

$$F = \frac{MSC}{MSE}$$

$$, MSC = \frac{SSC}{dfc}, MSE = \frac{SSE}{dfc}$$

$$dfc = C - 1$$

$$\text{d.f.e} = (C-1)(B-1) \quad dfc = (C-1)(B-1)$$

$$MSB = SSB / df_B(B-1)$$

$$df_T = N - 1$$

Understanding Statistics

- mean = $\sum x \cdot P(x)$ where x is the variable and $P(x)$ is the relative frequency (prob) of the variable.
- For continuous distribution, we do integration, but for discrete distribution we can do summation.
- Variance = $\sum (x - \bar{x})^2 \cdot P(x)$

How Variance is $\sum (x - \bar{x})^2 \cdot P(x)$

Mean of a binomial distribution is np .

pmf \rightarrow probability mass function.

stats. binom. pmf $(n, n, p) \rightarrow$ Binomial distribution.

cdf \rightarrow cumulative density function.

stats. binom. cdf will add all the pmf's.

Variance of the binomial distribution is npq . where $q = 1-p$.

stats. poisson. pmf $(n, \lambda) \rightarrow$ Poisson distribution.

where λ is the average number of events that happen in a period Δt .

$$P(x) = nCx p^n q^{n-x}$$

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

$$\sigma^2 = npq$$

$$\text{standard deviation } \sigma = \sqrt{npq}.$$

- The value you expect to get in statistical experiment is the mean.
It is the highest distributed value in the distribution.

- Binomial Distribution is used when,

we have categorical information, always we communicate this info in terms of proportion (%) or percentage.

whenever we have percentage information for categorical variables, we can generate binomial distribution.

- Poisson Distribution is used to generate counting statistics (counting information).

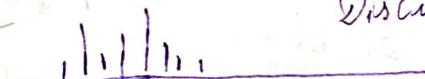
whenever we are dealing with discrete numerical quantity then,
Poisson distribution.

- whenever we are dealing with continuous numerical quantity, then
~~Poisson~~ normal distribution (continuous distribution).

In poisson distribution mean is λ .

	Binomial	Poisson
mean	np	λ
S.D	\sqrt{npq}	$\sqrt{\lambda}$

Discete distribution.



continuous distribution.



- Poisson distribution is independent of n .
(mean doesn't effect the Poisson distribution)
- Bernoulli is the collection of two binomials.

Characteristics of Standard Normal Distribution:

$\bar{x} = 0$ } The mean of the variable becomes 0, SD becomes 1.
 $SD = 1$ }

Inverse survival function:

stats.norm.isf function is used in python.
It gives the cut off value (range) for the required area.

Standard scalar, z scale, z score, (all are same).

Important z score values,

95% \rightarrow -1.96 to +1.96

90% \rightarrow -1.64 to +1.64

99% \rightarrow -2.58 to +2.58

In stats.norm.cdf function ($\rightarrow \mu, \sigma$)

Gaussian distribution / Normal distribution

This kind of distribution will form a bell curve



A random variable will form a Gaussian Distribution with ~~pop~~ some value of μ (mean) and ~~of standard deviation~~ ~~(variance)~~ σ (S.D.).

There is an empirical formula in Gaussian Distributions,

Probability of $(\mu - \sigma \leq x \leq \mu + \sigma)$ will be approx. equal to 68%.

It says that x which is a part of our random variable X .
That means 68% of the random variables will fall in this range.
Total number of elements from X present ~~S.D.~~ between first SD in the
left to first SD present in the right (i.e.) b/w the first standard
deviation is around 68% of the distribution.

68% of the elements of the random variable will be falling in the
first standard deviation. 68% is the percentage of the distribution of
data from the Random in the range from the random variable X .

so, within first standard deviation, 68% of the data exists. (68%). of the total distribution).

$$\begin{aligned}\mu - \sigma \leq x \leq \mu + \sigma &\Rightarrow 68\% \\ \mu - 2\sigma \leq x \leq \mu + 2\sigma &\Rightarrow 95\% \\ \mu - 3\sigma \leq x \leq \mu + 3\sigma &\Rightarrow 99.7\%\end{aligned}$$

This percentage of the total distribution exists in the given range.
(Confidence interval).

