

Principal Component Analysis:-

Sample \rightarrow

\rightarrow	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	9	10	2.5	1.3	2
Gene 4	5	7	6	2	5	7

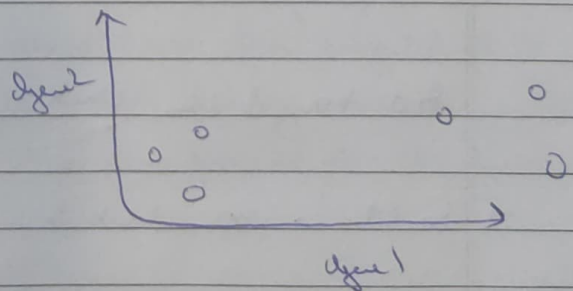
To represent this we require 4D which is not possible to visualize.

- So we are going to talk about how PCA can take 4 or more gene measurements (all the 4 or more dimensions of data) and make a 2-D PCA plot.
- We also talk about how PCA can tell which variable (gene) is most valuable for clustering data.

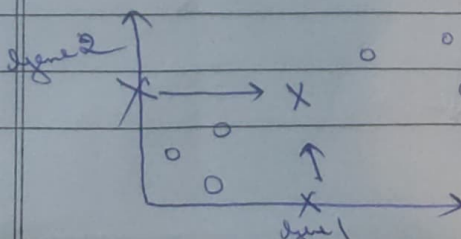
Working of PCA

\rightarrow	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

1) Plot the data



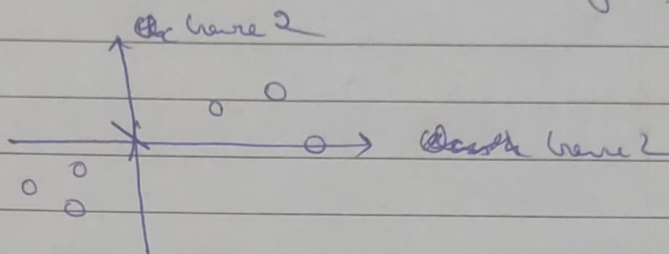
2) Calc the avg movement of Gene 1 and avg value of Gene 2



With avg values we can calculate the centre of data.

Now we will focus on graph, no longer need the original data

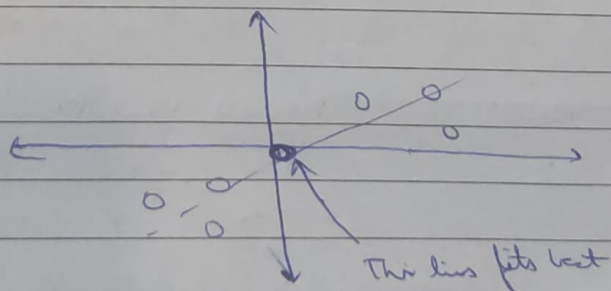
3) Now shift the data that centre is on origin



Note: Shifting the data did not change how the data pts are positioned relative to each other

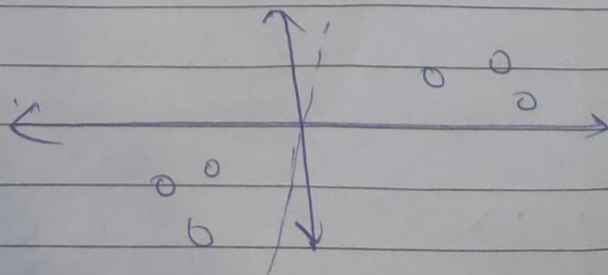
4) Try to fit a line

We start by drawing a random line that goes through origin. Then we rotate the line until it fits the data as well as it can, given that it has to go through origin.

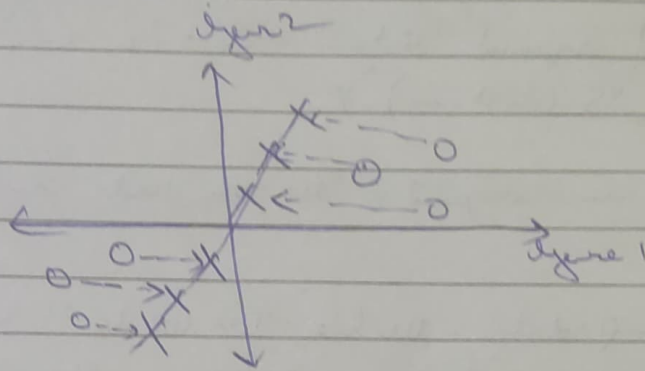


But how do we decide which line is the best?

