

# EDA (Exploratory Data Analytics)

①

✓ Correlation Matrix, heatmap.

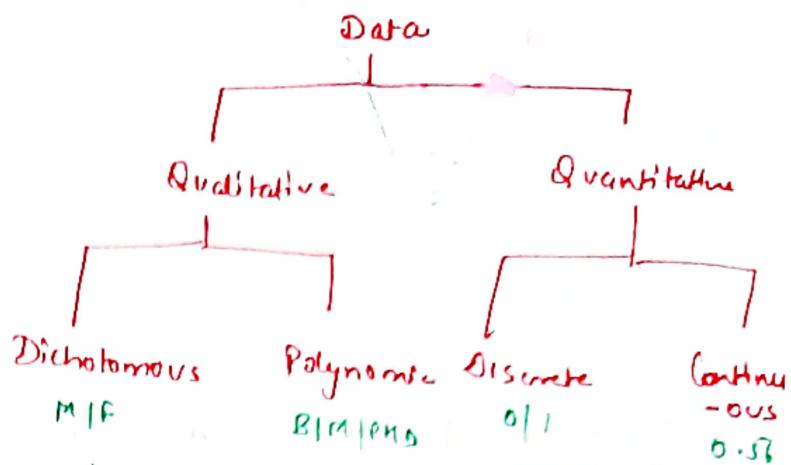
✓ Box Plot

✓ Histogram

✓ Kernel Density Estimate

✓ Central Tendency: Mean, Median, mode

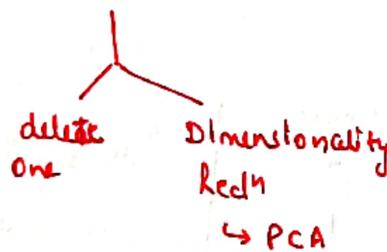
✓ Variance:- Range, std deviation, IQR ( $\theta_3 - \theta_1$ )



## Correlation Matrix

[Correlation matrix is useful to find highly correlated variables, these variables should be removed during the flow of analysis.]

↳ Multicollinearity Problem ] → When one variable is predictive enough of other variable misleading ML models.



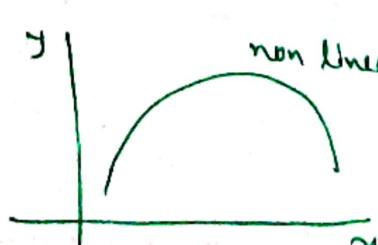
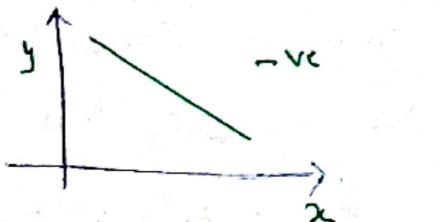
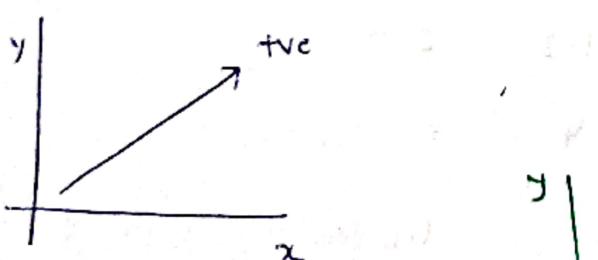
Correlation - Two variables

Scatter plot can tell us whether there

is a linear correlation b/w two variables

+ve linear, -ve linear

non linear, no correlation



(2)

Correlation Coeff - strength + direction of correlation.

(-1, 1)

$$r = \frac{n \sum xy - (\sum x)(\sum y)}{\sqrt{n \sum x^2 - (\sum x)^2} \sqrt{n \sum y^2 - (\sum y)^2}}$$

### Significance of Correlation Coefficient

$\alpha$  - level of significance  $\approx 0.05$

n - number of pairs of data in sample.

if  $|r| > \text{critical value at } \alpha \text{ in}$

in Pearson Correlation Coefficient table

Correlation Significant

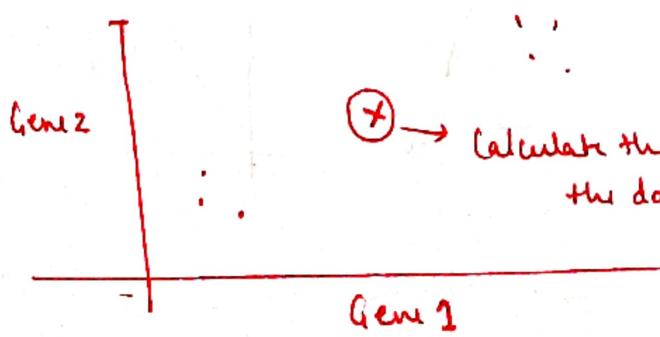
Correlation  $\neq$  Causation.

