

## CONFUSION MATRIX

		← Actual →	
		Positive	Negative
Predicted	Positive	TRUE POSITIVE (TP)	FALSE POSITIVE (FP)
	Negative	FALSE NEGATIVE (FN)	TRUE NEGATIVE (TN)

(\*)

The size of confusion matrix is determined by no. of thing we want to predict.

If we had to choose from 3 outcomes.

	thing 1	thing 2	thing 3.
thing 1	+	-	-
thing 2	-	+	-
thing 3	-	-	+

Confusion matrix tells us what our NL algorithm did right & what it did wrong.

$$\text{Sensitivity} = \frac{TP}{TP + FN}$$

(Recall)

Proportion of actual positives that are correctly identified.

$$\text{Specificity} = \frac{TN}{TN + FP}$$

Proportion of actual negatives that are correctly identified.

for 3 outcomes -

	OC1	OC2	OC3
OC1	2	4	6
OC2	8	10	12
OC3	14	16	18

OC1

$$\text{Sensitivity} = \frac{2}{8+14} ; \quad \text{Specificity} = \frac{(10+12+16+18)}{(10+12+16+18) + (4+6)}$$

OC2

$$\text{Sensitivity} = \frac{10}{4+16} ; \quad \text{Specificity} = \frac{(2+6+14+18)}{(2+6+14+18) + (8+12)}$$

OC3

$$\text{Sensitivity} = \frac{18}{14+16} ; \quad \text{Specificity} = \frac{(2+4+8+10)}{(2+4+8+10) + (6+12)}$$

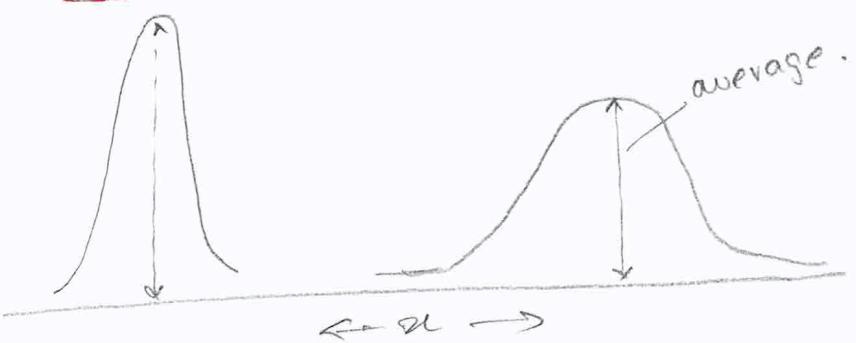
Precision =  $\frac{TP}{TP+FP}$

Recall =  $\frac{TP}{TP+\cancel{FN}}$

proportion of  
positive predictions  
were actually correct.

what proportion of actual  
positives was correctly  
identified.

## Normal Distribution:



mean,  $\bar{x} = \mu = \frac{x_1 + x_2 + \dots + x_n}{n}$

Variance = How data is spread around  $\mu$  =  $\frac{\sum (x_i - \mu)^2}{n}$

Std. Dev. =  $\sqrt{\frac{(x_i - \mu)^2}{n}}$

- 1) ND are always centred around average.
- 2) Width of Curve = Std. Dev
- Normal Curves are drawn such that 95% of measurements fall b/w  $\pm 2$  Std. dev. around the mean.

if data had been for apples, the variance = 100 (apple)<sup>2</sup>.  
 $\therefore$  we cannot plot this x-axis.  $\therefore$  Std. dev.

## Population vs Sample:

A population is the collection of all items of interest to our study & is denoted as 'N'. The mean, variance, std.dev... are population parameters

Sample is subset of population and its denoted = 'n'.

mean, variance...  $\Rightarrow$  are statistics.

