 + (53)^2 + (64)^2}{(2)(6)} + \frac{(177)^2}{36} = 17.16$$

$$\text{Cell 1} = 3+4+5+3+4+3 = 22, \text{Cell 4} = 5+4+3+5+5+4 = 26$$

$$\text{Cell 2} = 1+2+1+1+1+2 = 8, \text{Cell 5} = 7+8+7+8+7+7 = 44$$

$$\text{Cell 3} = 5+4+3+5+5+5 = 27, \text{Cell 6} = 9+8+9+8+7+9 = 56$$

$$SS_{\text{total}} = \sum Y^2 = \frac{T^2}{N} = \frac{1081 - (177)^2}{36} = 210.75$$

$$\begin{aligned} \sum Y^2 &= 8^2 + 4^2 + 5^2 + 3^2 + \dots + 1^2 + 2^2 + \dots \\ &\quad + 5^2 + 4^2 + \dots + 5^2 + 4^2 + \dots + 9^2 \\ &= 1081 \end{aligned}$$

All sum

$$SS_{\text{error}} = SS_{\text{total}} - SS_{\text{school}} - SS_{\text{dosage}} - SS_{\text{Interaction}}$$

$$= 16.17$$

$$MS_{\text{school}} = \frac{2.25}{1}$$

~~MS_{school}~~

$$MS = \frac{SS}{df}$$

$$MS_{\text{Interaction}} = \frac{17.16}{2}$$

$$F = \frac{MS_{\text{effect}}}{MS_{\text{error}}}$$

$$FA = \frac{2.25}{0.54} = 4.17$$

$$FB = \frac{57.59}{0.54} = 102.20, FA_{AB} = \frac{5.58}{0.54} = 15.89$$

[School] If F is greater than 4.17, reject Null hypothesis, Our F = 4.166 Retain Null hypothesis

[Dosage] If F is greater than 3.32, reject the Null hypothesis, Our F = 162.20, Reject Null hypothesis

[Interaction] If F is greater than 3.32, reject the Null hypothesis, Our F = 15.89, Reject Null hypothesis.

Conclusion

High school students & college students did not have significantly different anxiety levels. $F(1,30) = 4.166, p > 0.05$. There was a significant difference between 3 different levels of dosage, $F(2,30) = 162.20, P < 0.05$. An interaction effect was also present, $F(2,30) = 15.59, P < 0.05$.

Factorial ANOVA, Two Dependent Factors

Two factors with at least two levels each; Both factors are dependent

Factorial ANOVA (with two dependent factors) is most like an extension of a Repeated Measures ANOVA.

Researchers want to compare the anxiety levels of six individuals at two marital states after they have been divorced & then again after they have been gotten married. Anxiety is measured

at 3 times: week 1, week 2, week 3. Anxiety is rated on a scale of 1-10 with 10 being "high anxiety" & 1 being "low anxiety". $\alpha = 0.05$ to conduct Analysis of variance.

	subject	week 1	week 2	week 3
Divorced	1	3	5	7
	2	4	4	8
	3	5	3	7
	4	3	5	8
	5	4	5	7
	6	3	5	7
Married				
Married	1	1	5	9
	2	2	4	8
	3	1	3	9
	4	1	5	8
	5	1	5	7
	6	2	4	9

"Same person" across different weeks

$$H_0: \mu_{\text{divorced}} = \mu_{\text{married}} \quad H_0: \mu_{\text{week 1}} = \mu_{\text{week 2}} = \mu_{\text{week 3}}$$

$$H_1: \mu_{\text{divorced}} \neq \mu_{\text{married}} \quad H_1: \text{not all week means are equal}$$

$$H_0: \text{an interaction is absent}$$

$$H_1: \text{an interaction is present}$$

$$\alpha = 0.05$$

Degree of Freedom

Source	SS	df	MS	F
Marital Status (A)	2.25	1	2.25	9.00
Error (AxB)	1.25	5	0.25	
week (B)	175.17	2	87.59	93.18
Error (BxS)	9.4	10	0.94	
Interaction (AxB)	17.16	2	8.58	17.57
Error (AxBxS)	4.94	10	0.49	
Error (S)	0.58	5		
Total	210.75	35		