### Dimensionality Red

#### PCA: Principle Component Analysis

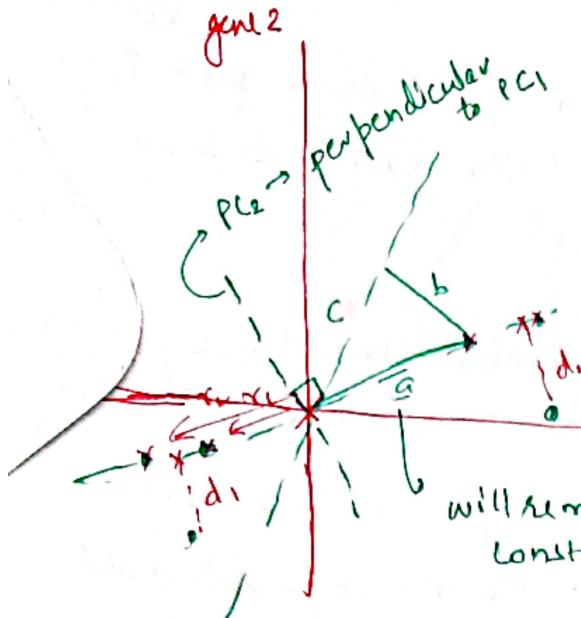
	mouse 1	mouse 2	mouse 3	mouse 4	mouse 5	mouse 6	
Gene 1	10	11	8	3	2	1	
Gene 2	6	4	5	3	2.8	1	
Gene 3	12	5	10	2.5	1.3	2	
Gene 4	5	7	6	2	4	7	

lets say this is the data



④ → calculate the centre of the data

lets find the centre of the data for the first two genes and then shift the graph such that centre comes to origin.



Now try fit a line to fit the data as well as it can.

- \* We find the line which minimize the distances from the data points  $d_1, d_2, \dots$

\* Or maximize the distance of projections from origin  $d_1, d_2, \dots$

$$\text{PCA} \quad a^2 = b^2 + c^2$$

N.B.  
This.

$$\begin{aligned} &\text{if } b \uparrow \text{ and } c \downarrow \\ &\text{if } c \uparrow \text{ and } b \downarrow \end{aligned}$$

- \* first best fit line is called PC<sub>1</sub>, Principle Component 1.

$\sqrt{b^2 + c^2}$  Linear combination of gene 1 and gene 2  
Linear combination of variables

- \* When we do PCA with SVD everything is scaled to 1, especially the hypotenuse, this unit vector along PCA is called a singular vector or eigenvector for PC<sub>1</sub>.

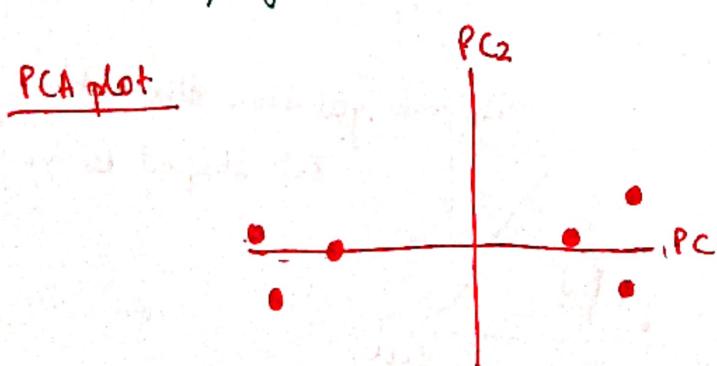
eigen value = sum of distance b/w projected point & the origin

projected distance from origin

$$\text{eigen value} = d_1^2 + d_2^2 + \dots$$

$$\text{variation} = \frac{d_1^2 + d_2^2 + \dots}{n-1}$$

We select a line with largest variations as PC<sub>1</sub>



(4)

