

Basics stats - Beginner

→ Probability

1) Intro to Basic Term,

2 Variables

3. Random Variables

4) Population, Sample, Population-mean, Sample mean.

5) Population Distribution, Sample Distribution & Sampling distri...

6) Mean, Median, Mode

7) Range

8) Measure of Dispersion

9) Variance

10) Std Deviation

11) Gaussian/Normal Distribution

Intermediate stats

12) Standard Normal Distribution

13) Z score

14) Probability Density function

15) Cumulative distribution fn

16) Hypothesis Testing

17) Many different plotting graphs

18) Kernel Density Estimator

19) Central Limit Theorem

20) Skewness of data

21) Covariance

22) Pearson Correlation Coefficient

23) Spearman Rank Correlation

24) Hypothesis Testing.

Advanced stats

25) Q-Q plot

26) Chebyshev's inequality

27) Discrete and cts Distribution

28) Bernoulli & Binomial dist.

29) Log Normal Dist.

30) Power law Distribution

31) Box Cox Transform

32) Poisson Distribution

33) Application of non-Gaussian Distribution

Population
Sample

Dyscriptive stats

- ① Analyzing data, Summarizing data
organizing data, in form of
numbers, & graphs

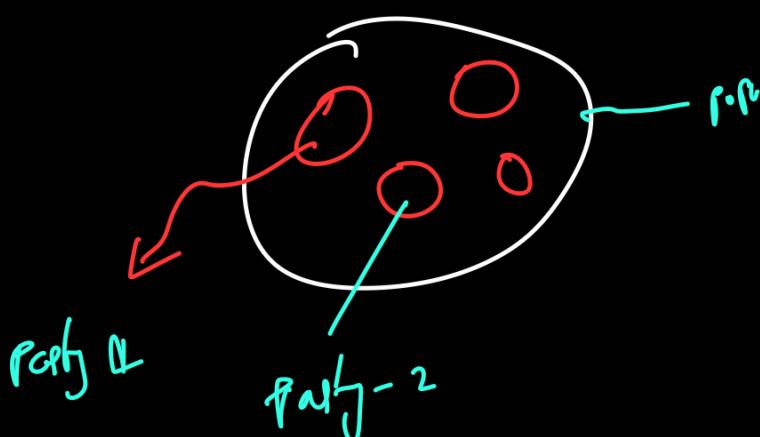
- ② Box plot, histogram, pie-chart,
PDF, CDF, Normal D.

- ③ Measures of Central Tendency
Mean, Median, Mode

Sample → inference & conclusions

↓
population
exist popls.

Inferential stats



- { ① Z Test → Hypothesis
② T test
③ Chi square Test

Statistics

- ① Population → population mean
② Sample → sample mean
③ Random variable → Discrete random variable
 cts random variable.

mean
median
mode } measures of central tendency

$$\rightarrow \text{Population mean} = \mu = \frac{1}{N} \sum_{i=1}^N h_i$$
$$\rightarrow \text{Sample mean } (\bar{x}) = \frac{1}{n} \sum_{i=1}^n h_i$$

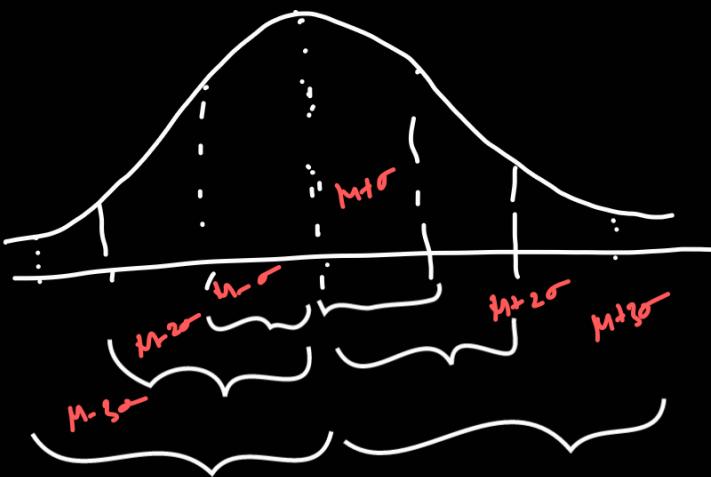
1 → Discrete → 2, 3, 4 [whole numbers]
2 → Cts → within a range of val.
10 — 15 [10.1, 10.2 ... 15]

Gaussian / normal distribution

$$x \sim \text{G.D}(\mu, \sigma)$$

mean $\mu = \frac{1}{n} \sum_{i=1}^n x_i$
 s_D

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



$$\Pr(\mu - \sigma \leq x \leq \mu + \sigma) \approx 68\%$$

$$\Pr(\mu - 2\sigma \leq x \leq \mu + 2\sigma) \approx 95\%$$

$$\Pr(\mu - 3\sigma \leq x \leq \mu + 3\sigma) \approx 99.7\%$$

Empirical formulae.

