

Lecture - 06: (2003) 2003 2003 2003 2003 2003 2003

Unsupervised ML:

\rightarrow k means Clustering.

→ Hierarchical Clustering.

→ Silhouette Score.

→ DB Scan Clustering.

118

20

11

$$G) \rightarrow \frac{f}{\partial f} \cdot \frac{\partial}{\partial f} =$$

Unsupervised ML? It means actually there is no output. We

and not as individual
need to form a clusters.

East - middle

101

f

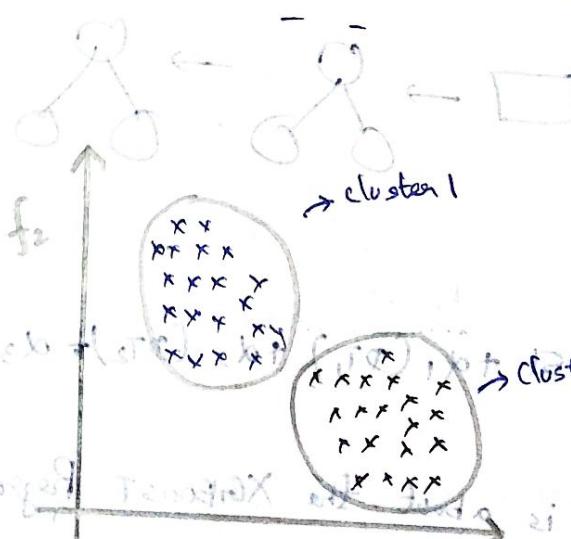
We form
Clusters to

make off

Category

Clusters

{ Similar kind of data }



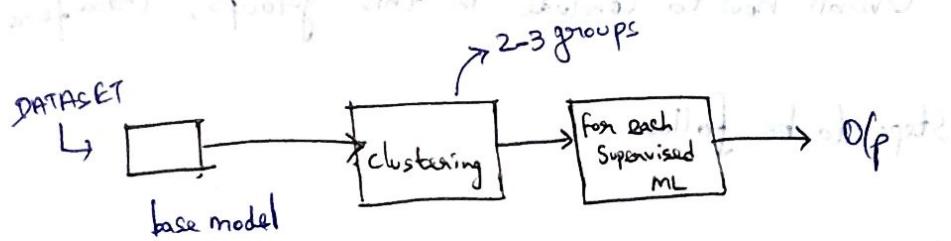
(+L) $\alpha_1 \dots \alpha_n$ + (+V) $\beta_1 \dots \beta_m$ $\vdash_{\text{L} \cup \text{V}} \gamma_1 \dots \gamma_p$

~~• Kassen-Post~~ ~~Faxbox~~ ~~abholen~~  si sdt

If it is also come under ⁴ labour code then it will be

We already know Custom Ensemble technique,

In this technique we create,



first clustering then can able to apply
Supervised ML algorithm.

Most of the time clustering is applied in ensemble
technique.

The "k means" \rightarrow Centroids.

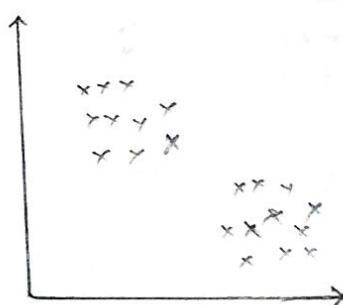
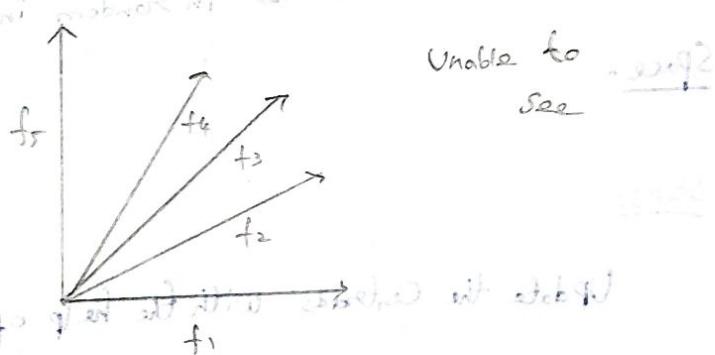
Suppose $k=2$, we get the two cluster groups.

Suppose $k=3$, we get the three cluster groups like that.

Suppose we have High Dimension data we cannot able to

See how the clusters will be formed.

Consider this,

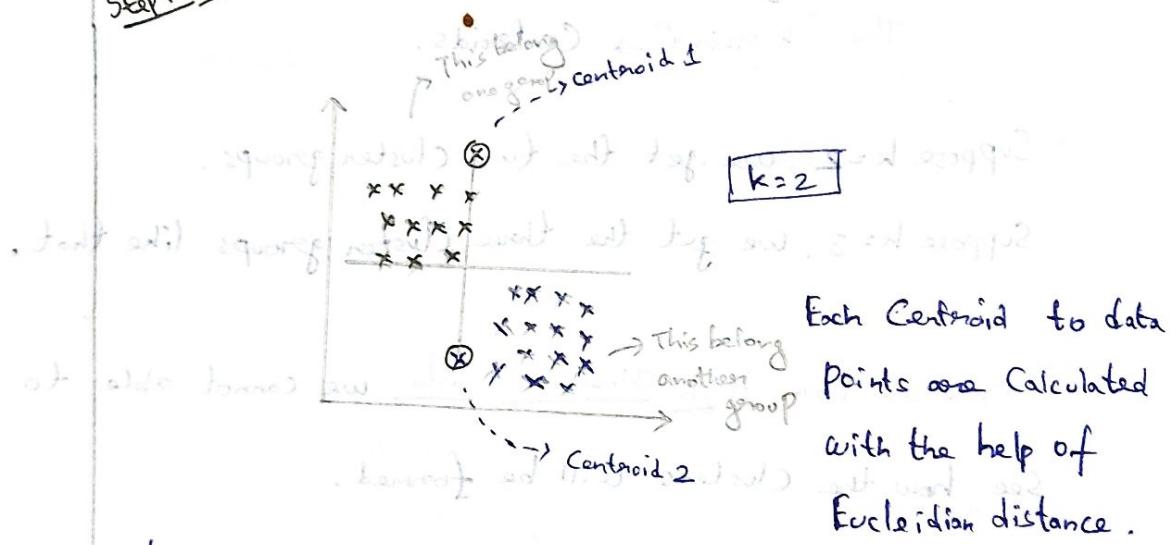


Overall how to conclude in this groups, Therefore Some
Steps to be follow.

Steps:

1. We try k values \Rightarrow Suitable $k=2$ just Consider.
2. Initialize the k number of Centroids.
3. Update the Centroids.

Step 1 and 2

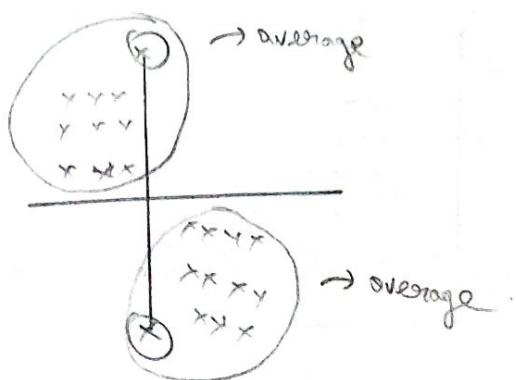


The Centroids are initialize in random in the graph.

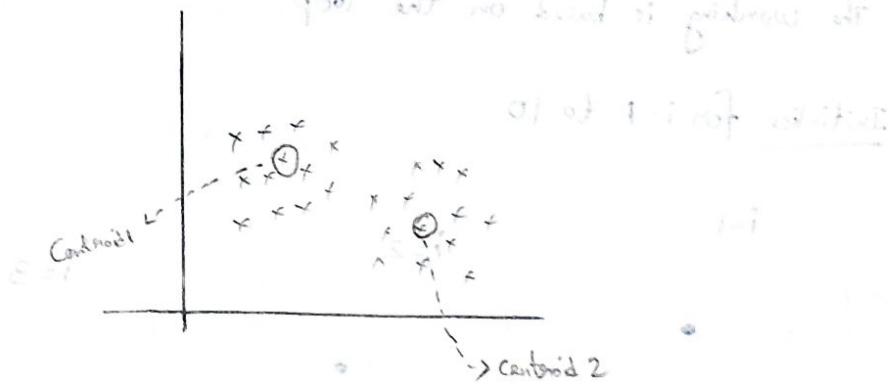
Space.

Step 3:

Update the Centroids with the help of taking the averages.



get all no. of point of following off

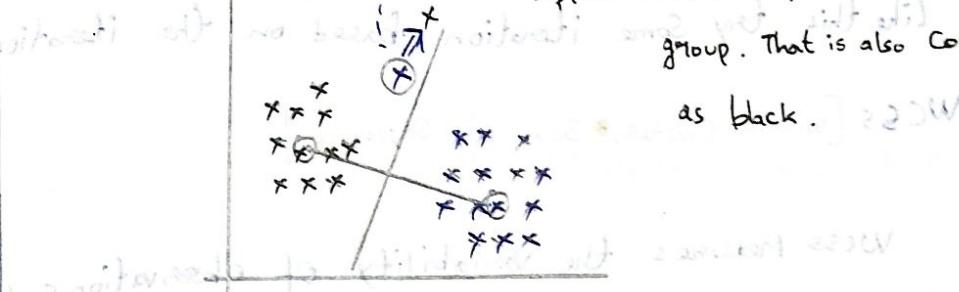


Again, the perpendicular line make, It means calculate Eucleidian distance.

blue charge
to black

Suppose some blue point in black cluster
group. That is also considered

as black.



Because, After updating all, still some blue points in

black group. It is also considered as black point.

How to decide the k-value?

With the help of Elbow method to find the optimized

k-value.

Elbow method:

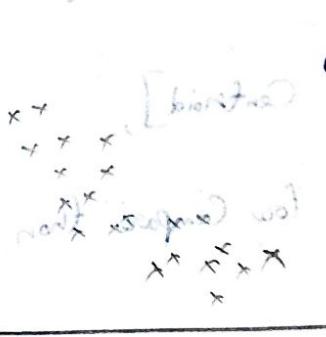
The elbow method is used to find the optimized k-value.

k-value.