Log Normal Distribution

A random variable x usually belongs to the log normal distribution, if $\ln(x)$ is normally distributed. ($\ln(x)$)

Generally this distribution is similar to bell curve, but right skewed.

- When Gaussian distribution is converted into Standard Normal distribution, then $(\mu=0)$ and $(\sigma=1)$. When we do this, values will be converted according to the standard scalar (Z score).

- Covariance of x and y is Variance of x .

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

• Lognormal distribution says that $x \sim \text{Log Normal}$ (x belongs to log normal dist) if $\ln(x) \sim N(\mu, \sigma)$.
(If $\log(x)$ belongs to Gaussian or normal distribution with some value of μ and σ).

Covariance

$x \uparrow y \uparrow$ then +ve \rightarrow It cannot say how much positive.

$x \uparrow y \downarrow$ then -ve \rightarrow It cannot say how much negative.

Pearson Correlation Coefficient

This is used to overcome the above disadvantage.

- The probability of a variable that falls in between $\mu - \sigma \leq x \leq \mu + \sigma$ which is basically the range of first standard deviation will approximately equal to 68%.
- 68% of the data points belonging to the random variable x fall ^{within} the range of above mentioned range.

- stats.binom.pmf(n, m, p) where This function gives the probability
 $n \rightarrow$ depends upon the question asked.
 $m \rightarrow$ number of sample size.
 $p \rightarrow$ probability given.
 - stats.poisson.pmf(n, λ) where. This function gives the probability
 $n \rightarrow$ depends upon the question.
 $\lambda \rightarrow$ mean
 - stats.poisson.cdf adds up everything
 - $\mu - \sigma \leq x \leq \mu + \sigma \rightarrow 0.6828 \rightarrow 68.28\%$
 $\mu - 2\sigma \leq x \leq \mu + 2\sigma \rightarrow 0.9545 \rightarrow 95.45\%$
 $\mu - 3\sigma \leq x \leq \mu + 3\sigma \rightarrow 0.9974 \rightarrow 99.74\%$
- Somewhat related to 2 SDs
- probability or the area
- This is the percentage of data points spread within the region.
-
- ### stats.norm.cdf(x, μ, σ)
- stats.norm.cdf(x, μ, σ) The function gives the probability
 $x \rightarrow$ depends upon the question
 $\mu \rightarrow$ mean (average)
 $\sigma \rightarrow$ standard deviation.
 - isf \rightarrow inverse survival function
 stats.norm.isf(0.025) # 95% area
 $\rightarrow 1.9599$
 - stats.norm.isf(0.005) # 99% area
 $\rightarrow 2.575$
 - stats.norm.isf(0.05) # 90% area
 $\rightarrow 1.64485$
- (The answer gives the range of the 95% area and 99% area)
- The isf actually gives us the -2 score to +2 score value.

Sampling error \rightarrow Difference b/w population mean and sample mean.

$$\begin{array}{r} x = 10 \\ 60 - 90 \\ 20 \\ 25 \\ 5 \\ 10 \end{array}$$

$$\begin{array}{l} 95\% \rightarrow 1.96 \\ 68\% \rightarrow \pi \end{array}$$

- confidence interval is sample mean \pm (margin of error)
- margin of error = $Z_{\text{critical value}} \times \frac{\text{population SD}}{\sqrt{n(\text{sample size})}}$
 $(= Z_{\text{score value}} \times \frac{\sigma}{\sqrt{n}})$
- standard error = $\frac{\sigma (\text{SD})}{\sqrt{n(\text{sample size})}}$

In theoretical approach we take many trials to calculate mean of the mean.
In practical approach we do only one round of sampling.
(margin of error and standard error were used for practical approach)
Standard error of mean is calculated in practical approach.

De-moivre's approximation of Discrete Distribution with Standard Normal Distribution.

generate Binomial distribution with $p=0.5$, $n=25$, and calculate $P(x \leq 14)$.