Population - Calculated

$$\mu = \frac{\text{sum of all } x_i}{n}$$

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

$$\text{Std. Dev.} = \sqrt{\frac{\sum (x - \mu)^2}{n}}$$

Sample - Estimate (since we don't have all data)

$$\bar{x} = \frac{\text{sum of } x_i}{n}$$

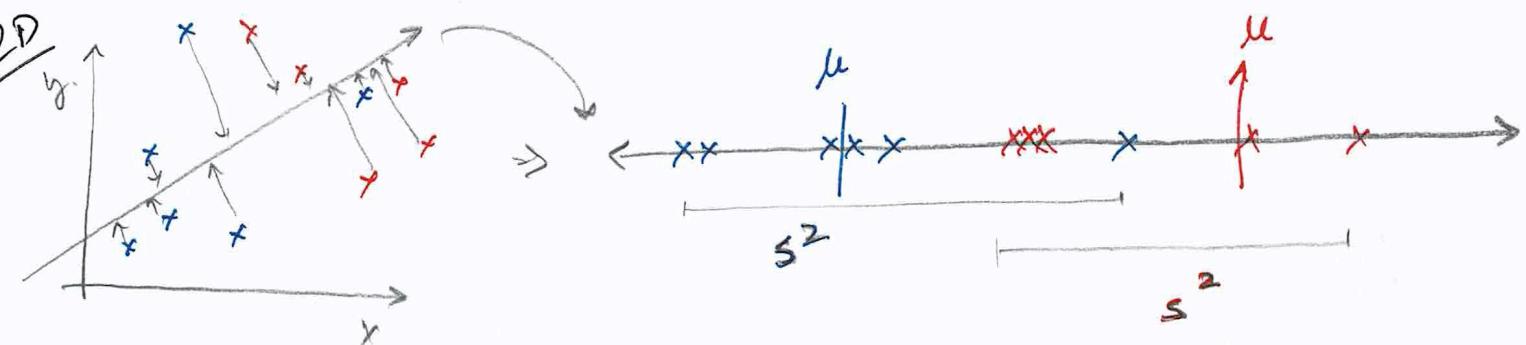
$$\text{Variance} = \frac{\sum (x - \bar{x})^2}{n-1}$$

$$\text{Std. Dev.} = \sqrt{\frac{\sum (x - \bar{x})^2}{n-1}}$$

## LINEAR DISCRIMINANT ANALYSIS.

PCA: Reduces dimensions by focusing on variables with most variation.

LDA: Maximizing the separability b/w 2 categories to make best decisions.



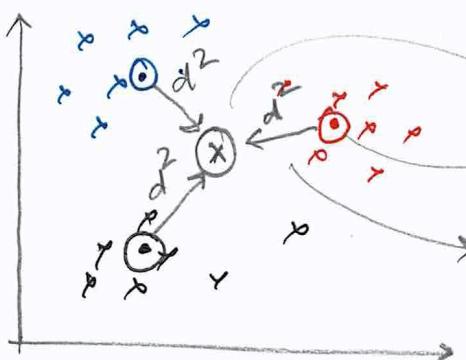
PCA creates a new axis & projects data on it.

LDA creates this axis -

- i) MAX. the dist. b/w  $\mu$  &  $\mu$
- ii) MIN. variation (scatter) b/w  $s^2$  &  $s^2$

$$\Rightarrow \frac{(\mu - \mu)^2}{s^2 + s^2} = \frac{d^2}{s^2 + s^2}$$

3D



Point that is central to all data  
point central to each category.

MAX. the distance b/w each category

while MIN. scatter

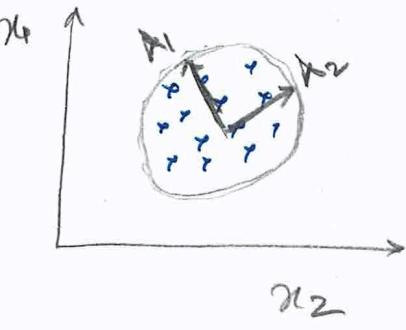
$$= \frac{d^2 + d^2 + d^2}{s^2 + s^2 + s^2}$$

## PCA

PC1 accounts for the most variation in data

④ PCA looks at variables with most variation

⑤ PCA - un supervised. ML technique

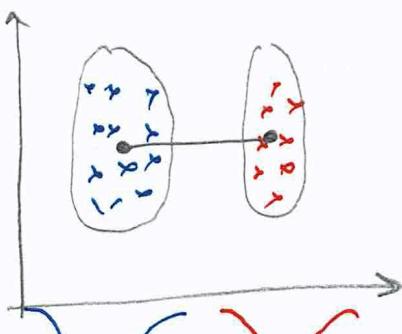


## LDA

⑥ LDI (first new axis) accounts for most variation b/w categories.

