

Statistics for Machine Learning :-

Learning pattern :-

Session 1 :-

(Binomial and Poisson)

Probability distribution \rightarrow discrete distribution & continuous distribution.

Session 2 :-

Continuous probability distribution (Normal distribution)

This concept will lead to CLT (Central Limit Theorem), it is the key topic to hypothesis testing | Range estimate of population parameter.

Session 3 :- Hypothesis Testing - start with one sample T-test.

Session 4 :- Two sample T-test (we have paired and unpaired).

Session 5 :- Proportion test \rightarrow Chi-Square Test (There are four types)

Session 6 :- ANOVA (It is a type of regression analysis).

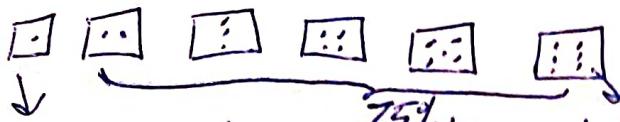
There is a strong connection between ANOVA & SLR.

Basic Probability Questions :-

1) A gambler cheats using a loaded die, with 'one' comes up 25% of time. What is the probability of getting '6' in that die?

Solution :-

It is a biased die here,



I comes 25% of the time $\xrightarrow{75\%}$ what is the prob of getting 6?

If it is a fair dice, answer is $1/6$.

Remaining 5 numbers have 75% chance, all are equally likely,
 $\text{so, } 75\% / 5 = \underline{\underline{15\%}} \text{ (answer)}$

2) what is the probability of getting a sum=3 while throwing two die (say white and black).

Solution:- Two dice are thrown, so the sample space will be $6^2 = 36$.

The list of Combinations are:-

$$\begin{array}{l} 1 + 2 = 3 \\ 2 + 1 = 3 \end{array} \quad \left. \begin{array}{l} \\ \end{array} \right\} \text{only two combinations.}$$

The probability is $2/36 = 1/18$.

33:30

3) For the above experiment, what is the conditional probability of getting a sum=3 when white die is already rolled and turns to be 'one'.

Solution:- white die is already rolled.

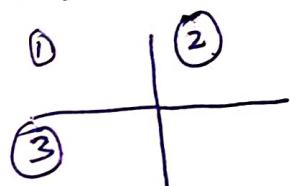
It is like a conditional probability.

$P(S=3 / w=1)$ [what is the probability that sum is 3 with]
 the condition of white die being 1]

One die is already fixed, then the chances left for the second die is just six.

The solution is :- $1/6$.

In probability we have four major types of calculations:-



① Mutually Exclusive events

Mutually exclusive event is going to calculate the probability of either A or B.

Example:- Head and Tail are mutually exclusive events.
when Head happens tail doesn't occur and vice versa.
then $P(H \text{ or } T)$ has a 100% chance.

$$\begin{aligned} P(H \text{ or } T) &= P(H) + P(T) \\ &= \frac{1}{2} + \frac{1}{2} = 1 \end{aligned}$$

Therefore,

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

what is the probability of getting king or queen in pack of 52 cards:-

$$P(K \text{ or } Q) \Rightarrow \frac{4}{52} + \frac{4}{52} = \frac{8}{52}$$

There is nothing common between King and Queen so,

$$P(A \cup B) = P(A) + P(B) - P(A \cap B)$$

where $P(A \cap B) = 0$

$$P(H \cap T) = 0$$

② Mutually Non-Exclusive Event :-

$$P(A \cap B) \neq 0 \quad [\text{need not be } 0]$$

Example:- What is the probability of getting a K or Q in a pack of 52 cards.

$$P(K \cup Q) = P(K) + P(Q) - P(K \cap Q)$$

$$= \frac{4}{52} + \frac{13}{52} - \frac{1}{52} = \frac{16}{52}$$

$$\boxed{P(A \cup B) = P(A) + P(B) - P(A \cap B) \text{ where } P(A \cap B) \neq 0}$$

The above two rules are addition rules (union rules).

③ Mutually Independent Event :-

- It means both the events can take place together.
- One event will not influence the other event.
- This is a product rule.

$$\boxed{P(A \cap B) = P(A) \cdot P(B)}$$

④ Mutually Non-Independent Event :-

• Here one event can influence the other event.

$$P(A \cap B) = P(A) \cdot P(B|A) \quad (P(B|A) \text{ means } A \text{ has already happened}).$$

We can rewrite it,

$$P(B|A) = \frac{P(A \cap B)}{P(A)} \quad \text{--- (1)}$$

Now, B is influencing A, (B has already happened).

$$P(A \cap B) = P(B) \cdot P(A|B)$$

$$P(A|B) = \frac{P(A \cap B)}{P(B)} \quad \text{--- (2)}$$

Equating (1) and (2), we get,

$$P(A) \cdot P(B|A) = P(B) \cdot P(A|B)$$

$$\begin{aligned} P(B|A) &= \frac{P(A|B) \cdot P(B)}{P(A)} \\ P(A|B) &= \frac{P(B|A) \cdot P(A)}{P(B)} \end{aligned}$$

This is Baye's rule.

Name Baye's rule works based on this philosophy.

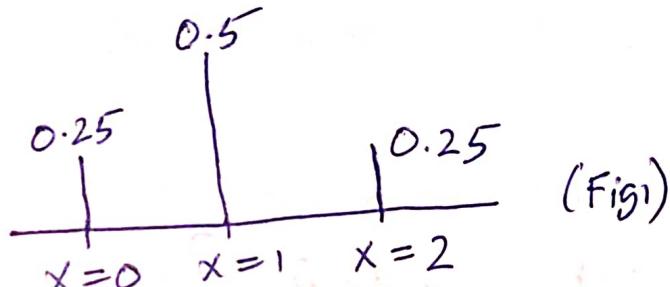
Probability Distribution :-

Let us take two coins, now the sample spaces are:-

$$\{ HT, TH, HH, TT \}$$

X is defined as a random variable getting head.

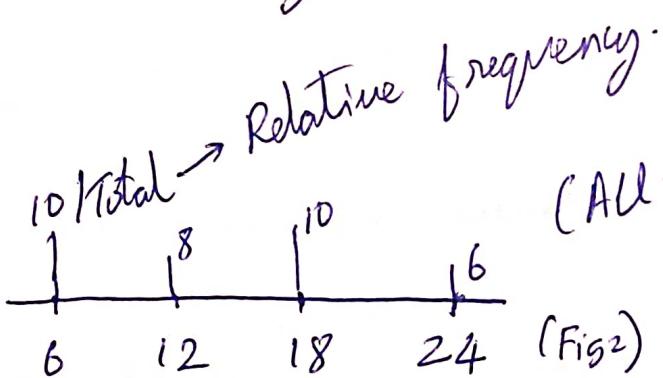
$$\left. \begin{array}{l} P(X=0) = \frac{1}{4} \\ P(X=1) = \frac{2}{4} = \frac{1}{2} \\ P(X=2) = \frac{1}{4} \\ P(X=3) = 0 \end{array} \right\} \text{Let us plot these values in a graph,}$$



This is the probability distribution

Let us take an example,

Consider the number of months in the x-axis and number of people in the y-axis.



(All these 10, 8, 10, 6 are frequencies)

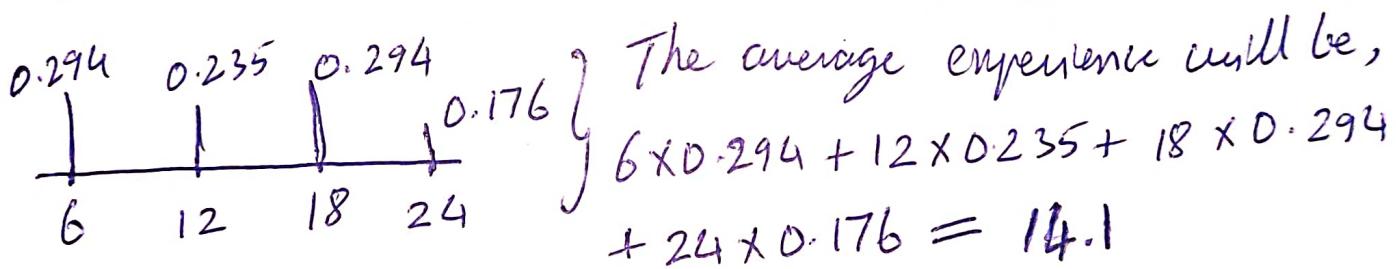
Your probability is nothing but the relative frequency.

Few observations about the above plots :-

- It's nothing more than a histogram (with relative frequency).
- The x-axis is discrete.
It is not continuous, because we cannot calculate what is the probability of 1.5 head or 0.8 head -
- So, it is an example of discrete distribution.

How can we calculate the average experience in Fig 2 :-

Let us replace the vertical lines with relative frequencies,



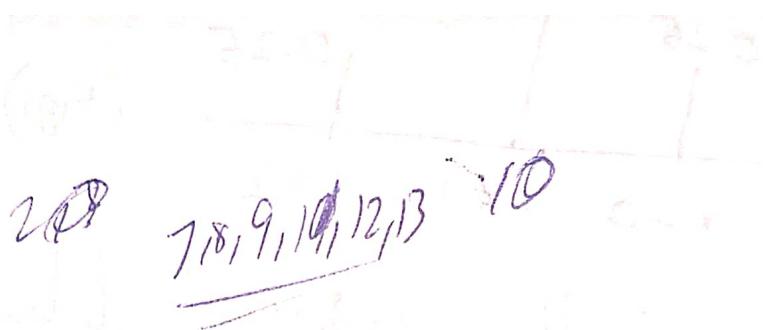
The average will be $\sum x \cdot p(x)$

This is the global representation for any distribution.