DF_r

$$df_{\text{marital status (A)}} = a - 1 = 2 - 1 = 1$$

$$df_{\text{error (A x S)}} = (a-1)(n-1) = (1)(5) = 5$$

$$df_{\text{week (B)}} = b - 1 = 3 - 1 = 2$$

$$df_{\text{error (B x S)}} = (b-1)(n-1) = (2)(5) = 10$$

$$df_{\text{error (A x B x S)}} = (a-1)(b-1)(n-1) = (1)(2)(5) = 10$$

$$df_{\text{marital status x week (A x B)}}$$

$$= (a-1)(b-1) = (1)(2) = 2$$

$$df_{\text{error (S)}} = n - 1 = 6 - 1 = 5$$

$$df_{\text{Total}} = N - 1 = 36 - 1 = 35$$

We have 3 Hypothesis so 3 decision rules

$$\text{Marital status} \left(\text{df}_{\text{marital status}}, \text{df}_{\text{error}} (A \times S) \right)$$

(A) \rightarrow column
 $= (1, 5)$ \rightarrow Row C.V = 6.61

$$\text{Week} \left(\text{df}_{\text{week}} (B), \text{df}_{\text{error}} (B \times S) \right)$$

Column \rightarrow
 $= (2, 10)$ C.V = 4.10
 \rightarrow Row

$$\text{Marital Status} \times \text{Week} \left(\text{df}_{\text{school} \times \text{week}} (A \times B), \text{df}_{\text{error}} (A \times B \times S) \right)$$

$\rightarrow (2, 10)$ C.V = 4.10

[Marital Status] If F is greater than 6.61 reject the Null hypothesis

[Week] If F is greater than 4.10 reject the Null hypothesis.

[Interaction] If F is greater than 4.10 - reject the Null hypothesis

Test statistic

$$SS_{\text{marital status}} (A), SS_{\text{week}} (B), SS_{\text{school} \times \text{week}} (A \times B)$$

$$SS_{\text{error}} (S), SS_{\text{error} (A \times S)}, SS_{\text{error} (B \times S)}$$

$$SS_{\text{error} (A \times B \times S)}, SS_{\text{Total}}$$

$$SS_{\text{marital status}} = \frac{\sum (\sum a_{ij})^2}{b n} - \frac{T^2}{N}$$

$$\text{Divorced sum} = 3+4+5+\dots+8+7+7 = 93$$

$$\text{Married sum} = 1+2+1+\dots+8+7+9 = 84$$

$$SS_{\text{marital status}} = \frac{(93)^2 + (84)^2 - (71)^2}{(3)(6)} = 2.25$$

$$SS_{\text{week}} = \frac{\sum (\sum b_{ij})^2}{n} - \frac{T^2}{N} = \frac{(30)^2 + (53)^2 + (64)^2 + (77)^2}{(2)(6)} = 175.1$$

$$\text{Week 1 sum} = 3+4+5+3+\dots+1+2 = 30$$

$$\text{Week 2 sum} = 8+4+3+\dots+5+5+4 = 53$$

$$\text{Week 3 sum} = 7+8+7+\dots+7+9 = 94$$

$$SS_{\text{marital status} \times \text{week}} = \frac{\sum (\sum a_{ij} b_{ij})^2}{n} - \frac{\sum (\sum a_{ij})^2}{b n} \cdot \frac{\sum (\sum b_{ij})^2}{n}$$

$$\begin{aligned}
 &= \frac{\sum (\sum a_i b_i)^2}{n} - \frac{(28)^2 + (4)^2}{(3)(6)} - \frac{(30)^2 + (53)^2 + (64)^2}{(2)(6)} + \frac{(177)^2}{36} \\
 &= \frac{(2^2 + 6^2 + 1^2 + 2^2 + 4^2 + 4^2 + 60^2)}{4} " " " = 17.16 \\
 \text{Cell}_1 &= 21.4 + 5.3 + 4.13 = 22, \quad \text{Cell}_4 = 5.4 + 3 + 4.13
 \end{aligned}$$