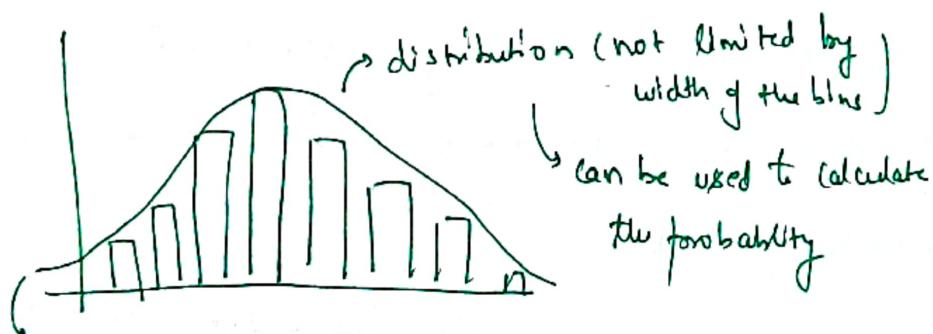
- \* PC<sub>1</sub> accounts for more variation than PC<sub>2</sub>. A score plot is a graphical representation for the variation each PC account for.

for 3D  $PC_1, PC_2, \dots, PC_3 \rightarrow$  It gives a new dimensional plane

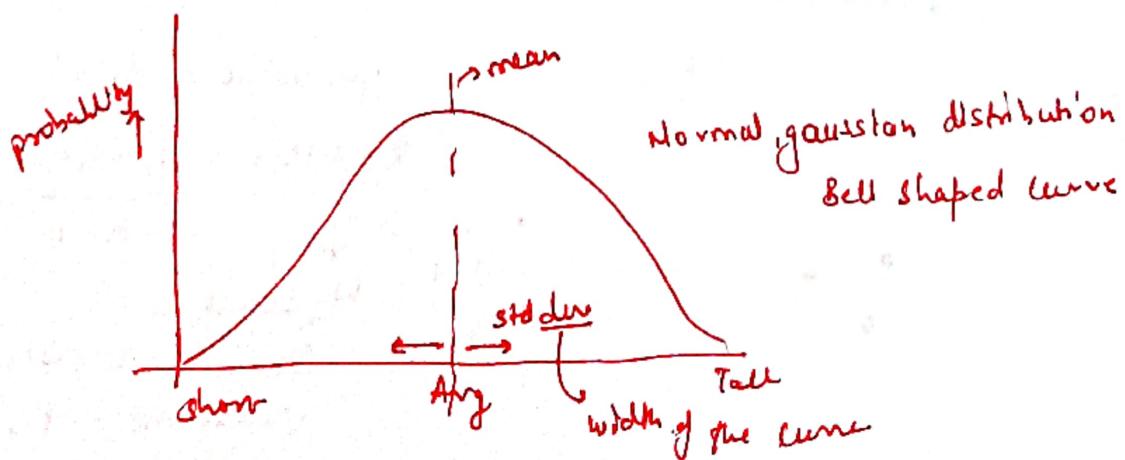
Now, how this reduces dimensionality, as we get the equal number of PCs as the dimensions. If the first two PCs account for substantial variation, we ignore other dimensions and take first two PCs regardless.

Histogram: we need to decide a good size of bin width to get a clear picture

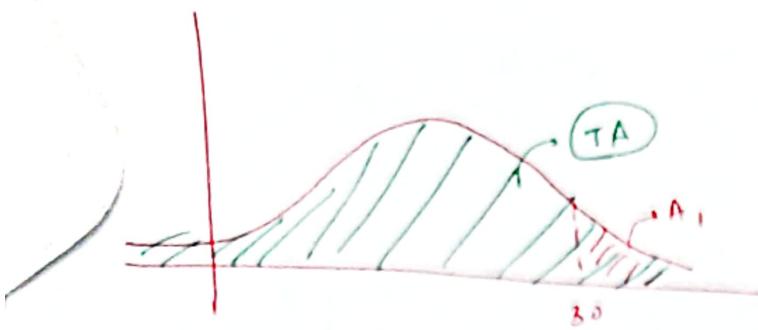
Curve over a Histogram is a distribution.