In discrete distribution, the histogram bars are centered on the numbers. This means $P(x \leq 14)$ in Discrete distribution is actually the area under the bars less than $x=14.5$, we need to account for that extra 0.5, while calculating the same area in continuous distribution.

With skewnorm function, we can generate a skewed distribution.

H_0 is null hypothesis (no effect hypothesis).

H_a is alternative hypothesis (effect hypothesis)

The equality goes to null hypothesis.

not equal to goes for alternate hypothesis.

~~The tail sample means one sided check.~~

standard error of the mean σ / \sqrt{n}

Central Limit Theorem

- CLT states that for a large sample drawn from a population with mean μ and standard deviation σ , the sampling distribution of mean, follows an approximate normal distribution with mean, μ and standard deviation σ/\sqrt{n} irrespective of the distribution of population for large sample size.
- As a general rule, statisticians have found that for any population distribution, when the sample size is at least 30 and above, the sampling distribution of the mean is approximately normal/Gaussian.

- Confidence interval is also known as acceptance zone for H_0 .
- t_{stat} is the distance between two means.

- A hypothesis is a testable statement made in support of a finding or a claim about something in the world around you.
- It should be capable of being tested, either by experiment (or) observation.

- Hypothesis testing is an objective scientific method that is used in making statistical decisions using observational or experimental data.

- Hypothesis testing is an assumption that we make about the population parameter, a premise or claim we want to test.

For population data, general formula for SD is:

$$\sigma_{\text{pop}} = \sqrt{\frac{\sum_{i=1}^N (x_i - \mu)^2}{N}}$$

For sample data, general formula for SD is:

$$S_{\text{samp}} = \sqrt{\frac{\sum_{i=1}^n (x_i - \bar{x})^2}{n-1}}$$

The numpy function np.std by default divides it by n .

ddof (delta degrees of freedom)

(ddof=1) means the denominator is $n-1$.

df['vol'].std() by default divides by $n-1$.

np.std[] by default divides by n .

ttest_1samp (samps-array, 300) will give us the t-stat score, p-value.

tstat is the hypothesis testing formula.

Every value has its own confidence interval.

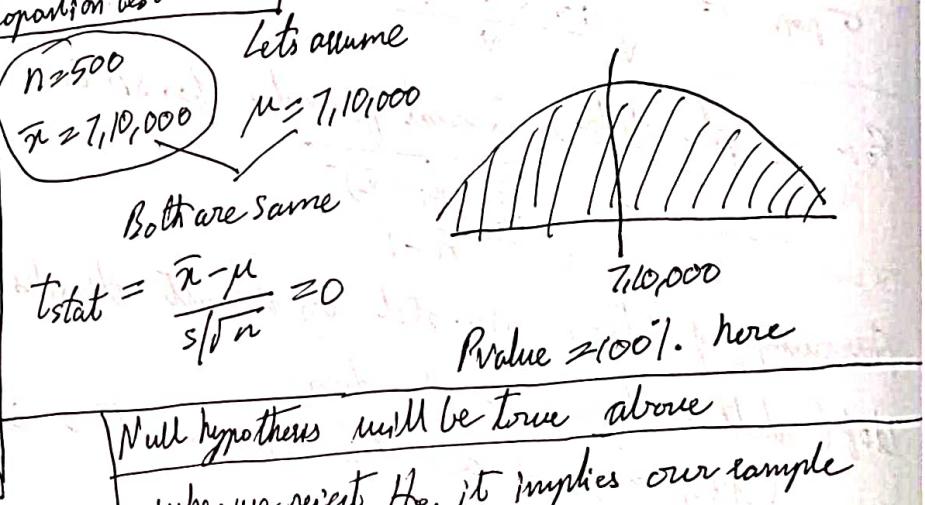
tstat is actually the distance between two means.

Practically point estimate is impossible, we can do a closer prediction with range estimate with different confidence interval levels.