$$\text{Cell 1} = 8+4+5+3+4+3=22, \text{Cell 4} = 5+4+3+5+5+4 = 26$$

$$\text{Cell 2} = 1+2+1+1+1+2 = 8 \quad \text{Cell 5} = 7+2+1 = 10$$

$$6813 - 511 = 178777777 = 44$$

$$\text{Cell 3} = 3+4+3+5+5+5 = 27, \text{ Cells} = 9.18 + 9.18 + 7 + 9 = 35.0$$

$$\begin{aligned} S^2_{\text{Total}} &= \frac{\sum Y^2 - \bar{T}^2}{N} = \frac{1081 - \bar{T}^2}{36} = \frac{1081 - (17)^2}{36} \\ \sum Y^2 &= 3^2 + 4^2 + \dots + 17^2 = 210.75 \\ &\quad + 1^2 + 2^2 = 1081 \end{aligned}$$

$$SS_{error(A \times s)} = \frac{\sum (\sum a_i \gamma_{ni})^2}{(b)} - \frac{\sum (\sum a_i)^2}{(b)(n)} - \frac{\sum (\sum n_i)^2}{(a)(b)} + \frac{1}{N}$$

$$= \frac{\sum (\sum a_i \gamma_{ni})^2}{b} - \frac{(93)^2 + (84)^2}{(3)(6)} - \frac{\sum (\sum n_i)^2}{(a)(b)} + \frac{177^2}{36}$$

$$= \frac{\sum (\sum a_i \gamma_i)^2}{b} - \frac{(a_3)^2 + (a_4)^2}{(3)(6)} - \frac{\sum (\sum \gamma_i)^2}{(a)(b)} + \frac{(177)^2}{36}$$

$$\text{Subject 1} = 3+5+7+1+5+9 = 30 \quad \text{Subject 2} = 8+5+1+1+1 = 16$$

$$\text{Subject 2} = 4+4+8+2+4+8=30, \text{ Subject 3}=6, \dots$$

$$\text{Subject 3} = 5+3+7+1+3+9 = 28, \text{ Subject 1} = 8+5+1+2+1 = 17$$

$$a_1 n_1 = 3 + 5 + 7 = 15$$

$$\alpha_2 \cdot n_2 = 4+4+8 = 16$$

$$a_3 n_3 = 5 + 3 + 7 = 15$$

$$\begin{matrix} & 7 \\ \therefore & 2+4+9=15 \end{matrix}$$

$$= \frac{(15)^2 + (16)^2 + (15)^2 + (16)^2 + (16)^2 + (15)^2 + (15)^2 + (15)^2 + (15)^2 + (15)^2}{10(15)^2 + (15)^2}$$

$$= \frac{9(3+7)(84)^2}{(3)(6)} - \frac{(30)^2 + (30)^2 + (28)^2 + (30)^2 + (29)^2 + (30)^2}{(2)(3)}$$

$$+ \frac{(177)}{36} = 1.25$$

ss
error ($B \times s$)

$$= \frac{\sum (\sum b_i \eta_i)^2}{a} - \frac{\sum (\sum b_i)^2}{(a)(b)} - \frac{\sum (\sum \eta_i)^2}{(a)(b)} + \frac{T^2}{N}$$

$$= \frac{\sum b_{ini}}{a} - \frac{(30)^2 + (53)^2 + (94)^2}{(2)(6)} - \frac{(30)^2 + (89)^2 + (21)^2 + (34)^2 + (29)^2 + (50)^2}{(6)(7)}$$