Consider this, we cannot

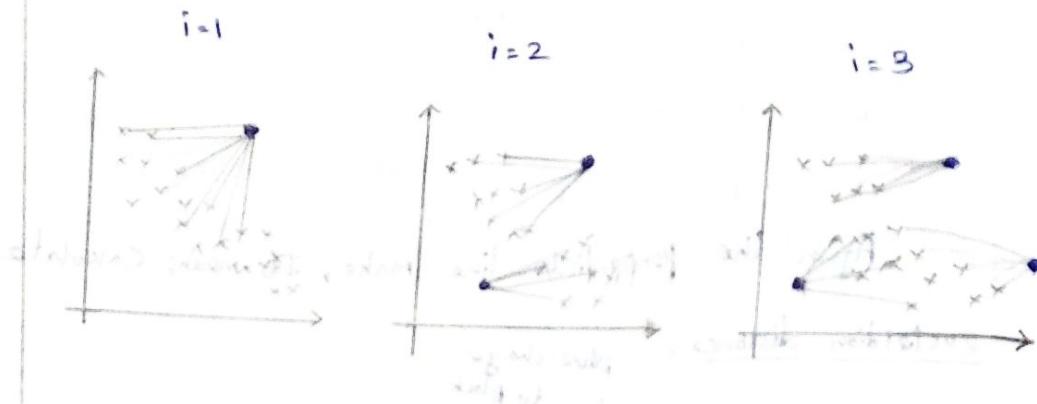
directly split clusters. In this

case we can't use elbow method.



The working is based on the loop.

Initiate for $i=1$ to 10



Like this try some iteration based on the iteration and WCSS [within clusters sum of squares].

WCSS measures the variability of observations within each cluster using Euclidean distance.

BCSS [between clusters Sum of Squares] which measures the squared average distance between all Centroids.

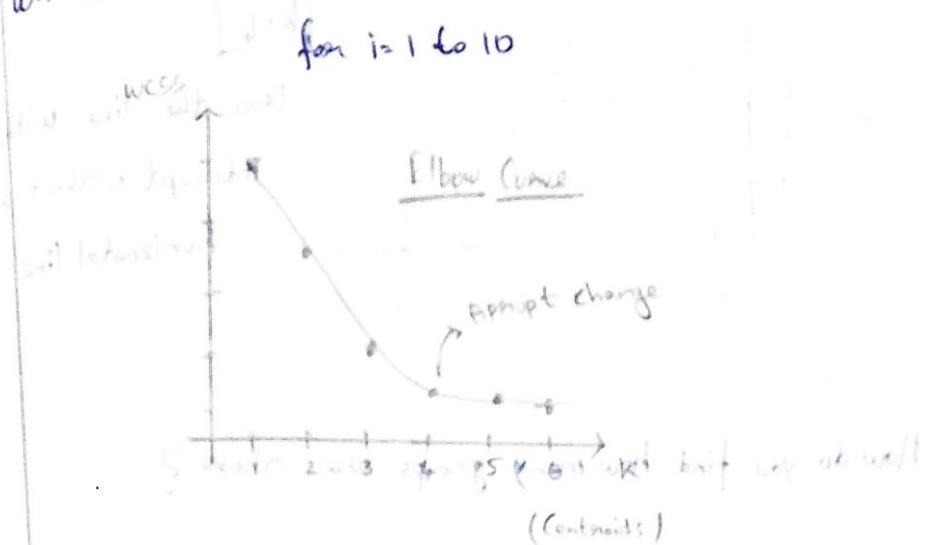
To calculate BCSS, you find the Euclidean distance from a given Cluster Centroid to all other Cluster Centroids.

Obviously $i=1$ [It means one centroid],
The WCSS is quite large.

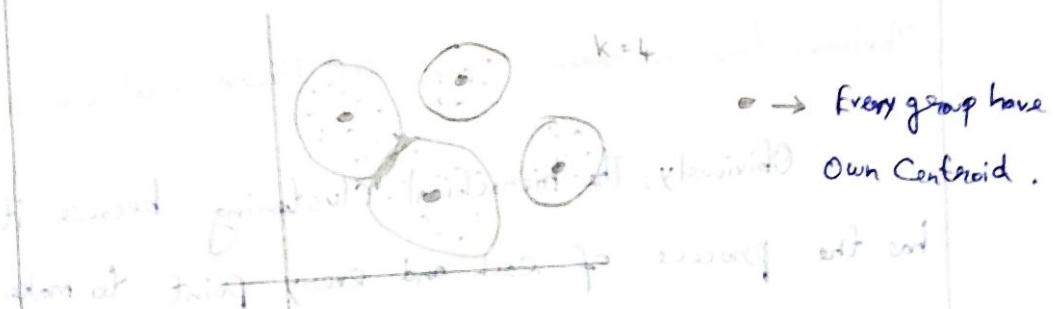
$i=2$ [It means two Centroid],

The WCSS is low Compare than $i=1$.

Because, Based on Centroids the whole data points split to calculate WCSS. Obviously the more Centroids the distance will be low.



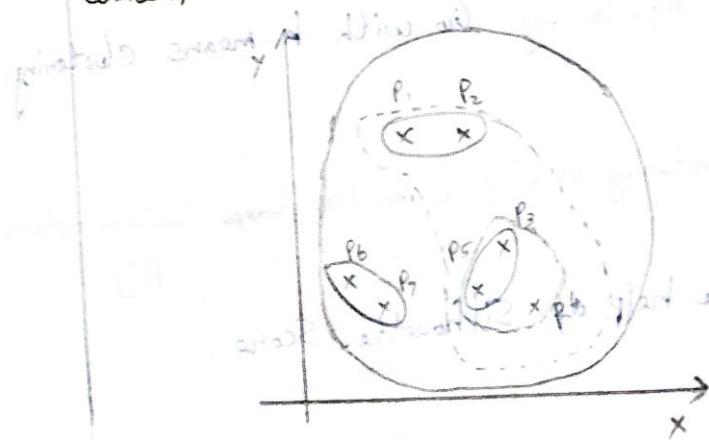
The abrupt change happened at 4 Centroid. It means actually 4 groups are available in the dataset.



This is all about the k-means Clustering.

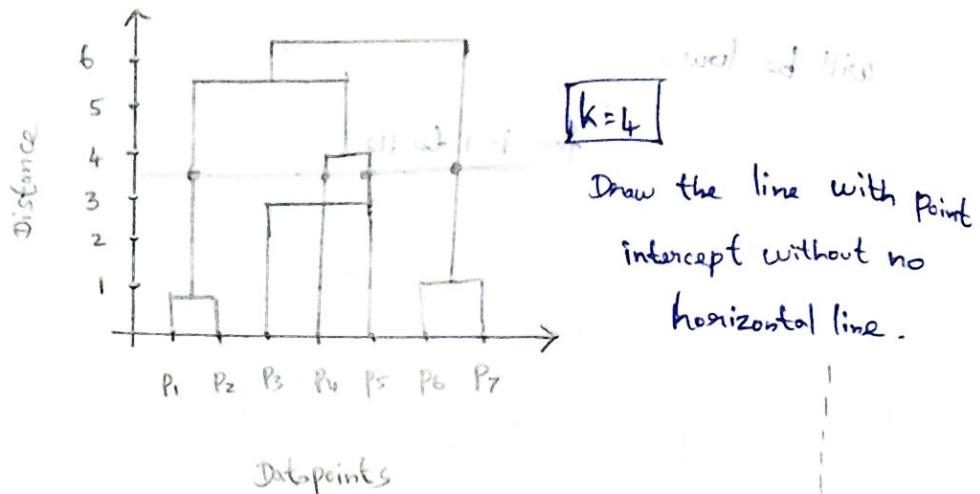
Hierarchical Clustering:

Consider,



This clustering make clusters based on minimum distance to maximum distance point.

After clustering state which will be cluster no 1 and no 2
so on till we obtained seven clusters. This is called as Dendrogram.



How do you find how many groups are above?

Answer: You need to find the longest vertical line that has no horizontal line passed through it.

Maximum time is taken k-means (or) Hierarchical Clustering?

Obviously, The hierarchical clustering because it has the process of each and every point to make dendrogram is high time complexity.

If the dataset is small, Go with hierarchical Clustering.

If the dataset is high (or) big, Go with k-means Clustering.

How to Validate Clustering model? [Like Performance Confusion metrics, R^2].

Ans: With the help of Silhouette Score.

Validating Clustering Models:

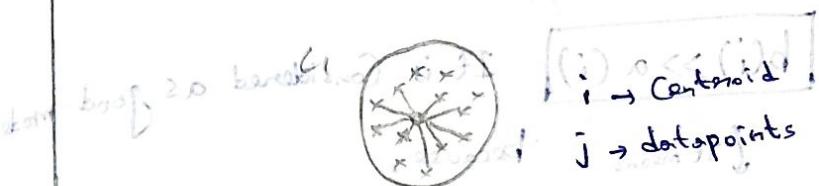
Final notes

The Silhouette Score is used to validate the cluster models.

Some steps to be followed,

→ To find out the $a(i)$ and $b(i)$.

Consider $a(i)$ is nothing but one group.



The Centroid (i) calculates the differences and squares of all data points with summation.

$$\sum d(i,j)$$

Then calculate the overall averages of each with the formula,

$$a(i) = \frac{1}{|C_i|-1} \sum_{j \in C_i, i \neq j} d(i,j)$$

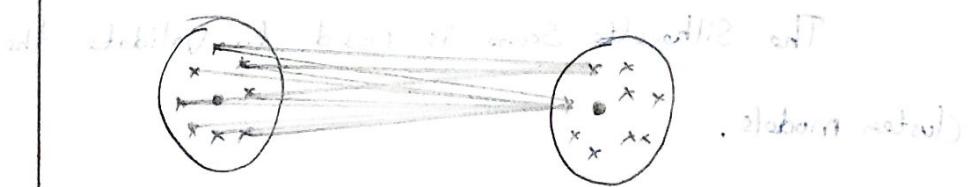
Similarly calculate $b(i)$ with the help of Considering

Other group.



