

## # Normal Distribution

(1SD) (2SD) (3SD)

- 68% - 95% - 99.7% rule (use z-table)

$Z \rightarrow$  score away from mean in terms of SD

$Z=0 \rightarrow$  on mean  $\rightarrow$  divides population ~~into two~~ in half.

- Dist' of height with mean at 67in & SD 3in. what % are b/w 63 & 67.

$$Z = \frac{x - \mu}{\sigma} = \frac{63 - 67}{3} = \frac{-4}{3} = -1.33$$

↓

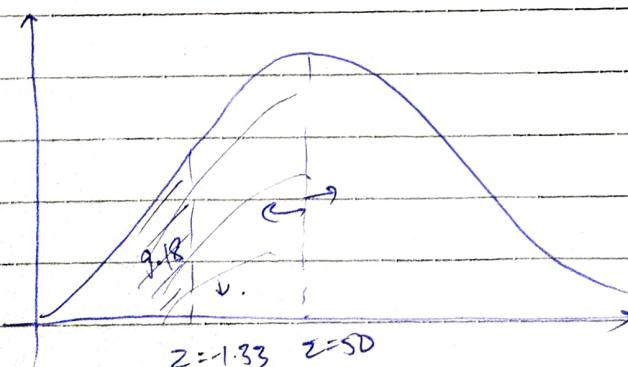
$\rightarrow 0.03$

$\rightarrow 0.0918$

$\times 100 = 9.18\% \text{ to the left of } 63.$   
(less than 63)

$0.0 \rightarrow 50\%$

$$\therefore 50 - 9.18 = 40.82 \quad \text{answer}$$



$9.18 \quad 90.82$

$100$   
-  
 $9.18$   
-  
 $90.82$

- Q: What is the 40<sup>th</sup> percentile of heights?  
 → 40% ppl following you!

find 0.4000 in table ← that z-score is answer.

- Monthly rents in neighbourhood

$$\text{Mean} = \$900$$

$$\text{SD} = \$600 \rightarrow$$

→ No bound on monthly rent

→ Is it a normal dist<sup>n</sup>? → No. (Reason)

→ SD is very large, order of mean

→ ∵ Dispersion of data is huge.

→ Normal dist<sup>n</sup>? atleast 3 SD to one side of mean  
should be in range

→ ∵ For this, t-table / stats is used.

→ In these cases, we can use chebychev inequality  
for bound

→ Good bound is better than a very weak approximation.

- Binomial approximation with large values of  $n$

- approxed by a Normal dist<sup>n</sup>

- Mean =  $np$ , Variance =  $np(1-p)$ , draw normal curve

## # Central Limit Theorem

- Aggregate of large no. of independent random vars,  $\xrightarrow{\text{sum}}$  leads to a no random var. that is approximately normal.

- o If sample size  $> 30$ , whatever be the dist', sample mean distr will be a normal dist'.

↓

Sample mean dist': Take the same-size sample, randomly, many times, and plot the distribution of all the means in each sample-set.

- o As sample size ↑, we get closer to normal mean
- o sample size = 1  $\rightarrow$  mean = Actual mean.
- o Mean of sample mean dist' = mean of actual popul.  
What abt their SDs?

→ Not the same

→ SD of sample < SD of popul.

Reason: In samples, we are evening out the values closely.

$$\text{sample mean dist}' = \frac{\sigma}{\sqrt{n}} \text{ sample size}$$

(Standard Error)

## # Inferential statistics

Population: collection of same units

Parameter: A number we're interested in about popul.

Sample: Subset of popul.

Estimate: guess for parameter, calculated from sample

- Estimate is good only if sample is good

- Sample should avoid Bias

~~Bias~~

- To avoid bias: cover larger area, select elements for sample randomly.
- sample should be representative of all complexities of the population

# Selection Bias: literary digest used PhoneBook.

- phones were a luxury at that time
- Biased towards ppl with phone

# Non-response Bias:

- only 2.4 mn responded out of 10 mn.
- covering almost 75% of popl
- o Sample should not be a sample of convenience
- o Should draw uniformly at random with or w/o replacement from population.

o Simple Random Sample:

- Random sample
- all units equally likely
- o Random Sample
- all units need not be equally likely

Q: A population has 4 ppl A, B, C, D.  
 A is chosen in sample & other 3 are chosen only if heads on coin toss  $\rightarrow$  Biased towards A.