Log normal Distribution

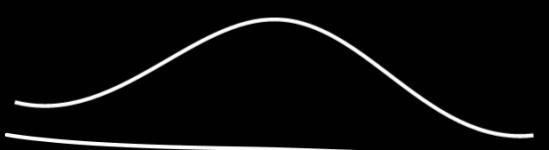
$X \sim \text{Log normal distribution}$

if $\ln(x)$ is normally distributed

$$\rightarrow X = \{x_1, x_2, x_3, \dots, x_n\}$$

$$\log(x_1), \log(x_2), \dots, \log(x_n)$$

normally distributed



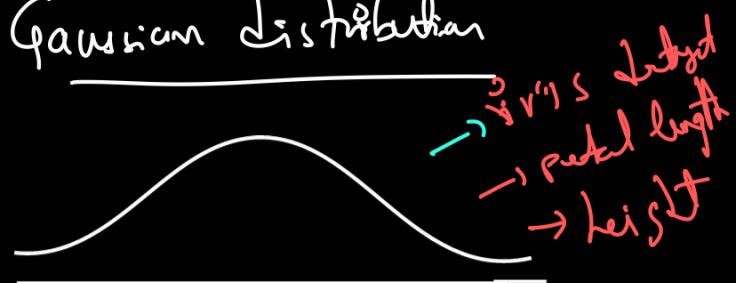
→ lognormal
↓

$X \sim \text{log normal}$

If $\log(x)$ is normally distributed.

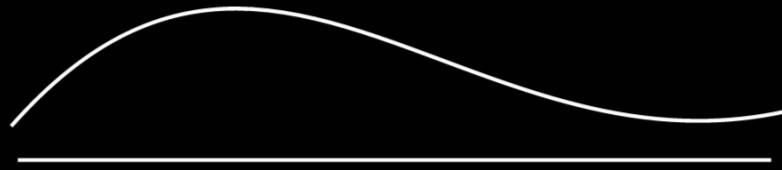
↳ $\log(x)$ normal distribution \Rightarrow a log normal distribution

Gaussian distribution



Log normal

→ right skewed



↳ income of people.

↳ feedback of product

↳ comments length in amazon.

→ why distributions are 'im'?

→ to convert all into same std scale.

Σ_{α} :

Eqn

market

theo

prct

$$\Rightarrow \text{Cov}(\alpha, \alpha) = \text{Var}(\alpha)$$

$$\rightarrow \text{Cov}(x, y) < \begin{array}{l} \text{if } y \uparrow \text{ -ve cov} \\ \text{if } y \downarrow \text{ -ve cov} \end{array}$$

\rightarrow Covariance only says +ve or -ve

\hookrightarrow but don't say how much
true or -ve, so we use
Pearson correlation coefficient!

mean, median, mode

$$[1, 2, 3, 4, 5] = \text{mean} = \frac{1+2+3+4+5}{5}$$

$$M=3$$

If 1 add an outlier 50

$$[1, 2, 3, 4, 5, 50] = \text{mean} = \frac{1+2+3+4+5+50}{6}$$

$$M_{\text{new}} = 13$$

\rightarrow So, mean is highly influenced by outliers

So, we use median.

$$[1, 2, 3, 4, 5, 50]$$

$$\frac{3+4}{2} = \frac{7}{2} = \underline{\underline{3.5}}$$

Covariance

$$\Sigma_{\alpha} = \frac{\text{sum of hours}}{\text{sum of price}}$$

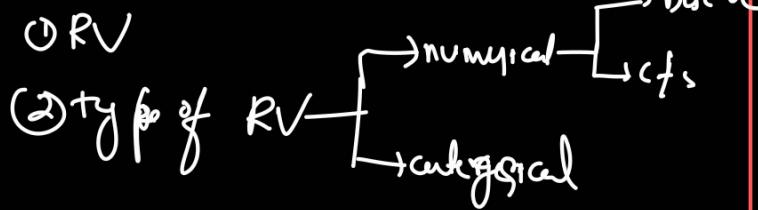
$S \uparrow$	$P \uparrow$
$S \downarrow$	$P \downarrow$

$$\left[\text{Cov}(S_{\text{theo}}, P_{\text{theo}}) = \frac{1}{n} \sum_{i=1}^n (\alpha_i - \bar{\alpha}) \times (y_i - \bar{y}) \right]$$

$$\text{Var}(\alpha) = \frac{1}{n} \sum_{i=1}^n (\alpha_i - \bar{\alpha})^2$$

\rightarrow If there are lot of outliers,
use median & mode
 \downarrow
otherwise mean.

Random variables



Central Limit Theorem

$$X \not\sim \text{G.D.}(\mu, \sigma^2)$$

$$\boxed{n \geq 30}$$

$$x_1, x_2, \dots, x_{30} = \bar{x}_1,$$

$$\begin{aligned} x_2 &= \bar{x}_2 \\ x_3 &= \bar{x}_3 \\ &\vdots \\ x_{100} &= \bar{x}_{100} \end{aligned}$$

$$\bar{X}(\underbrace{\bar{x}_1, \bar{x}_2, \bar{x}_3, \dots, \bar{x}_{100}}_{\text{mean}}) \approx \text{G.D.}(\mu, \frac{\sigma^2}{n})$$

$$\bar{X} \approx \mu$$

CAE SYSTEM'S INEQUALITY

$$X \sim \text{G.D.}(\mu, \sigma^2)$$

$$\Pr(\mu - \sigma \leq X \leq \mu + \sigma) \approx 68\%$$

$$\Pr(\mu - 2\sigma \leq X \leq \mu + 2\sigma) \approx 95\%$$

$$\Pr(\mu - 3\sigma \leq X \leq \mu + 3\sigma) \approx 99.7\%$$

But if not follow G.D.