⑦ LDA MAX. the separation of categories.

⑧ LDA - supervised ML technique

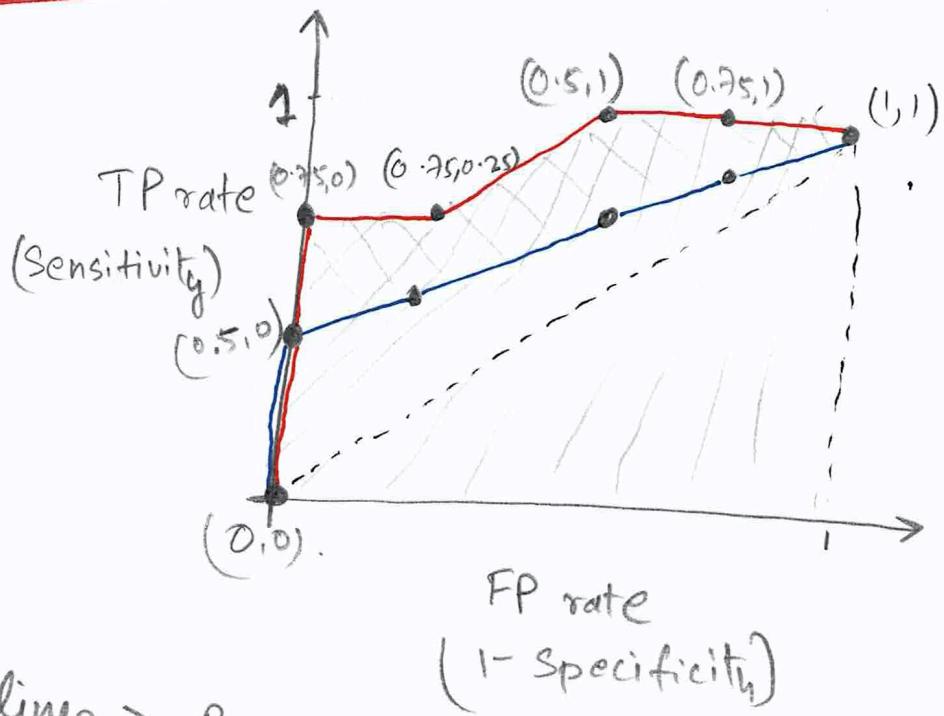


Good projection separates the classes well.

## ROC & AUC

ROC - Receiver Operator Characteristics.

AUC: Area Under Curve.



the red and blue lines  $\Rightarrow$  ROC graph.

Area under ROC graph  $\Rightarrow$  AUC.

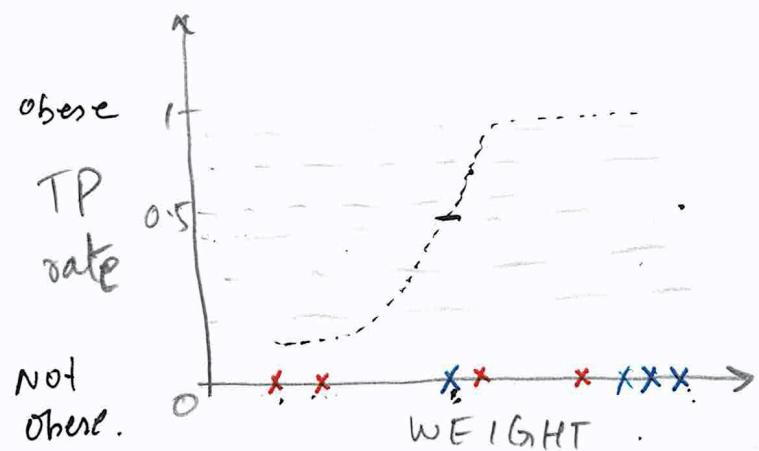
$$TP\text{ rate} = \frac{TP}{TP + FN}$$

Proportion of 'T' samples are correctly classified.