Q: In a pop. avg age = 37, SD = 15.  
 sample size = 200 taken from pop  
 Avg. age for sample = ?

Sol<sup>n</sup>: Expected Age = 37  
 Standard Error =  $\frac{15}{\sqrt{200}} = 1.06$

$$\frac{\text{SD}(\bar{x})}{\sqrt{n}}$$

Q: Sample size = 100, Avg = 38, SD = 14.  
 Avg age of population?

Sol<sup>n</sup>: Estimated age = 38  
 Std. error =  $\frac{SD}{\sqrt{100}} = 1.4$ .

- SD is for sample here
- we want SD of population
- $\therefore$  contradiction

But it is correct coz of assumption:

Bootstrap: A large sample size is like a population,  
 then SD of the sample is an approximation  
 to the SD of population.

$$\begin{array}{c} 7.42 \\ 8.7 \xrightarrow{S} 8.8 \xrightarrow{S} 8.5 \xrightarrow{S} 7.8 \end{array}$$

## Confidence Intervals

DATE: / /  
PAGE NO.:

- 8: Sample of 625 households,  
avg. income = \$ 63000,  
SD = \$ 40,000

Find an approximate 95% confidence interval for the avg. income of households.

Sol<sup>n</sup>: std error =  $\frac{40000}{\sqrt{625}} = 1600$  (previous method)

For 95% accuracy : need to see  $\textcircled{2}$  SD away from mean

$$\therefore 63000 \pm 2 \times \frac{40000}{\sqrt{625}}$$

$$= 63000 \pm 2 \times 1600$$

$$= 59800 \text{ to } 66200 \text{ (confidence interval)}$$

~~- 95% of time, it will fall into this range,  
if a house is picked at random.  
The mean will fall into this range~~

Given dist<sup>r</sup> is not normal, but how did we use normal dist<sup>n</sup> approximation for a normal dist<sup>r</sup>

- Sample mean dist<sup>n</sup> by central mean th. says that it is a normal dist<sup>r</sup>.

# On confidence Intervals

DATE: / /  
PAGE NO.:

- Good confidence Interval
  - on which decisions can be made
- Bad : when claim does not match the actual position of mean.

## Testing hypothesis

- Toss coin 15 times
  - 10 heads
  - 5 tails

- Coin is fair :  $H_0$  → null hypothesis
- Coin is unfair : biased towards head -  $H_A$   
 $H_A$  : alternate hypothesis

This is a 1 side test (Don't see if biased towards tails)

- Assume  $H_0$  is true & calculate change of  $H_A$ ,  
i.e. above obs. of  $\frac{10 \text{ heads}}{15 \text{ toss}}$  = p-value

### Method

If p-value < 5% (small) → choose  $H_A$   
we call it statistically significant  
o/w (p-value > 5%), no change in our assumption  
that  $H_0$  is true.

- $P(10 \text{ heads} / 15 \text{ tosses}) < 5\%$  → means prob was less than 5%, but it still happened, that means coin is unfair.

Q: Red-flowering = 25% in each plant  
Data: 400 plants, 88 are red-flowering  
Is data in sync with theory?

Soln:  $H_0 = p = 0.25$ ,  $H_A = p < 0.25$

P-value: binomial with  $p = 0.25$  &  $k \leq 88$   
 $\Rightarrow P = 9.08\%$

of 400 plants, each  $\begin{cases} \text{red } (0.25) \\ \text{not red } (0.75) \end{cases}$   
 $k \leq 88$

since  $p(9.08) > 5\%$ .  $\rightarrow$  hence theory (25%) is fine.

This is called exact binomial test.