$$+ \underbrace{(177)}_{36}^2$$

$$\begin{aligned}
 & \text{SS}_{\text{error}} = \frac{(4^2 + 5^2 + 6^2 + 8^2 + 10^2 + 16^2 + 15^2 + 9^2 - 16^2)}{8} \\
 & = \frac{(30) + (25) + (36) + (64) + (100) + (256) + (225) + (81) - (256)}{8} \\
 & = \frac{677}{8} = 84.625
 \end{aligned}$$

$$\begin{aligned} \text{Subject 1} &= 3+1=4 ; \text{subject 2} = \frac{36}{4}=9 \\ \text{Subject 1} &= 5+5=10 ; \text{subject 2} = 4+4=8 \\ \text{Subject 1} &= 7+9=16 ; \text{subject 2} = 8+8=16 \end{aligned}$$

Subject 6 = 5

Subject 6 = 9

Subject 6 = 16

$$SS_{\text{error}(S)} = \frac{\sum (\varepsilon \eta_i)^2}{(a)(b)} - \frac{T^2}{N}$$

$$= \frac{(30)^2 + (30)^2 + (28)^2 + (30)^2 + (29)^2 - (147)^2}{36} = 0.58$$

$$\text{Error}(AXBXS) = \text{Total} - \text{rest we found} (CSS)$$

$$MS = \frac{SS}{dt} ; F = \frac{MS_{\text{effect}}}{MS_{\text{error}}}$$

$$F_A = \frac{2.25}{0.25} = 9.00; F_B = \frac{87.59}{0.94} = 93.19$$

$$F_{AXB} = \frac{8.58}{0.49} = 17.51$$

[mantal status]: For greater than 6.61
our $F = 9.00$, reject H_0

[Week] : F is greater than 4.10, reject null
 $\text{One } F = 93.18$,

[Interaction] : F_{11} is greater than 4.10 reject
Our $F = 17.50$, Reject H_0

Conclusion: Anxiety levels differed significantly for divorced & other remained individuals, $F(1,5) = 9.00, p < 0.05$. There was a significant difference between three different weeks, $F(2,10) = 93.18, p < 0.05$. An interaction effect was also present, $F(3,10) = 17.51, p < 0.05$

Factorial ANOVA, Two Mixed Factors:

Two factors with at least two levels such One factor is independent, while the other factor is dependent

This Factorial ANOVA is combination of One way ANOVA & Repeated Measures ANOVA.

	Week 1	Week 2	Week 3
High school students	3	5	7
	4	4	8
	5	3	7
	3	5	8
College students	3	5	7
	2	4	9
	1	3	8
	2	4	7

$H_0: \mu_{\text{highschool}} = \mu_{\text{college}} ; H_0: \mu_{\text{week}_1} = \mu_{\text{week}_2} = \mu_{\text{week}_3}$
 $H_1: \mu_{\text{highschool}} \neq \mu_{\text{college}} ; H_1: \text{not all means are equal}$

$H_0: \text{an interaction is absent}$

$H_1: \text{an interaction is present}$

$$\alpha = 0.05$$

	SS	df	MS	F
School (A)	2.25	1	2.25	12.5
Error (S/A)	1.83	10	0.18	
Week (B)	175.17	2	87.59	121.65
Interaction (AxB)	17.16	2	8.58	11.92
Error (BxS/A)	14.34	20	0.72	
Total	210.75	35		

$$df_{\text{School}}(A) = a - 1 = 2 - 1 = 1$$

$$df_{\text{Error}}(S/A) = (a)(n-1) = (2)(5-1) = 10$$

$$df_{\text{Week}}(B) = b - 1 = 3 - 1 = 2$$

$$df_{\text{School} \times \text{Week}}(AxB) = (a-1)(b-1) = (2)(2) = 4$$

$$df_{\text{Error}}(BxS/A) = (a)(b-1)(n-1) = (2)(2)(5-1) = 20$$

$$df_{\text{Total}} = N - 1 = 36 - 1 = 35$$

3 hypothesis, F table, 3 degrees of freedom

School

$$\left(\text{df}_{\text{School}(A)}, \text{df}_{\text{Error}(\frac{S}{A})} \right) = (1, 10) \quad \text{CV} = 4.35$$