- For continuous distribution we replace \sum with \int symbol.
- \int is the extended version of 'sum'.

$$\boxed{\text{Variance} = \sum (x - \bar{x})^2 \cdot p(x)}$$

Continuous, we do Integration.
Discrete, we do summation.



$$\text{variance} = \frac{\sum (x_i - \bar{x})^2}{n}$$

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

$$\sum (x_i - \bar{x})^2 p(x)$$

$$7(7-10)^2 p(x) + 8(8-10)^2 p(x)$$

$$\frac{(7-10)^2}{n} + \frac{(8-10)^2}{n}$$

Binomial Distribution & Poisson Distribution of Distributions

-Discrete Distributions.

Assume there is a T-shirt manufacturing company which is operating for last 5 years, and manufactured 2,00,000 T-shirts. The QA always finds the defective pieces, they observed,

- ① 5% of the T-shirts are defective.
- ② Out of the defective T-shirts, (10,000), they are counting the number of defects in each T-shirt.

$D \rightarrow 3$
 $D \rightarrow 2$
 $D \rightarrow 4$
 $P \rightarrow 5$
⋮
 $10,000 D \rightarrow 4$

How can we communicate these defects.
 $\frac{\text{Sum of Defects}}{10,000}$ will give the average number of defects.

The mean of the binomial distribution is given as np .

pmf \rightarrow probability mass function

There is a function in python,

`stats.binom.pmf(0, n, p)` $\boxed{\text{Binomial Distribution}}$

(This function is used to obtain the probability mass function for a certain value of n, n, p)

`stats.binom.cdf` will add everything
(cumulative density function)

- Mean of the Binomial distribution is np and the variance of the binomial distribution is npq .
 we have another function stats.binom.cdf (cumulative distribution function)
- Formula for the Binomial Distribution is :-

$$P(X) = {}^n C_x p^x q^{n-x}$$
 where np (mean)
 stats.binom.pmf(x, n, p)
- Formula for the Poisson Distribution :-

$$P(X) = \frac{e^{-\lambda} \lambda^x}{x!}$$
 where λ (mean)
 stats.poisson.pmf(n, lambda)
 $\lambda \rightarrow$ The average number of events that happen in a period Δt .
 (mean of the binomial distribution)
- Mean = $n p$ (where n is the number of events and p is the probability)
 $\mu = np$
- The value you expect to get in a statistical experiment is the mean. It is the highest distributed value in the distribution
- Variance = npq where $q = 1 - p$
 $\sigma^2 = npq$
- Standard Deviation $\sigma = \sqrt{npq}$

If there are categorical variables then it is binomial distribution.

Defect/no defect → categorical column

men/women → categorical column

- whenever we have a categorical information, always we communicate this information in terms of proportion (or) percentage.

Let us look at three statements :-

(Counting)

① In a toll gate, during 8AM-10:30 AM, average 78 cars are passing.

→ numerical quantity (discrete)

② A customer care unit receiving on an average 150 calls.

(Counting)
→ numerical quantity
(discrete)

③ A bank data tells that 15% of customers default their loan EMI.

All the above examples come under categorical variables with numeric information that is discrete.

Default or Not Default
(Categorical Variable)

- whenever we have the percentage information for categorical variables, we can generate binomial distributions.
- For counting statistics (Counting info) we can generate Poisson distribution.

- whenever we express a categorical variable, we express in terms of proportions or percentages. (Binomial)
- whenever we address a numerical quantity (continuous) then continuous distribution (normal distribution).
- If we are dealing with numerical quantity (discrete) then poisson distribution.

• In Poisson distribution mean is λ

General Rule:-

$$\sum x \cdot p(x) \quad (\text{mean})$$

$$\sum (x - \bar{x})^2 \cdot p(x) \quad (\text{Variance})$$

$$\sqrt{\sum (x - \bar{x})^2 \cdot p(x)} \quad (\text{Standard deviation}).$$

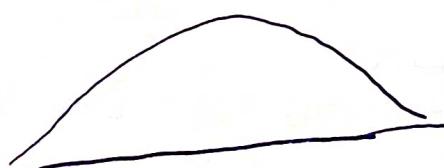
Binomial	Poisson
mean np	λ
s.d. \sqrt{npq}	$\sqrt{\lambda}$

$p \rightarrow$ probability of the occurrence of event happening.

$n \rightarrow$ sample size

$\lambda \rightarrow$ mean

Discrete Distribution



Continuous (normal) distribution.

If we want to add multiple pmf's we can go for cdf's.

- Poisson Distribution is independent of n . (meaning mean doesn't effect the Poisson Distribution.)
- Bernoulli \rightarrow Collection of Two Binomials.

- whenever discrete, we use the terminology mass, whenever it is continuous, we use the terminology density.
- Only in the normal distribution, we use the word standard normal distribution.
- In discrete distribution, we just talk about the counts.
(no units mentioned)
- In continuous distribution, units are mentioned.
(every variable has its own units).

The characteristic of Standard Normal Distribution is :-

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

Inverse survival function :-

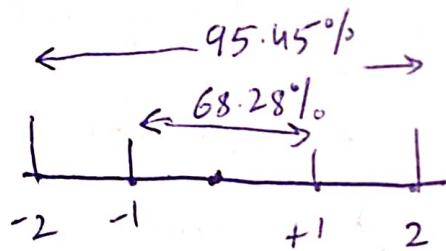
`stats.norm.isf` function is used in python.

It gives the cut off values (range) for the required area.

Central Limit Theorem is the foundation for inferential statistics.

We need to understand the population parameter in the inferential statistics.

- Standard Scores, Z scale, Z-score (same)
- In stats.norm.cdf function ($-$, μ , σ).
- I&f (Inverse Survival Function)



`stats.norm.isf(0.025)`

we use this if we want to calculate 95% of the area. Then 5% will be distributed as 2.5% on the left and right each.
Give this value in the `stats.norm.isf` function.

95%	-1.96 to +1.96
90%	-1.64 to +1.64
99%	-2.58 to +2.58

Objective:- To estimate the population μ .

CLT will give the following rules in order to achieve the objective. (Estimate the population mean or proportion (categories))

- Takes the samples randomly from everywhere.
- Take the mean values.
- Do one more round and take the mean values.
- Do for a multiple times (100 trials).
- So we got total 100 means.
- If we plot the distributions of 100 means, these 100 means will form the normal distribution.
- It will form the normal distribution irrespective of the population distribution.

$\bar{x} - 1SD$ to $\bar{x} + 1SD$ } 68% CI
(68% of data will lie within this range) 32% Risk

$\bar{x} - 2SD$ to $\bar{x} + 2SD$ } 95% CI
(95% of data)

$\bar{x} - 3SD$ to $\bar{x} + 3SD$ } 99% CI
(99% of data will lie within this range) 1% Risk

SD \rightarrow Standard Deviation
CI \rightarrow Confidence Interval
and Risk is the Error Risk.

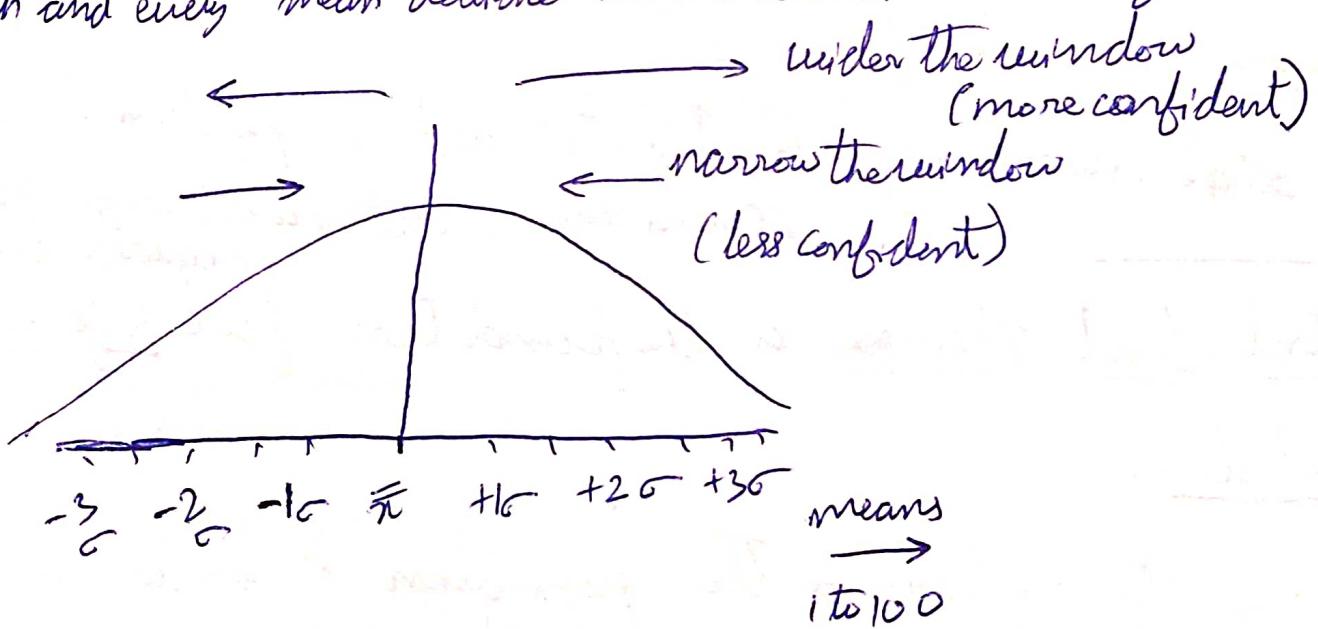
Wider the Window
more confident.