→ So let's go back to 'random line that goes through origin'.

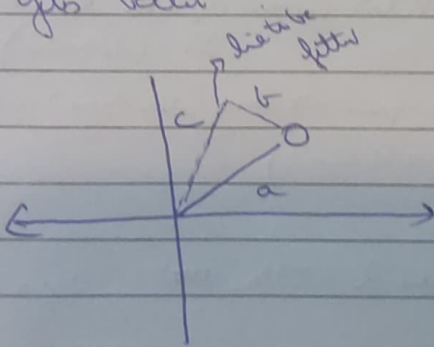


To quantify how good this line fits the data, PCA projects data onto it



Then it can either measure the distances from the data to the line or try to find the line that minimizes those distances or it can try to find the line that minimizes the distance from the projected pts to the origin.

The projected pts distance on line gets longer when the line fits better

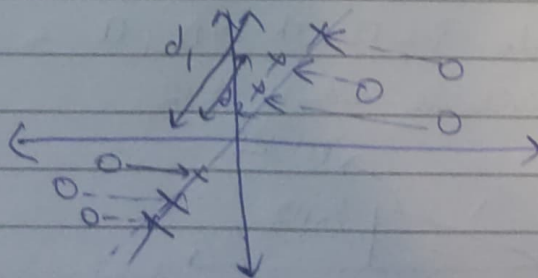


a is the distance from pt to origin (fixed)

$$a^2 = b^2 + c^2 \quad (\text{Pythagorean theorem})$$

In order to minimize b , we need to increase c .

Minimize b , the distance from pt to line but it's easier to calc c , the distance from projected pt to origin, so PCA finds best line by minimizing the sum of squared distances from the projected pts to the origin.



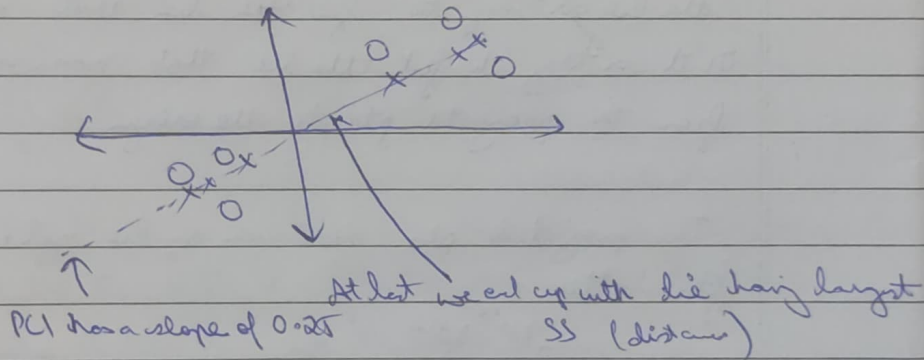
So 6 distances $d_1, d_2, d_3, d_4, d_5, d_6$

$$\text{Sum of squared distances} = d_1^2 + d_2^2 + d_3^2 + d_4^2 + d_5^2 + d_6^2$$

$$(\text{SS (distances)}) =$$

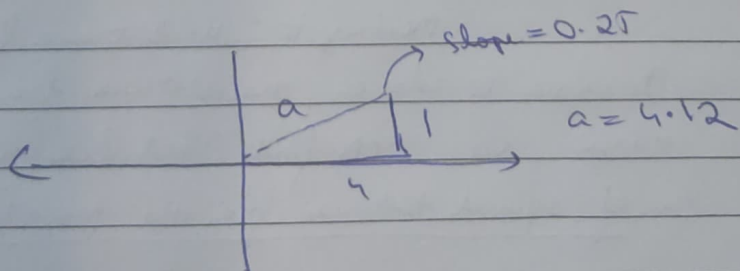
Distance is squared so the -ve doesn't cancel +ve.

