

## P-hacking:

- > Exploit the statistical analysis.
- > If it is false they reject null hypothesis.

## P values:

evaluate the data the null hypothesis

is true

---

> Manipulating data or analyses to artificially get significant P values.

In null hypothesis testing it is either reject or

fail

The binary decision process leads us to 4 possible scenarios.

	$H_0$ is True	$H_0$ is false
Fail to reject $H_0$	✓	X
Reject $H_0$	X	✓

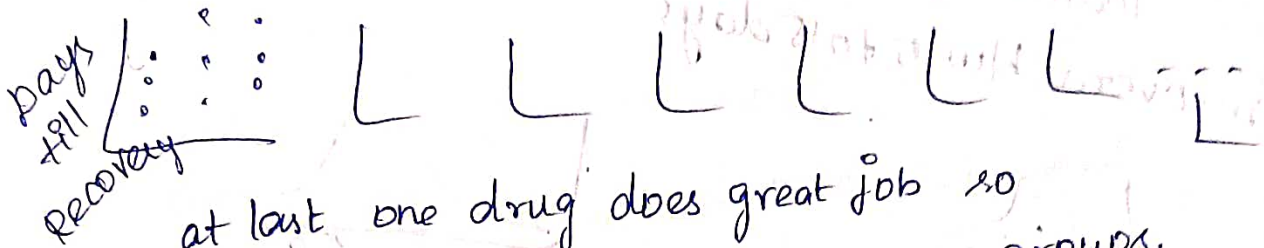
P hacking is analyses are being chosen based on the P values significant, not the best analysis Plan.

Eg.:-

Some people affect with virus we need to develop a drug in this



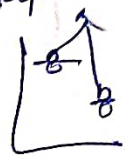
check with people by this drug & measure



at last one drug does great job so

we calculate the means of the two groups.

then we get P value = 0.02



$0.02 < 0.05$  We reject the null hypothesis there is no difference b/w not taking a drug & taking drug z.

Use P. hacking

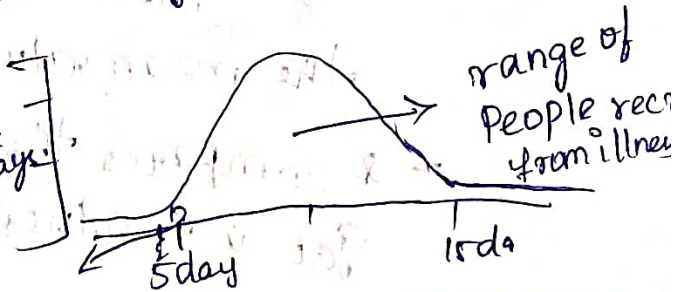
↓  
refer to: misuse & abuse of analysis

techniques & results in being false positives.

we measured: recovery times for a whole

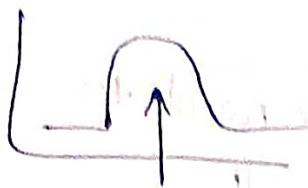
lot of people did not take any drugs to fight virus.

2.5-1. of area under the curve is for duration less than 5 days.



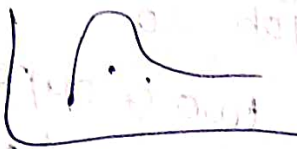


indicating that 95% of people recovered in less than 5 days.



95% of area under the curve is b/w 5 to 15 days

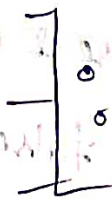
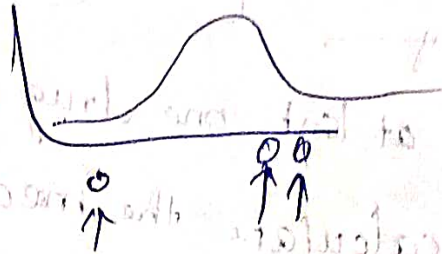
indicating that 95% of people recovered b/w 5 to 15 days



3 people recover

from illness there's good chance but say

something b/w 5 to 15d



Two group of people.