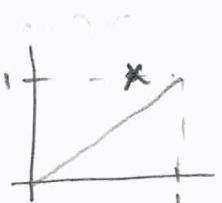
$$FP\text{ rate} = \frac{FP}{FP + TN}$$

Proportion of 'F' samples that were incorrectly classified.

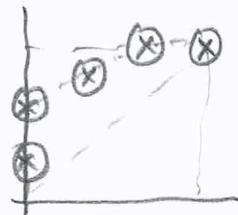
## Roc Graphs



To test the effectiveness of the Logistic Reg. model - we can set a classification threshold = 0.5  
ie,  $> 0.5$  - obese  
 $< 0.5$  - not obese.

Threshold	Confusion Matrix	TP rate	FP rate	Graph.
0	$\begin{array}{c c} 4 & 4 \\ \hline 0 & 0 \end{array}$	1	1	
0.1	$\begin{array}{c c} 4 & 3 \\ \hline 0 & 1 \end{array}$	1	0.75	
0.2	$\begin{array}{c c} 4 & 2 \\ \hline 0 & 2 \end{array}$	1	0.5	
0.3	$\begin{array}{c c} 3 & 1 \\ \hline 1 & 3 \end{array}$	0.75	0.25	
0.7	$\begin{array}{c c} 3 & 0 \\ \hline 1 & 4 \end{array}$	0.75	0	

This how the ROC Graph is obtained



AUC can be used compare two ROC graphs.

From graph on previous page - RED ROC is greater than BLUE ROC

∴ Red ROC is BETTER.

red ROC - logistic regression  
blue ROC - random forest

ROC can be used to compare 2 ML's.

ODDS, LOG(ODDS), ODDS RATIO, LOG(ODDS RATIO).

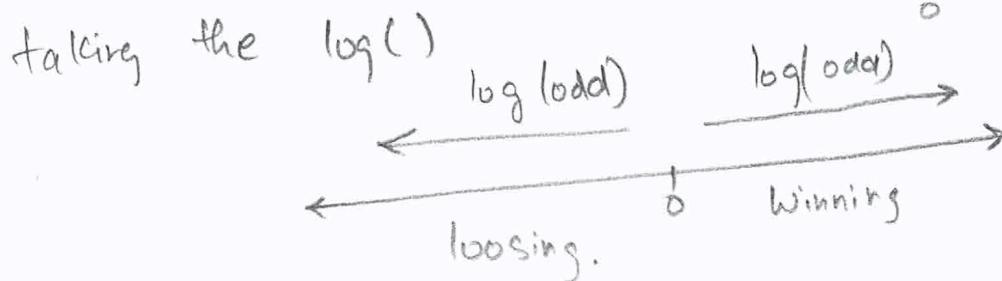
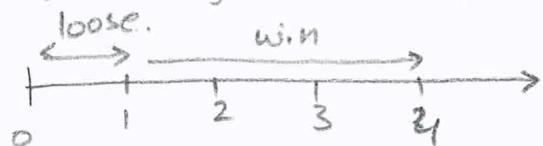
$$\text{odds} = \frac{\text{Some thing happening.}}{\text{Some thing NOT happening.}} = \frac{\times \times \times \times \times}{\times \times \times}$$

$$\text{Probability} = \frac{\text{Some thing happening.}}{\text{Everything that could happen}} = \frac{\times \times \times \times \times}{\times \times \times \times \times \times \times}$$

$$\text{odds} = \frac{\text{probability}}{1 - \text{probability}}$$

$\log(\text{odds})$  - makes values symmetrical & easier to interpret.