Then calculate $b(i) = \sum d(i,j)$

After that,



Each and every points are calculated with each every point of another group.

Based on above,
 Group associated with point i is (i) (silhouette)

$b(i) > a(i)$ It is Considered as good model

It means because,

The distance between $b(i)$ is greater than the distance of $a(i)$.

The more distance, the clusters are correct

group without any point.

Bad model

$a(i) > b(i)$ → It means the greatest distance

of $a(i)$ is greater than $b(i)$.

The value of the Silhouette Score its ranges

between (-1 to 1).

George will

The value towards +1 ⇒ It Considered as good model.

The value towards -1 ⇒ It Considered as Bad Model.

The formula will be, at which is given as

Silhouette Score,

$$sc(i) = \frac{b(i) - a(i)}{\max\{a(i), b(i)\}}, \text{ if } |c_i| > 1$$

$$\max\{a(i), b(i)\}$$



By this formula gives the value (-1 to +1).

DBScan Clustering:

Density-Based Spatial clustering of Applications with

Noise (DBSCAN)

The four components of DBScan Clustering,

→ Epsilon (ϵ) → Epsilon, defined as radius ϵ

→ Min Points → Minimum points required to form a cluster

→ Core Points

Epsilon refers to the radius
of the circle.

→ Border Points

→ Noise Points

Min points:

Core Points:

At least present 4 points.

[min points = 4]

Border Points:

In this case,

core
border points.

At least one point

present, we call it is a core border

border points.

Noise Point:



without any points, we can say Noise Points.

[It means the points are present outside]

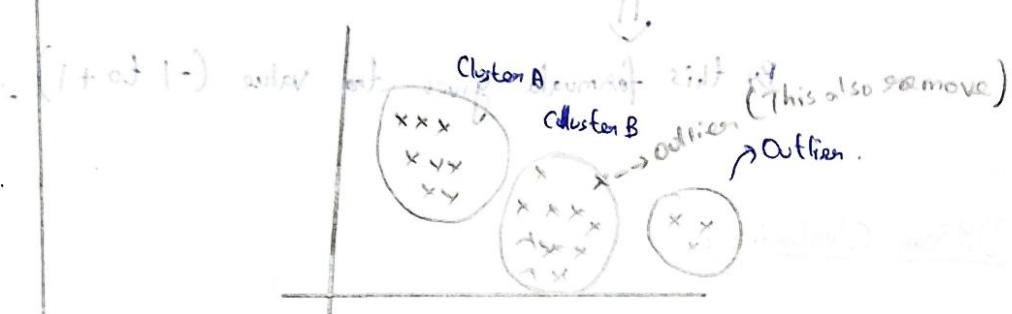
↳ Outliers.

The Min pts is decided by hyperparameter.

The ϵ -epsilon value also decided in some specific way.

Why Use? $(1) + (2) = (3)$ $\rightarrow 130$

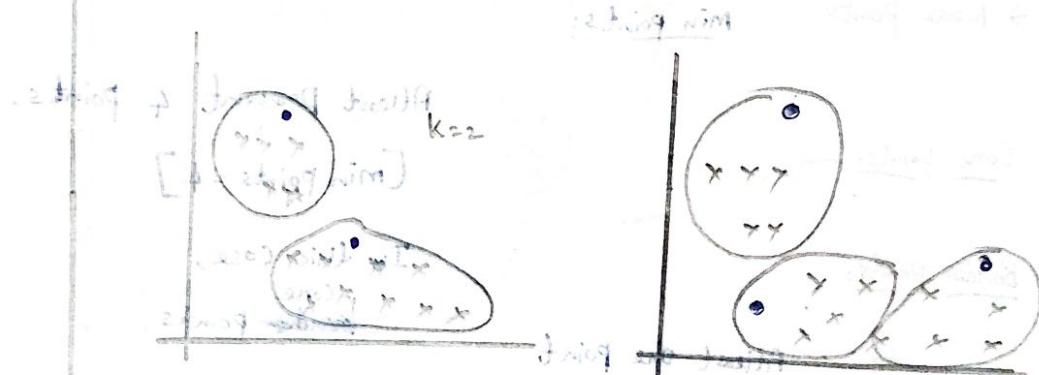
To detect the outliers efficient way,



Once find that Outlier, we don't need to take the group for analyze, Just neglecting whole Outlied group.

One issue with k-means:

If initialize the centroids in k-means it should be Very Very far. Otherwise, Suppose initialize Centroids in short distance it may be take two groups of two centroid instead of three Centroid of three groups.



Centroid use in short distance | Centroid use in long distance.

short distance (bad) vs. long distance (good)
[at start during one step of learning it]
initialization

Efficient One

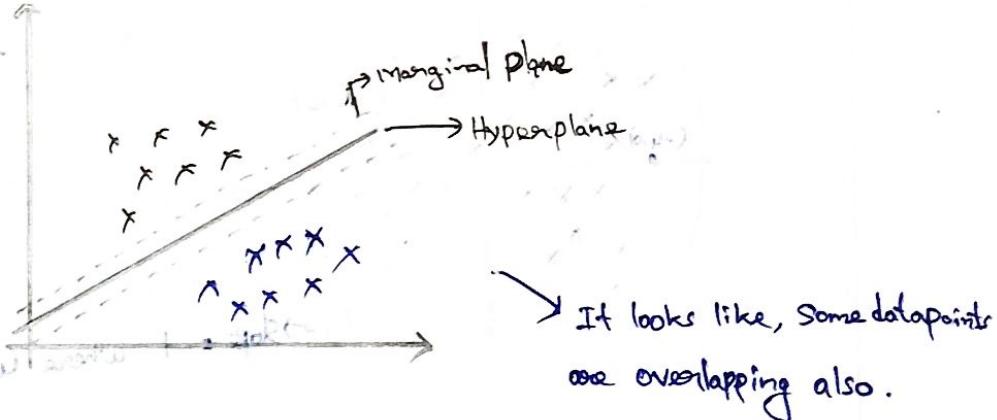
To set the centroid very very far in random space, something algorithm is used called as k means ++.

Support Vector Machine: (SVM)

The main aim of the Support vector machine create best fit line with marginal plane. In this way, we create or decide the output.

It is supervised machine learning algorithm for classification.

looks like,

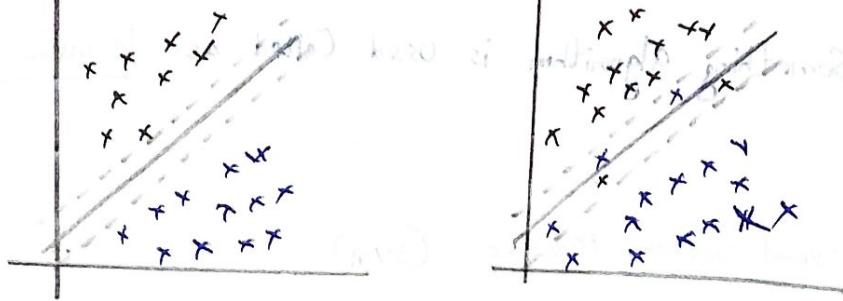


Hard margin and Soft margin:

The data is linearly separable and we don't want to have

any misclassification, we use SVM with a hard margin.

• If we want to classify data with hard margin with loss of

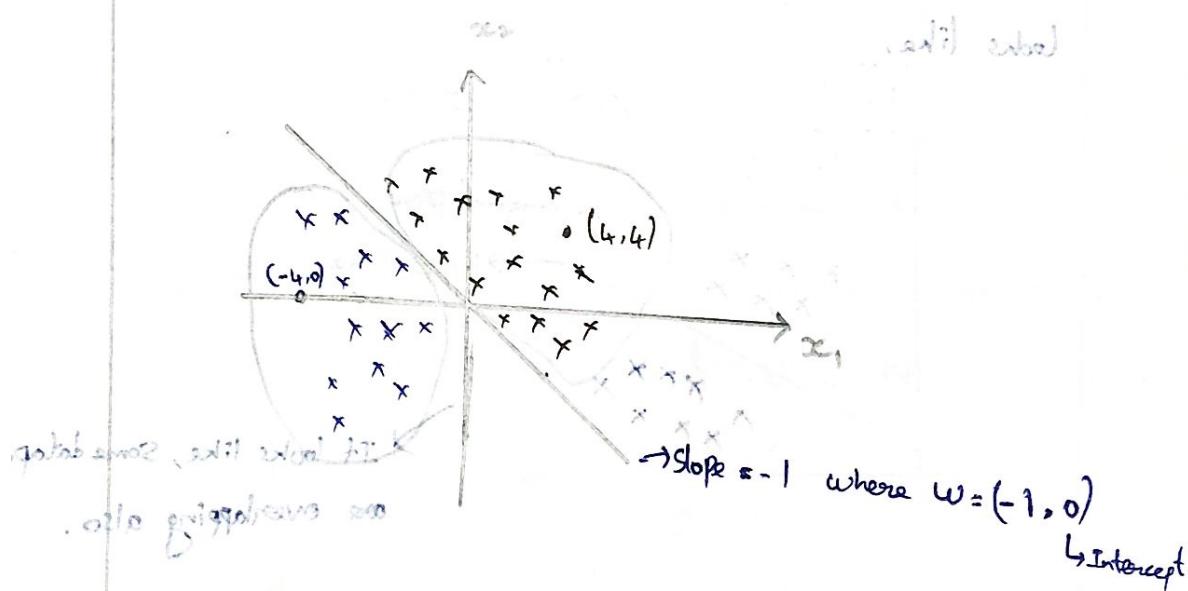


hard margin, Hard margin happens with no misclassification

However, when a linear boundary is not feasible, OR we want some misclassification is called soft margin.

Generally, In normal dataset only contains soft margin.