↑ column
↓ Row

$$\text{week} \left(\text{df}_{\text{Week}(B)}, \text{df}_{\text{Error}(\frac{B \times S}{A})} \right) : (2, 20) \quad \text{CV} = 3.49$$

$$\text{School} \times \text{Week} \left(\text{df}_{\text{School} \times \text{Week}(A \times B)}, \text{df}_{\text{Error}(\frac{B \times S}{A})} \right) : (2, 20) \quad \text{CV} = 3.49$$

[School]: If F is greater than 4.35, reject null hypothesis

[Week]: If F is greater than 3.49, reject H₀

[Interaction]: If F is greater than 3.49, reject H₀

$SS_{\text{School}(A)}$, $SS_{\text{Error}(\frac{S}{A})}$, $SS_{\text{Week}(B)}$,

$SS_{\text{School} \times \text{Week}(A \times B)}$, $SS_{\text{Error}(\frac{B \times S}{A})}$

SS_{Total}

$$SS_{\text{School}} = \frac{\sum (\sum a_i)^2}{(b)(n)} - \frac{T^2}{N} = \frac{(93)^2 + (84)^2 - (177)^2}{(3)(6)} = 36 = 2.25$$

$$\text{High school sum} = 3+4+5+\dots+7+7 = 93$$

$$\text{College sum} = 1+2+1+1+1+\dots+7+9 = 84$$

$$SS_{\text{Week}} = \frac{\sum (\sum b_i)^2}{(a)(n)} - \frac{T^2}{N} = \frac{(20)^2 + (53)^2 + (44)^2 - (177)^2}{(2)(6)} = \frac{36}{36} = 175.17$$

$$\text{Week 1 sum} = 3+4+5+3+\dots+1+2+\dots+2 = 30$$

$$\text{Week 2 sum} = 5+4+3+5+\dots+8+5+4 = 53$$

$$\text{Week 3 sum} = 7+8+7+8+\dots+8+7+9 = 94$$

$$SS_{\text{School} \times \text{Week}} = \frac{\sum (\sum a_i b_i)^2}{(n)} - \frac{\sum (\sum a_i)^2}{(b)(n)} - \frac{\sum (\sum b_i)^2}{(a)(n)} + \frac{T^2}{N}$$

$$\frac{\sum (\sum a_i b_i)^2}{n} - \frac{(93)^2 + (84)^2}{(3)(6)} - \frac{(30)^2 + (63)^2 + (64)^2 + (77)^2}{(2)(6)} - \frac{177^2}{36}$$

$$\frac{(22)^2 + (5)^2 + (27)^2 + (26)^2 + (64)^2 + (50)^2}{(6)} - \frac{(93)^2 + (64)^2}{(3)(6)} - \frac{(30)^2 + (63)^2}{(2)(6)} - \frac{177^2}{36} = 17.16$$

$$\text{Cell 1} = 3+4+5+3+4+3 = 22$$

$$\text{Cell 2} = 1+2+1+2+1+1 = 8$$

$$\text{Cell 3} = 27$$

$$SS_{\text{Total}} = \sum y^2 - \bar{y}^2$$

$$\text{Cell 4} = 26$$

$$\sum y^2 = 3^2 + 4^2 + \dots + 7^2 + 9^2$$

$$\text{Cell 5} = 44$$

$$= 1081$$

$$\text{Cell 6} = 50$$

$$= 1081 - \frac{177^2}{36} = 210.75$$

$$SS_{\text{error}(B)} = \frac{\sum (\sum a_i n_i)^2}{(b)} = \frac{\sum (\sum a_i)^2}{(b)(n)}$$