Eg: When a team is playing very bad, the odds of losing may fall b/w 0 & 1 always ; the odds of winning may be  $> 1$ .



$$\text{ODDS RATIO} = \frac{\overline{xx}}{\overline{xxx}} = \frac{2/4}{3/1} = 0.17.$$

When denominator > numerator  $\Rightarrow$  value b/w 0 & 1

$\therefore \log(\text{odds ratio}) \Rightarrow$  nice & symmetrical.

Eg:

		has Cancer	
		Y	N
Has mutated gene	Y	23	117
	N	6	210

Use odds ratio to determine relation b/w gene & cancer.

Q: If someone has mutated gene, are the odds of them having cancer higher?

if Gene = 'Y', odds they can have cancer =  $\frac{23}{117}$

Gene = 'N', " " " " " " =  $\frac{6}{210}$

(odds are 6.88 times greater than someone with muted gene also has cancer.)

odds ratio =  $\frac{\frac{23}{117}}{\frac{6}{210}} = 6.88$

$\log(6.88) = 1.93$

Relation b/w gene/can

\* odds ratio - larger value means gene is good predictor of cancer

To find if this relation is statistically significant -  
i.e. odds ratio,  $\log(\text{odds ratio})$

1) Fishers Exact Test.

2) Chi-Square test

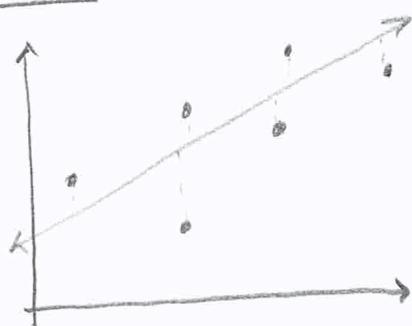
3) the Wald test.

## Bias & Variance Trade-off

Bias: the inability for a ML algorithm to capture the true relationship is called Bias.

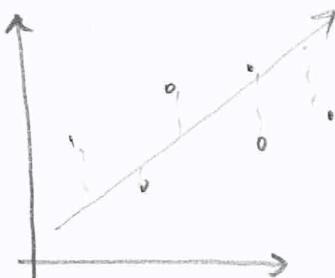
Variance: the difference in fits b/w datasets (test & train d'ts) is Variance.

training set

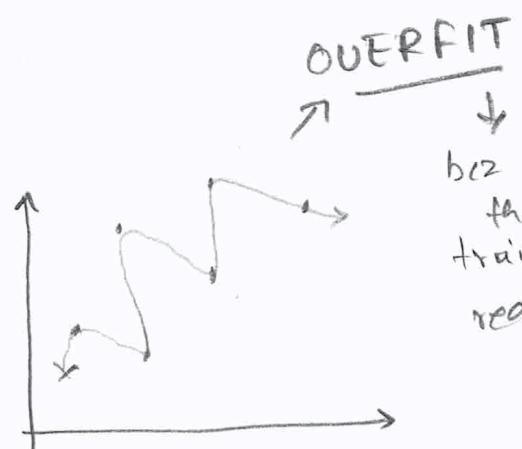


High Bias

testing set

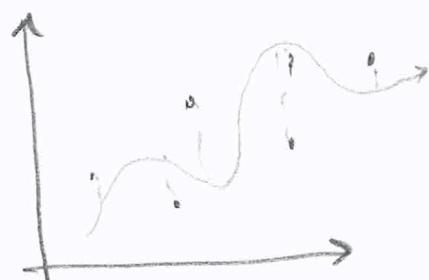


low variance



Low Bias

OVERFIT  
↓  
b'z it fits  
the training d's  
really well



high variance

Best ML model: low bias (accurately model true relationship)  
low variability (consistent prediction across various d'ts)