We keep rotating the line till we don't get a line with the largest sum of squared distances b/w the projected pts and the origin



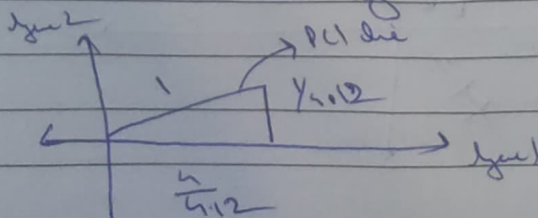
(This means that data mostly spread out along the lye1 axis and only a little bit spread out along lye2 axis)

The ratio of lye1 to lye2 tells that lye1 is more important



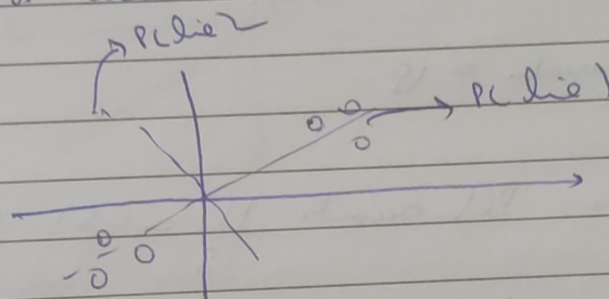
When you do PCA with SVD, a must be scaled to 1.

So \therefore we divide all sides by 4.12.



- ~~Eigen~~ → This unit long vector is called "Singular vector" or "Eigenvector" for PC1
- $\frac{SS(\text{distances for PC1})}{n-1} = \text{Eigenvalue for PC1}$
- $\sqrt{SS(\text{distances for PC1})} = \text{Singular value for PC1}$

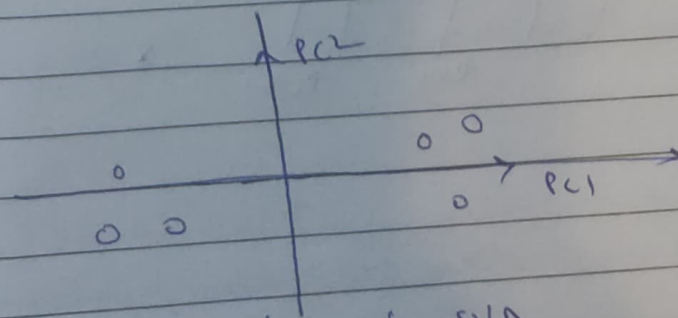
→ PC2 line is line \perp to PC1



According to PC2 gene 2 is more up than gene 1.

$$\frac{SS(\text{distances for PC2})}{n-1} = \text{Eigenvalue for PC2}$$

- To get final plot,
we rotate everything so we get PC1 as horizontal



That is how PCA is done using SVD

Eigen values are just measure of variation

$$\frac{SS(\text{distance for PC1})}{n-1} = \text{variance for PC1}$$

$$\frac{SS(\text{distance for PC2})}{n-1} = \text{variance for PC2}$$

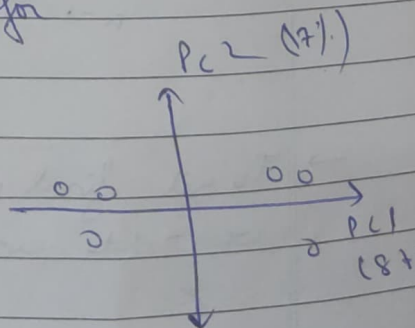
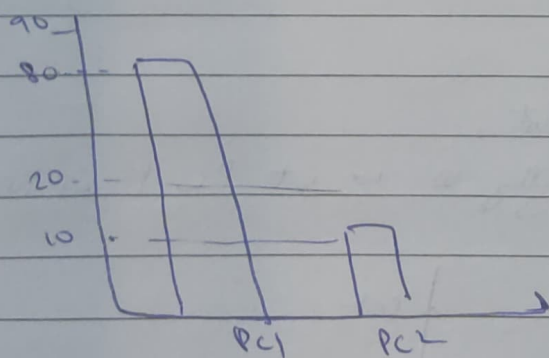
Eg. variance PC1 = 15
variance PC2 = 3

$$\text{Total variance} = 18$$

That means PC1 accounts $\frac{15}{18} = 0.83 = 83\%$ of total variance around PCs