Consider,



The equation of straight line $y = mx + c$

It is also written as,

and it is known that $\boxed{ax + by + c = 0}$ (and it is also written as)

∴ $y = \frac{-c}{b} - \frac{ax}{b}$ if $a \neq 0$, where $m = -\frac{a}{b}$, $c = \frac{-c}{b}$

$$y = mx + c$$

slope Intercept

Same like written as,
and with $f(x)$ as a linear function of x .

$$y = w_1 x_1 + w_2 x_2 + w_3 x_3 + \dots + w_n x_n + b$$

(why write this? because

based on feature (x)

$$\boxed{y = w^T x + b}$$

Consider $(-4, 0)$,

Apply in y ,

$$y = \begin{bmatrix} -1 \\ 0 \end{bmatrix} \begin{bmatrix} -4 \\ 0 \end{bmatrix} \Rightarrow -4 + 0 = 4 \quad \boxed{y=4}$$

$$\cancel{\begin{bmatrix} -4 + 0 \\ 0 + 0 \end{bmatrix}} = \begin{bmatrix} 4 \\ 0 \end{bmatrix} \quad \cancel{4} \quad \text{(Positive Value)}$$

and Consider $(4, 4)$,

$$y = \begin{bmatrix} -1 \\ 0 \end{bmatrix} \begin{bmatrix} 4 \\ 4 \end{bmatrix} \quad 0 = 4 + 4 \cdot 0$$

$$\cancel{\begin{bmatrix} 4 + 4 \\ 0 + 4 \end{bmatrix}} = \begin{bmatrix} 8 \\ 4 \end{bmatrix} \quad \cancel{8}$$

$$y = -4 + 0$$

first minimum at $b=4$

$$\boxed{y = -4}$$

$0 = d + x^T w$ (Negative Value)

$$1 = d + x^T w$$

So, we conclude that the points in below the linear line gives positive value. So that the each line. It gives positive value. So that the each

data points below are give positive value.

The value above the points in the linear line.

It gives y as negative value. So all the points are

Considered as negative.

How to Create best Marginal plane?

The answer is very simple. According to above

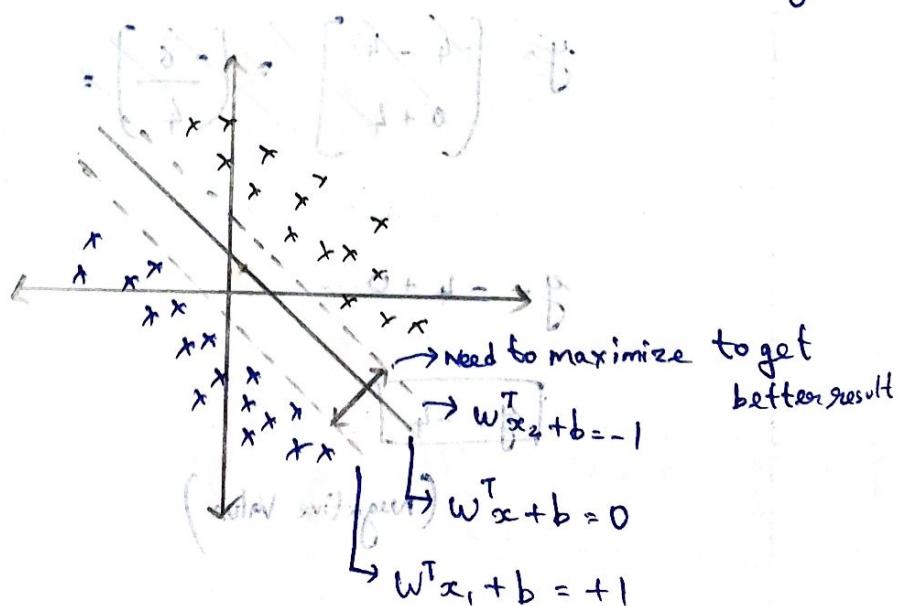
(Calculation)

$$ax + by + c = 0$$

↓

$$w^T x + b = 0 \quad [\text{It indicates the } mx + c = y]$$

straight line



The line equation is written with the help of

previous calculation.

To find out the distance between two marginal planes,

[writing getting out] so what's the margin width?

we do,

$$w^T x_1 + b = 1$$

$$w^T x_2 + b = -1$$

so we can do [writing getting out]

$$w^T(x_1 - x_2) = 2$$

In Generally, The marginal plane are the Vectors. we
don't need any mag. magnitude. So, remove with the

$$\frac{1}{\|w\|} \text{ with } \Leftrightarrow \frac{2}{\|w\|} \text{ maximization}$$

[writing] $\frac{w^T(x_1 - x_2)}{\|w\|} = \frac{2}{\|w\|}$, [Just dividing]

so we can divide with $\|w\|$

$$\frac{\|w\|}{\|w\|} \text{ with}$$

So, Our aim is to maximize the marginal, It means that

$$\text{Maximize } \left[\frac{2}{\|w\|} \right].$$

[writing] $\frac{2}{\|w\|}$ with respect to (w, b) .

so it will be maximized

$$\text{Maximize}_{(w, b)} \frac{2}{\|w\|} \text{ with respect to } (w, b)$$

Such that,

so that all the condition of margins will not

$$y_i \begin{cases} +1 & w^T x + b \geq 1 \\ -1 & w^T x + b \leq -1 \end{cases} \quad \text{(by writing with the help of diagram)}$$

and margin can be noted with the help of

So, this overall written as [for getting positive].

$$\boxed{y_i * (w^T x_i + b) \geq 1} \rightarrow \text{constraint}$$



for correct point

for correct point it means we

With the help of

get the positive value.

above equation to get

the positive value.

So,

$$\underset{(w,b)}{\text{Maximize}} \frac{2}{\|w\|} \Leftrightarrow \underset{(w,b)}{\text{Min}} \frac{\|w\|}{2}$$

The both are equal, so find minimizer the function,

$$\underset{(w,b)}{\text{Min}} \frac{\|w\|}{2}$$

with adding two more parameters.

So,

$$\underset{(w,b)}{\text{Min}} \frac{\|w\|}{2} + C_i \sum_{i=1}^n \xi_i$$

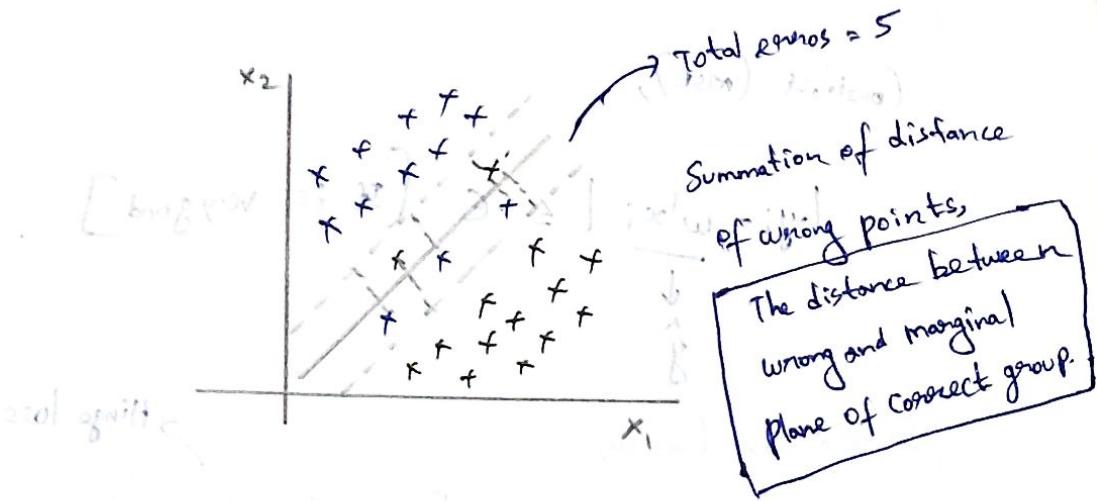
Now we have $\frac{\|w\|}{2}$ + $\sum_{i=1}^n \xi_i$ to minimize at least one

How many

errors we

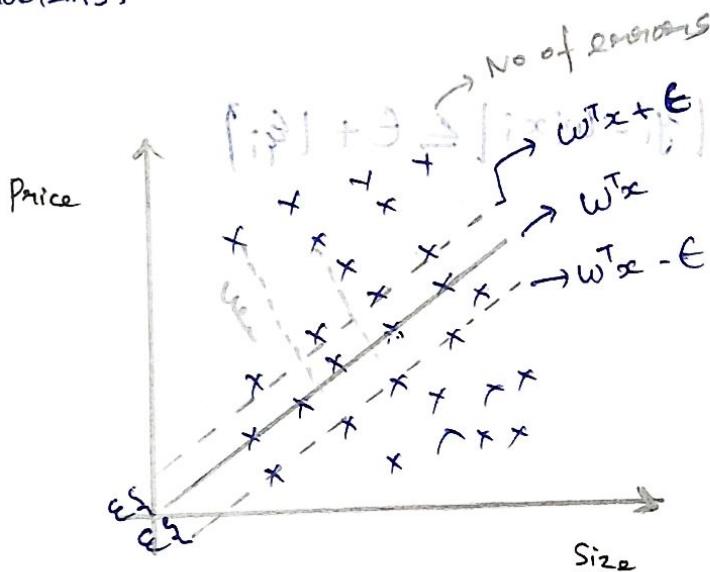
can have.

Summation of the distance
of the wrong data points.



Support Vector Regression (SVR): $\frac{\text{Minimize}}{C} \|w\|$

The SVR is basically tells about the Regression Problems.



ϵ -epsilon denotes the distance

So Cost function,

$$\text{Minimize}_{(w,b)} \frac{\|w\|}{2}$$

Constraint,

$$|y - w^T x_i| \leq \epsilon$$

Constraint (mse),

$$|y_i - \hat{w}^T x_i| \leq \epsilon \quad [\text{it is very good}]$$

\downarrow
 \hat{y}

So total cost function,

$$\underset{(\omega, b)}{\text{Minimize}} \quad \frac{\|\omega\|}{2} + C \sum_{i=1}^n |\xi_i|$$

Hinge loss

So, Constraint would be

$$|y_i - \hat{w}^T x_i| \leq \epsilon + \xi_i$$

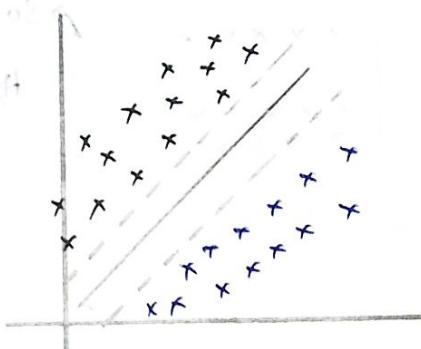
so solve w.r.t. weight vector \hat{w}

$\|\omega\|$ minimization

SOME CONCEPTS IN MACHINE LEARNING :

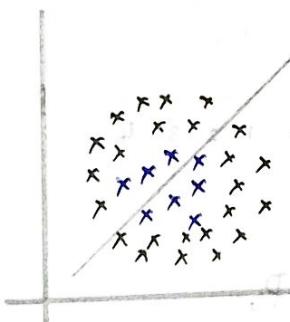
SVM kernels:

Consider a SVM for datapoints,



The datapoints may be an
Soft margin (or) Hard
Margin.