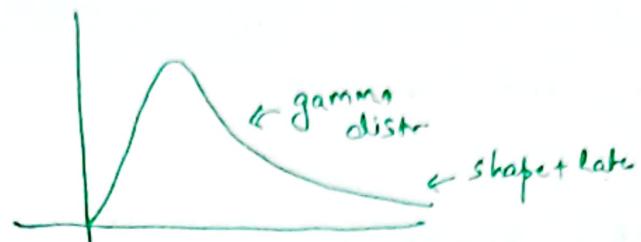
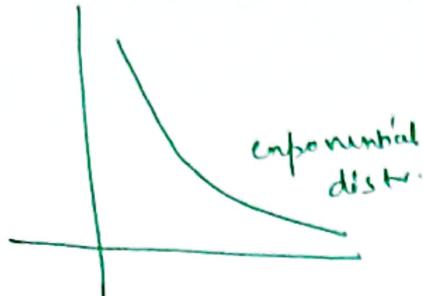
Curve with mean and std deviations-



(5)



Probability for values &gt; 3.0

 $\frac{A_1}{TA}$   
Total Area

We use data to generate these curves and it's very important to determine how much confidence we have in these estimates. Statisticians use p-values & confidence intervals for that.

mean =  $\bar{x} = \frac{\sum x}{n}$  → this mean will be an estimate of population  
 not exactly the population mean because it's  
 not always the case that we have enough  
 data.  
 Estimated mean

 $\mu \rightarrow$  population mean

$$\text{Population Variance} = \frac{\sum (x-\mu)^2}{n}$$

$$\text{std dev} = \sqrt{\frac{\sum (x-\mu)^2}{n}}$$

$$\text{estimated population variance} = \frac{\sum (x-\bar{x})^2}{n-1}$$

$$\text{estimated std dev} = \sqrt{\frac{\sum (x-\bar{x})^2}{n-1}}$$

so that we  
 do not underestimate  
 the variance

6

## Central Tendency Measures - mean, median, mode

### Box plot | Box & Whiskers Plot

$Q_1$  - 25% below, 75% above.

$Q_2$  - median (50% above, 50% below) - Not affected by the outliers.

$Q_3$  - 75% below 25% above

$$IQR = Q_3 - Q_1$$

$$Q_1 - 1.5(IQR) \rightarrow \text{outlier}$$

$$Q_3 + 1.5(IQR) \rightarrow \text{outlier}$$

### Hypothesis Testing

If we generate a hypothesis we can clearly test the hypothesis by repeating the experiments for long, but we do not always have this option of repeating experiments multiple times. Hence we need a quick method to check the hypothesis.

$H_0$  / null hypothesis - negative hypothesis. - We assume a basic null hypothesis. In order to avoid the large number of possible hypotheses we can share with us. We reject the null hypothesis then. or fail to reject it.

$H_0 \rightarrow$  Drug A is not different from Drug B

↳ p value is a number b/w 0-1 which talks about the confidence in

$H_0$ . Closer the pvalue to 0 more confidence that drug A is ~~not~~ different from drug B.

P-values

- One sided - rarely used.
- two sided

Example 1:- Null hypothesis: ~~whether~~ My coin <sup>not</sup> is different from other coins. It is not special. It does not always gives a head.

Let's validate this hypothesis using P-value.

We toss the coin twice

<u>Outcome</u>	<u>Probability</u>
2 heads	0.25
1 h 1 t	0.25
1 t 1 h	0.25
2 tails	0.25

$P\text{-value} = \frac{\text{Probability of getting something ran}}{\text{Probability of getting something more extreme}}$