Narrow the window
less confident.

If there is no Inferential Statistics, there is no ML, then no DL.
 CLT says that, if we plot the distribution of 100 means, these 100 means will form a normal distribution. irrespective of your population distribution. (whatever the shape is).

- The mean of the 100 means is $\bar{\bar{x}}$.
- Standard Deviation = 0 for the true mean
- $SD = \text{standard error of the mean}$

Each and every mean deviate with a certain amount from $\bar{\bar{x}}$



Examples for the type of distribution:-

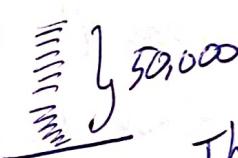
• It is expected 20% of covid19 cases in India effecting children.

Binomial

• It is also observed the avg age of covid affected is 38.45 years
Continuous (Normal Distribution)

• The avg-number of symptoms for the +ve cases is found to be 6.
Poisson

• From the before objective:

To collect the information of 50,000 Data Scientists, 

• True mean is unknown.

• We can only infer(OR) predict the true mean
Take Average :- $\mu \rightarrow$ This mean
is known as true mean.

mean.

• Point estimate is not at all possible in statistics.

• CLT is used).

• We can do a Range Estimate.

• Objective is to Estimate the population mean (OR) proportion

if there are Categorical variables.

Example :-

Assume $n = 40$, true mean is 24.6 years.

Let us say the range estimate is :-

[18 years to 50 years] Wider Interval, more confident.

-0.5 to +0.5 is the acceptable range for skewness.

Normal Distribution

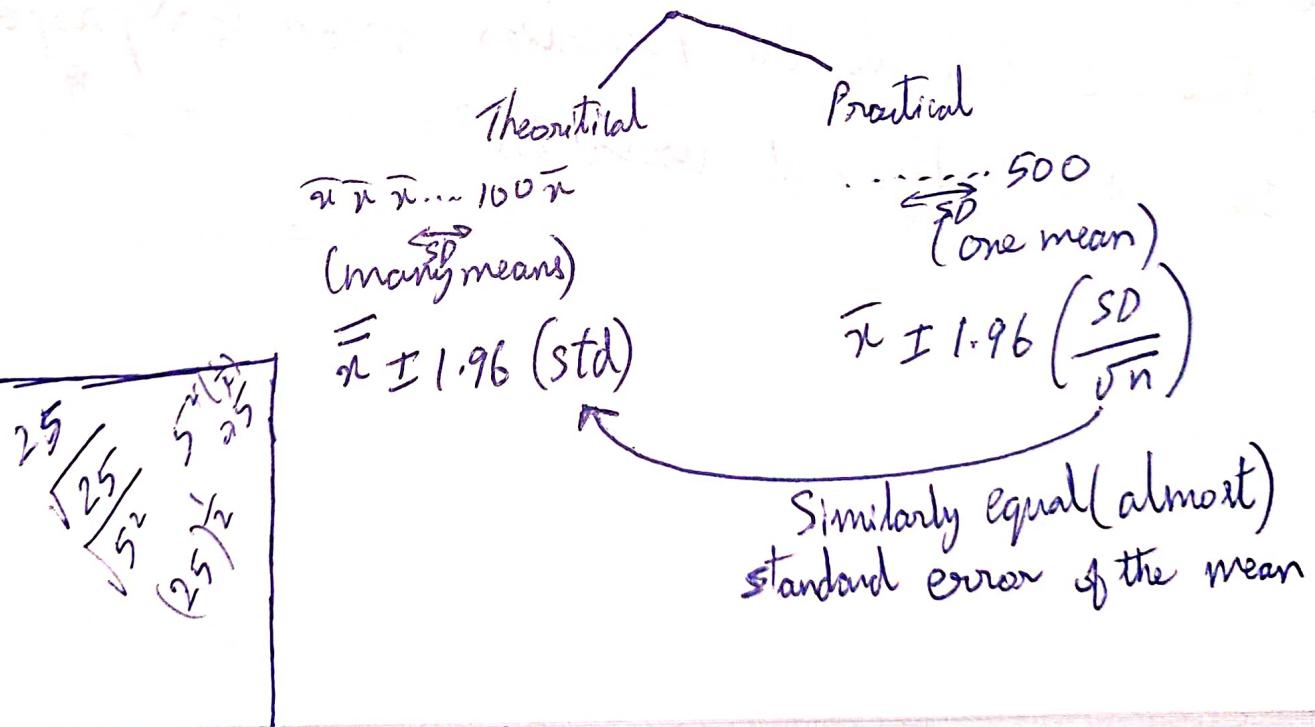
Area Under the normal Distribution

$$\mu - \sigma \leq x \leq \mu + \sigma \rightarrow 0.6828$$

$$\mu - 2\sigma \leq x \leq \mu + 2\sigma \rightarrow 0.9545$$

$$\mu - 3\sigma \leq x \leq \mu + 3\sigma \rightarrow 0.9974$$

- As we increase the sample size the error will reduce as the sample size will almost be same as population size.
- The CLT states that, if we take a number of means from different samples and plot them, they will be normally distributed, irrespective of the population distribution.
- With np.random.seed() same random numbers will be generated.



Lab Session 1 :-

- Poisson Distribution $\Rightarrow P(x) = \frac{e^{-\lambda} \lambda^x}{x!}$ where λ is the average number of events that happen in a period Δt . Events in each of these Δt 's are independent of each other.
- Binomial and Poisson is for Discrete Probability.
- Normal Distribution is for Continuous Probability.
- Normal Distribution is a symmetrical bell curve, whenever the data is put together, it takes the shape of a symmetrical bell curve, where mean, median and mode lie at the same point. (which is the center of the curve). (distribution takes shape)

For normal distribution, we perform `stats.norm.cdf(n, mu, sigma)` where μ is the mean and σ is the standard deviation.

- Sampling Error :- Difference b/w Population mean and Sample mean.
- Standard Error :- $\sigma(SD) / \sqrt{n}$ (sample size)
standard Deviation (σ) of the sample divided by root over the n .

`stats.poisson.rvs` is used to generate random variables.

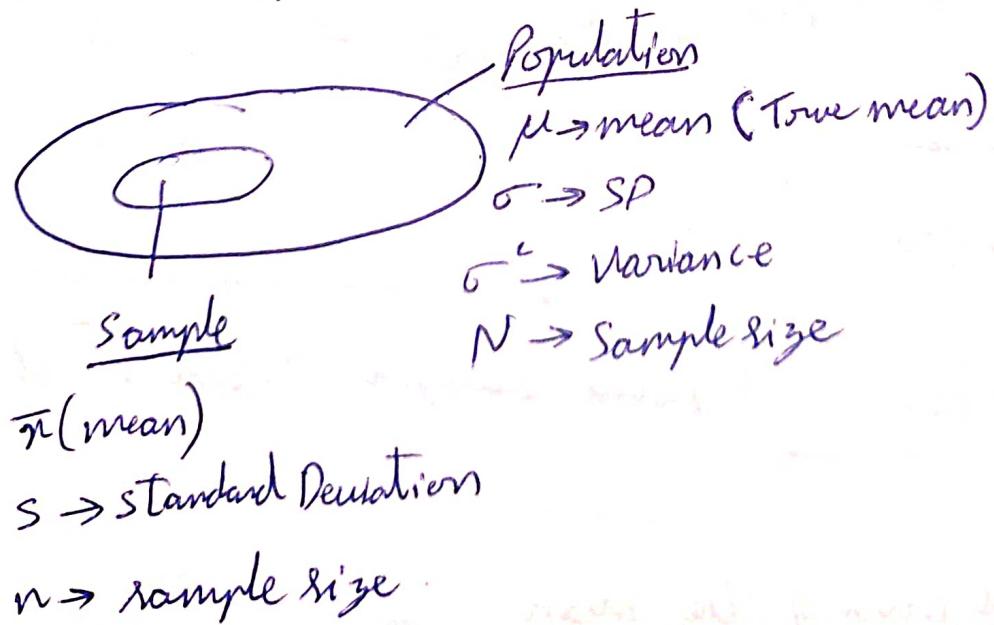
- `stats.poisson.rvs(loc, mu, size)`.
 loc is the lowest x value, mu is the middle of distribution.
- `np.random.seed()` makes the random variables predictable.

- Confidence interval is sample mean \pm (margin of error).

statistics (Session 3)

1:18:40

- Range Estimate leads to one sample test and other tests.
- With skewnorm function, we can generate a skewed distribution.
- In inferential statistics we are going to infer whether our sample is a representation of population. (we have to verify this with hypothesis testing)
- In inferential statistics, we can only infer about μ (or) proportion.
- In theoretical approach we take many trials to calculate mean of the mean.
- In practical approach we do only one round of sampling.



- Z critical value, margin of error value

$$\text{margin of error} = \text{Z critical Value} \times \frac{\text{Population SD}}{\sqrt{\text{Sample Size}(n)}}$$

Z critical Value \Rightarrow we can find using $ISF(0.025)$ or anyone of them.

Stats.norm.ppf(q), pass the q, value where $q = 1 - P$

Example:- The interval from 2.634 to 2.966 forms a 95% confidence interval for μ . In other words, we are 95% confident that the average delivery time lies between 2.634 and 2.966 hours.

• In both (theoretical and practical) the range contains the true mean.