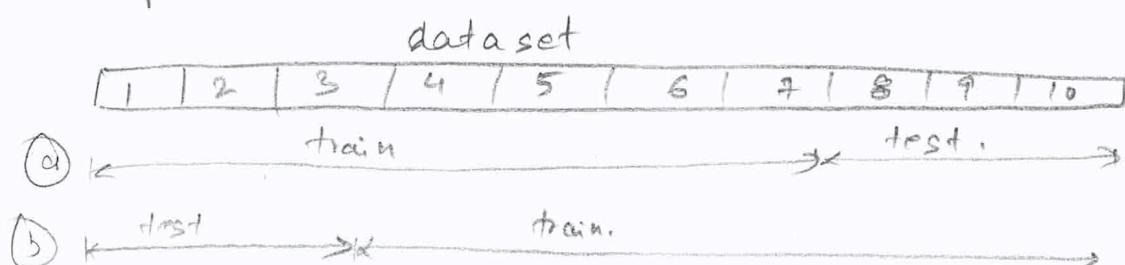
Regularization } methods to find the sweet spot  
Boosting } b/w simple & complicated models.  
Bagging } - to find ML model with low bias/variance

## CROSS VALIDATION

How do we decide which machine learning mode to use?

 Cross Validation allows us to compare different ML models & get a sense of how well they will perform.

In supervised ML, dataset is divided - test & train.



How do we know that first consider 80% of the data for training & last 20% of data for testing is the best option.

∴ Cross validation chooses all of them, one at a time & then summarizes the result.

Since the data set is divided into 10-blocks  
it is called 10-fold cross validation

## K-fold Cross validation

K - no. of blocks the df  
is divided into.

## Leave one out cross validation-

when the ' $K$ ' (no of fold) = no. of observations

- this is only Recomd. when df is small.

## Definition:

K-fold - we split df into  $k$  different subsets. we use

' $k-1$ ' subsets to train df & leave last subset for test.

then Average agnts each fold & then fine tune  
our model.

## HYPOTHESIS TESTING

Hypothesis: A claim that we want to test. The current fact =

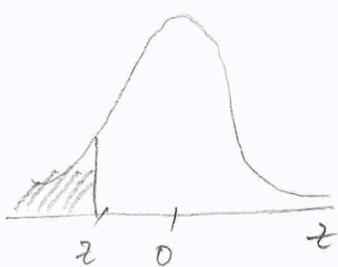
Null Hypothesis,  $(H_0)$ : Currently accepted value of a parameter.

Alternative Hypothesis,  $(H_1/H_a)$ : Research Hyp; Involves the claim to be tested. (what we are proposing to be true)

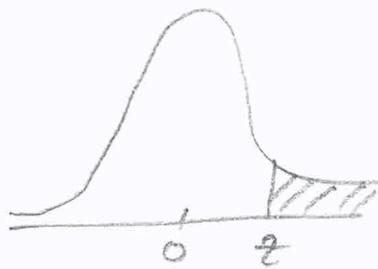
Possible outcomes: Reject  $H_0$ : p-value  $< 0.05$

Fail to Reject  $H_0$ :

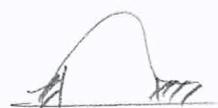
### left-tail test



### right-tail test



$$\text{eg: } z = -1.34, \quad H_0: \mu \geq 0.5, \quad H_a: \mu < 0.5 \quad | \quad H_a: \neq \\ \therefore \text{left tail test.}$$