$$\frac{\sum (\sum a_i n_i)^2}{(b)} - \frac{(93)^2 + (84)^2}{(3)(6)}$$

$$a_1 n_1 = 3+5+7 = 15$$

$$a_2 n_2 = 4+4+8 = 16$$

$$a_3 n_3 = 15$$

$$\therefore a_1 n_2 = 18$$

$$\begin{aligned} & (15)^2 + (6)^2 + (15)^2 + (6)^2 + (16)^2 + (5)^2 + (15)^2 + (4)^2 + (13)^2 \\ & - (14)^2 + (13)^2 + (15)^2 \\ & \frac{1}{3} - \frac{(93)^2 + (84)^2}{(3)(6)} = 1.83 \end{aligned}$$

$$SS_{\text{Total}} - SS_{\text{School}} - SS_{\text{Error}(S/A)} = SS_{\text{Error}(B \times S/A)}$$

$$SS_{\text{Error}(B \times S/A)} = 14.34$$

$$MS = \frac{SS}{df} ; F = \frac{MS_{\text{effect}}}{MS_{\text{error}}}$$

$$F_A = \frac{2.25}{0.18} = 12.5 ; F_B = \frac{87.59}{0.72} = 121.65$$

$$F_{AXB} = \frac{8.58}{0.72} = 11.92$$

[School] F is greater than 4.35, reject H_0
Our $F = 12.5$

[Week] F is greater than 3.49, reject H_0
Our $F = 121.65$

[Interaction] F is greater than 3.49, reject H_0
Our $F = 11.92$

Conclusion:

High school students & college students had significantly different anxiety levels. $F(1,10) = 12.5$, $p < 0.05$. There was a significant difference between 3 different weeks, $F(2,20) = 121.65$, $p < 0.05$. An interaction effect was also present, $F(2,20) = 11.92$, $p < 0.05$.

Chi-square Test for Goodness of fit

χ^2 Test for Goodness of fit tests claims about population proportion.

It is non parametric test that is performed on categorical (nominal or ordinal) data.

Ex: In the 2000 U.S. Census, the ages of individuals in a small town were found to be following.

Less than 18	18 - 35	Greater than 35
20%	30%	50%

In 2010 ages of $n = 500$ individuals were sampled. Below are the results:

Less than 18	18 - 35	Greater than 35
121	288	91

Use $\alpha = 0.05$, would you conclude that the population distribution of ages has changed in the last 10 years?

	Less than 18	18 - 35	Greater than 35
Expected	20%	30%	50%

If $n = 500$

	Less than 18	18 - 35	Greater than 35
Observed	121	288	91
Expected	$500 \times 20\% = 100$	$500 \times 30\% = 150$	$500 \times 50\% = 250$

H_0 : Data meets expected distribution

H_1 : Data does not meet expected distribution.

$k = \text{no of groups}$

$$df = k - 1 = 3 - 1 = 2$$

In χ^2 table $\alpha = 0.05$ & $df = 2$

If χ^2 is greater than 5.99, reject H₀.

$$\chi^2 = \sum \frac{(f_o - f_e)^2}{f_e}$$

$$\chi^2 = \frac{(121-100)^2}{100} + \frac{(288-200)^2}{150} + \frac{(91-250)^2}{250}$$

$$\chi^2 = 232.494$$

So we reject Null hypothesis.

Conclusion:

Age of 2010 population are different than those expected based on 2000 population.

Chi-Square Test for Independence

- * It evaluates relationship between two variables.
- * It is a non parametric test that is performed on categorical (nominal or ordinal) data.

Ex: 500 elementary school boys & girls are asked which is their favorite color: Blue, green or pink. Results are shown below:

Observed	Blue	Green	Pink	Column Total
Boys	100	150	20	300
Girls	20	30	180	200
Row Total	120	180	200	N=500

Since $\alpha = 0.05$, would you conclude that there is relationship between gender & colour favorite.