Q: 25% chance of red flowering  
Data: 400 plants, 88 → red

DATE: 7/9/22

PAGE NO.:

Ans.

$$z = \frac{88.5 - 100}{8.66} = -1.328$$

Wednesday

- Get area under curve
- $P = 9.21\% > 5\% \rightarrow$  theory looks good
- This is one-sample z-test

Q: 16000 jurors → 26% of them are black  
Data: 100 men on jury → 8 are black.  
Is data % in sync with theory %?

Sol: Selection of 100 men is like a trial conducted 100 times with selecting black as 1 & white as 0.

$$\mu = 100 * 0.26 = 26$$

$$\sigma = \sqrt{100 * 0.26 * 0.74} = 4.39 \quad (O = \sqrt{np(1-p)})$$

Calculate z.

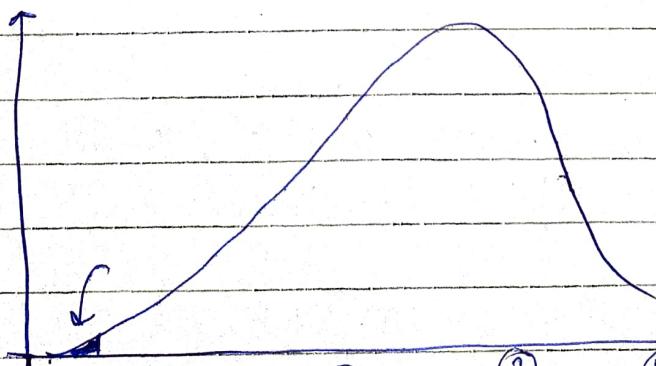
$x = 8$ , how far is it from 26. in terms of SD.

$$z = \frac{8.5 - 26}{4.39} = -3.99 \approx -4$$

$\therefore z$  is almost 4 std error away from the mean.  $\therefore \leq 5\% \rightarrow$  Not in sync

→ Extreme cond?

with given  
(16000 - 26%)  
hypothesis



① 68% ② 95% ③ 99.7% ④ SD → very less

- # Significance value: level of tolerance for error.
- # Here, by default, we consider significance value = 5%.

$\therefore p < 5\%$ .  $\rightarrow$  alternate hypothesis considered.

W.R.T. the given sample, & the significance val., we conclude that theory is unfair.

- Q: coin tossed 20 times -  $H_0 : p = 0.5$  &  $H_A : p \neq 0.5$   
 If no. of head is more than 14, then select  $H_A$ ,  $\therefore / H_0$

Sol: Possible Scenario

		$p = 0.5$	$p = 0.8$
<del>Reality</del>	<del>Observation</del>	correct	Type-I error $\text{prob} = 5.0\%$
	$p = 0.5$	Type-II error	correct $\text{prob} = 91.3\%$

Type 1: Significant level

Type 2: Power.

- # Calculate Type I error

Binomial,  $n = 20$ ,  $p = 0.5$ ,  $k \geq 14$ .

Prob = 5.8%.  $\rightarrow$  significance level

- # Type II

Binomial,  $n = 20$ ,  $p = 0.8$ ,  $k \geq 14$   
 Prob = 91.3%.  $\rightarrow$  Power

$\therefore$  Type II error =  $100 - 91.3\% = 8.7\%$ .

- Significance & Power are inter-related
- Want to minimize error
- If error  $\downarrow$ , power  $\uparrow$  & vice-versa.

### Neyman-Pearson Lemma

Design the test such that we maximize power, given the fixed significance level.

- See slides -

Eg: Obs: 400 coin toss, 225 heads seen

$H_0$ : coin is fair,  $p = 0.5$

Under  $H_0$ , expect 200 heads with std. error of 10 (Jnplif)

- 1<sup>st</sup> alternative : coin is biased towards  $\rightarrow p > 0.5$
- 2<sup>nd</sup> : coin is not fair  $\rightarrow p \neq 0.5$
- 3<sup>rd</sup> : coin is biased to tail  $\rightarrow p < 0.5$

Find P-value of all above cases.