Suppose, we have a data look like below,



In this data, How do
Split, Because these
data are no correlation.

The above linear line split is wrong, so therefore my problem is that create the linear line for the problem. So, SVM kernel which comes in place.

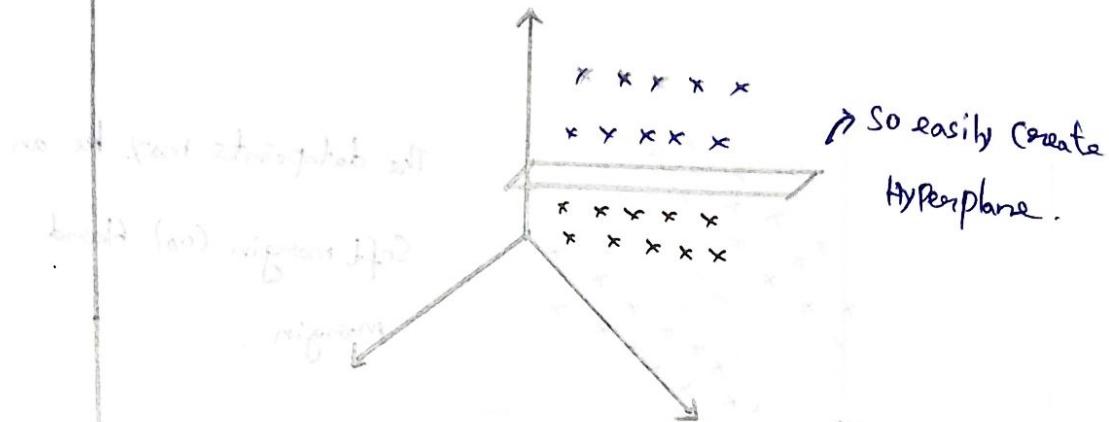
What Can I do for above Problem?

probably we need to apply some transformation. It means we need to change to lower dimension to higher dimension

The transformation is done with the help of

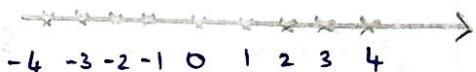
Mathematical formula .

After applying transformation, The diagram would be



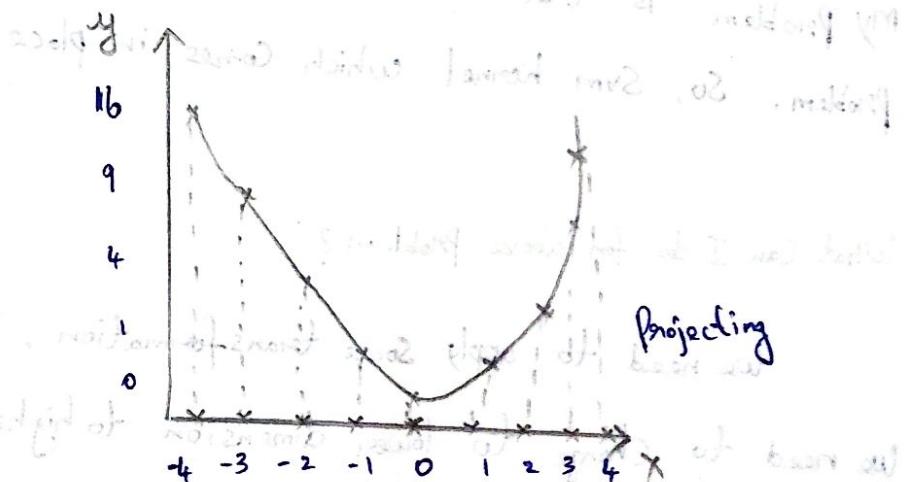
Suppose we have One dimensional,

One dimension (f_1)



We need to change 1D - 2D,

Applying $y = f(x)$, $f(x) = x^2$



$$y = (-4)^2 = 16$$

$$y = (-3)^2 = 9$$

$$y(4) = 16$$

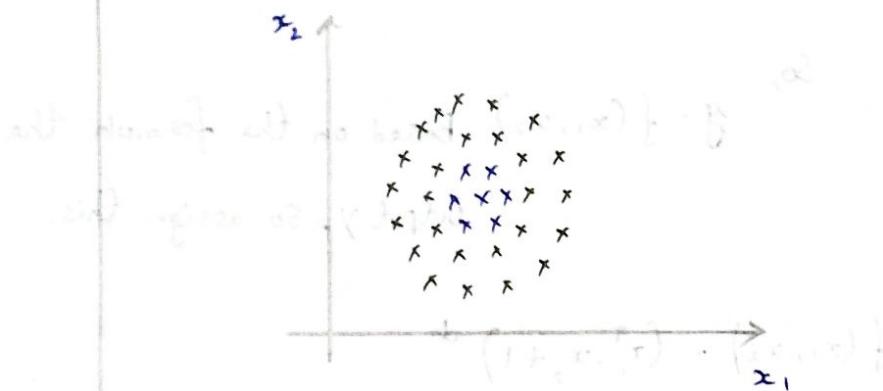
After plotting the value, we need to project the

One dimensional data.

Same like, the 1D - 2D, 2D - 3D, 3D - 4D like this

the line is fitted properly and make split.

Suppose we have,



We have this type of the data, This is 2D. We need to change 3D.

Some types of kernels are,

→ Polynomial Kernel.

→ RBF Kernel.

→ Sigmoid Kernel.

Based on situation of data, we use the kernel types.

In this, we use the polynomial kernels.

The above dataset, we have the features like

x_1, x_2 and O/p as y .

The formula for Polynomial kernel,

$$f(x_1, x_2) = (x_1^T \cdot x_2 + 1)^d \quad d \rightarrow \text{dimensions}$$

So,

$y = f(x_1, x_2)$ Based on the formula the output y . So assign this.

So,

$$f(x_1, x_2) = (x_1^T \cdot x_2 + 1)^d$$

Apply feature,

$$[x_1 \cdot x_2]^T * [x_1 \cdot x_2]$$

$$\begin{bmatrix} x_1 \\ x_2 \end{bmatrix} * \begin{bmatrix} x_1 & x_2 \end{bmatrix}$$

$$= \begin{bmatrix} x_1^2 & x_1 \cdot x_2 \\ x_1 \cdot x_2 & x_2^2 \end{bmatrix}$$

Take the Unique elements for creating 3D.

So, After taking this now we have five features,

$$x_1 \quad x_2 \quad y \quad x_1^2 \quad x_2^2 \quad x_1 x_2$$

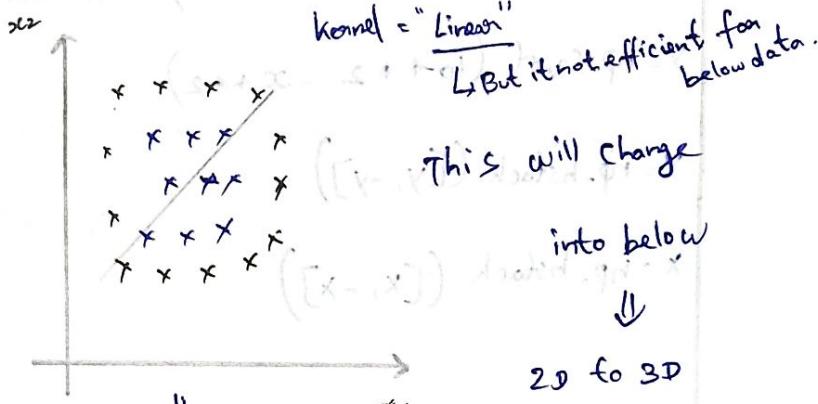
$$2 \quad 4 \quad 1 \quad 4 \quad 16 \quad 8$$

$$3 \quad 3 \quad 0 \quad 9 \quad 9 \quad 9$$

$$4 \quad 5 \quad 1 \quad 16 \quad 25 \quad 20$$

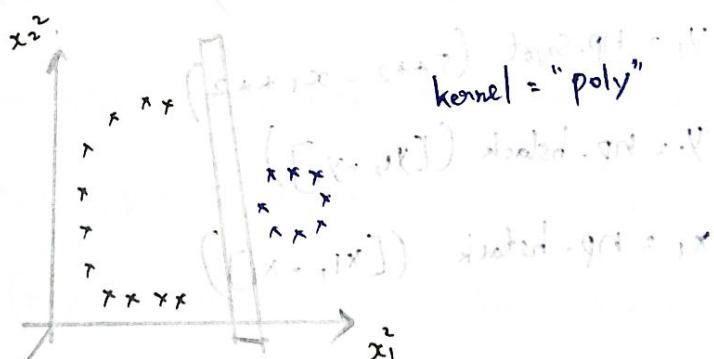
$$\dots \quad \dots \quad \dots \quad \dots \quad \dots \quad \dots$$

Dataset A ((x_1, x_2, y)) original form



2D to 3D

Dataset B ((x_1, x_2, y)) expanded form



↳ Easily split and take the decision.

The kernel type "poly", "rbf", "Sigmoid" are taking with the hyperparameter tuning.

This is all about the SVM kernels.