Independent (unpaired)	
dependent (paired)	2 samples 2 groups
Mean	Two Sample 't' Test ANOVA
Proportion	Two Sample proportion test Chi Square Test

full information to check	Two sample test Dependent / Independent	One sample mean test → one row of data two sample mean test → two rows of data
One sample Test Two tail left tail Right tail	one sample T test (or) Proportion test	One sided check means one tailed sample

Formula for One sample T test
$t_{\text{stat}} = \frac{\bar{x} - \mu}{s/\sqrt{n}}$
$t_{\text{stat}} = \frac{\bar{x}_1 - \bar{x}_2}{\sqrt{\frac{s_1^2}{n_1} + \frac{s_2^2}{n_2}}}$



- Whenever $P < 5\%$ we reject H_0 .
- When $P > 5\%$ H_0 becomes true

To reject H_0	t_{stat}	P value (area)
95%	> 1.96	< 0.05
99%	> 2.58	< 0.01
90%	> 1.64	< 0.1

$H_a > -$ (Right tailed)
 $H_a < -$ (Left tailed)

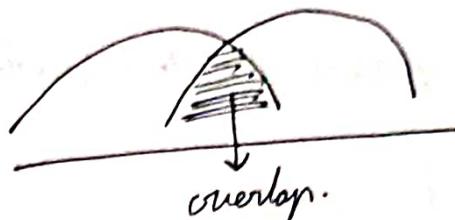
Two tail problem statements in One Sample T test actually have equal or not equal.
 Left tail (or) right tail generally have less than or greater than in the problem statement.

Two Sample T Test:

Two sample groups are taken (different groups) in the problem statement.
If H_0 must be true, two means must be 1.96 SD's away.

- Overlap happens, then p-value increases.

Diagrammatic representation of overlap:-



- Two sample t-test
 - Independent (unpaired)
 $t\text{test-ind}(g_1, g_2)$
 - Dependent (paired)
 $t\text{test-rel}(g_1, g_2)$

g_1 and g_2 are two groups of arrays.

rel means relative.

- When the population standard deviation is not known to us, then we use T-test.

- If it is known, we go for z test.

- Whenever we are calculating mean for two samples (or) groups, we call it Two Sample t test (or) Two sample mean test.

- t test and z test are same for $n=500$.

- Considering one sample,

$$t\text{test} \Rightarrow \frac{\bar{x} - \mu}{s/\sqrt{n}}, \quad z\text{test} \Rightarrow \frac{\bar{x} - \mu}{\sigma/\sqrt{n}}$$

- According to statistics, when $n \geq 30$, $s \approx \sigma$ and ~~tstat = zstat~~.

$t_{\text{stat}} = z_{\text{stat}}$.

where n_1 and n_2 are the sample sizes

x_1 and x_2 are the sample means.

s_1^2 and s_2^2 are sample variances

$$\frac{\bar{x}_1 - \bar{x}_2}{\sqrt{\frac{s_1^2}{n_1} + \frac{s_2^2}{n_2}}} \leftarrow t_{\text{stat}}$$

Rg(2017) Problem statement

Speedy delivery is the effect we wanted to study.

$M_d time < 3 \text{ hrs} \rightarrow H_a$ (speedy delivery is the effect we wanted to study)

$M_d time \geq 3 \text{ hrs} \rightarrow H_0$ (we don't want to study late delivery)

(another example)