• Tolerance level of skewness is ± 0.5 .

- Null Hypothesis \rightarrow No effect Hypothesis (H_0)
- Alternative Hypothesis \rightarrow effect hypothesis (H_a)
- whenever we are checking equal or not equal, the equality goes to null hypothesis. (H_0)
- Something not equal to, will go for effected hypothesis. (H_a)
- If it is a one sided check, then it is a one tail sample.

- How the sample mean \bar{x} is away from hypothesized mean μ . This formula $\Rightarrow \frac{\bar{x} - \mu}{S/\sqrt{n}}$. (This is generally used for using this formulas, we get the SDs) (One sample T test)
- (It is going to tell how many standard deviations \bar{x} is away from μ). (These are called Z-scores)

$$\frac{\bar{x}_1 - \bar{x}_2}{\sqrt{\frac{s_1^2}{n_1} + \frac{s_2^2}{n_2}}} \quad (\text{This formula is used for Two sample T Test).}$$

Independent groups

We can find t_{stat} by how many SDs away with these formulas. (t_{stat})

$$\frac{\bar{x} - \mu}{S/\sqrt{n}} \quad (\text{One sample T test})$$

- Confidence interval is also known as acceptance zone for H_0 .

Let us take an example: Swiggy claims that it can deliver the food within 50min, $\mu_{\text{avg}} < 50\text{min}$ $\rightarrow H_a$. The avg delivery time of hyderabad and chennai is equal and not equal.

But, the agencies find out it is more than 50min.

$$\mu_{\text{avg}} \geq 50\text{min} \rightarrow H_o$$

(The above is H_o because the delivery speed is effected.)

One Sample Mean Test

Two Sample Mean Test

$$\mu_{\text{avg(Hyd)}} = \mu_{\text{avg(Ch)}}$$
$$\rightarrow H_o$$

$$\mu_{\text{avg(Hyd)}} \neq \mu_{\text{avg(Ch)}}$$
$$\rightarrow H_a$$

$\cdot \frac{s}{\sqrt{n}}$ is the standard error of the mean

- The ~~minimal~~ probability required for the null hypothesis being true is
- Null hypothesis being true is less than 5%
- Whenever $P < 5\%$, we reject H_0 .
- H_0 becomes true once it is greater than 5%.
-

To reject H_0

	t_{stat} (dist. b/w two means)	P-Value (area)
95%	> 1.96	< 0.05
99%	> 2.58	< 0.01
90%	> 1.64	< 0.1

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

Two Tail Vs Two Sample

One Tail Vs One Sample

These terminologies are completely different.

Sample
One tail Test

Two Tail
left Tail

Right Tail

Two sample Test

Dependent Independent

$$n = 500$$

$$\bar{x} = 7,10,000$$

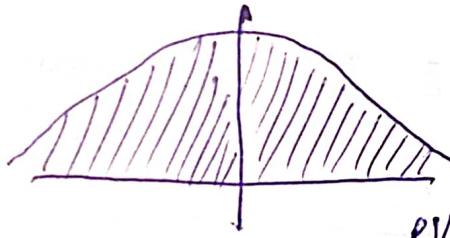
SD

let us assume

$$\mu = 7,10,000$$

both are same.

$$t_{\text{stat}} = \frac{\bar{x} - \mu}{s/\sqrt{n}} = 0.$$



$$p\text{Value} = 100\%.$$

only for one sample T Test or proportion test we have:-

Two tail, left tail and right tail.

Two Tailed One-Sample t test :-

Company claims they fill their cans with 300ml. We need to verify whether the expected population mean is 300ml.

First we need to do hypothesis,

We need to frame null hypothesis H_0 and alternate hypothesis H_a .

$$H_0 = \mu_{\text{Pop}} = 300 \text{ ml}$$

favourable for H_0 being true

$$H_a = \mu_{\text{Pop}} \neq 300 \text{ ml}$$

H_0 favourable zone

↳ what are the effects

1.96

95% confidence

300ml

+1.96

→ They can overfill the can (spillage)

→ Under filling.

(Expected)

When the value falls within the range, $t_{\text{stat}} = 0$ (if sample \bar{x} is 300ml) then the P-value is 100%.
 where null hypothesis is true, outside Meaning it correctly represents the μ .
 this range it is not true. When Null hypothesis is true (100%)
 the value falls within the range probability percentage increase by 5%. When outside it is less than
 5%.

$$n = 400 \quad (\text{Randomly Taken 400 samples})$$

$$\bar{x} = 298.56$$

$$s = 1.5$$

In jupyter notebook, we got,

$$\bar{x} = 298.51$$

$$s = 1.54$$

} Information given

$$t_{\text{stat}} = -19.18$$

(Let us draw a diagram in the next page).

General formula for SD is, $\sigma_{\text{pop}} = \sqrt{\frac{\sum_{i=1}^N (x_i - \mu)^2}{N}}$

For sample data, SD is, $s_{\text{samp}} = \sqrt{\frac{\sum_{i=1}^n (x_i - \bar{x})^2}{n-1}}$

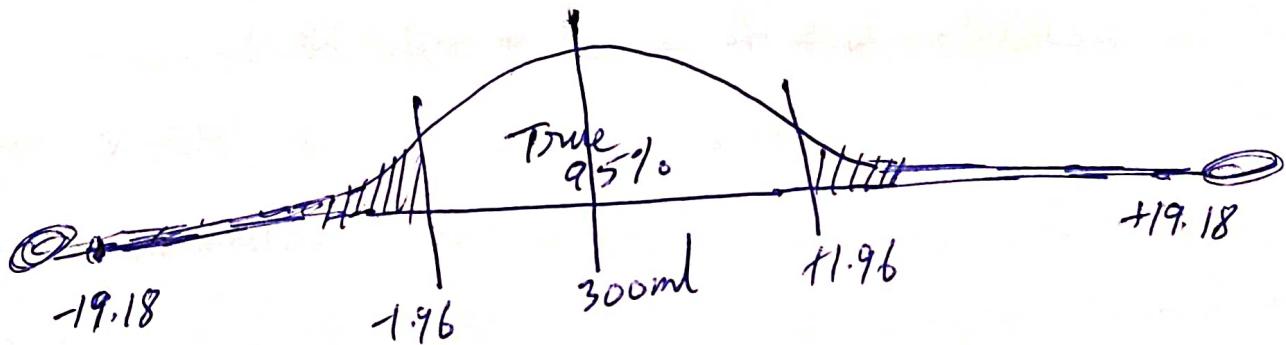
- The numpy function np.std by default it divides by n .
- ddof=1 means, the denominator is $n-1$.

ddof (delta degrees of freedom)

ttest_1(ramp(ramp=array, 300)) will give us the t-stat score, p-value

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

np.std[] by default divides by n .



only 0.5 data lies beyond 30.

The area will be obviously less than 5%.

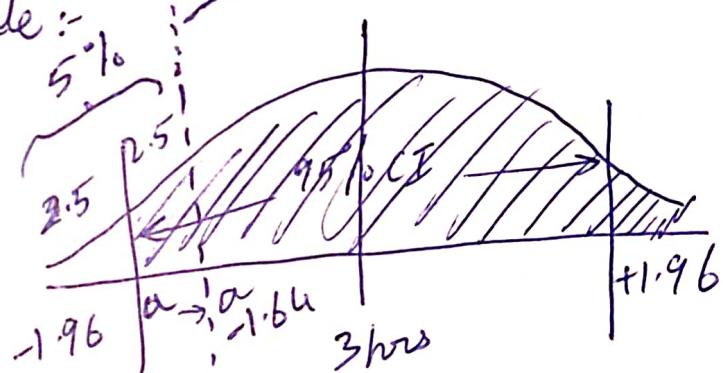
Since the P value is $P < 5\%$, we reject the null hypothesis.
 Let's do the one tail test.

From the Newyork Courier example:

$$\mu_{\text{dtime}} < 3 \text{ hrs} \rightarrow H_a$$

$$\mu_{\text{dtime}} \geq 3 \text{ hrs} \rightarrow H_0$$

(This is left tail example)
 (Right tail is similar)



This 95% will represent 3 hrs

The whole shaded part is $95 + 2.5 = 97.5\%$. We need to move α to the right by 2.5%. Since it is a 1 tail test 5% will be on the left side. We use if function to get -1.64.

If null hypothesis is being rejected, then H_a is true.

- Every value has its own confidence interval.
- t_{stat} is the hypothesis testing formula.

- Let us consider, a JEE college is claiming average marks more as :-
- $\mu_{avg} > 98.5$ we always take equality sign as an advantage to the H_0 .
- $\mu_{avg} \leq 98.5$
- we should always think what is H_a and H_0 .
- H_a always is ~~rejective~~ and tells about the efficiency.

If H_a is supportive

$H_a > \underline{\quad}$ (Right tailed)

$H_a < \underline{\quad}$ (Left tailed)

Two Sample T Test :-

Let us understand the type of testing in each statement:-

- Avg age of DSE at GL is > 25 years.

(Here the population is entire India and age is the variable. This is the test of mean. This is one sample because we are considering the entire country. This is right tail test.)

