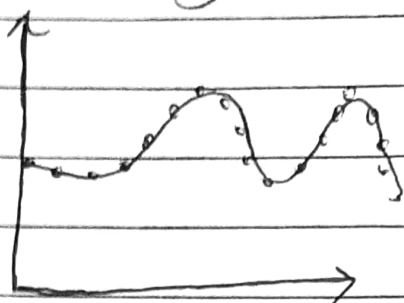


FDA pre-processing technique

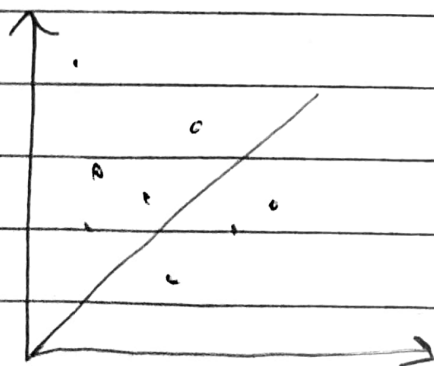
PCA (this principal component Analysis)
(~~old~~ modern data Analysis)

before starting PCA concept lets understand what is overfitting and underfitting
its modeling error



overfitting modeling error

it happens when excess knowledge of attributes where model tries to use all attributes (even the least important ones);



underfitting modeling error

Due to lack of enough attributes model doesn't have enough knowledge and gives incoherent result

and what is best fitting
it is when model make the prediction with 0 error.

Now we come to the PCA

to overcome the overfitting problem. we need PCA

what it do
it reduces the overfitting problem.

why overfitting happens when there is lot of dimensions (columns) in our data set. Now in order to reduce dimensions of our dataset in order to get rid of overfitting problem, we need PCA (Principal Component Analysis).

→ Some important point.

PCA change the orientation of dataset

No. of PCA principal component can be less than or equal No. of attributes

PCA have job to compress the columns and find out unnecessary columns.

in PCA Accuracy is just going to hurt but at best for simplicity.

DATE	/	/
------	---	---

* whenever No. of PC is generated whom should we give priority?
 \Rightarrow to PC, in ascending order.

* No. of PC should be independent to each other.
and PC have to capture essence of column.

* let take example in order to understand PCA.

① Eigen values :- it's going to capture the essence of the dataset.

② Eigen _{vector} :- it's going to capture direction for us.

Use of PCA

① identify relation between columns

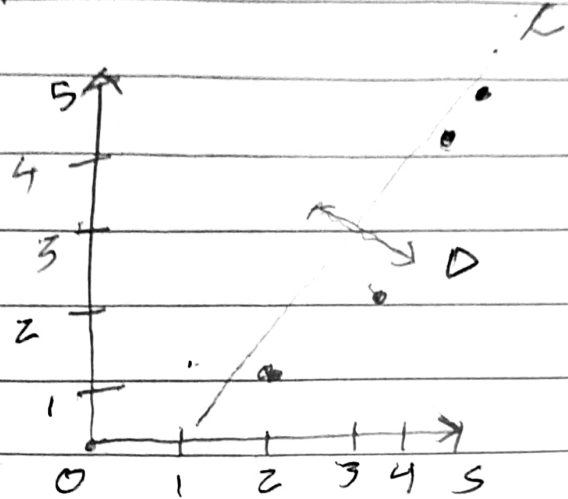
② reduce columns

③ visualize data in 2D

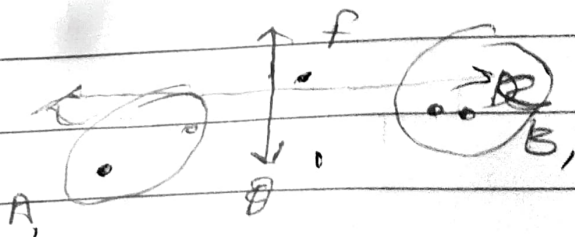
* data in PCA should be in linear fashion among the columns then only PCA will perform well otherwise not.

Suppose we dataset

A	B
1	2
5	6
2	3
4	5



So we have dataset in linear fashion
 So we can fit line to calculate variance (line C) and since we have 2nd column too we can put more line (note:- the line we are putting we assume as axis) and we can calculate more variances now lets assume C as PC_1 and D as PC_2 and if we have more columns and so on.

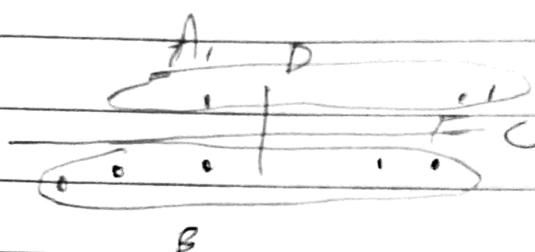


So the most data on circle is deleted
for line C and datapoint with them.

~~PC₁~~

<u>AB</u>	<u>Influence</u>	<u>un numbers</u>
A ₁	high	10
B ₁	medium	5
F	low	0.2

Similar for D



~~PC₃~~

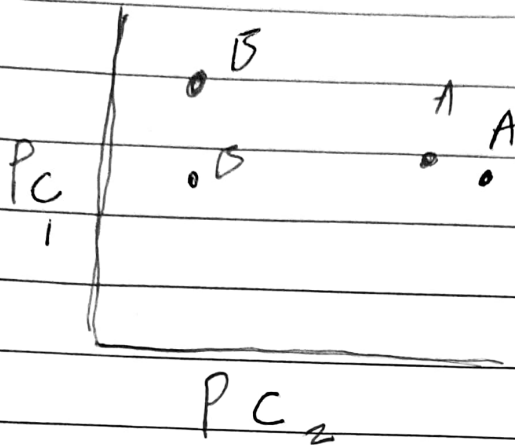
<u>AB</u>	<u>influence</u>	<u>un number</u>
A	medium	3
B ₁	high	10
F	low	0.5

$$A \text{ PC}_1 \text{ Score} = (1 \times 10) + (5 \times 0.5)$$

$$P_{C_1} = 12$$

Similar for PC_2

$$= 6$$



conclusion:- this is basic understanding for PCA we will be using mathematical formula for this calculating PCA.