$y \not\sim \text{G.D.}$

$$\Pr(\mu - k\sigma \leq Y \leq \mu + k\sigma) > 1 - \frac{1}{k^2}$$

$$k=1 \quad \Pr(\mu - \sigma \leq Y \leq \mu + \sigma) > 1 - \frac{1}{1^2} = 0$$

$$k=2 \quad \Pr(\mu - 2\sigma \leq Y \leq \mu + 2\sigma) > 1 - \frac{1}{4} = 75\%$$

$$k=3 \quad \Pr(\mu - 3\sigma \leq Y \leq \mu + 3\sigma) > 1 - \frac{1}{9} = 88\%$$

PEARSON CORRELATION COEFFICIENT

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

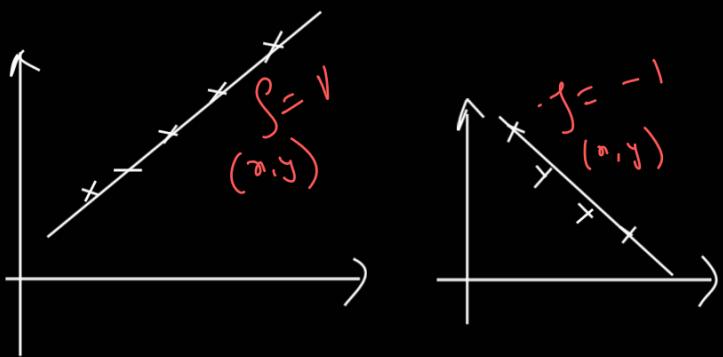
$$\text{② PEARSON CC} = \rho_{(x,y)} = \frac{\text{cov}(x,y)}{\sigma_x \sigma_y}$$

① SIGNIFICANT

② DIRECTION OF RELATIONSHIP

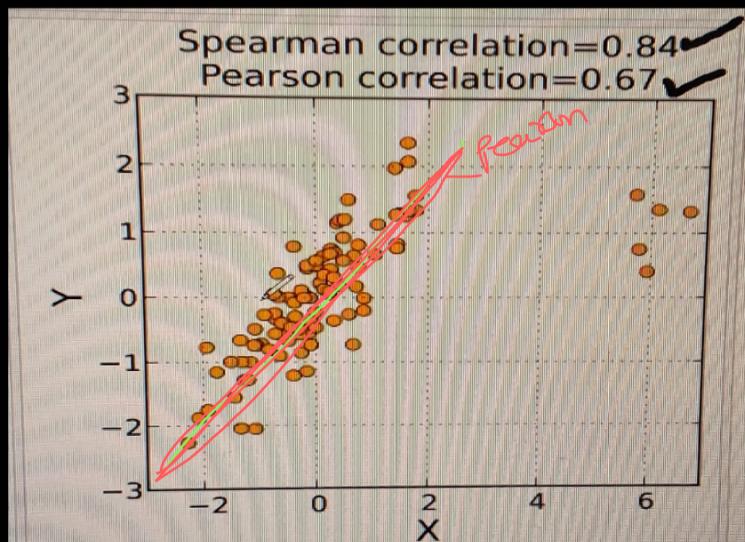
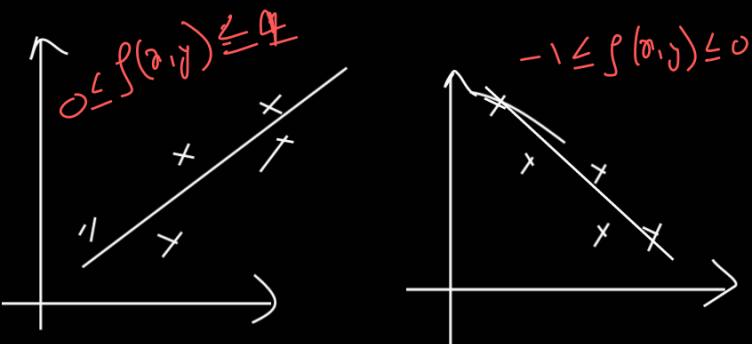
$$\boxed{-1 \leq \rho \leq 1}$$

③ Spearman rank CC



$$\textcircled{5} \quad \rho = 1 - \frac{6 \sum d_i^2}{n(n^2 - 1)}$$

If ρ is close to 0 \rightarrow no corr
 1 \rightarrow highly corr
 -1 \rightarrow negative corr



Spearman Correlation

x	y	rank x	rank y	d_i	d_i^2
1	2	1	2	0	0
2	1	2	1	-1	1
3	3	3	3	0	0
4	4	4	4	0	0
5	5	5	5	0	0
6	6	6	6	0	0
7	7	7	7	0	0
8	8	8	8	0	0
9	9	9	9	0	0
10	10	10	10	0	0

① Sort x in ascending order

② keep corresponding values in y .

③ assign 1, 2, 3... rank for x .

↳ corresponding ranks to y .

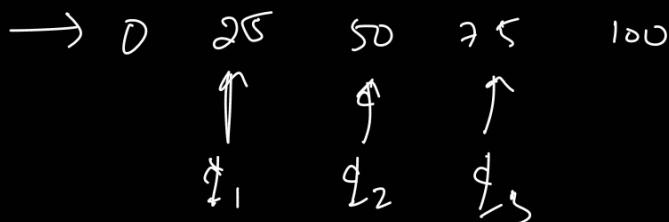
④ find diff of rank of x and y .