①

— chance of getting 225 heads or more, assuming coin is fair.

$$z = \frac{225 - 200}{10} = 2.5 \cancel{\text{not}} \leftarrow \text{significant}$$

$\therefore$  can say biased towards head.

This is 1-tailed test

- (2) no. of heads as  $\geq 25$  more or less than the expected value of 200, assuming fair

$$z = \pm 2.5, p = 2 * 0.71\% = 1.42\%$$

$p < 5\% \rightarrow$  choose 2<sup>nd</sup> alt.

called - 2-tailed test.

- (3) chance of getting heads = 225 or fewer, assuming coin is fair

$$z = \frac{225 - 200}{10} = 2.5$$

$p = 99.4\% > 5\%$ , so choose  $H_0$

→ This is another 1-tailed test.

#

### Z-test

Q: Avg. height of Indian men = 175 cm

Data: simple random sample of 100 studs taken from univer. in India with 10kg str.

Avg height of sample = 174 cm  
Std. Dev = 5 cm

Q: Is avg height of students shorter than Indian men?

t-test

For P-value, find t & look at t-table.

df      : sample size = 5

: see row of df = 4.

Now see,  $t = 1.79$

it is b/w 1.533 & 2.132 (in that row)

# why  $\div$  by  $n-1$ , not  $n$ .

# Two-Sample Test

Eg  $\Rightarrow$  Take 2 samples of egg, 1 before  $\uparrow$ ing temp<sup>r</sup>, 1 after temp<sup>r</sup>  $\downarrow$ e, compare their qualities

Method:

-  $H_0: \mu_1 = \mu_2$  &  $H_A: \mu_1 \neq \mu_2$

- assume  $H_0$  & calc. z.

$$z = \frac{\text{mean of sample } 1 - \text{mean of sample } 2}{\sqrt{\frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}}}$$

$\text{do: } \mu_1 - \mu_2 \text{ in } H_0$

$\boxed{\text{Here do} = 0.}$

& most cases

sample size of  
1st sample

• Look at z-table, find P-value. If  $P < 5\%$   $\rightarrow$  discard the null hypothesis

• Unpaired Test: different sample items, diff. sample size  
 Paired Test: Same sample items in both situations

If sample size too small, use t-test

- If variance of both samples are same

$$t = \frac{(\bar{x}_1 - \bar{x}_2) - d_0}{S_p \sqrt{\frac{1}{n_1} + \frac{1}{n_2}}}$$

$$S_p^2 = \frac{(n_1-1)s_1^2 + (n_2-1)s_2^2}{n_1+n_2-2}$$

$$\rightarrow df = n_1 + n_2 - 2$$

degree of freedom

- If variance diff

$$t = \frac{(\bar{x}_1 - \bar{x}_2) - d_0}{\sqrt{\frac{s_1^2}{n_1} + \frac{s_2^2}{n_2}}}$$

$$df = \frac{\left( \frac{s_1^2}{n_1} + \frac{s_2^2}{n_2} \right)^2}{\frac{s_1^2}{n_1} + \frac{s_2^2}{n_2}}$$

## ANOVA (More than 2 samples)

Mean of all  $i$  samples =  $\bar{y}$

mean of  $i^{\text{th}}$  sample =  $\bar{y}_i$

$n \rightarrow$  sample size (~~total~~)

$y_{ij} \rightarrow$  individual value

$\bar{y}_i \rightarrow$  mean of that ~~is~~ particular sample.

$N \rightarrow$  total items, including  
= na.

$$\text{F-score} = \text{MSB} / \text{MSE}$$

$\chi^2$  - chi-square distribution

By Stoke Dutta

DATE 12/9/22

PAGE NO.:

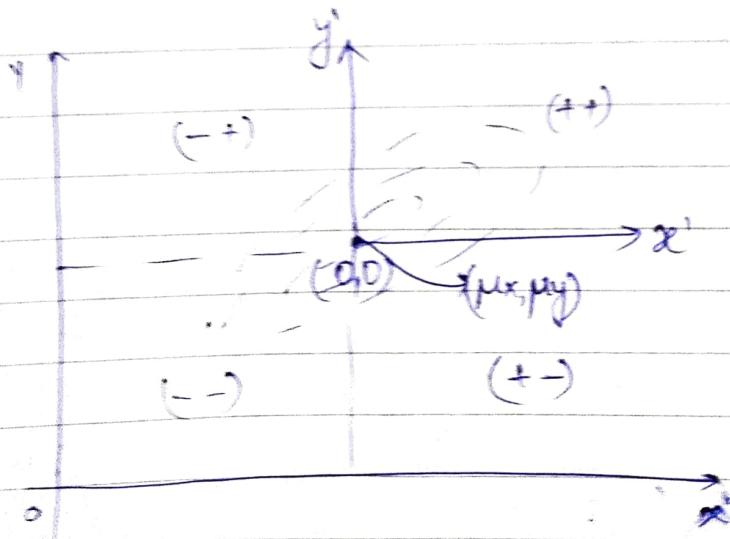
New

$\mu_x = 130$ (f <sub>1</sub> )	$\mu_y = 60$ (f <sub>2</sub> )	$x - \mu_x$	$y - \mu_y$
height = x	weight = y		
130	55	0	-5
135	60	5	0
125	65	-5	5
140	70	10	10
120	50	-10	-10

$$\mu_x = \frac{1}{n} \sum_{i=1}^n x_i = \sum_{i=1}^n p_i x_i$$

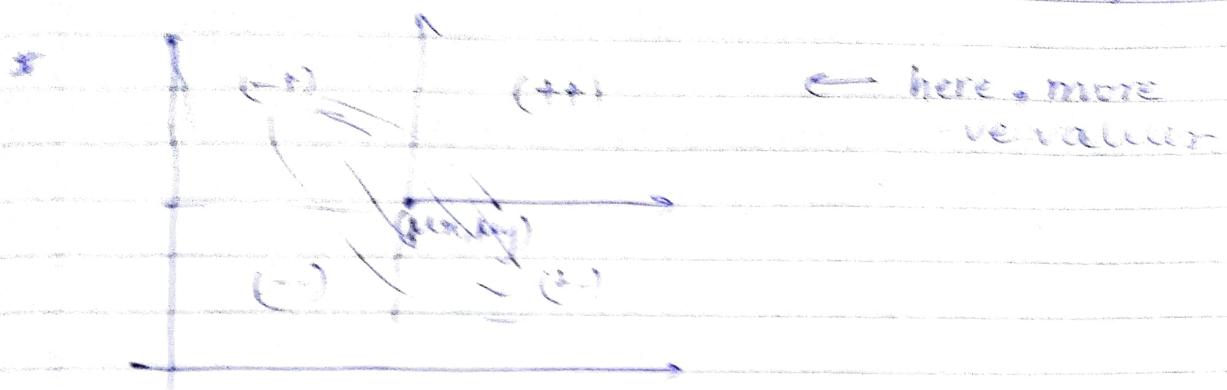
$$\begin{aligned} \text{var}(x) &= \frac{1}{n} \sum_{i=1}^n (x_i - \mu_x)^2 \\ &= \frac{1}{n} \sum_{i=1}^n (x_i - \mu_x)(x_i - \mu_x) \end{aligned}$$