Here, H_0 is $\mu > 25$ years

H_0 is $\mu \leq 25$ years)

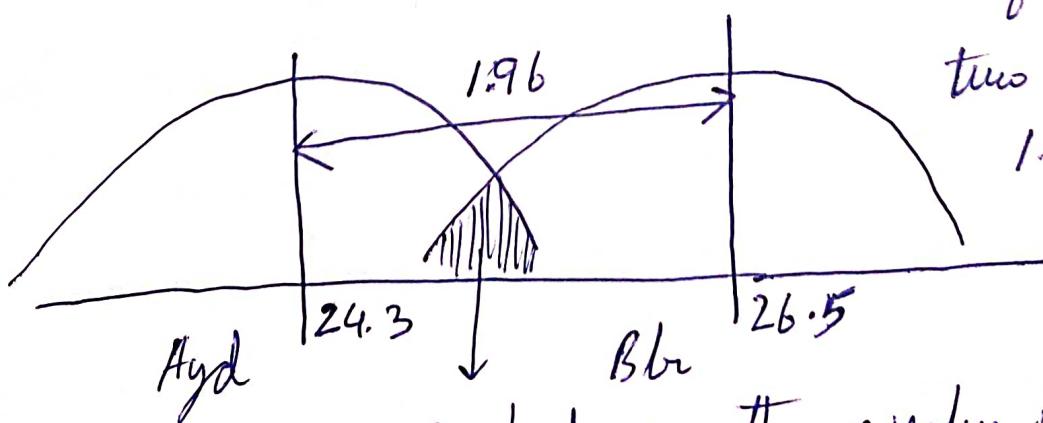
- Average age of DSE student at Bangalore is more compared to Hyderabad.

$$\begin{array}{ll}
 \text{Bbr} & \text{Hyd} \\
 \text{---} & \text{---} \\
 n_1 = 1000 & n_2 = 800 \\
 \bar{x}_1 = 26.5 & \bar{x}_2 = 24.3 \\
 s_1 = 2.3 & s_2 = 1.8
 \end{array}$$

It becomes two sample mean test.

Here,
 $H_0 = \bar{x}_{\text{Hyd}} = \bar{x}_{\text{Bbr}}$
 $H_a = \neq$

If H_a must be true
 two means must be
 1.96 SDs away -



Overlap happens, then p value increases.

Overlap doesn't actually happen for 1.96 standard deviations.

If we calculate, the two means are significantly different.
The t.stat score we got is 22.7640.

- There are two types of Two Sample Tests :-

- Independent
- Dependent

Independent (unpaired)

t-test-ind(g_1, g_2)

Dependent

t-test-rel(g_1, g_2) (paired)
(Here rel means relative)

g_1 and g_2 are two group arrays. Just upload the generated sample arrays into the function.

- Let us understand about Dependent (Relative groups)
- Consider $n=100$ people are taken randomly.
- The number of antibodies (mean) before vaccination is μ_{ab} before vaccination.
- After Vaccination it will be μ_{ab} after vaccination.
- The company claims that the antibodies will increase after vaccination.

$$H_0 : \mu_{ab} \underset{\text{before Vaccine}}{=} \underset{\text{after Vaccine}}{\mu_{ab}}$$

This is
Dependent
type.

$$H_a : \neq$$

Let us check about the independent :-

- Randomly we are giving covarin to 15000 Indians.
- Randomly we have taken 10,000 African samples covarin.
- After generating the covarin shot, the average num of antibodies in the Indians and Africans is μ_{In} , μ_{Af} .

$$H_0 : \mu_{In} = \mu_{Af} \quad (\text{Both are same})$$

$$H_a : \neq \quad (\text{Both have a difference}).$$

This is the example of an independent test.

Lab Session :-

- Our research question is always an 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.
- There are Type 1 and Type 2 errors.

Type 1 error :- denoted by α , rejecting the null hypothesis incorrectly.

Type 2 error :- Incorrectly fail to reject the null hypothesis.

which means that we should have rejected the null hypothesis but we failed to reject the null hypothesis. Denoted by β .

- When the population standard deviation is not known to us, then we use T-test.
- If it is known we go for Z Test.



Sessions 5 :-

Summarizing the key statistical test :-

It is useful in machine learning model building.

Two cross two Quadrant :-

Whenever we talk about population parameter, either we check mean (or) proportion. When it is a continuous distribution we check whether the means are statistically equal (or) different.

• Whenever we are calculating mean for two samples (or) groups, we call it Two Sample t' Test. (or) Two Sample Mean Test.

- For more than two groups we call it ANOVA.
- t test and Z test are same for $n = 500$.

$$t\text{-test} \Rightarrow \frac{\bar{x} - \mu}{s/\sqrt{n}} \quad \text{To the context of one sample,}$$

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

According to statistics, when $n=30$, $s \stackrel{\text{and above}}{\approx} \sigma$ and $t_{\text{stat}} = z_{\text{stat}}$.

In paired groups we have before and after effects -
For eg:- Before vaccination and after vaccination.

we are going to work with the same samples.

Let us consider we have $n=50$ and $n=50$ as g_1 and g_2 .

First we will calculate $g_1 - g_2$ for all the 50 records.
 Each and every record you take the difference. ΣD will be the sum of total differences.

Another, column, take a squared difference $(g_1 - g_2)^2$.
 ΣD^2 will be the sum of squared differences.

Unlike independent events n_1 and n_2 , we have only one n in dependent groups.

		Independent (unpaired)
(Paired) Dependent	Mean	(Samples)
Proportion	2 groups	Two Sample 't' test
	>2 groups	ANOVA Chi Square Test

• Tstat for paired test \Rightarrow ~~t_{stat}~~ $\frac{\bar{ED}/n}{\sqrt{\frac{\sum ED^2 - \bar{ED}^2}{(n-1)n}}}$.

Tstat $> 1.96 \Rightarrow P\text{value} < 5\% \Rightarrow \text{Reject } H_0$

• Y is dependent variable.

In prediction model, we have combination of dependent and independent variables.

$$Y = f[X]$$

Dependent \downarrow $\begin{array}{l} \text{analog} \\ \text{body space} \\ \text{Colour} \end{array}$ Independent variables.

The price of car is continuous numerical variable.

Regression Analysis.

- $Y = f[X]$ is y as a function of x.
- This kind of prediction model, we call it regression model.
- The outcome whatever we are predicting is continuous numerical model.
- ANOVA will form a foundation for regression.

$$Y = f[X] \text{ (age, bmi etc.)}$$

↓
Blood glucose level

[Instead of telling the blood glucose level,
the model will tell Diabetic / Healthy.]

- Dependent variable Y will have only categorical column which is Diabetic / Healthy.
- Model will predict whether diabetic / healthy. (Classification model).

— Assume an example, our dependent variable ' y ' is categorized into H (Healthy) $\rightarrow 0$ and C (Cancer) $\rightarrow 1$.
The independent variables are x_1 (Age), x_2 (Gender),
 x_3 (Drinking)

— Assume \rightarrow

Age has continuous values $\rightarrow 42.3, 60.36, 38.42$. (x_1)

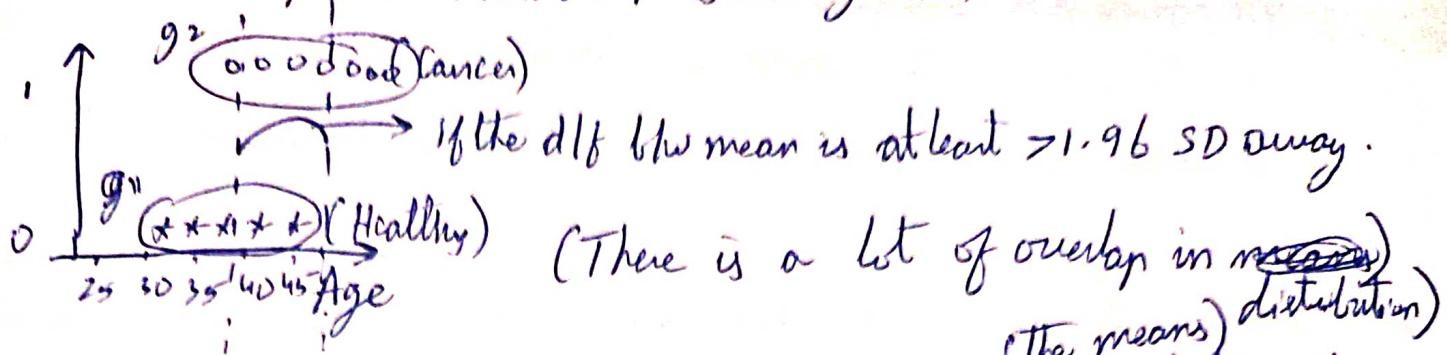
Gender has categories like M and F (x_2)

Drinking has categories like 0 (No drinking), 1 (occasionally)
and 2 (Severely).

— We need to plot the independent and dependent variable now.

X_1	X_2	X_3	y'
Age	Gr	Drinker	
42.3	M	0	H → 0
60.36	M	1	C → 1
38.42	F	1	
:	F	2	
:	M	2	
:	M	2	
!	M	1	
!	M	0	
!	F	2	

Let us consider Dependent variable Y is age in years.



(There is a lot of overlap in ~~means~~ distribution)
(the means)

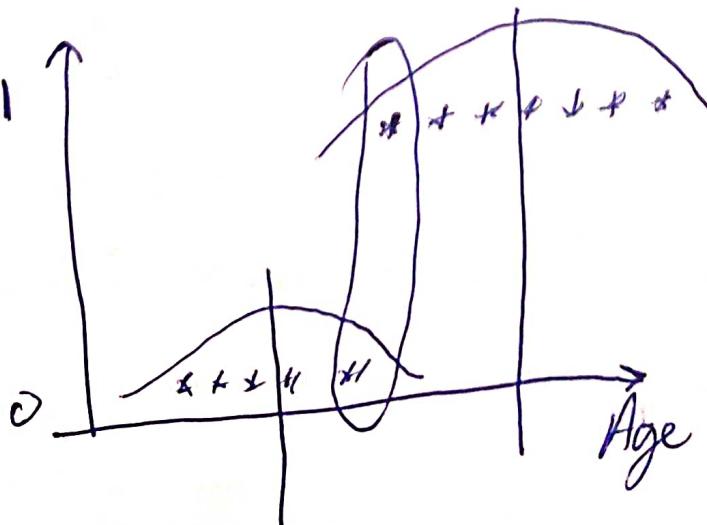
If the means difference is above 1.96, they are statistically different.