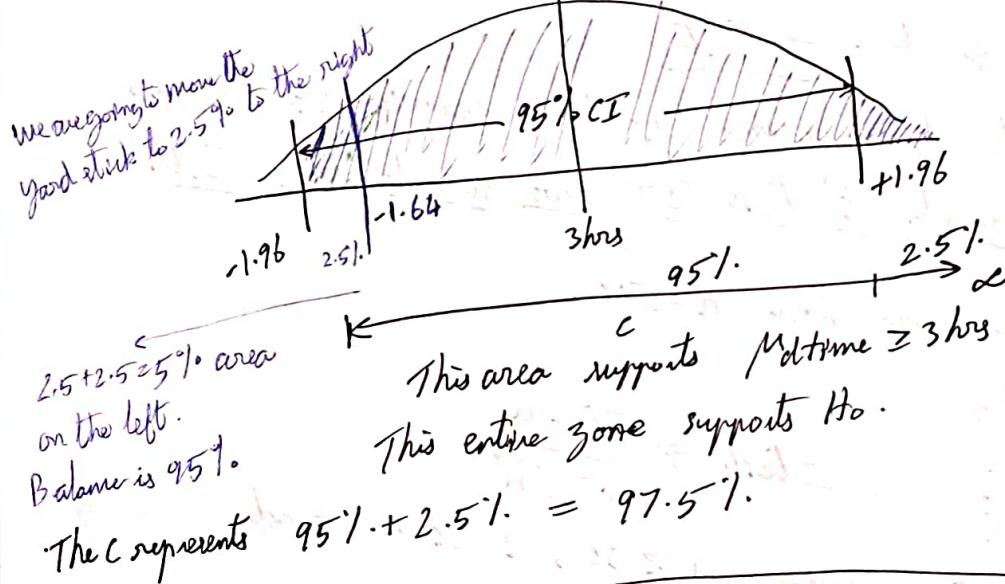
Pharmaceutical company is claiming their medicine works within 10min.

$M_d time < 10 \text{ min} \rightarrow H_a$ (we want to study this less time effect)

$M_d time > 10 \text{ min} \rightarrow H_0$

(we always take the contradiction statement as H_0)

This new portion is favourable for H_0 True



In the jupyter lab
has created a bad score
in the last. He has
taken only 50 samples

How to decide H_a and H_0

H_a is alternate hypothesis

H_0 is null hypothesis

Our research question is always on alternate hypothesis whereas the null hypothesis is the status quo.

Null hypothesis is the -ve version of our research question.

Alternate hypothesis is what we are trying to question.

One sample mean test

It has three different flavours,
to, to, to, and

- Two tail, left tail and right tail.

Two tasks one sample I test
in the full

- Two tailed one sample I test
 company claims they fill cans with 300ml. we need to verify whether
 the expected population mean is 300ml.

Left tail and Right tail

In these problems we use greater than or less than ($<$)

- In these problems we are given:
- Let us take the delivery time of a courier company.
- Delivery Company claims $\mu < 3$ hrs.
- + / the effect

Let us take the delivery time of a μ \sim hrs.
 Delivery Company claims $\mu < 3$ hrs.
 To study the effect

Let us take the delivery time of Dilmay Company claims $\mu < 3$ hrs.
 $H_0: \mu \geq 3$ hrs. $\rightarrow H_a: \mu < 3$ hrs. want to study the effect hypothesis

$M_d t_{\text{time}} < 3 \text{ hrs} \rightarrow H_0$ we want to
 $M_d t_{\text{time}} \geq 3 \text{ hrs} \rightarrow H_0$ we don't want to study late delivery

- $Md\text{time} \geq 3 \text{ hrs} \rightarrow H_0$ we don't want to study late delivery
- We always take the contradiction statement as H_0 which is the negative version of our research question.

negative version of our results is claiming average marks above 98.5.

negative
Let us consider a college is claiming average marks.

$$\text{Mang} > 98.5 \rightarrow \text{Ha}$$

$$\text{Mang} \leq 98.5 \rightarrow H_0$$

Brauerian example (previous page)

Pharmaceutical example (previous page)
right tail, for the Ha statements,

In the left tail and right tail, for the Ha statement
+ +) inside Ha statement then (Right Tailed)

If \rightarrow (greater than) inside Ha statement
It) inside Ha statement then (left Tailed).

If $>$ (greater than) inside Ha
If $<$ (less than) inside Ha statement
(we can remember which side we keep the 5%)

If $<$ (less than)
we can remember which side we keep
write out H_a and H_0 in One sample mean

If $<$ (less than)
we can remember which side we keep
This is how we figure out H_a and H_0 in One sample mean test.

- Only one sample t test (or) one sample proportion test (or) one sample mean test has three flavours (two tail, left tail and right tail).
- Going further we will be doing only two tail.
- ~~random~~ random will generate values in standard normal scale.

KRISH NAIK STATISTICS PLAYLIST (statistics for machine learning)

Statistics is the discipline that concerns the collection, organization, analysis, interpretation and presentation of data.