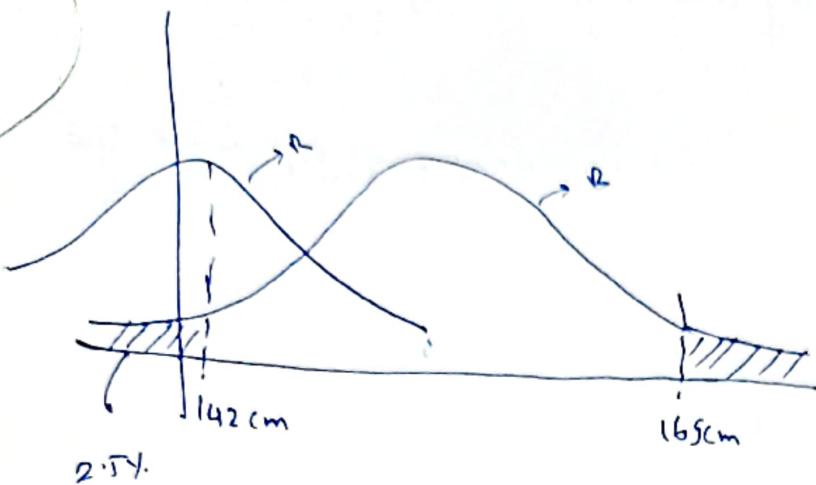
(as only two extreme cases are 2 head/2 tail)

$$P\text{-value} = 0.5$$

Hypothesis - If I get 2 heads in a row, my <sup>coin</sup> ~~head~~ is no different from the normal coin.

We only reject this hypothesis if  $P\text{-value} < 0.05$

Now it's not easy to calculate probabilities for other cases like p-values  
hence we use some distribution of probability to calculate those



Now let's say the height is 162 cm. Is this height comes from blue distribution or from red distribution.

~~Let's test~~

$H_0 \rightarrow$  It does not come from ~~blue~~ distribution

$$p\text{-value} = 0.0025 + 0.025 = 0.0275$$

↳

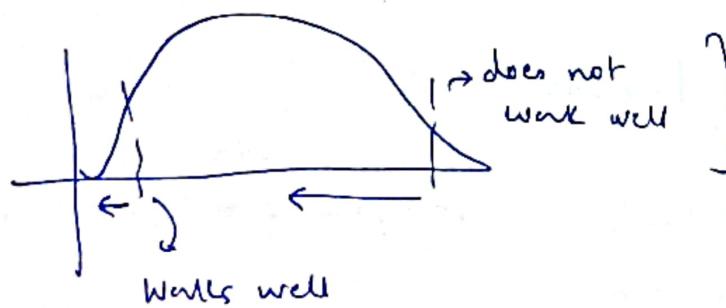
P value of the point that this height comes from blue distribution is 0.025 → Now this is at border line of significance hence we say that this may or may not come from blue distribution.

Now this is ~~not~~ inconclusive.

#### \* One sided p value

We calculate the area only in the direction we want to change

see the



Caution don't use this.

✓ Confidence Interval

✓ ANOVA + F test

✓ Chi-Square test

✓ Z test

✓ T test

✓ Cramers V

✓ Point Biserial Correlation

- Probability

- Bayes

- Naive Bayes

- Gaussian Naive Bayes.

(g)

\* Categorical and Categorical Variables Relations.  
goodness of fit

→ Chi square test

Day	M	T	W	Th	F	S
Expected	10	10	15	20	30	15
of						
Observed	30	14	34	45	57	20
Expected	20	20	30	40	60	30

$H_0$ : Owner's distribution is correct

$H_1$ : Owner's distribution is not correct

Significance level = 0.05 (5%)

Expected amount of customer is told by the shop owner and observed is the actual number of customers visiting the restaurant. We need to see how good is our hypothesis with regards to the observed customers.

$$\text{chi square statistic } (\chi^2) = \sum \frac{(O - E)^2}{E} = 11.44$$

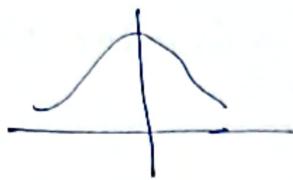
Degree of freedom = 5 (there are six values in datapoint)  
 $6-1$

Critical chi-square value from the chart =  $11.07 = (f_c)^2$

$11.44 > 11.07$  hence null hypothesis is rejected  
i.e. Expected  $\chi^2$  are less likely

What is chi-square distribution?