⑤ square d_i^2

$$\rho = \frac{\text{cov}(R_x, R_y)}{\sigma_x \times \sigma_y}$$

→ Various ways to find outliers

- ① Using scatter plot
- ② Box plot
- ③ Using z-score
- ④ Using the IQR interquartile range



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

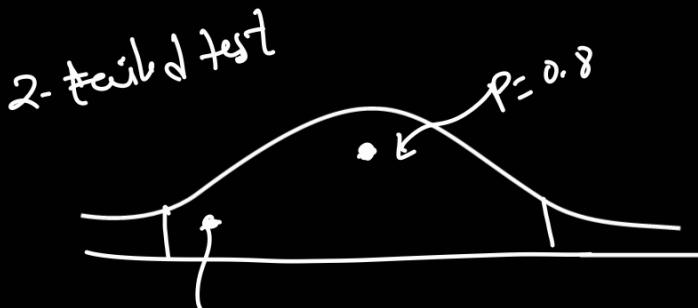
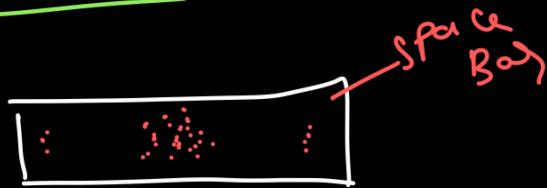
lower bound = $Q_1 - 1.5 \times \text{IQR}$
 upper bound = $Q_3 + 1.5 \times \text{IQR}$

Normalization (min-max normalization)

↪ b/w (0, 1)

$$X_{\text{norm}} = \frac{X - X_{\text{min}}}{X_{\text{max}} - X_{\text{min}}}$$

P-value



P=0.01
 (if we repeat 100 times
 we will touch this place
 1 time.)

Probability

Penalty

$P=0.8 \rightarrow$ if we repeat 100 times
 we will touch 80 times at the place

→ null hypothesis:

Treats everything same or equal.

→ p-value:

If $\hat{\mu}$ is the probability for the "null hypothesis" to be true.

Confidence Intervals

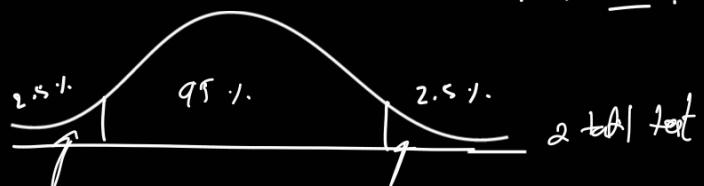
① Point estimate $\hat{\mu} \rightarrow \bar{x}$



② Population std given

③ Population std not given

④ Average size of sharks in sea. with 95% CI.



lets assume $\sigma = 100$ $n = 30$ $\bar{x} = 500$
 if given $\mu = 386 - 613$

$$\text{CI} = \text{Point Estimate} \pm \text{margin Error}$$

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

$$= 500 \pm 2.00 \frac{100}{\sqrt{30}}$$

$$\text{lower limit} = 500 - 1.96 \frac{100}{\sqrt{30}} = 386$$

$$\text{higher limit} = 500 + 1.96 \frac{100}{\sqrt{30}} = 613$$

$$\begin{aligned} z_{0.025} &= 1.96 \\ &= 0.975 \\ \therefore 1.96 & \end{aligned}$$

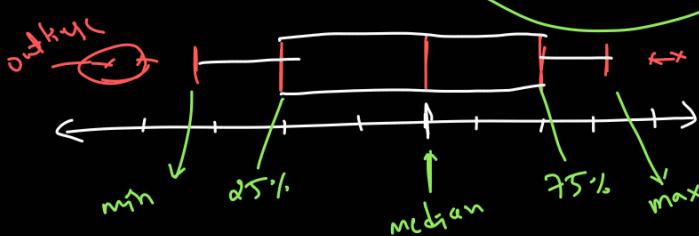
- # 5 Number Summary
- ① Minimum
 - ② 25% Percentile Q_1
 - ③ Median
 - ④ 75% Percentile Q_3
 - ⑤ Max

$$IQR = Q_3 - Q_1$$

[lower bracket higher bracket]

$$Q_1 - 1.5 IQR$$

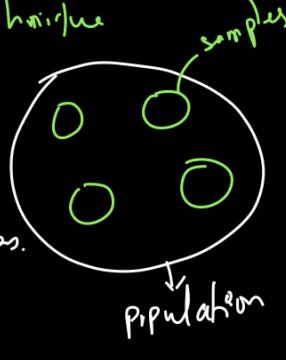
$$Q_3 + 1.5 IQR$$



Q) What are different sampling Techniques?

① Random Sampling Technique

→ for exit polls, we are picking some random people from different areas. whether boy/girls

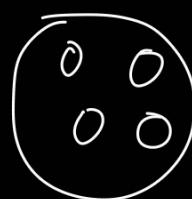


② Stratified Sampling Technique

→ if we pick random samples such that we take males : females

$$= \underline{1:1}$$

↳ no bias



③ Systematic Sampling

→ if we select a particular location and if we select every "n" person in that place.



④ Cluster Sampling

↳ if I want to do a survey of AI ML

↳ only take survey of domain people.