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

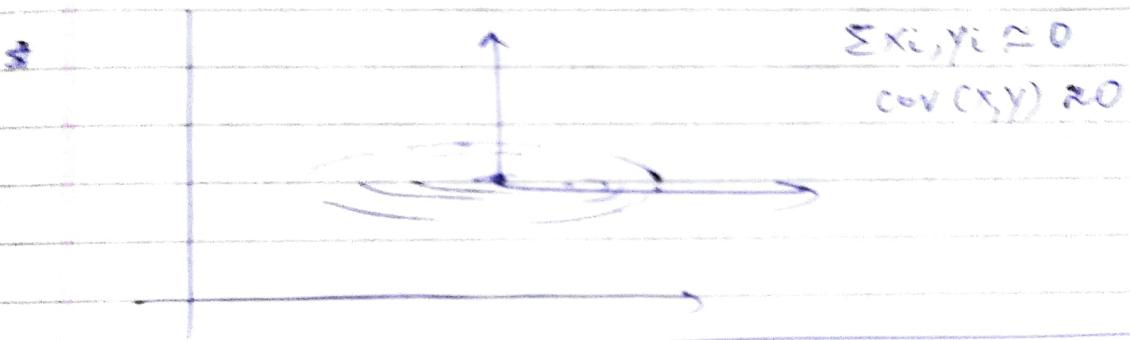


+ve values are larger (more values in 1<sup>st</sup> & 3<sup>rd</sup> quadrant)

$\Sigma x_i y_i > 0 \rightarrow \text{cov}(x,y) > 0$  ( $x, y \rightarrow$  directly prop)



$\Sigma x_i y_i < 0 \rightarrow \text{cov}(x,y) < 0$  ( $x, y \rightarrow$  inverse prop)



correlation  $r_{x,y} = \frac{\text{cov}(x,y)}{\sqrt{\text{var}(x) \text{var}(y)}}$

$r_{x,y} \in [-1, 1]$

let there are D attributes

$f_1(\text{height}), f_2(\text{weight}), \dots, f_D(\text{age})$

- Need to calculate rel' b/w all pairs of attributes

$$\Sigma = \begin{bmatrix} f_1 & f_2 \\ f_2 & f_3 \end{bmatrix}_{D \times D}$$

no. of attributes

each cell = covariance of ( $f_i$  &  $f_j$ )

$$\rightarrow \underbrace{\text{cov}(f_i, f_i)}_{\text{diagonal elements}} = \text{var}(f_i)$$

$$\rightarrow \text{cov}(f_i, f_j) = \text{cov}(f_j, f_i)$$

$\Sigma \rightarrow$  dispersion matrix  $\rightarrow$  symmetric wrt  
diagonal elements

$\rightarrow$  Also called, covariance matrix.

#  $\Sigma$  is real symmetric Positive Definite Matrix

Real: all elements are real values (never imaginary)

Symmetric:  $\text{cov}(f_i, f_j) = \text{cov}(f_j, f_i)$

upper & lower are sym wrt diagonal

Positive Definite:

$A$  is +ve definite if  $\forall$

$$\vec{x}' A \vec{x} > 0, (\forall x)$$

Eg:  $A = \begin{bmatrix} 2 & -1 \\ -1 & 3 \end{bmatrix} \rightarrow$  Is it +ve definite

Let  $\underline{x} = \begin{bmatrix} a \\ b \end{bmatrix}$   $\{a, b \in \mathbb{R}\}$

$$\underline{x}' = [a \ b]$$

$$\therefore \underline{x}' A \underline{x} = [a \ b] \begin{bmatrix} 2 & -1 \\ -1 & 3 \end{bmatrix} \begin{bmatrix} a \\ b \end{bmatrix}$$

$$= \begin{bmatrix} 2a-b \\ 3b-a \end{bmatrix} \begin{bmatrix} a \\ b \end{bmatrix}$$

$$= [(2a-b) \quad (3b-a)] \begin{bmatrix} a \\ b \end{bmatrix}$$

$$= \boxed{2a^2 - 2ab + 3b^2}$$

$$= 2\left(a^2 - \frac{2ab}{2} + \frac{b^2}{4}\right) + \frac{5}{2}b^2$$

$$= 2\left(a - \frac{b}{2}\right)^2 + \left(\frac{5}{2}\right)b^2 \geq 0$$

$\rightarrow |\Sigma| > 0$  (Determinant of  $\Sigma > 0$ )

$\rightarrow \Sigma^{-1}$  also be +ve definite matrix

$\rightarrow$  All the eigen values of  $\Sigma$  are  $> 0$ .

$$A \underline{x} = \lambda \underline{x}$$

$$[A - \lambda I] = 0$$

gives values of  $\lambda \rightarrow$  eigen values.