• Descriptive statistics

- Analyzing data, summarizing data, organizing data in the form of number of graphs.
- BAR PLOT, HISTOGRAM, Piechart, PDF, CDF, Normal distribution.
- Measure of Central Tendency. (mean, median, mode)
- Variance, standard deviation.

• Inferential statistics

Kind of predictive analysis

- The gaussian/normal distribution can be converted into standard normal distribution by using standard scalar (z score).
- If we have a column with a few values and we already know that these values are forming a log normal distribution. we need to convert these values into standard normal distribution.
 Convert the values into log → Scale them using standard scalar for standard normal distribution.

Cov (x, y) (Covariance)

$$\text{Cov}(x, y) = \text{cov}(x, y) = \text{var}(x)(\text{Variance})$$

- Covariance helps us to quantify the relationship between features, between random variable in a particular data set.

Central Limit Theorem :-

It specifies that, if the sample size is $n \geq 30$ selected randomly. let us consider this sample as S_1 . Similarly S_2 has 30 more sample. similarly we take till S_{100} (samples).

We calculate all the 100 means. If we take all the means and try to plot them, it will follow normal distribution and the sample mean will be approx. equal to population mean.

$$\bar{X} \approx \text{Gr.D} \left(\mu, \frac{\sigma^2}{n} \right)$$

Chebychev's Inequality

- When a random variable n belongs to Gaussian distribution, we can tell the percentage of the distribution like 68%, 95% and 99.7%.
- When a random variable n doesn't belong to gaussian distribution, to tell the distribution percentage, Chebychev's Inequality is used.

$$P_n(\mu - \sigma < n < \mu + \sigma) \geq 1 - \frac{1}{k^2}$$

$k \rightarrow$ the number of standard deviations

$$P_n(\mu - k\sigma < n < \mu + k\sigma) \geq 1 - \frac{1}{k^2}$$

If K is 2, then,

$$\begin{aligned} P_n(\mu - 2\sigma \leq n \leq \mu + 2\sigma) &\geq 1 - \frac{1}{4} \\ &\geq \frac{3}{4} \quad (\text{which is } 75\%) \end{aligned}$$

Gaussian

So, if the y doesn't follow any ~~random~~ distribution, ~~more than~~ then ~~less than~~ 75% of the data points ~~will be~~ belonging to y random variable will be falling within the range of second standard deviation.

Pearson Correlation Coefficient

This concept will also be used for feature selection process.

- Covariance helps us to find the direction of the relationship.

Pearson Correlation Coefficient

$$P(x,y) = \frac{\text{Cov}(x,y)}{\sigma_x \sigma_y}$$

- In covariance we don't know how much positive or negative the variation is.

- Pearson Correlation Coefficient states that, based on the variance of x and y , it will be able to tell that strength (how strong it is correlated) and Direction of the relationship.
- The value ranges from -1 to $+1$. ($-1 \leq P \leq 1$)

Spearman's rank Correlation Coefficient

It is somewhat better than Pearson Correlation Coefficient.

- Instead of x and y , we use rank of x and rank of y .

$$\Rightarrow \frac{\text{Cov}(\text{rank}_x, \text{rank}_y)}{\sigma_{\text{rank } x} \sigma_{\text{rank } y}}$$

Even if we have non-linear data points, we can use this method.

$$r_s = 1 - \frac{6 \sum d_i^2}{n(n^2 - 1)}$$

Outlier: An outlier is a data point in a data set that is distant from all other observations. A data point that lies outside the overall distribution of the data set.

- Reasons for an outlier to exist are:
- Variability in the data
 - An experimental measurement error.

• Impacts of having outliers in a data set:-

→ It causes various problems during our statistical analysis.

→ It may cause a significant impact on the mean and the standard deviation.

Interview Questions :-

. Tell us a few examples about left skewed and right skewed distribution.

. What is the relationship b/w mean, median and mode in skewed dist.

Right skewed :- Wealth distribution, length of comment

Normal Distribution :- Age, weight, height.

Left skewed :- Life span of a human being.