Normal distribution  $(0, 1)$



$\sigma^2$  (var)  
mean

$$\theta_3 = \sigma_1, 2\sigma, 3\sigma$$

$Q_1 = \chi^2$  → Chi-square distn

DOF = 1

$$Q_2 = \frac{\chi_1^2}{\sigma} + \frac{\chi_2^2}{2\sigma}$$

Chi-square  
distribution

DOF = 2

and like this the distribution  
are there.

upto  $\sigma$

Sample random values  
from this normal distn

-square and square  
them

take two samples of  
values from this normal  
distribution and square  
them → upto  
two

### Categorical & Categorical Variable relations

	Herb1	Herb2	Placebo	Total
Sick	20	30	30	80
	25.3	29.4	25.3	214.
not Sick	100	110	90	300
	94.7	110.6	94.7	294.
Total	120	140	120	380

These are observed value for herb1, herb2 and Placebo impact on  
person being sick & not sick.

Assum

$H_0$  = Herbs do nothing → If herbs do nothing then % of

$H_1$  = Herbs have an impact

→ Hence  $80/380 = 21\%$ .

If herbs do nothing then % of  
sick / total and not sick / total  
will remain the same irrespective  
of herb.

$$\chi^2 = \frac{\sum (O - E)^2}{E} = 2.53$$

$$\alpha = 0.1\gamma.$$

$$DOF = \underbrace{(C-1)}_{\text{number of columns}} \times \underbrace{(R-1)}_{\text{number of rows}} = (3-1) \times (2-1) = 2$$

~~H<sub>c</sub>~~  $\Rightarrow \boxed{\chi^2 < \chi_{c,0}^2}$  Hence we do not reject the null hypothesis.

### ANOVA (Analysis of Variance)

$m$		
1	2	3
3	5	5
2	3	6
1	4	7
$\bar{x}$	2	4
	4	6

$$\begin{aligned} \text{Sum of Squares} &= \frac{(3-4)^2 + (2-4)^2 + (1-4)^2 + \dots}{n} \\ &= 30 \\ \text{grand mean} \rightarrow \text{mean of all of stuff} \\ \bar{x} &= \frac{3+2+1+5+3+4+5+6+7}{9} = \frac{36}{9} = 4 \end{aligned}$$

$$\text{Variance} = \frac{SST}{(m-1)(n-1)} = \frac{SST}{dof}$$

$$dof = m \cdot n - 1$$

$$= \frac{30}{8} = 3.75$$

$$\begin{aligned} SSW \quad (\text{sum of squares within}) &= (3-2)^2 + (2-2)^2 + (1-2)^2 \\ &\quad + (5-4)^2 + (4-4)^2 + (3-4)^2 + \\ &\quad (5-1)^2 + (6-6)^2 + (7-6)^2 \\ &= 6 \end{aligned}$$

$\therefore$ , 6 out of total 30 variation comes from within these groups.

$m \cdot (n-1)$  d.f  $\rightarrow$  6 d.f.

column mean  
total mean

(12)

$$\begin{aligned} \underline{SS_B} \text{ (variation b/w samples)} &= \frac{(2-4)^2 + (2-4)^2 + (2-4)^2}{m-1} \text{ or } 3(2-4)^2 \\ &\quad + 3(4-6)^2 + 3(4-6)^2 \\ &= 12 + 0 + 12 = 24. \end{aligned}$$

$24/30$  variation b/w samples

$$dof = m-1 = 2$$

∴ Inference from ANOVA is regarding the actual difference in the population means or distributions or just a difference in the means.

$H_0 \rightarrow$  means are the same or food doesn't make a difference

$H_1 \rightarrow$  food does make a difference.

Assume  $H_0$ :

$$F\text{-statistic.} = \left[ \frac{\left( \frac{SS_B}{m-1} \right)}{\left( \frac{SS_W}{m(n-1)} \right)} \right]$$

↓

if this is small

chi square

distr

2 dof

$$\left[ \frac{SS_W}{m(n-1)} \right]$$

chi square

distr

6 dof

If this number is big  
means there will be a  
difference b/w the population  
mean of the sample.

harder to reject null hypothesis

difficult to reject

null hypothesis.

$$= 12$$

Calculate critical f value

$$= 3.46$$

Significance level = ~~0.05~~ 0.1

f value > critical f value

reject null hypothesis  
probably a difference in  
population means.

T-test - Applied to one numerical variable.

↓

check mean of the sample reflects good with population

$$\bar{x} \pm t^* \cdot \frac{s}{\sqrt{n}} \quad \begin{matrix} s \rightarrow \text{std dev of sample} \\ \sqrt{n} \end{matrix}$$

t statistics help us get 95% confidence intervals where the mean of various samples of populations will fall with this range. This number improves highly with t statistics and z statistics fails to improve the results on this number.