PRACTICAL IMPLEMENTATION:

SVM kernels and Practical

```
import numpy as np
```

```
import matplotlib.pyplot as plt
```

```
x = np.linspace (-5.0, 5.0, 100) # Feature 1
```

```
y = np.sqrt (10 ** 2 - x ** 2)
```

```
y = np.hstack ([y, -y])
```

```
x = np.hstack ([x, -x])
```

```
x1 = np.linspace (-5.0, 5.0, 100) # Feature 2
```

```
y1 = np.sqrt (5 ** 2 - x1 ** 2)
```

```
y1 = np.hstack ([y1, -y1])
```

```
x1 = np.hstack ([x1, -x1])
```

```
plt.scatter (y, x)
```

```
plt.scatter (y1, x1)
```

Create Dataframes with Pandas

```
import Pandas as pd
```

```
df1 = pd.DataFrame(np.vstack([y, x]).T, columns=['x1', 'x2'])
```

```
df1['y'] = 0
```

```
df2 = pd.DataFrame(np.vstack([y1, x1]).T, columns=['x1', 'x2'])
```

```
df2['y'] = 1
```

```
df = df1.append(df2)
```

```
df.head()
```

Independent and dependent features

```
X = df.iloc[:, :-2]
```

```
y = df.y
```

Split the dataset into train and test.

```
from sklearn.model_selection import train_test_split
```

```
X_train, X_test, y_train, y_test = train_test_split
```

```
(X, y, test_size=0.25, random_state=0)
```

y_train

Applying in Algorithm

```
from sklearn.svm import SVC
```

```
Classifier = SVC(kernel="linear")
```

If linear for above datasets.

Check the accuracy.

```
from sklearn.metrics import accuracy_score
```

```
y_pred = classifier.predict(x_test)
```

```
accuracy_score(y_test, y_pred)
```

Output: 0.45

Therefore Use Polynomial kernel

we need to find Components for the polynomial kernel.

```
# x1, x2, x1-Square, x2-Square, x1*x2
```

```
df['x1-Square'] = df['x1'] ** 2
```

```
df['x2-Square'] = df['x2'] ** 2
```

```
df['x1*x2'] = (df['x1'] + df['x2'])
```

```
df.head()
```

Independent and dependent feature

```
x = df[['x1', 'x2', 'x1-Square', 'x2-Square', 'x1*x2']]
```

```
y = df['y']
```

y

SVC benign and malignant

("malignant") SVC + classified

X-train, X-test, y-train, y-test = train-test-split

(x, y, test_size=0.25, random_state=0)

X-train

#(visualize the 3D Graph)

import plotly-express as px

fig = px.Scatter_3d (df, x='x1', y='x2', z='x1+x2',

color='y')

.fig.show()

Classifier = SVC (kernel="linear")

Classifier.fit (X-train, y-train)

y-pred = Classifier.predict (X-test)

accuracy-score (y-test, y-pred)

Output: 1.0

It denotes 100% accuracy with Correct linear

Split.

Principal Component Analysis: [PCA]

Dimensionality Reduction

Why we use PCA?

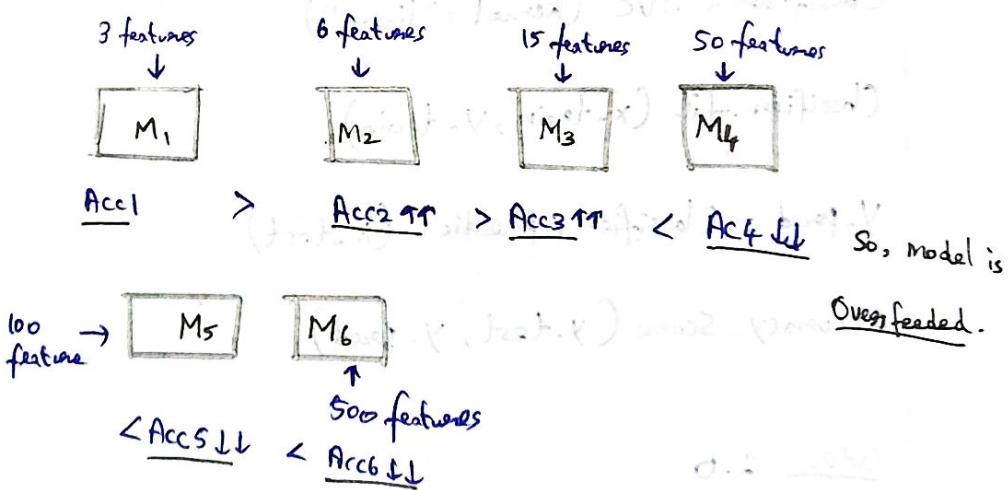
Some reasons are used to perform PCA.

①. Curse of Dimensionality:

Suppose we have 500 features in the dataset.

We want to train ML models with these features,

Suppose,



Note that, The increase in features means. Some features are not important at all. So, the accuracy will decrease;

The features would be,

- House Size
- No of bedrooms.
- No of bathrooms. like that.

② Model Performance Degrade:

Suppose Consider a human being, The person thinks about not only the prices of the room apartment and also thinks about the facilities which comes along with the room.

feature 1

Loc A

Price

2 Bhk \leftarrow 450k - 500k

3 Bhk \leftarrow 500k - 600k

Near Beach $\leftarrow \uparrow\uparrow\uparrow$

Near to celebrity \leftarrow TPP

house

Grocery shop \leftarrow Its is not more important but machine also train this.

School \leftarrow TPP

This says curse of dimensionality.

Sometime, The domain expert also Confused to give

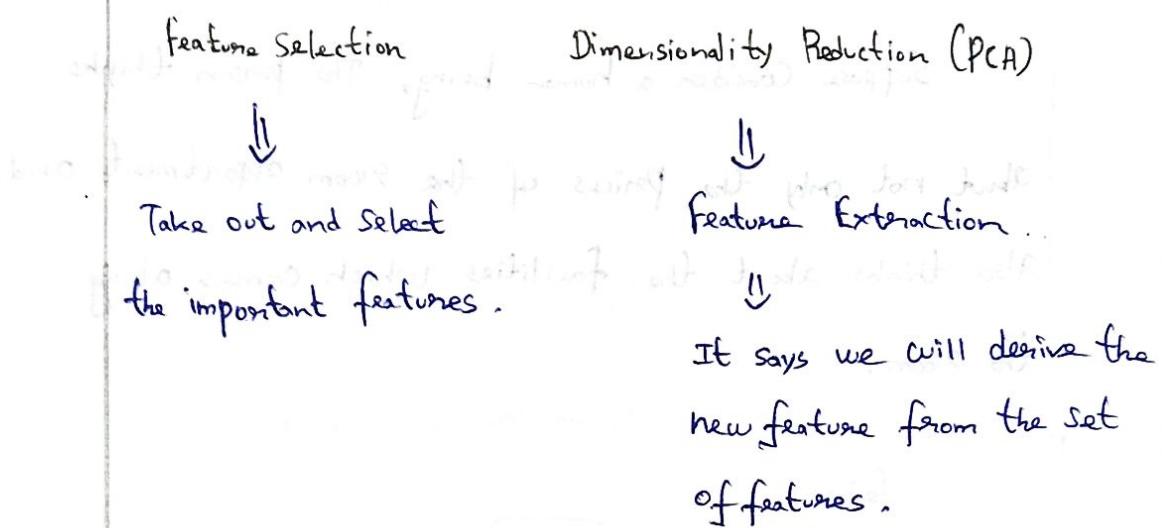
accurate data. This will shows the many features are feeded into the machines. So, the model is Overfitted.

How to remove and prevent this Curse of dimensionality?

Two different ways are,

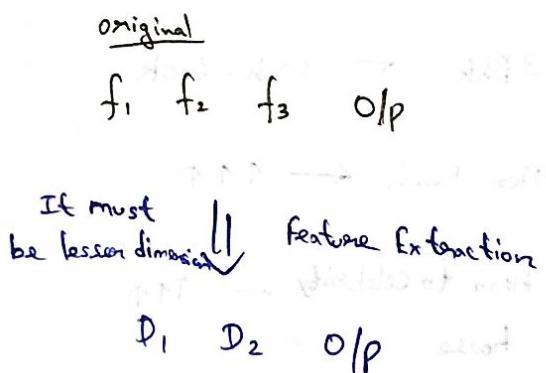
→ Feature Selection.

→ Feature Dimensionality Reduction (PCA).



The feature extraction would

be an,



This is all about the feature engineering on PCA.

We will see about the Feature Selection and

Feature Extraction.

End of lesson 12 (PCA) of machine learning

Feature Selection Vs Feature Extraction:

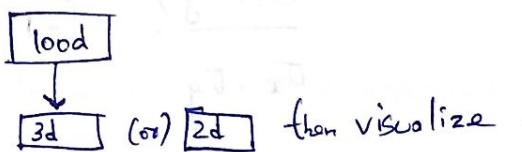
Why Dimensionality Reduction?

→ Prevent → Curse of dimensionality.

→ Improve the performance of model.

→ Visualise the data → Understand the data.

3d 2d visualization.



Feature Selection:

feature selection is the technique to select the most important feature from set of features.

Consider,

I/p O/p

x y

- -

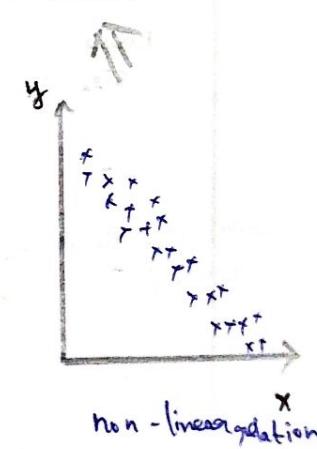
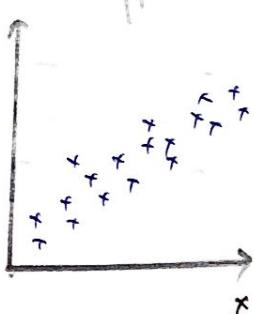
- -

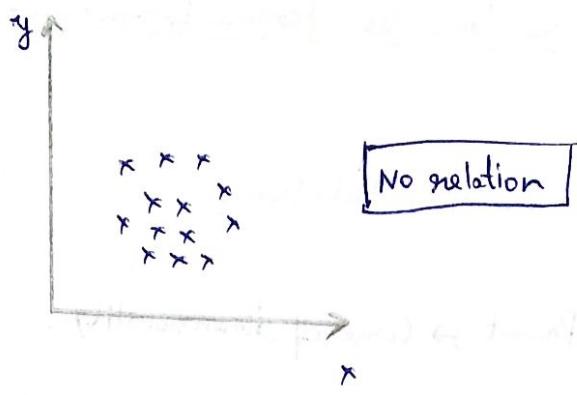
- -

- -

$x \uparrow y \uparrow$
 $x \downarrow y \downarrow$

$x \downarrow y \uparrow$
 $x \uparrow y \downarrow$





$$\text{Cov}(x, y) = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{N-1} = \text{+ve (or) 0 (or) -ve}$$

Pearson Correlation = $\frac{\text{Cov}(x, y)}{\sigma_x \cdot \sigma_y} = [-1 \text{ to } 1]$

More towards +1 the feature x and y are more relationship.