Bernoulli Distribution

We only have two outcomes in Bernoulli distribution.

Example of a Ctrn:

$X = 1 \rightarrow$ Head Success
 $X = 0 \rightarrow$ Tail Fail

$$P(X=x) = P^x (1-P)^{1-x}$$
 (Probability mass function) (outcome is discrete in pmf)

Probability density function (outcome is continuous value)

It will usually have zero and one.

$$P(X=0) = P^0 (1-P)^{1-0} = (1-P) \Rightarrow q$$

$$P(X=1) = P^1 (1-P)^{0} = P$$

$$\begin{cases} P(X=0) = 0.4 = 1-P \\ P(X=1) = 0.6 \end{cases}$$

$$\begin{aligned} \text{Mean } E(n) &= \sum_{i=1}^x i \cdot P(i) \\ &= 0(0.4) + 1(0.6) = 0.6 \Rightarrow P \end{aligned}$$

The mean of the Bernoulli distribution is P .

$$\text{Variance} = P(1-P) = pq$$

$$\text{S.D is } \sqrt{pq}$$

What is p-value :-

It is the probability for the "null hypothesis" to be true.

$p=0.05 \rightarrow$ significance value

Different Sampling Techniques :-

- Suppose we have to do an exit poll, we collect the data randomly from random places. This is **Random Sampling Technique**.
- **Stratified Sampling Technique :-**

what if the exit polls are wrong. The sampling might be wrong.

May be when 1000 samples were taken, there were 700 men and 300 women. It is completely biased. we may not get proper output.

50:50 ratio is better.

- **Systematic Sampling :-**

we select the n^{th} person from every group and try to record the information.

- **Cluster Sampling Technique :-**

It is more domain related. If we are researching medical domain. The person we select \oplus randomly are familiar with this domain knowledge.

- Continuous Variable means mean test and others is proportion Test.
- whenever we are building a prediction model to predict a numerical variable, that is called **Regression model**.

Model is going to say Diabetic / Healthy which is a classification model.

Two Sample Z-Test for Difference in proportions

$$Z_{\text{data}} = \sqrt{\frac{P_1 - P_2}{P_{\text{pooled}} (1 - P_{\text{pooled}}) \left(\frac{1}{n_1} + \frac{1}{n_2} \right)}} \quad \text{where } P_{\text{pooled}} = \frac{x_1 + x_2}{n_1 + n_2}$$

ANOVA :

Consider $\bar{A}, \bar{B}, \bar{C}$ are means of ages.

$$H_0: \bar{A} = \bar{B} = \bar{C}$$

If null hypothesis H_0 holds good (meaning all the means are overlapping), then age is an useless feature.

If there is a significant difference of mean in the feature, then we can consider it. H_0 should hold good for the feature to be considered.

For a strong feature, MSTR must be very high

- magnitude and direction

↳ Correlation (positive or negative)

- Covariance and Coefficient of Correlation

↳ $-1 \rightarrow +1$

Random Covariation between two variables is so called covariation.

Covariance :-

$$\text{Cov}(x, y) = \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})$$

$$\text{Cov}(x, x) = \text{Var}(x)$$

Covariance helps us to quantify the relationship between features, between random variable in a particular data set.

ANOVA (Session 6) points :-

- whenever we reject the null hypothesis, whatever the test it might be, we have some interesting story to tell from H_a . This story is known as post hoc analysis.
- In $mod = ols('Income ~ Age_Group', data = data)$: fit()
 \downarrow
(First parameter must be numeric variable)
- Tukey Test will do the pair wise analysis.
It performs a pair wise two sample T test.
- Let us look at the example of a confusion matrix:-
PREDICTION

		Healthy (-ve)	Covid 19 (+ve)
		Healthy (-ve)	True Negative
		Covid 19 (+ve)	False Positive
A	C	Healthy (-ve)	False positive \rightarrow Type I error (α) error
U	A	Covid 19 (+ve)	True Positive \downarrow Type II error (β) error
L			

- Predicting Healthy as Healthy is a true negative
- Predicting positive case as a positive is true positive
- Type II error is the costly error in the above example.