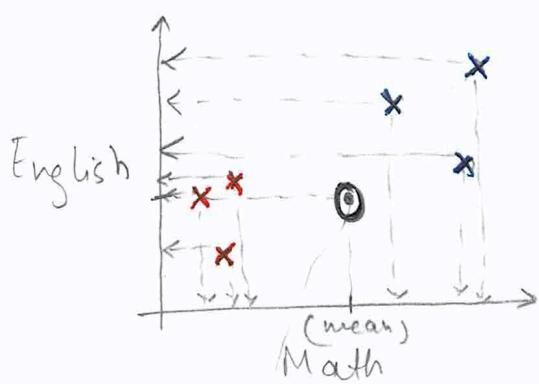


# PRINCIPAL COMPONENT ANALYSIS (PCA)

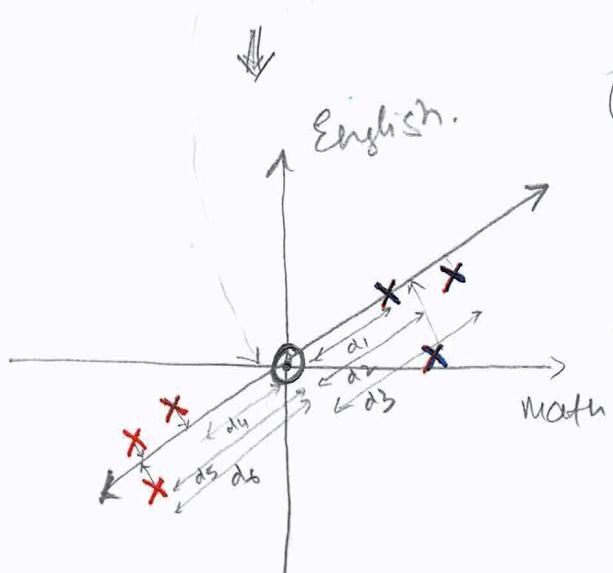
SVD: Singular Value Decomposition

Q:

	Math	English	Science	History	3 variables
Student 1	10	11	8	3	
Student 2	6	4	5	3.	
Student 3	12	9	10	2.5	
Student 4	5	7	6	2	



- ① Find the mean of Math & English & project the points on why axis

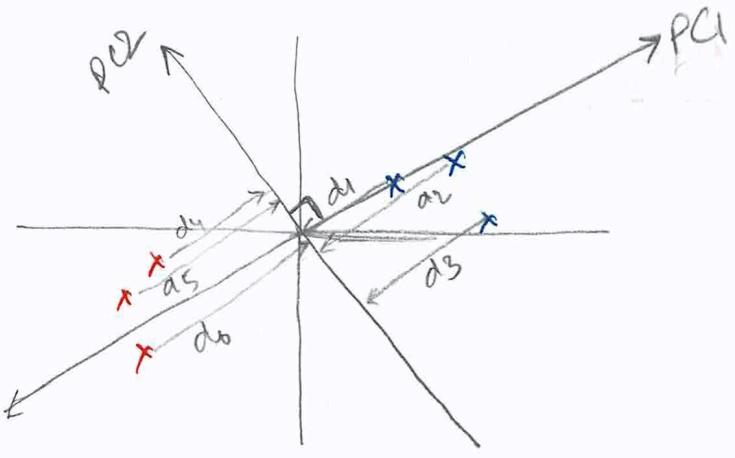


- ② Shift the data so the the mean is in center.  
 ③ we try to fit a line that MAXIMIZES the distance b/w origin (mean) to the points.

$$\text{ie } d_1^2 + d_2^2 + d_3^2 + d_4^2 + d_5^2 + d_6^2 = \text{sum of sq. of distances} \\ = SS(\text{distance})$$

$$= PC_1$$

= Eigenvalue of  $PC_1$



④ PC2 is at  $90^\circ$  of PC1.  
Perpendicular

$$\therefore d_1^2 + d_2^2 + d_3^2 + d_4^2 + d_5^2 + d_6^2 = \text{PL2}$$

= Eigenvalue for PC2

PC1 & PC2 now becomes x/y axis & its rotated.  
 We convert into variation around origin  $(0,0) \Rightarrow$

$$\frac{\text{SS}(\text{dist. PC1})}{n-1} = \text{Variation of PC1} \quad (\text{eg. } = 15)$$

$$\frac{\text{SS}(\text{dist. PC2})}{n-1} = \text{variation of PC2.} \quad (\text{eg. } = 3)$$

$\Rightarrow$  PC1 accounts for  $15/18 = 0.83 = 83\%$  of total variation &  
 PC2 accounts for  $3/18 = 0.17 = 17\%$  of total variation.

SCREE PLOT is a graphical rep. of % of variations in PCs.