Conditions for t-statistics  $\left\{ \begin{array}{l} \text{Random} \\ \text{normal central limit theorem } n > 30 \\ \text{Independent} \\ \text{observation} \end{array} \right.$  - Normal distribution

T statistics tells you about the plausibility of ~~one~~ a value in a distribution considering sample mean.

e.g. if t stats give a range of -95, 105 and we expect the average thickness of 110 then our sample and population might not meet the standards.

Cramers V  $\rightarrow$  two categorical variables.

$$\text{Cramers V} = \sqrt{\chi^2 / [n(v-1)]}$$

$$v = \min(\text{rows}, \text{columns})$$

- 0 : not associated (values)
- 0.1 : perfectly associated
- 0.45 : weakly associated
- 0.75 : strongly associated

## Point Biserial Correlation (Continuous + Categorical)

(14)

continuous variable      binary variable  
 $y$                            $x$

### Z-score + Z-test

#  $\#$  std dev away from pop's mean for a particular datapoint.

Length of wilyed turtles

2	-0.59
2	-0.59
3	0
2	-0.59
5	1.3
1	-1.3
6	1.77

$$\underline{Z\text{-score}} \quad \frac{x - \mu}{\sigma}$$

Why? → how usual or unusual a certain data point is  
 Outliers.

$$\mu = 3 \quad \sigma = 1.59$$

Is there a diff b/w dist A & dist B

Support new law	District A		District B		Total
	Yes	No	Yes	No	
Yes	58	42	52	48	110
No					90
Total	100	100	100	100	200

$$\alpha = 0.05$$

- Random
- Normal
- Independence

$H_0 \rightarrow$  There is no difference b/w Dist A & Dist B

$H_1 \rightarrow$  There is a difference b/w two districts

(15)

$$\hat{\sigma}_{\hat{P}_A - \hat{P}_B}^2 = \hat{\sigma}_{\hat{P}_A}^2 + \hat{\sigma}_{\hat{P}_B}^2$$

$$\approx \frac{0.55(1-0.55)}{100} + \frac{0.55(1-0.55)}{100} \quad (\text{Considering } H_0 \text{ as true})$$

$$\hat{\sigma}_{\hat{P}}^2 = \frac{p(1-p)}{n}$$

$$\hat{\sigma}_{\hat{P}_A - \hat{P}_B} \approx \sqrt{\frac{0.55(0.45)}{50}} \approx 0.07$$

size of A      size of B

In this case they are equal but may not be always.

$$z_{\text{score}} = \frac{\hat{P}_A - \hat{P}_B}{\hat{\sigma}_{\hat{P}_A - \hat{P}_B}} = \frac{0.05}{0.07} \approx 0.86$$

$P(\text{at } z\text{-score}) \rightarrow$  find using the table value

if  $P\text{value} > \text{significance level}$  fail to reject our null hypothesis. There is not enough evidence to suggest that myopia becoming more common over time

### Confidence Interval

When we put a confidence interval of 95% i.e.  $\mu \pm 3\sigma$  this means that when we sample random data from the population true mean of the pop will be within this range of  $\mu \pm 3\sigma$  95% of the time.

### Significance level

It talks about the comparison of conditional probability with regards to the significance level ( $\alpha$ ). If  $P\text{value} < \alpha$  reject  $H_0$  otherwise do not reject  $H_0$ .

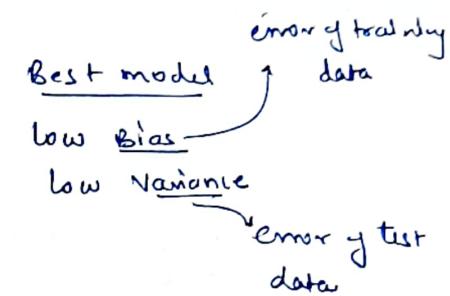
## Bias vs Variance

error of fit  $\uparrow\uparrow \rightarrow$  under fitting

error of fit  $\downarrow\downarrow \rightarrow$  over fitting

Low Bias  
High Variance

High Bias  
High Variance



## Data skewness