Confusion matrix

		Predicted	
		P	N
Actual	P	TP	FN
	N	FP	TN

Type II Error $\rightarrow FN$

Type I Error $\rightarrow FP$

① Whether the person has a disease or not.

→ for positive [cancer], if detect -ve, big problem.
So, we need to decrease FN

↳ there is no problem, if for -ve, gives +ve.

↳ so we adopt $\frac{TP}{TP+FP} \rightarrow \text{Precision}$

↳ $TP \uparrow, TN \uparrow, FN \downarrow, FP \downarrow$

② Market will crash or not

↳ same, if market falls [true], & predict -ve, so lose money.

↳ Also, if market don't fall, also loose money by selling.

↳ so, we need to focus on both $FN \downarrow, FP \downarrow$

$$\frac{TP+TN}{TP+TN+FP+FN} = \text{accuracy}$$

↳ f score

↳ In this problem, both reducing & focusing on both FP & FN is imp.

③ Vaccination side effect.

→ if +ve, predicts -ve big problem, so $FP \downarrow$

→ also, we loose vaccination if -ve, predict +ve

→ So, both $FN \downarrow, FP \downarrow$

→ $FP \downarrow$ is imp → adopt $\frac{TP}{TP+FP} \rightarrow \text{Recall}$

↳ spam mail detection.

↳ it's ok, if some (-ve) mails not detected & predict (FN), [we can remove]

↳ it shouldn't remove imp mails $[FP \downarrow \downarrow]$

Z-score importance

- z_1 is how many std away from mean
- helps in standardization.
- helps to find the best performance b/w the different groups of data.

→ Σ : compare the scores b/w different tests

2020

$$\text{Avg} = 181$$

$$\sigma = 12$$

$$\text{final score} = 187$$

$$z_{2020} = \frac{187 - 181}{12} = \frac{6}{12} = 0.5$$

2021

$$\text{Avg} = 182$$

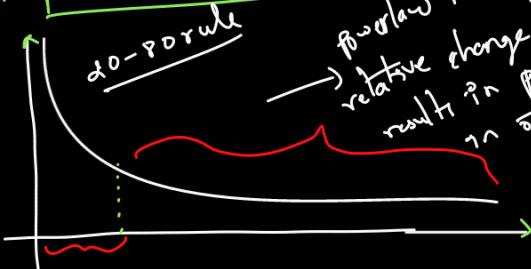
$$\sigma = 5$$

$$\text{final} = 185$$

$$z_{2021} = \frac{185 - 182}{5} = \frac{3}{5} = 0.6$$

→ So, 2021 performance is better.

Power law distribution

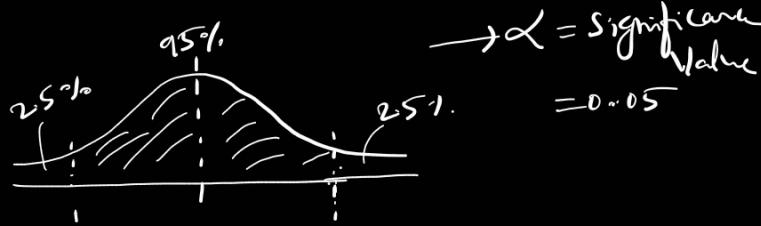


- ① 80% of sales are coming from 20% of points.
- ② 80% of windows crash by 20% of bugs.

Hypothesis Testing - confidence Interval



→ out of 100 touches, I will touch 80 times then



→ α = significance value
 $= 0.05$

Hypothesis testing

- ① Null hypothesis → Alternate hypothesis
- ② Perform Experiments $\rightarrow \alpha = 0.05$ { t-test, z-test, chi-square, ANOVA }
- ③ P-value falls in extreme End \rightarrow Reject Null

→ Average height of people in India

↳ population (X) → samples → conduct exp

↙ do statistical tests

- ### # Statistical tests
- ① t-test [Comparison of mean] → sample (s)
 - ② z-Test [Comparison of mean] → a mean
 - ③ Anova Test (Analysis of Variance)
 - ④ F Test (Comparison of variance)
 - ⑤ Chi square (Comparison Categorical variables)

① One Sample Z-test

→ In a population, the avg IP is 100 & std is 15. The doctor tested the new medicine whether it \uparrow s or \downarrow s the IQ. After 1 month sample of 30 people were taken, had a mean of $\bar{x} = 140$. Did medicine effect intelligence? ($\alpha = 0.05$).

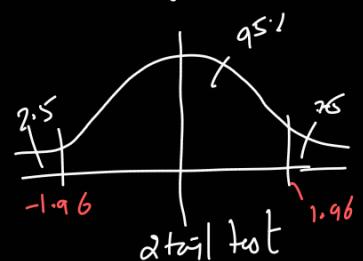
- ① Null hypothesis $H_0 \Rightarrow \mu = 100$ [no effect]
- ② Alt hypothesis $\Rightarrow \mu \neq 100$ [effect].

$$\text{③ significance level } \alpha = 0.05$$

④ state decision rule →

$$\text{if b/w } -1.96 \text{ to } +1.96 \text{ accept}$$

$$\text{if } < -1.96 \text{ or } > +1.96 \text{ reject Null}$$



$$\text{④ } z\text{-test: } z = \frac{\bar{x} - \mu}{\sigma/\sqrt{n}}$$

$$= \frac{140 - 100}{\frac{15}{\sqrt{30}}} = \frac{40}{2.74} = 14.60 > 1.96 \rightarrow \text{so reject } H_0$$

* So, change in IQ

One Sample Z test with proportion
 Q) A survey claims that 9 out of 10 doctors recommend aspirin for patients with headache.
 To test this claim, random sample of 52 doctors taken out of 82 recommends this aspirin. Is this claim accurate?

$$\alpha = 0.05.$$

$\rightarrow 9 \text{ out of } 10 \rightarrow \text{proportion.}$

$$\rightarrow H_0: \text{Null Hypothesis} = p = 0.90 \quad [9/10]$$

$$H_1 \neq 0.90.$$

(2) Significance $\alpha = 0.05$

(3) Decision Rule.