$$A_{D \times D} \underline{x}_{D \times 1} = \lambda_{1x1} \underline{x}_{D \times 1}$$

(Scalar)

D values of  $\lambda$ , D eigen values ( $\lambda_1, \lambda_2, \dots, \lambda_D$ )

corresponding eig. vectors ( $v_1, v_2, \dots, v_D$ )

$$d_2^2(\mu_1, \bar{x})$$

If distance b/w  $\bar{x}$  &  $\mu_1$   
is very small

$$f_2$$

(1)

Exemplar

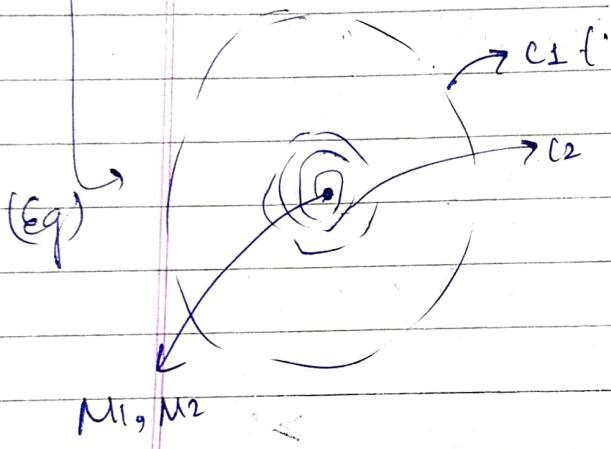
(2) say

$$\begin{pmatrix} 5 \\ 5 \end{pmatrix}$$

Exemplar

$$f_1$$

→ If  $\mu_1$  is not representing the class sufficiently,  
cultivate 2<sup>nd</sup> order stats of class,  
e.g.: variance



$$d_2^2(\mu_1, \bar{x}) = (\bar{x} - \mu_1)' (\bar{x} - \mu_1)$$

$$(C_1, \mu_1, \Sigma) d_M^2(\mu_1, \bar{x}) = (\bar{x} - \mu_1)' \left( \frac{\Sigma}{2} \right)' (\bar{x} - \mu_1)$$

(Mahalanobis distance)

$$d_M^2(\mu_1, \mu_2) = (\mu_1 - \mu_2)' \Sigma^{-1} (\mu_1 - \mu_2)$$

where

$$\Sigma = \frac{\Sigma_1 + \Sigma_2}{2}$$

$f_1$	$f_2$
2.5	2.4
0.5	0.7
2.2	2.9
1.9	2.2
3.1	3
2.3	2.7
2	1.6
1	1.1
1.5	1.6
1.1	0.9

Cos difference

$$\mathbf{x} = \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} \quad \mathbf{y} = \begin{bmatrix} y_1 \\ y_2 \end{bmatrix}$$

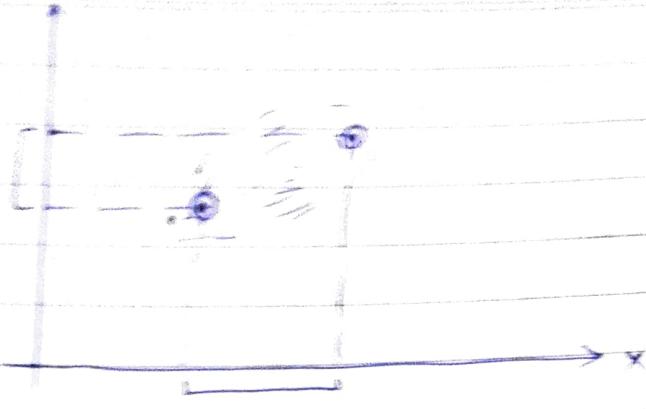
$$\mathbf{x} \cdot \mathbf{y} = |\mathbf{x}| |\mathbf{y}| \cos \theta. \quad (\text{Dot product})$$

$$\therefore \cos \theta = \frac{\mathbf{x} \cdot \mathbf{y}}{|\mathbf{x}| |\mathbf{y}|}$$