If the total score is 1.12, then there is no relation between the age and cancer.

If the means are within 1.12 standard deviations, there is high P-value $P > 5\%$, null hypothesis holds good.

$$H_0 : \bar{age}_0 = \bar{age}_1$$

$$H_a : \neq$$



In this graph there is very less overlap.

$$P < 5\%$$

We need to understand between x_1 and y what statistical test is performed, two sample T test is done to understand whether x_1 is the potential variable in predicting the disease or not.

- whenever we get $P > 5\%$, we will not use the feature in the model.
- If at all we need to use the feature, $P < 5\%$ which means $P < 0.05$, they are statistically significant, we can use in the prediction model.

Now let us consider with x_2 (Gender). We will use the cross tab function now,

`pd.crosstab(df['G'], df['Y'])`.

	H	C
100 M	50	50
80 F	40	40

Will the gender effect cancer?
Answer :- No

$$\frac{50}{100} = 50\%, \quad \frac{40}{80} = 50\%.$$

$$H_0 : P_{MC} = P_{FC}$$

(Proportion of male having cancer is equal to proportion of female having cancer).

$$H_a : \neq$$

- Let us consider,

	C	
n_1	100M	45
n_2	80F	30
		x_1
		x_2

then $P_1 = 0.55$
 $P_2 = 0.63$

P_{pooled} is total proportion of cancer irrespective of the gender.

$$P_{\text{pooled}} = \frac{105}{180}$$

Proportion-Z test formula is $\sqrt{\frac{P_1 - P_2}{P_{\text{pooled}}(1 - P_{\text{pooled}}) \cdot \left(\frac{1}{n_1} + \frac{1}{n_2}\right)}}$

We have a python function,

`proportions_ztest([x1, x2], [n1, n2])`.

- A travel agency wanted to study the effect of holiday choice among men and women
- | | | Dependent | Cruise |
|-------------|-------|-----------|--------|
| | | Beach | |
| Independent | men | 209 | 280 |
| | women | 225 | 248 |
- Travel agency wanted to study whether the gender effects preferred holiday.
- We can formulate the hypothesis in two ways, to the content of beach or to the content of cruise. If any one side of the story is satisfied, automatically the other side is satisfied.

We will formulate hypothesis w.r.t beach first,

$$H_0 : P_{MB} = P_{WB}$$

$$H_a : \neq$$

Then, w.r.t ① cruise holiday,

$$H_0 : P_{MC} = P_{WC}$$

$$H_a : \neq$$

Let's solve w.r.t beach first.

$$P_1 = \frac{209}{489} \quad P_2 = \frac{225}{473} \quad n_1 = 489 \\ n_2 = 473$$

$$P_{\text{pooled}} = \frac{434}{(489+473)}$$

In the interview, there are two contents, t-stats and z-stats,

For continuous variable we say,

tstat (s) In real life we calculate $n \geq 30$, when n

zstat (σ) approaches 500, $n = 500$, then

$$\text{tstat} = \text{zstat}$$

• Two sample proportion test, popularly called as Zstat (or)

Zdata. Proportion-ztest (two sample proportion test).

Let us look at a problem statement:-

The number of Covid19 positive cases in India which is effected for women, expectation is 60%.

$$\begin{aligned} \text{Pop}_{\text{prop CovidW}} &= 60\% \rightarrow H_0 \quad \rightarrow z_{\text{stat}} = 0 \text{ for this} \\ &\neq 60\% \rightarrow H_a \quad P\text{-value} = 100\%. \end{aligned}$$

To verify this we randomly collect 1500 samples, Male are 650 and Female are 850.

What proportion of females are effected? - $850/1500, P_f = 56\%$.

Sample proportion is 56% but expected is 60%.

This above problem comes under one-sample proportion test.

Proportion-Z test ($[x_1, x_2], [n_1, n_2]$)

Whether it is one sample or two sample proportion test we have common function.

We have only one set of information here, then The arguments are:-

Prop-Z test (850, 1500, 0.60)

1:09:31
Chi-square test

- one sample mean test
- two sample mean test
- one sample proportion test
- two sample proportion test

Test of Mean:

one-sample mean

Test

$t\text{test_one sample}(x \text{amp}, \mu)$

Two sample mean test

paired
(dependent)

$t\text{test_rel}(g_1, g_2)$

unpaired
(independent)

$t\text{test_ind}(g_1, g_2)$

Chi-square Test for Goodness of fit of multinomial data :-

$$\chi^2_{\text{data}} = \sum \frac{(O - E)^2}{E} \quad E \text{ is expected, } O \text{ is observed}$$

$$\text{Expected frequency} = \frac{(\text{Row total})(\text{Column total})}{\text{Grand total}}$$

	Healthy	Mild	Severe	✓ This is observed count
Male	10	5	5	
Female	12	6	6	

• 25. If .. .
Similar pattern is observed between the male and female.
then only the null hypothesis holds true

Let us change the observations :-

(Expected count is male proportion
should match the female proportion.

How to arrive the count logically.

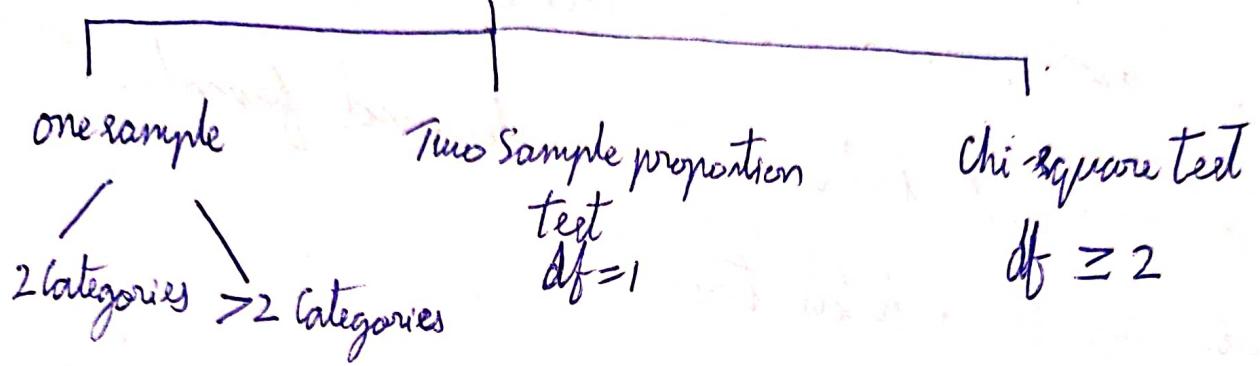
	H	M	S
20 M	8 40%	8 40%	4 20%
24 F	12 50%	6 25%	6 25%

Male healthy expected count is $\frac{(\text{Total male})^{\text{20}} \times (\text{Total healthy people})}{(\text{Total male and female})^{\text{44}}}$

$$\Rightarrow \frac{20 \times 20}{44} = 9.09$$

We have the four flavours of proportion test :-

→ Test of proportions



• Chi square :- we have two groups with degrees of freedom 2 and above
chi²-Contingency() is used.

• Two Sample proportion :- prop-ztest([x₁, x₂], [n₁, n₂])

• One Sample (2 categories) :- prop-ztest(x, n, expected proportion)

• One Sample (>2 categories) :- chisquare([], [])

we should put either numerical list or numpy array

Similarly, the expected count table will be,

Then the proportion percentage will be

9.1	6.3	4.5
10.9	7.63	5.45

45.5%	31.5%	22.5%
45%	31.7%	22.7%

Let us draw the two tables once again:-

Observations Table

	Healthy	Mild	Severe
20 M	8 (40%)	8 (40%)	4 (20%)
24 F	12 (50%)	6 (25%)	6 (25%)

Expected Count

9.1 (45%)	6.3 (31.5%)	4.5 (22.5%)
10.9 (45%)	7.63 (31.7%)	5.45 (22.7%)

Chi-square value is defined as:-

$$\frac{[(\text{Observed Count}) - (\text{Expected Count})]^2}{\text{Expected Count}}$$

Expected Count

+ (for all the six entries).

$$\sum \frac{(8-9.1)^2}{9.1} + \frac{(8-6.3)^2}{6.3} + \frac{(4-4.5)^2}{4.5} + \dots$$

• There is an effect if chi-square value is large.

• If there is correct proportion, then the observations match with the expected count.

Then the chi-square value becomes zero.

So the lower the chi-square value, it ends up with high p-value, the null hypothesis holds good.

Large chi-square value, we end up with small p-value, we reject the null hypothesis.

The chi-square distribution is square of the standard normal distribution.

For Z score the cutoff points are

95%	1.96
99%	2.58
90%	1.67

If p-value is less than 0.05 we reject H_0 .

- Chi₂-contingency is the function for chi-square.
Chi₂-contingency (CT) is the parameter where CT is the name of the table.
- post hoc analysis :- we are going to convey story about the H_a statement.

- If a dependent variable is continuous, it is a regression problem.
- If a dependent variable is categorical, it is a classification problem.