(4) Z-statistics

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

$$Z_0 = \frac{0.82 - 0.90}{\sqrt{\frac{0.90(0.10)}{100}}} = \frac{-0.08}{\sqrt{0.009}} = \frac{-0.08}{0.0949} = -0.84$$

= [-2.667] → reject H_0

(5) Make Conclusion :-

$$-2.667 \notin [-1.96, 1.96]$$

Reject the Null hypothesis
 Accept the Alt hypothesis

* CLAIM IS INACCURATE

Probability

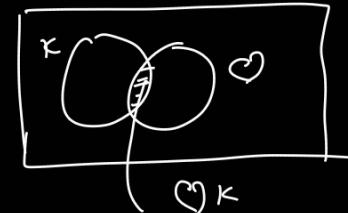
{It is the measure of likelihood of an event}

→ Ex: Tossing a coin

$\rightarrow \text{Probability} = \frac{\# \text{ of ways an event can occur}}{\# \text{ of possible outcomes.}}$

(1) Mutual Exclusive → tossing a coin +/T

(2) Not mutual exclusive → taking out a card from deck (X & C)



$$\rightarrow \text{Probability of } (+ \text{ or } T) = \frac{1}{2} + \frac{1}{2} = 1$$

↪ Additive rule for mutual exclusive

$$\rightarrow P(K \text{ or } C) = P(K) + P(C) - P(K \cap C)$$

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

Additive

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

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

Multiplicative Rule

(1) Dependent Event

$$+ + + \\ \frac{1}{2} \times \frac{1}{2} \times \frac{1}{2}$$

(2) Independent event

$$\text{deck of cards} \\ \frac{1}{52} \times \frac{1}{51} \times \frac{1}{50}$$

$$\downarrow \text{conditional probability} \\ P(K \text{ and } C) \Rightarrow P(K) \times P(C|K) \\ = \frac{1}{52} \times \frac{1}{51}$$

Permutation of combination

6 reads.

Rock	Raj
Run	Kum
chin	cut

$$\frac{6}{1} \times \frac{5}{2} \times \frac{4}{3} = 120$$

↪ order matters

Permutation

$$n_p = \frac{n!}{(n-r)!} = \frac{6!}{(6-3)!} = \frac{6 \times 5 \times 4 \times 3 \times 2 \times 1}{3 \times 2 \times 1} = 120$$

Combinations → order don't matter.

$$n_c = \frac{n!}{(n-r)!r!} = \frac{6!}{(6-3)!3!} = \frac{6 \times 5 \times 4 \times 3 \times 2 \times 1}{3 \times 2 \times 1} = 20$$

Why Sample Variance is divided by $n-1$ # for cts data → Histogram.

Population (N)	Sample (n)
$M = \frac{1}{N} \sum_{i=1}^N x_i$	$\bar{x} = \frac{1}{n} \sum_{i=1}^n x_i$
Variance $\sigma^2 = \frac{1}{N} \sum_{i=1}^N (x_i - M)^2$	$s^2 = \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})^2$

→ if we select only sample which is extreme so, our (s^2) will vary much.

so, scratch fact $n-1, n-2, \dots$ & get good result with $(n-1) \rightarrow \sigma^2 \approx s^2$

↳ skewed data → unbiased estimation.

Variables

Random variables → Numerical → Discrete cts
→ Categorical

* Variable Measurement Scales :

4 types of measurement Variable

- ① Nominal → {categorical data} color, gender, type of flower.
- ② Ordinal → order of data matters, value doesn't
- ③ Interval → order & value matters, '0' not present
- ④ Ratio →

Frequency distribution

→ flower class {Rose, lily, Rose, lily, sunflower, }.

flower	freqn	Cumulative
Rose	3	3
lily	4	7
sunfl.	2	9

\bar{x} = 1.866

$\#$ for cts data → Histogram.

$Ages = \{10, 12, 14, 18, 24, 26, 30, 35, 36, 37, 40, 41, 42, 43, 50, 51, 52\}$

$\hookrightarrow Bins = 10$

→ pdf → smoothing of histogram.

When the variance is more?

Variance → spread of data.