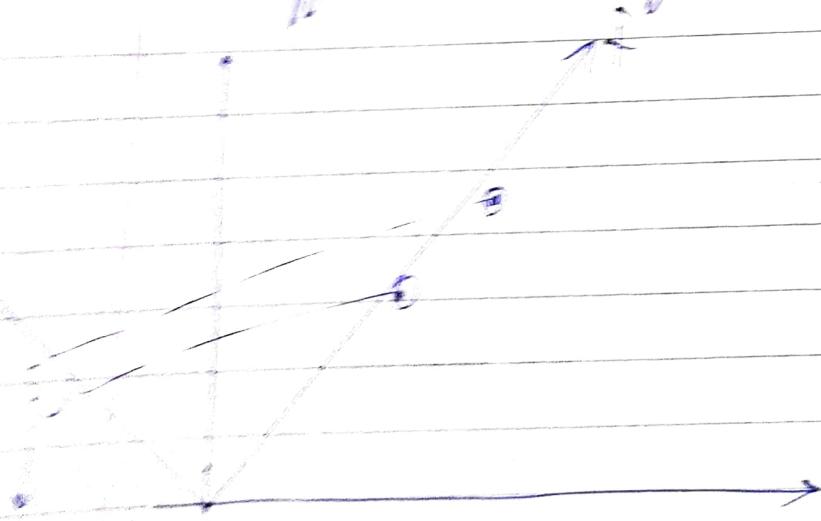
$$\frac{\mathbf{x}^T \cdot \mathbf{y}}{\sqrt{\sum_{i=1}^D x_i^2} \sqrt{\sum_{i=1}^D y_i^2}}$$

# Eigen Vectors are perpendicular to each other.  
 $\therefore v_i \cdot v_j = 0$ .

- All dimensions are 1<sup>r</sup> cos effect of element of 1D should be zero on element of another dimension



- If variance is large wrt a dimension, it becomes easier to differentiate b/w the elements i.e. spreadness ↑.
- we transform dataset into other dimensions to differentiate among patterns easily.



projection of these points in this dimension

1. calculate Dispersion Matrix of  $\Sigma$  of given datasets.

$X \rightarrow \text{dataset} \rightarrow n \times D^e$  attributes  
entries ,  $X^T = D \times n$

- a) \* Subtract mean from each data in the data set. ( $\mu$ )  
 (making mean-subtracted data)  
 (Mean of new dataset = 0)

b)  $\Sigma = X^T \cdot X$ .

c) calculate eigen values and eigen vector of  $\Sigma$

- \* trace ( $\Sigma$ ): sum of diagonal elements  
 : sum of variance wrt each attribute

trace ( $\Sigma$ ) =  $\sum_i \lambda_i$  , all  $\lambda_i \geq 0$   
 (sum)

\* If arranged,

$$\lambda_1 > \lambda_2 > \dots > \lambda_D$$

$v_1 v_2 \dots v_D$  ← corresponding eigen vectors  
 ↳ all are  $1^T$  to each other

$$\frac{\sum_i^D \lambda_i}{\sum_{i=1}^D \lambda_i} > 99\%$$

∴ We can avoid 'd to D'

→ Not important in terms of variance

$$\text{Determin. } |\Sigma| = \lambda_1 \times \lambda_2 \times \lambda_3 \dots \times \lambda_d$$

- inand

$$\begin{matrix} \lambda_1 & \lambda_2 & \lambda_d \\ \downarrow & & \\ v_1(Dx) & v_2 & \xrightarrow{d} v_d \end{matrix} \dots \dots \dots v_D$$

Transformation Matrix ( $W$ ) contains all selected eigen vectors

$$W_{D \times d} = [v_1 \ v_2 \ \dots \ v_d] \quad \boxed{D \times d}$$

$$x_i = \boxed{\quad}_{D \times 1}$$

$$y_i = W^T \cdot x_i$$

$\hookrightarrow$  smaller dataset

• Variance ( $y_i$ ) > Variance ( $x_i$ )

$\uparrow$   
Transformed