H₀: Gender & color not related

H₁: Gender & color are related.

$$\alpha = 0.05$$

$$df = (\text{rows}-1)(\text{columns}-1)$$

$$= (2-1)(3-1)$$

$$= 2$$

$$\alpha = 0.05 \text{ & } df = 2$$

If χ^2 is greater than 5.99
using χ^2 table *reject H₀*

$$\chi^2 = \sum \frac{(f_o - f_e)^2}{f_e}; f_e = \frac{f_c}{n}$$

Expected	Blue	Green	Pink	
Boys	72	108	120	300
Girls	48	72	80	200
	120	180	200	N=500

$$(Boys, Blue) = \frac{(120-72)^2}{500} = 72$$

Similarly all values

$$\begin{aligned} \chi^2 &= \frac{(100-72)^2}{72} + \frac{(20-48)^2}{48} + \frac{(150-108)^2}{108} + \frac{(80-72)^2}{72} \\ &+ \frac{(20-120)^2}{120} + \frac{(180-80)^2}{80} = 276.389 \end{aligned}$$

If χ^2 is greater than 5.99, reject H_0
 $\chi^2 = 276.39$

so we reject Null Hypothesis.

Conclusion: In population there is a relationship between gender & favorite colour.

Mann Whitney U Test:

It is a version of the independent samples t-Test that can be performed on ordinal (ranked) data

Ordinal data is table, $\alpha = 0.05$, Is there a difference between treatment A & treatment B

Treatment - A	28	31	36	35	32	33
Treatment - B	12	18	19	14	20	19

H_0 : No difference between marks of two treatment

H_1 : There is difference " " " "

$\alpha = 0.05$

It follows 'Z' distribution if sample size is greater than 20

Two tailed Test

If Z is less than -1.96 or > 1.96 reject

H_0

Treatment A	Treatment B	Rank	Score	Sample Points	
				1	2
28	12	1	12	(B)	6
31	18	2	14	(B)	6
36	19	3	18	(B)	6
35	14	4.5	19	(B)	6
32	20	6	20	(B)	6
33	19	7	28	(A)	0
		8	31	(A)	0
		9	32	(A)	0
		10	33	(A)	0
		11	35	(A)	0
		12	36	(A)	0

$$U_A = 0 + 0 + 0 + 0 + 0 + 0 = 0$$

$$U_B = 6 + 6 + 6 + 6 + 6 + 6 = 36$$

$U = 0$ (smaller value among U_A & U_B is 0)

Usually, if the more different groups are

$U = 0$

$$Z = \frac{U - \frac{n_A n_B}{2}}{\sqrt{\frac{n_A n_B (n_A + n_B + 1)}{12}}} = \frac{0 - \frac{(6)(6)}{2}}{\sqrt{\frac{(6)(6)(6+6+1)}{12}}} = \boxed{Z = -2.88}$$

so we reject H_0

Conclusion

There is difference between ranks of the two treatments $Z = -2.88, P < 0.05$

Wilcoxon Signed-Ranks Test :- (U-Test)

It is a version of dependent sample t-test performed on ordinal (ranked) data.

$\alpha = 0.05$

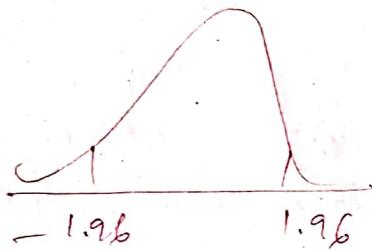
Before	28	17	36	35	32	33
After	12	31	19	14	20	19

H_0 : No difference between treatments

H_1 : Difference $\neq 0$

$\alpha = 0.05$

- At least 20 samples were applying
Z distribution.