Then we calculate the mean values for a group

& compares these two means & get a P-values = 0.86.

$0.86 > 0.05$  would fail to see difference  
blw the 2 groups of observation

## False Discovery Rate

All  $P$  value



$$P = 0.65$$

$$= 0.05$$

$$= 0.02$$

Prblm to adjust all the  $P$  value by using  
mechanism called False Discovery Rate (FDR)

↓  
Mathematical adjustment of  $P$  values  
increases them by some values & in the end the  
 $P$  values are incorrectly come lower adjusted  
to get higher values  $0.05$ .

## Multiple Testing Problem.

↑ no. of tests.

maximum no. of tests are resulting

in rejection of null.  
More test will mean that more false

Positives

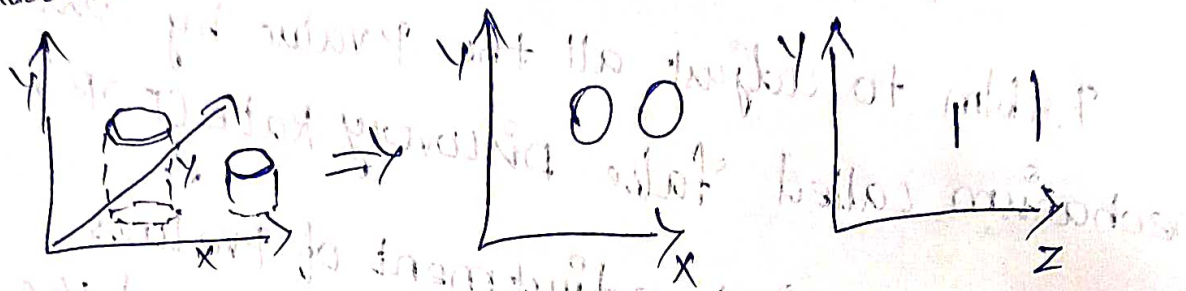
(5% of total tests in our case)  
5 out of 100, 50 out of 1000, 500 out of 10,000

Called multiple Testing problem.



# Dimensionality Reduction:

Process of reducing the no. of random variables or feature called DR  
No. of IP features, variables or columns present in given data set known as dimension.  
Technique  $\rightarrow$  way of converting higher dimension dataset into lower dimension data set.



Components DR..

Feature selection - find subset of original set of variable or features to get a smaller subset.  
Used to model the problem.

- 3 ways
- Filter
  - Wrapper
  - Embedded.

Feature extraction: reduces the data in a high dimensional space to lower dimension space with lesser no. of dimensions.

Methods

PCA

Principal Component Analysis

LDA

Linear Discriminant

GDA

Generalized Di

"

"

---

is may both linear or non linear dep

## PCA:

Finds a linear projection of data into orthogonal basis sys.

that has minimum redundancy & preserves the variance in data.

Application PCA/SVD.

Dimensionality reduction

Latent Semantic Indexing

Kleinberg / Hits algorithm

Google / PageRank algorithm

Image - compression

Data visualization (Project data in 2D).

## Eigen vector.

$A$  is square matrix

nonzero vector  $v$  is eigenvector of  $A$

if there is a scalar  $\lambda$  (eigenvalue)

$$Av = \lambda v$$

$$\text{eg: } \begin{pmatrix} 2 & 3 \\ 2 & 1 \end{pmatrix} \begin{pmatrix} 8 \\ 2 \end{pmatrix} = \begin{pmatrix} 12 \\ 8 \end{pmatrix} = 4 \begin{pmatrix} 3 \\ 2 \end{pmatrix}$$

Covariance of matrix  $X$  of  $x$

$$C_x = \frac{1}{n-1} X^T X$$

symmetric

$$\frac{1}{n-1} (X^T X) = C_x$$



Goal & assumptions of PCA.