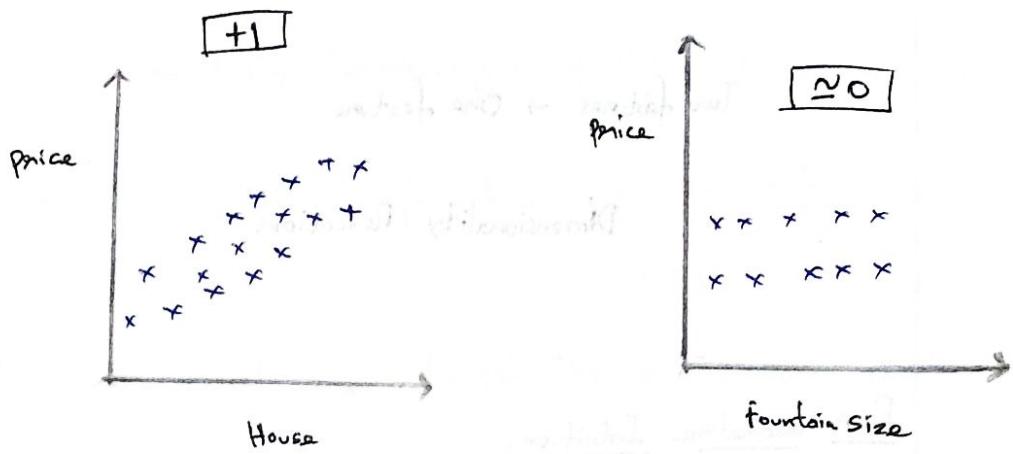
More towards -1 the feature x and y are no relationship.

Is equal to 0. The two features are no-correlated.

Example,

Dataset housing

House size	fountain size	price
100	100	100
200	200	200
300	300	300
400	400	400
500	500	500



\Downarrow This is linear and more correlated.
 \Downarrow This is not an linear and correlated.

Therefore, The house size is very important feature. Select that house size and neglect the fountain size for an analysis.

Feature Extraction:

Feature Extraction is the technique to extract the new feature from Set of features.

Consider that,

Rooms Size	No of Rooms	Price
-	-	-
-	-	-
-	-	-

\Downarrow *Principal Transformation applied to extract new feature*

House Size	Price
-	-
-	-

Two features \rightarrow One feature

Dimensionality Reduction

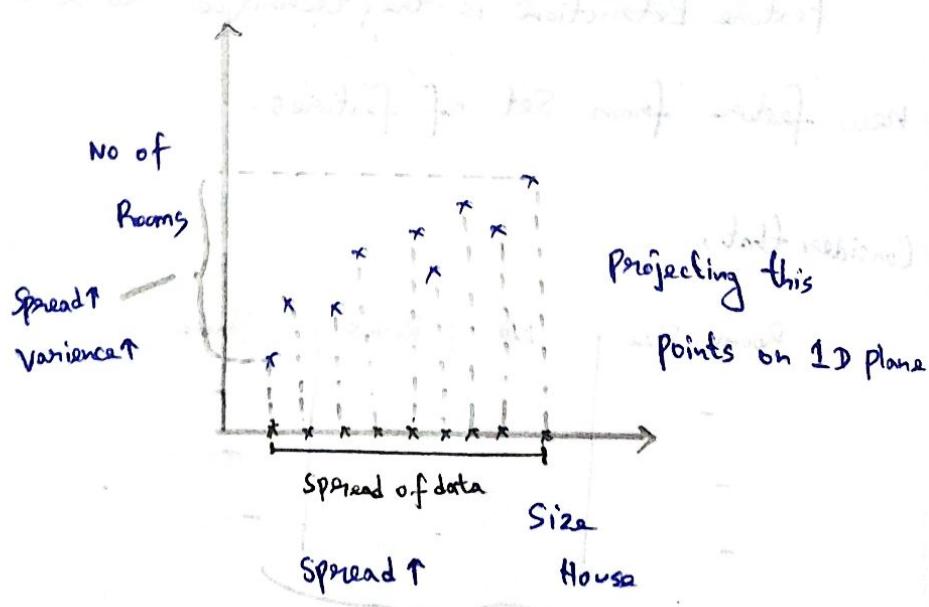
PCA: Geometric Intuition:

Consider the housing dataset,

House	Size	No of Room	Price (o/p)	PCA, 2 dimensions \rightarrow 1 dimensions
-	-	-	-	-
-	-	-	-	-
-	-	-	-	-
-	-	-	-	-

Obviously, we know the size house and no of room are the linear relationship.

So, plot the points,



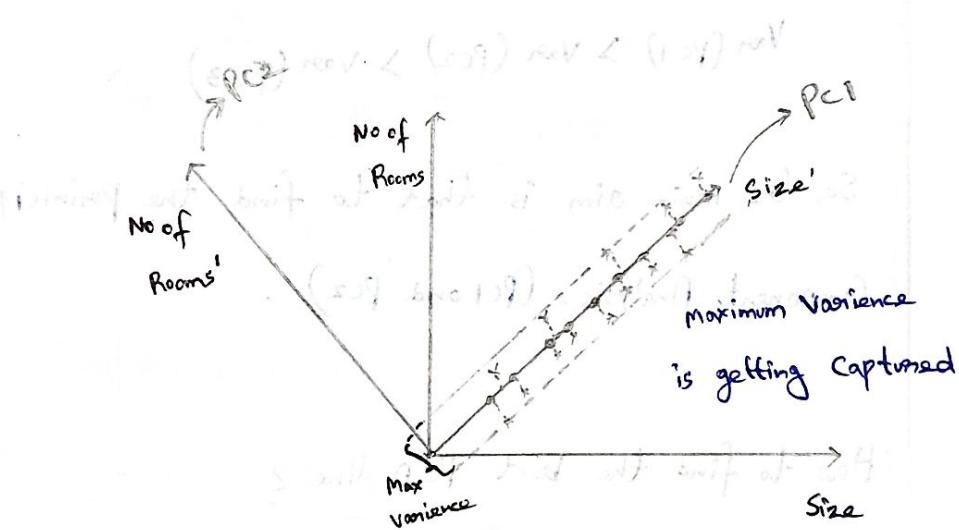
So, the main disadvantages of this approach is that
the No of rooms information will be lost.

So, the model unable to predict and perform well.

→ Loss of information (No of Rooms)

Therefore, the different approach was used in PCA.

That different approach is nothing but the Some transformation
is applied to Capture maximum Variance of data.



How the axis line Created?

It is Created with the help of Some
transformation is nothing but Eigen decomposition on
Matrix.

So, in this approach we tried $2D \rightarrow 1D$.

Therefore much information is not lost.

2 Dimensions

PC₁, PC₂

$$\text{Var}(\text{PC}_1) > \text{Var}(\text{PC}_2)$$

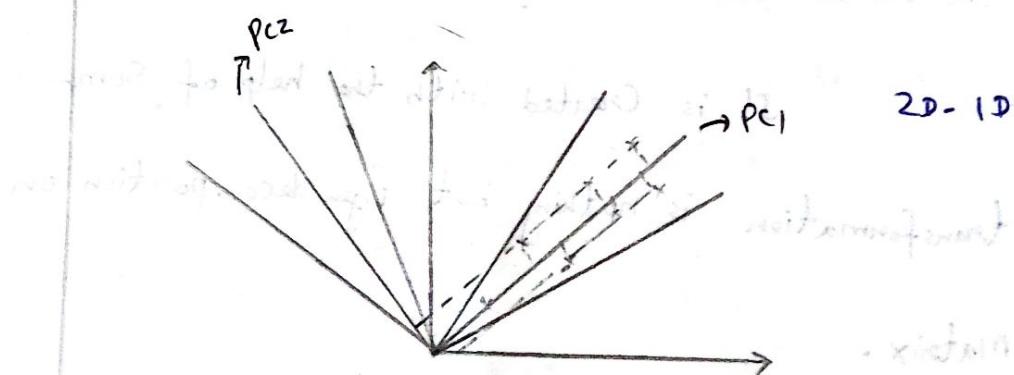
3 dimensions

PC₁, PC₂, PC₃

$$\text{Var}(\text{PC}_1) > \text{Var}(\text{PC}_2) > \text{Var}(\text{PC}_3)$$

So, the main aim is that to find the Principal Component Analysis. (PC₁ and PC₂) .

How to find the best PCA line?



2 Best PCA

The finding of PCA is very simple to check the
Which Principal Component lines Captures more Variance of
data points.

That data PCA lines are selected.

To get the best Principal Component which Captures
Maximum Variance.

Suppose we have,

$$3D \rightarrow 1D$$



We know; $\underbrace{PC_1, PC_2, PC_3}_{\curvearrowright 1D}$

$$\text{Var}(PC_1) > \text{Var}(PC_2) > \text{Var}(PC_3)$$

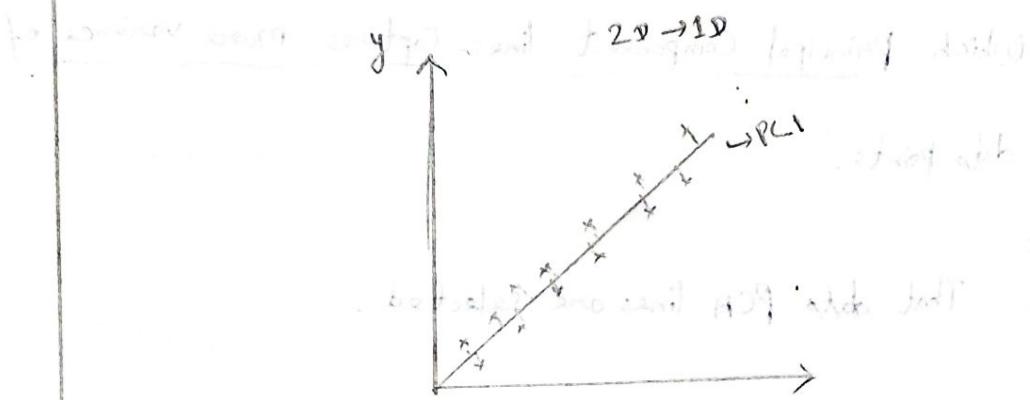
Suppose, $3D \rightarrow 2D$,

$$\underbrace{PC_1, PC_2, PC_3}_{\curvearrowright}$$

\hookrightarrow This two taken and convert to 2D.