Before After Difference Rank

Before	After	Difference	Rank
28	12	16	4
17	31	-14	1
36	19	17	5
35	14	21	6
32	20	12	2
33	19	14	3

Lower \rightarrow Lower number, Higher \rightarrow Higher number rank

(Pointing) Rank addition

$$\Sigma R_+ = 4 + 5 + 6 + 2 + 3 = 20$$

$$\Sigma R_- = 1$$

T is smaller among

$$T = 1$$

$$\Sigma R_+ \text{ & } \Sigma R_-$$

$$Z = T - \frac{n(n+1)}{4} = 1 - \frac{(6)(6+1)}{4}$$

$$\sqrt{\frac{n(n+1)(2n+1)}{24}}$$

$$\sqrt{\frac{(6)(6+1)(2(6)+1)}{24}}$$

$$Z = -1.99$$

so reject H_0

Conclusion: there is difference between before & after groups $Z = -1.99$, $p < 0.05$

Kruskal Wallis Test (H-Test)

Is a version of an independent measure (One way) ANOVA, performed on

Ordinal (ranked) Data

Group 1 Group 2 Group 3

27	20	34
2	8	31
4	14	3
18	36	23
7	21	30
9	22	6
T	39	65
n	6	6

H_0 : No difference between treatments

H_1 : Difference

$\alpha = 0.05$

$k = \text{no of groups}$

$$df = k - 1 = 3 - 1 = 2$$

We use Chi-square Table

If χ^2 is greater than 5.99,
reject H_0

Original score	Rank
2	1
3	2
4	3
6	4
7	5
8	6
9	7
14	8
18	9
20	10
21	11
22	12
23	13
27	14
30	15
31	16
34	17
36	18

With small number lower rank, higher large rank

Test Statistic (H)

$$H = \frac{12}{N(N+1)} \left(\sum \frac{T_i^2}{n} \right) - 3(N+1)$$

$$= \frac{12}{(18)(18+1)} \left(\frac{(39)^2}{6} + \frac{(65)^2}{6} + \frac{(67)^2}{6} \right) - 3(18+1) = 2.854$$

So we reject H_1

Conclusion: significant
There is no difference among 3 groups

$$H = 2.854 (2, N=18), P < 0.05$$

The Friedman Test

As a version of Repeated Measures Anova, performed on ordinal (Ranked) data.

$$\alpha = 0.05$$

Week 1	Week 2	Week 3	Week 1	Week 2	Week 3
27	20	34	2	1	3
2	8	31	1	2	3
14	14	3	2	3	1
18	36	23	1	3	2
7	21	30	1	2	3
9	22	6	2	3	1

R 9 14 13

H_0 : No difference in 3 conditions

H_1 : Rea difference in 3 conditions

$$df = k-1 \quad k: \text{no of groups}$$

$$23-1 = 2 \quad \text{we are comparing}$$

using χ^2 square table

if χ^2 is greater than 5.99

Rank the values in Rows

$$\chi^2_r = \frac{12}{nk(k+1)} \sum R^2 - 3n(k+1)$$

$$= \frac{12}{(6)(3)(3+1)} (9^2 + 13^2 + (4)^2) - (3)(6)(3+1)$$

$$\chi^2_r = 2.33$$

as critical value $\chi^2_{\text{critical}} < \chi^2_{\text{statistic}}$

we reject H_0

Conclusion

there is no difference in 3 groups, $\chi^2 = 2.33$
(2, n=6), $p > 0.05$

Calculating Required Sample size to Estimate Population Proportion:

In household of 200 ex: previous

we found 95% confidence, at least one computer is between 0.701 & 0.819

Confidence = Point estimate \pm Margin error

$$= 0.76 \pm 0.059$$

now what sample size would I need to change margin of error from 0.059 to 0.030 in 95% confidence interval

If we have prior estimate of population proportion exists:

$$n = \hat{P}(1-\hat{P}) \left(\frac{2\alpha/2}{E} \right)^2 = (0.76)(1-0.76) \times \left(\frac{1.96}{0.030} \right)^2 = 778.81$$

If prior estimate of population proportion does not exist,

$$n = 0.25 \left(\frac{2\alpha/2}{E} \right)^2 = (0.25) \left(\frac{1.96}{0.030} \right)^2 = 1067.111 \approx 1068$$

Alpha levels:

If H_0 is rejected when $\alpha = 0.03$,

* What could you say if $\alpha = 0.05$?

* If $\alpha = 0.01$?