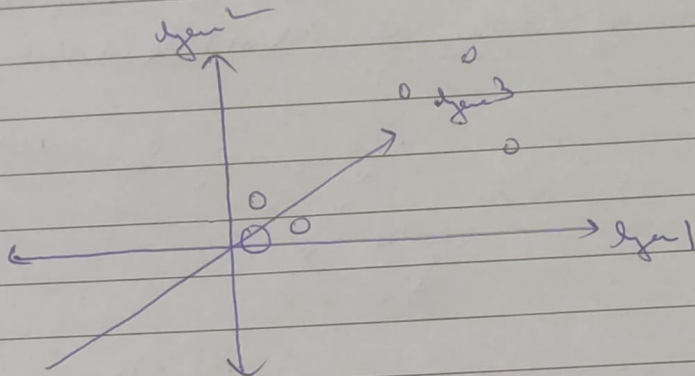
PC2 accounts for $\frac{3}{18} = 0.17 = 17\%$ of total variance around PCs

→ A Scree Plot - graphical representation of % of variance that each PC accounts for.



→ PCA with 3 variables

	Mouse1	Mouse2	Mouse3	Mouse4	Mouse5	Mouse6
Case 1	10	4	8	3	2	1
Case 2	6	4	5	3	2.8	1
Case 3	12	9	10	2.5	1.3	2



- You enter the data
- Fit best fitting line PC1

0.62 Ratio Case 1

0.15 Ratio Case 2

0.77 Ratio Case 3 ∴ Case 3 - most up for PC1

- Fit PC2 in to PC1

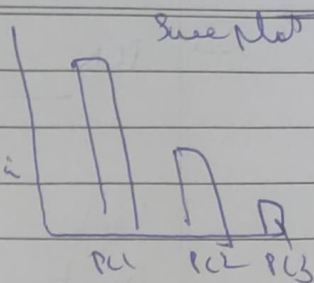
- Fit PC3 in to PC1 & PC2 as goes through origin

* In theory, No of PC = No of variables / Samples (Whitaker's rule)

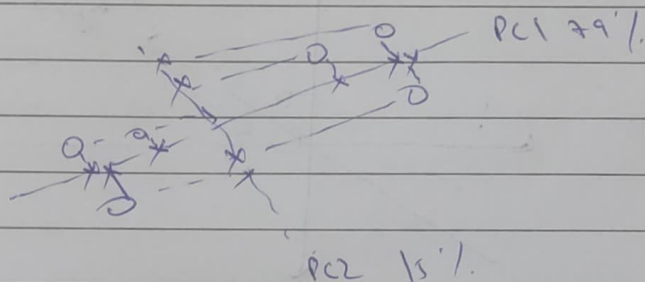
- Once you have all PCs, you can use eigenvalues (ie SS divided) to determine proportion of variance that each PC accounts for

∴ PC1 - 79% variance, PC2 - 15%, PC3 - 6%.

- ∴ It can a 2D graph using just PC1 and PC2 good approximation of 3D graph since it would account for 94% of variation in data.



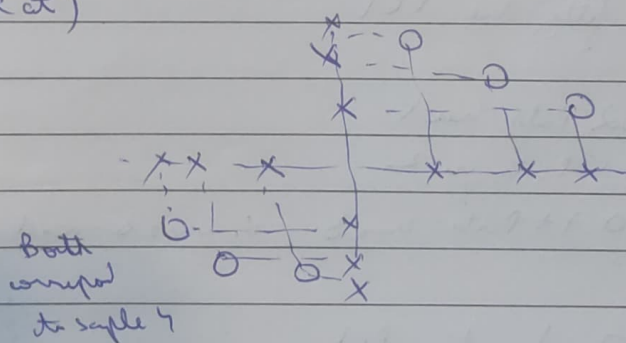
- To convert the 3D graph into a 2D PCA, we just strip away everything but PC1 and PC2 of data.



Project the sample
onto PC1
and PC2

push test

The we rotate, PC1 - horizontal, PC2 - vertical (easier to look at)



→ If we have higher dimension, we get scree plot and depend on the value from scree plot we choose the PC