Great learning (Statistics Course) session 6 (ANOVA) :-

Randomized 50 countries
Randomized 50 countries
Randomized 50 countries
Randomized 50 countries

- ANOVA means we are checking the mean for more than two groups.
- If it is two groups, the degree of example is 1, we do two sample mean test.
- But in the example, the degree of freedom is $2(\text{no. of groups} - 1)$, whenever we have two degrees of freedom and we are checking the mean, it is an ANOVA test.
- The below example is age feature w.r.t three categories.
- Sample ages for Groups A, B and C

Group A	Group B	Group C
30	25	25
40	30	30
50	50	40
60	55	45

Sample ages for Groups D, E and F

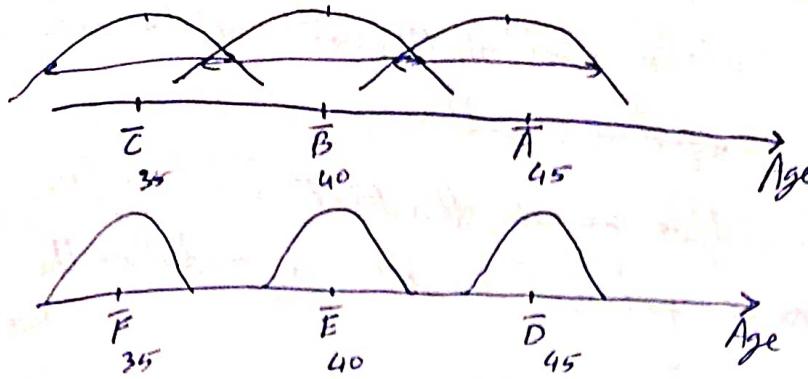
Group D	Group E	Group F
43	37	34
45	40	35
45	40	35
47	43	36

- The groups A, B, C has large variations in their values.
- The groups D, E and F has less variations in their values.

• $\bar{A} = 45$ } The means are same but the spread is large.
 $\bar{D} = 45$

• $\bar{B} = 40$ • $\bar{C} = 35$
 $\bar{E} = 40$ • $\bar{F} = 35$

Let us understand these values diagrammatically.



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

$$H_a: \neq \neq$$

Hence we are checking

$$45 = 40 = 35$$

The null hypothesis will be, $H_0: \bar{D} = \bar{E} = \bar{F}$

$$H_a: \neq \neq$$

Hence also we are checking $45 = 40 = 35$

In the first case H_0 holds true.

In the second case H_a holds true.

In the first case we can see many overlaps, due to overlaps, p value will be high ($P > 0.05$)

If P value is high, H_0 holds good.

In the second case, statistically, one confidence interval will not overlap other confidence interval. So, the p value will be near to zero.

~~As~~ $P \approx 0$, so we reject the null hypothesis, alternate hypothesis

H_a holds good.

Whenever we get P value high, the feature is a useless feature.

There is no significant difference in the mean (first case), hence we have useless features.

In the second one, we have significant difference in the mean and the features are useful.

Now we are going to prove this mathematically.

Let us take another scatterplot representation.

In this example, we have taken continuous variable age as an independent variable, and the categories are SD (severe diabetes), MD (mild diabetes) and H (Healthy). It's something like a classification problem.

For classification problem, definitely we can do this ANOVA test to, infer whether we can infer age as an potential variable in predicting the category of the disease.

Example: Suppose we three different rice grain variety.

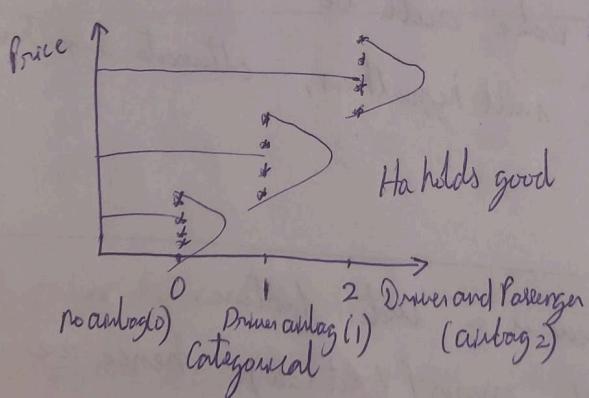
Rice1	16.3 mm
Rice2	18.6 mm
Rice3	22.6 mm

This is a classification problem

Different rice varieties have different shapes. (eccentricity)

To predict the type of rice, whether the $\frac{\text{mm}}{\text{mm}}$ feature is useful or not, If means are significantly different, we can use the feature to categorise.

Assume the price of the car.



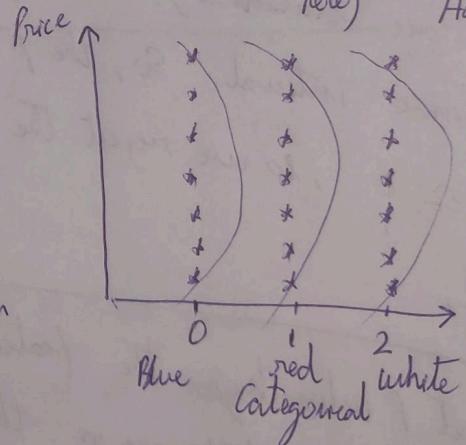
$$\text{Price} = f(\text{airbags})$$

$$\text{Price} = f(\text{color of the car})$$

For ANOVA we have F_{stat} formula,

$$F_{\text{stat}} = \frac{\text{MSTR} (\text{mean square treatment})}{\text{MSE} (\text{mean square error})}$$

(The means are very close here) H_0 holds good



In this example, we have taken continuous variable age as an independent variable, and the categories are SD (severe diabetes), MD (mild diabetes) and H (Healthy). It's something like a classification problem.

For classification problem, definitely we can do this ANOVA test to, infer whether we can infer age as an potential variable in predicting the category of the disease.

Example : Suppose we three different rice grain variety.

Rice1	16.3 mm
Rice2	18.6 mm
Rice3	22.6 mm

This is a classification problem

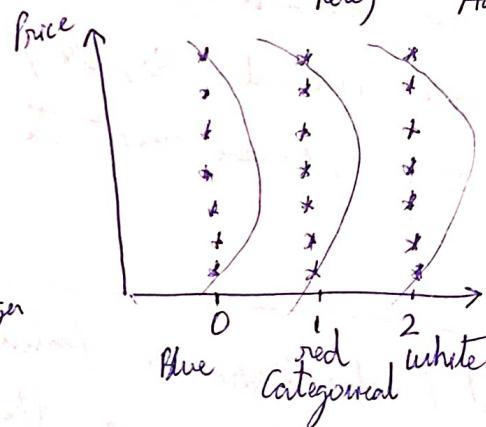
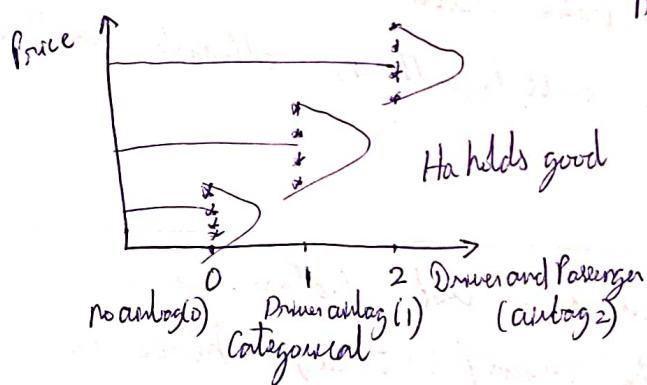
Different rice varieties have different shapes. (eccentricity)

To predict the type of rice, whether the ^{mm} feature is useful or not,

If means are significantly different, we can use the feature to categorise.

Assume the price of the car.

(The means are very close here) H_0 holds good



For ANOVA we have F_{stat} formula,

$$F_{\text{stat}} = \frac{\text{MSTR} (\text{mean square treatment})}{\text{MSE} (\text{mean square error})}$$

MSTR is going to talk about between group variability. (Between sample variability)

MSE means within group variability (within sample variability)

In order the feature to be significant, H_0 holds good. P should be less than 0.05.

If H_0 holds good, F_{stat} should be high.

F_{stat} will be high when Numerator is large and denominator is less.

Higher the value, the means are far apart.

Smaller the values, the means are closer

Whenever the F_{stat} is very high, the p value is very small and vice versa.

Between sample variability: Finding the variance of k-sample means, weighted by sample size and expressed as Mean Square Treatment (MSTR).

$$\text{MSTR} = \frac{\sum n_i (\bar{x}_i - \bar{\bar{x}})^2}{k-1}$$

where n_i is the number of samples in each group. (Here

Numerator of MSTR is the sum of the squares treatment (SSTR) and denominator term $k-1$ is the degree of freedom.

Here,

n_i → number of records in each group.

\bar{x}_i → one group mean value.

$\bar{\bar{x}}$ → global mean

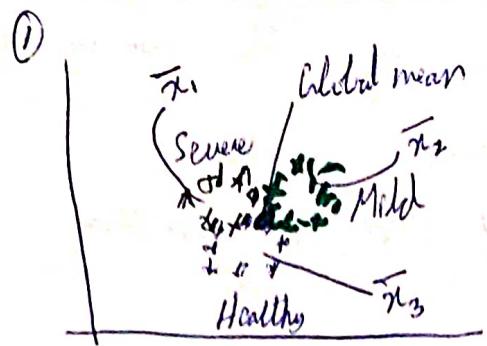
Larger the MSTR value, good the separation is.