$Y = XA$  → covariance matrices capture information

minimal noise & redundancy. → best transformation (orthogonal basis vector)

Derivation PCA

$C_Y$ : covari of  $Y$  expressed in terms of  $A$ .

$$C_Y = \frac{1}{n-1} Y^T Y$$

$$= \frac{1}{n-1} (XA)^T (XA)$$

$$= \frac{1}{n-1} A^T X^T X A$$

$$= \frac{1}{n-1} A^T (X^T X) A$$

Assume  $A = V$  i.e. each column is an eigen vector of  $X^T X$ .  $X = VD$

$$C_Y = \frac{1}{n-1} V^T (X^T X) V$$

$$= \frac{1}{n-1} V^T (V D V^T) V$$

$$= \frac{1}{n-1} V^T V D V^T V$$

$$= \frac{1}{n-1} V^{-1} V D V^{-1} V$$

$$= \frac{1}{n-1} D$$



# Loss function:

Measure estimated value & true value.

used in Linear Loss.

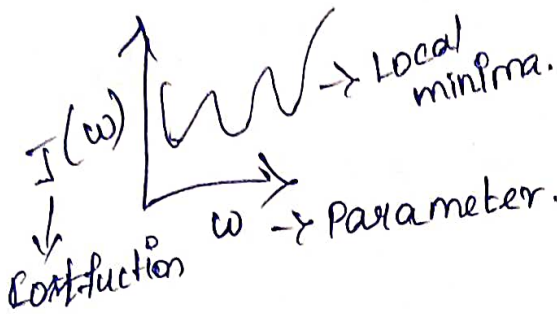
$$Y = 1 \text{ or } 0$$

cannot use this

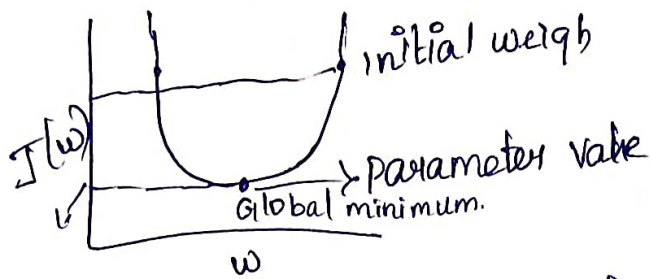
$$Loss = \frac{1}{n} \sum_{i=1}^n (Y_i - \hat{Y}_i)^2$$

↑ function

↓  
 $\hat{Y}_i$  Predicted value given in data set  
 ↓  
 Salary of person based on the experience about 20,000 but set model 30,000



in which Parameter value the loss func is very less. then model is work accurate.



$\hat{Y}$  lies b/w 0-1

Binary cross entropy loss func or log loss

$$L(Y, \hat{Y}) = -(Y \log \hat{Y} + (1-Y) \log (1-\hat{Y}))$$

where  $Y=1$   $L(1, \hat{Y}) = -(\underbrace{1 \log \hat{Y}}_0 + \underbrace{(1-1) \log (1-\hat{Y})}_0)$

$$L(1, \hat{Y}) = -\log \hat{Y}$$

$Y=0$   $L(0, \hat{Y}) = -(\underbrace{0 \log \hat{Y}}_0 + \underbrace{(1-0) \log (1-\hat{Y})}_1)$

$$L(0, \hat{Y}) = -\log (1-\hat{Y})$$

we need minimum loss func (L)

↓  
slight difference between value & true value

$y=1$   
↑

value is large so  $-\log \hat{y}$  is large negative num

$y=0$   
↑

value is small so  $-\log \hat{y}$  is large - value.

loss function (J)

$$J(w, b) = \frac{1}{m} \sum H(y^i, \hat{y}^i) = \frac{1}{m} \sum (y^i \log \hat{y}^i + (1 - y^i) \log (1 - \hat{y}^i))$$

m denotes no. of data point in set