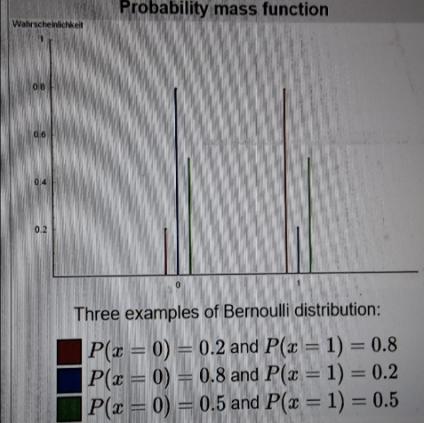
how much? $\mu = 100$, $\sigma = 15$

$Z = \frac{85 - 100}{15} = \frac{-15}{15} = -1$

$Z_{@-1} = 0.15 \times 15 = 2.25$

Bernoulli distribution

Bernoulli distribution
Probability mass function



→ Outcome $\rightarrow \{0, 1\}$
 Eg: Tossing a coin
 $x=1 \rightarrow \text{success}$
 $x=0 \rightarrow \text{fail}$
 $f(x=x) = P^x (1-P)^{1-x}$

This is called {pmf}
 ↓
 discrete discrete

$$P(\text{success}) = P, P(\text{fail}) = 1-P = q$$

$$P(x=0) = P^0 (1-P)^{1-0} = (1-P) = q$$

$$P(x=1) = P^1 (1-P)^{1-1} = P$$

$$\text{PMF} = \begin{cases} q = 1-p & \text{if } x=0 \\ p & \text{if } x=1 \end{cases}$$

$$P^x (1-P)^{1-x} \Rightarrow \underline{\text{pmf}}$$

$$\begin{aligned} \text{mean} \quad \bar{x}(x) &= \sum_{x=1}^2 x \cdot p(x) & x=0 \quad \text{or} \quad 1 \\ & P(x=0) = 0.4 = 1-p & P(x=1) = 0.6 \\ & \downarrow & \downarrow p \\ & = 0 \cdot x (0.4) + 1 \cdot x (0.6) & \\ & = 0.6 \rightarrow P \end{aligned}$$

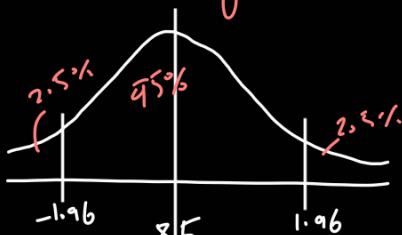
$$\Rightarrow \text{Variance: } P(1-P) = pq$$

$$\sigma = SD = \sqrt{pq}$$

mean

mode.

Eg: Colleges in Karnataka have 85% placement rate. A new college was opened and it was found that a sample of 150 students had a placement rate of 88% with a std 4%. Does this clg has a diff placement rate? $\alpha = 0.05$



2 tail test.

- ① Null hypothesis (H_0): It has diff placement rate
- ② Alt Hypothesis (H_1): It has some placement rate
- ③ Significance value = 0.05
- ④ Decision Rule: $Z = \frac{\bar{x} - \mu}{\sigma/\sqrt{n}} = \frac{88 - 85}{4/\sqrt{150}} = \frac{3}{\sqrt{150}} = 9.18$

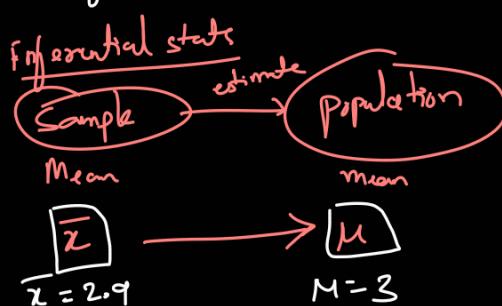
→ 9.18 falls outside of interval (-1.96, 1.96)

→ we reject the null hypothesis

* So, the clg don't have diff placement rate

Point Estimate:

The value of any statistic that estimates the value of a parameter.

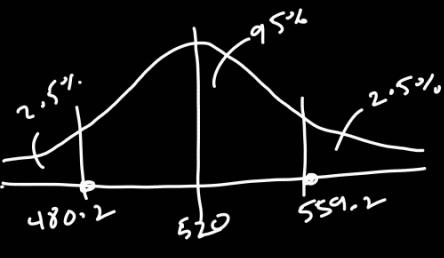


→ Confidence Intervals

$$\text{Point Estimate} \pm \text{Margin of error}$$

- ⑤ On the quant test of CAT Exam, the std deviation is known to be 100. A sample of 25 test has a mean of 520 score. Construct a 95% CI about the mean?

$$\rightarrow \sigma = 100, \bar{x} = 520, n = 25, \alpha = 0.05$$



$\rightarrow CI =$
Point Estimate + margin
of error.
 \rightarrow population std given
 $\rightarrow \bar{x} \pm Z_{\frac{\alpha}{2}} \frac{\sigma}{\sqrt{n}}$
std error.

\rightarrow this formula usually use when $n \geq 30$

$$\text{upper bound} = \bar{x} + Z_{0.05/2} \frac{100}{\sqrt{25}}$$

$$\text{lower bound} = \bar{x} - Z_{0.025} \frac{100}{\sqrt{25}}$$

$$Z_{0.025} = 1.96$$

$$UB = 520 + 1.96(20) = 559.2$$

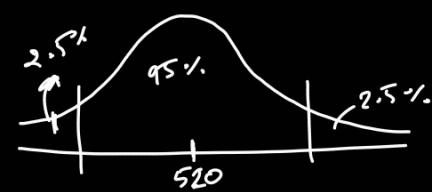
$$LB = 520 - 1.96(20) = 480.2$$

Q) On the quant test of CAT exam, a sample of 25 test takers has a mean of 520 with a std deviation of 80. Construct 95% CI about the mean?

$$n = 25, \bar{x} = 520, s = 80$$

If population std is not given $\rightarrow t$ -test

$$CI = \bar{x} \pm t_{\frac{\alpha}{2}} \frac{s}{\sqrt{n}}$$



\rightarrow degree of freedom = $n - 1 = 25 - 1 = 24$.

$$t_{0.025} = 2.064$$

$$UB = 520 + 2.064 \left(\frac{80}{\sqrt{25}} \right) = 553.024$$

$$LB = 520 - 2.064 \left(\frac{80}{\sqrt{25}} \right) = 486.97$$

$$[486.97 \longleftrightarrow 553.024]$$

Chi Square Test

Q) Chi square Test claims abt population proportions.