Substituting values into the formula,

$$4 * (45-40)^2 + 4 * \frac{(40-40)^2 + 4 * (35-40)^2}{2} = \frac{100+100}{2} = 100$$

$$\text{MSTR} = 100$$

If we plot the records into a scatter plot,



$$\bar{A} = \bar{D} = 45$$

$$\bar{B} = \bar{E} = 60$$

$$\bar{C} = \bar{F} = 35$$

- For (2) MSTR will be large and for (1) it is small.

- Within sample variability :- Finding the weighted mean of sample variances, expressed as Mean Square Error (MSE).

$$MSE = \frac{\sum (n_i - 1) s_i^2}{n_f - k}$$

$$F_{\text{data}} = \frac{MSTR}{MSE}$$

- Numerator of MSE is the sum of the squares error (SSE) and denominator term $n_f - k$ is the degrees of freedom.
- The total number of squares (SST) is the sum of SSTR + SSE. ANOVA table is the convenient way to display all these parameters.

Here,

$n_i - 1$ is each group number of records - 1

s_i^2 is the variance of groups. (s_1^2, s_2^2 and s_3^2)

n_f means the total number of records

k is the total number of groups.

Whenever we reject H_0 , we have some interesting story ~~to tell~~ to tell from H_a .

That is the post hoc analysis.

In the reality:

stat. Test

Reject H_0

Accept H_0

H_0 : True

This is error (\times)
Type I (or) α

Whenever H_0 is true, we
are accepting H_0

Right ✓

H_0 : False

Whenever H_0 is false
we reject H_0 .

Right decision

This is error (\times)

Type II or β

If no one knows there
in the reality.

one sample test

Test of mean test - 1 sample (sample, μ)

Test of prop

prop-ztest (x, n, prop)

Z-test

Chi-square

Summary of Learning

standard

To check the normality of the data :-

Acceptable range of skewness is -0.5 to 0.5.

from scipy.stats import shapiro.

shapiro (df['age']), it is going to return the p value.

Data = Normal

H_0 :

$\boxed{\text{Data} \neq \text{Normal}}$

H_a :

This is skewness effect

If ~~shapiro~~ when we do shapiro of our numerical column, it returns, for example :-

$P \rightarrow 0.000138$

(Data is skewed, ~~it~~ because we reject the null hypothesis)

H_0 holds good, the data is normal.

If $P = 0.348$, then

Summary of Learning:

- Understanding probability distribution

- Discrete

- Continuous

- Normal

- Central Limit Theorem } Theoretical and practical approach
In practical approach what approximation we do, $\frac{s}{\sqrt{n}}$, ie standard error of mean.

- Range estimate for the confidence interval

- Test of mean t-test-one-sample(xamp, μ)

- One sample Test

- Test of proportion prop-ztest(x, n, prop)

- Two sample Test

- Test of mean
 - Dependent
 - Independent

- Test of proportion prop-ztest($[x_1, x_2], [n_1, n_2]$)

one sample proportion test with more than two categories.

- One sample Test / > 2 categories.

- Chi-square (observed count, expected count)

- Two sample proportion test doesn't have any flavours.

- If we are checking more than two group test.

- Test of mean

- ANOVA

- Test of proportion

- chi². Contingency

Lab statistics

what are the applications of chi-square.

→ Goodness of fit test

Example: Suppose you have a data, and you want to test whether the data supports poisson distribution. Such kind of tests are called goodness of fit.

→ Independence Test

The relationship between two features, for example x and y are independent or dependent.

→ Homogeneity test

To test how uniform the data is.

If the calculated chi-square is exceeding the chi-square critical, then we reject the null hypothesis. (like any other test).

Contingency Table

means the cross tabs

Example: A dice is rolled 132 times, find whether it is biased or unbiased.

chi-square critical value $\chi^2_{\text{crit}} = 11.07$.

1	2	3	4	5	6
16	20	25	14	29	28

	f_o (frequency observed)	f_e (expected)
1	16	22
2	20	22
3	25	22
4	14	22
5	29	22
6	28	22

$O-E$	$(O-E)^2$	$(O-E)^2/E$
-6	36	36/22
-2	4	4/22
3	9	9/22
-8	64	64/22
7	49	49/22
6	36	36/22

$$\chi^2 = \frac{\sum (O-E)^2}{E} = \frac{36+49+64+9+49+36}{22} = 9.$$

chi-square calculated < χ^2_{crit} , so we will fail to reject the null hypothesis.

This is the goodness of fit example.

Let's talk about Test of independence,
How do we calculate the expected frequency when we have two features
and you are trying to look at the proportions of the features w.r.t each other.

$$\text{Expected} = \frac{\text{Column Total} \times \text{Row Total}}{\text{Overall Total (N)}}$$

Example:- A total of 10,000 people ~~voted~~, out of which some voted and some did not vote. Out of the people who voted 2792 men and 3591 women and out of the people who didn't vote 1486 were men & 2131 women. We need to check if the gender and voting are independent to each other or not.

Solution:-

$$\chi^2_{\text{calc}} = 3.841$$

Observed data

	Men	Women
vote	2792	3591
no vote	1486	2131
	↓	↓
	4278	5722

$$\rightarrow 6383$$

$$\rightarrow 3617$$

Expected data

	Men	Women
vote	$\frac{4278 \times 6383}{10,000} = 2731$	$\frac{6383 \times 5722}{10,000} = 3652$
no vote	$\frac{4278 \times 3617}{10,000} = 1547$	$\frac{5722 \times 3617}{10,000} = 2070$

$$\text{Chi-square } \chi^2_{\text{cal}} = 6.6.$$

Since calculated value is exceeding the critical value, we reject the null hypothesis which means the gender and voting are dependent on each other.

Example: A sample of $n=60$ size is taken, and the observations that are seen are

0	1	2	3
32	15	9	4

Does the data follow poisson distribution or not.

Solution:-

Hint: To calculate the expected frequencies, we should do the (probability distribution function) \times (sample size) for each case.

$$f_e = \text{PDF} \times n \quad \hat{n}_{\text{cnt}} = 7.814 \\ = f(n) \times n$$

Because, this is a poisson distribution, we need to know, what is lambda.

λ is the average $\therefore \lambda = \frac{(0 \times 32) + (1 \times 15) + (2 \times 9) + (3 \times 4)}{60} = 0.75$

Exp data Calculation

$$f(n) = \frac{e^{-\lambda} \cdot \lambda^n}{n!}$$

$$f(0) = \frac{e^{-0.75} \times (0.75)^0}{0!} = 0.47$$

$$f(1) = \frac{e^{-0.75} \times (0.75)^1}{1!}$$

$$\bar{f}_e = 0.47 \times n = 28.34$$

The Chi-square calculated = 4.47.

So, it is lesser than \hat{n}_{cnt} . We fail to reject the null hypothesis.

Data follows Poisson distribution.

This is a goodness of fit example.

ANOVA \rightarrow Analysis of Variance

We are analysing the variance between the samples.

Why ANOVA?

To this point we have been comparing two populations.

\rightarrow Independent samples t-test (random).

\rightarrow Matched sample t-test (paired).

- Of course limiting ourselves to the comparison of two populations is well ... limiting.
- what if we wish to compare the means of more than two populations?
- what if we wish to compare populations each containing several levels or subgroups?
- Enter... ANOVA; Analysis Of Variance.

Things solved during ANOVA test,

- Suppose we want to compare three sample means to see if a difference exists somewhere among them.
- Is one mean so far away from the other two that it is likely not from the same population?
- What we are asking is:
 - Do all three of these means come from a common population?
 - Or are all three so far apart that they all likely come from unique populations?

Null hypothesis

$$H_0: \mu_1 = \mu_2 = \mu_3$$

We have one way ANOVA

Two way ANOVA with repetition

Two Way ANOVA without repetition

Suppose we have three samples, S_1, S_2 and S_3 and their respective means are \bar{x}_1, \bar{x}_2 and \bar{x}_3 .

How would you calculate the f ratio.

- Let's talk about one-way ANOVA first,
- we need F value and compare it to the Fcritical
- If $F_{\text{value}} > F_{\text{critical}}$ then reject the null hypothesis.

$$F_{\text{cal}} = \frac{MSC}{MSE}$$

$$\text{SSC} = (\bar{x}_1 - \mu)^2 + (\bar{x}_2 - \mu)^2 + (\bar{x}_3 - \mu)^2$$

(between)

$$MSC = \frac{\text{SSC}}{dfc}$$

$$dfc = \text{num of columns} - 1$$

F_{cal} is $F_{\text{calculated}}$.

dfc is degree of freedom of columns.

SSC is sum of squares of the columns.

$$MSE = \frac{SSE}{dfc}$$

dfc is degree of freedom of error
 $= \text{total population (N)} - \text{num of columns (C)}$

, SSE is within and SSC is between

. Within means Variance within each sample

Basically the variance of each data point in a sample from its respective mean

$$SSE = SST - SSC$$

where SST is total sum of squares
 $\sum (x - \mu)^2$

Problem statement

Twenty-one students at the autonomous university of madrid (UAM) in spain were selected for an informal study about student study skills; 7 first year, 7 second year and 7 third year undergraduates were randomly selected.

The students were given a study-skills assessment having a maximum score of 100. As researchers we are interested in whether or not a difference exists somewhere between the three different year levels. We will conduct this analysis using a one-way ANOVA technique.

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

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

$$dfc = c - 1$$

$$\cancel{dfe} = \cancel{(c-1)(B-1)} dfc = (c-1)(B-1)$$

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

$$df_T = N - 1$$