\rightarrow It is a parametric test that is performed on categorical, [nominal or ordinal] data.

Q) In the 2000 Indian Census, the age of the individual in a small town were found to be the following:

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

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

<18	18-35	>35
121	288	91

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

\rightarrow Here proportion of population is given chi² test

population 2000		
<18	18-35	>35
90%	30%	50%

n = 500		
<18	18-35	>35
121	288	91
$\frac{121 \times 500}{100}$	$\frac{30 \times 500}{150}$	$\frac{50 \times 500}{250}$

① $H_0 =$ The data meets the distribution 2000 in

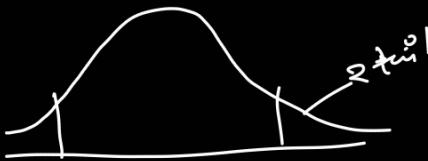
$H_1 =$ The data doesn't meet ...

② $\alpha = 0.05$ [95% (I)]

③ Degree of freedom = $n - 1 = 3 - 1 = 2$

$\rightarrow df = 2$, $\alpha = 0.05$, Chi² table

→ ② Decision Boundary



$\rightarrow \chi^2$ is greater than 5.99, reject H₀

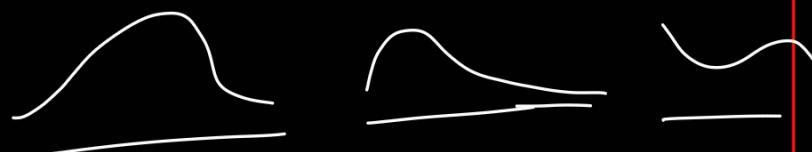
⑤ Calculate Test statistics

$$\begin{aligned} \chi^2 &= \sum \frac{(f_0 - f_c)^2}{f_c} \\ &= \frac{(121 - 100)^2}{100} + \frac{(289 - 150)^2}{150} + \frac{(91 - 250)^2}{250} \\ &= 232.494 \end{aligned}$$

$\chi^2 = 232.494 > 5.99 \quad \{ \text{reject null Hyp}\}$

Central Limit theorem

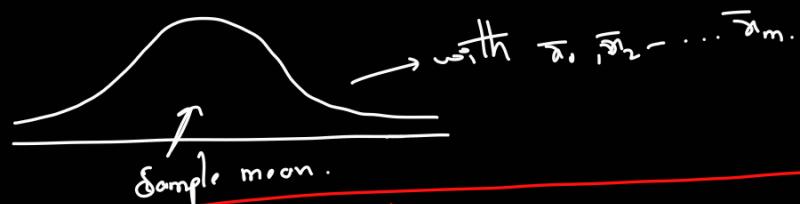
\rightarrow Any kind of distribution



$\rightarrow n \geq 30$ [n ≥ 30] sample size.

$$\begin{array}{l} s_1 \xrightarrow{\text{prob}} \bar{x}_1 \\ s_2 \xrightarrow{\text{prob}} \bar{x}_2 \\ s_3 \xrightarrow{\text{prob}} \bar{x}_3 \\ \vdots \end{array} \left. \begin{array}{c} \text{sample mean} \end{array} \right\}$$

$s_m \dots \bar{x}_m$ more bigger in more CLT holds good



Poisson distribution

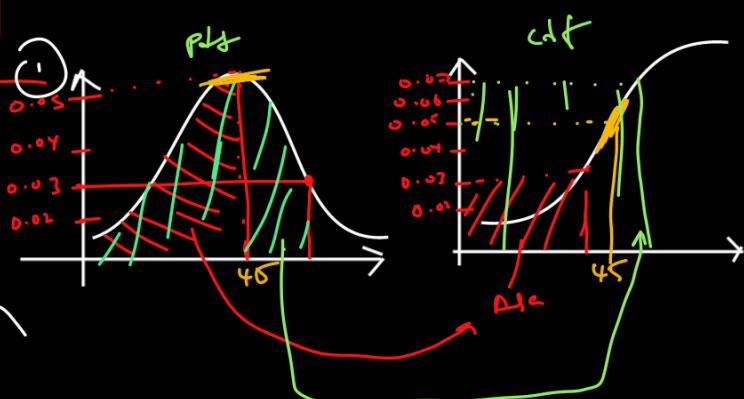
Imp Stats Qn

\rightarrow PDF vs CDF vs PMF

\rightarrow PDF \rightarrow Probability Density Function

\rightarrow CDF \rightarrow Cumulative Distribution Function

\rightarrow PMF \rightarrow Probability Mass Function



PDF \Rightarrow is the derivative of a CDF

② PMF \Rightarrow Discrete Random Variable.

Eg: Tossing a coin
Rolling a dice