Maths Intuition behind PCA algorithm.

at least at square root of λ to get half add.

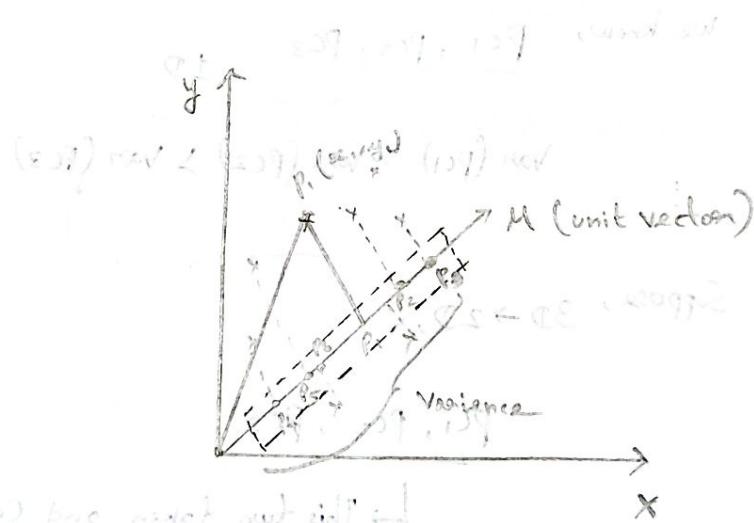


and also having following two steps of
The two important things are,

1 - Projection.

2. Cost function related to Variance.

Consider that,



Q.E.D. I found this point and will go!

Suppose I want to project $p_i(x_i, y_i)$ on unit vector,

and trying to build gradient descent

$$\text{Proj } p_i \text{, } M = \underline{p_i \cdot M} \quad \text{where,}$$

$$\|u\| \quad \|u\| = 1$$

Therefore, $\boxed{\text{Proj } p_i, u = p_i \cdot u}$ \Rightarrow Give Scalar value.

\downarrow
 p_i' refers to only
Scalar refers to only

So, In our case we want to find Magnitude.

the variance. (i.e) variance which

comes under the distance. Therefore

We say magnitude. (Scalar).

So, Computing every points we need to project. It looks

like,

$\boxed{p_0', p_1', p_2', p_3', p_4', \dots, p_n'}$ \rightarrow Scalar value
 \downarrow
Variance (distance)
spread.

The projecting point on the Unit vector we say p_i' .

Let's take,

$\boxed{p_0', p_1', p_2', p_3', p_4', \dots, p_n'}$

We use different notation as,

$x_0, x_1, x_2, x_3, x_4, \dots, x_n$

Goal:

Find the best Unit vector which captures maximum Variance.

$$\text{So, Max Variance} = \sum_{i=1}^n \frac{(x_i - \bar{x})^2}{n-1}$$

Want to minimize \downarrow

So, this is the Cost function of PCA.

Then, we cannot obviously try to find which drift

Vector line is capturing maximum variance.

In order to find that, we something used called as Eigen Values and Eigen Vectors.

Eigen Vectors and Eigen Values:

To calculate this,

→ Firstly, Covariance matrix between features need to be find.

→ Then the eigen vectors and eigen values will found (or) find out from this Covariance Matrix.

→ Eigen Vector → Eigen Value → It shows the Magnitude of eigen Vector. Using this to capture the maximum variance.

The Second Point shows some mathematical equation to find eigen value and eigen vector. Something are the,

$$AV = \lambda V$$

→ It is nothing but linear transformation of matrix.

Using these, to find out the maximum variance of the datapoints captured.

Eigen Vectors and Eigen Values: [Linear Transformation]

[Eigen decomposition of Covariance matrix for finding which is Matrix]



Eigen Vector and Eigen Values.

$$\begin{bmatrix} v & \lambda \end{bmatrix}$$

$(v, \lambda) = (x, x) \rightarrow (x, x) \cdot v$ The linear transformation would

$$\begin{bmatrix} v \\ \lambda \end{bmatrix} \rightarrow \text{vector}$$



$$[] * [v] = \lambda * v$$

See the difference in
some website..

Eigen Value

So,

$$A * V = \lambda * V$$

↓ So this give

Eigen Vector

which gives

Maximum Magnitude

↳ from this we get principal components

and get the Maximum Variance Captured.

The graphical visualization of this available in,
<https://shad.io/MatVis/>

Steps to calculate Eigen value and Vectors:

1. Firstly, Covariance of features.

Consider,

$$\begin{matrix} x & y \\ \downarrow & \downarrow \\ x & y \end{matrix} \quad z \quad \xrightarrow{\text{dependent feature}}$$

$$\text{Cov}(x, y) = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{N-1}$$

Suppose the two feature, represented as 2×2 ,

	x	y
x	Vaar(x)	Cov(x, y)
y	Cov(y, x)	Vaar(y)

why Vaar(x), because $\text{Cov}(x, x)$,

$$\text{Cov}(x, x) = \frac{\sum_{i=1}^n (x_i - \bar{x})(x_i - \bar{x})}{N-1}$$

$$= \frac{\sum_{i=1}^n x_i^2 - \bar{x}\bar{x}}{N-1}$$

$$= \frac{\sum_{i=1}^n (x_i - \bar{x})^2}{N-1} \Rightarrow \text{Vaar}(x)$$

That's why we put $\text{Cov}(x, x)$ and $\text{Cov}(y, y)$ are mentioned at $\text{Var}(x)$ and $\text{Var}(y)$.

Suppose, the same as 3×3 matrix

	x	y	z
x	$\text{Var}(x)$	$\text{Cov}(x, y)$	$\text{Cov}(x, z)$
y	$\text{Cov}(y, x)$	$\text{Var}(y)$	$\text{Cov}(y, z)$
z	$\text{Cov}(z, x)$	$\text{Cov}(z, y)$	$\text{Var}(z)$

Then Consider,

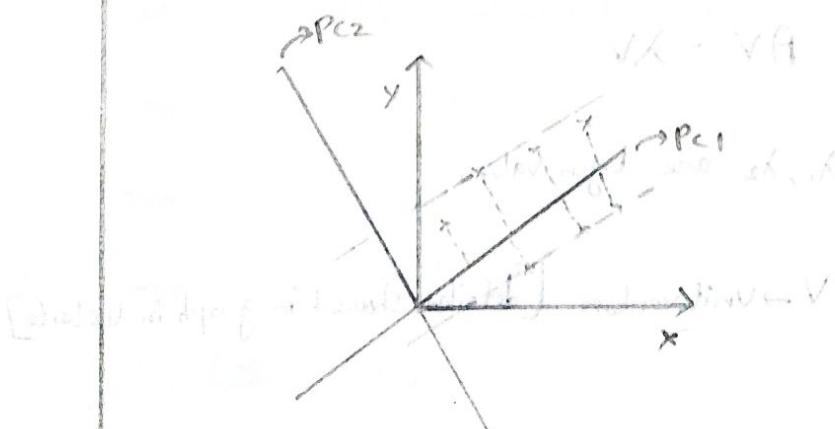
$$A = \begin{bmatrix} \text{Var}(x) & \text{Cov}(x, y) \\ \text{Cov}(y, x) & \text{Var}(y) \end{bmatrix}$$

Also $A \cdot V = \lambda \cdot V$

Where λ denotes λ_1, λ_2 i.e. (f_1, f_2)

$$\begin{matrix} \downarrow & \downarrow \\ p_{c1} & p_{c2} \end{matrix}$$

Then the $\boxed{\lambda_1, \lambda_2}$ are the Eigen values.



Overall Steps for Calculations:

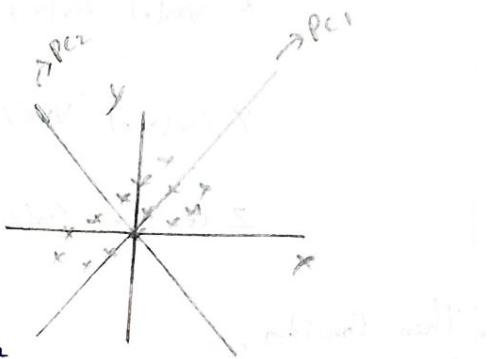
1) Change dimension



2) Standardize the data.

Once, Standardization is applied. It comes to the

Center.



3) Covariance matrix X and Y

$$\begin{matrix} X & Y \\ X & \text{Var}(x) \end{matrix} \quad \boxed{\sqrt{6} = \sqrt{A_1}, \sqrt{3}}$$

$$A = \begin{pmatrix} \text{Var}(x) & \text{Cov}(x, y) \\ \text{Cov}(y, x) & \text{Var}(y) \end{pmatrix}$$

+ standard deviation of each variable

4) Find out Eigen vectors and Value

$$AV = \lambda V$$

λ_1, λ_2 are Eigen values.

$V \rightarrow$ Unit vector [detailed showed in graph in website]

λ_1, λ_2

$\lambda_1 - PC_1$

$\lambda_2 - PC_2$

↳ First PC has maximum variance.

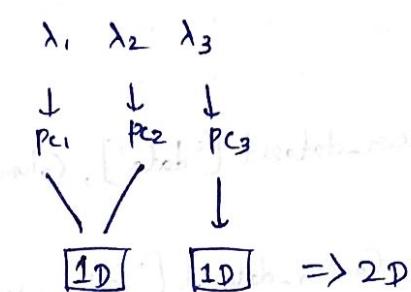
λ_1 denotes the magnitude of eigen vector.

(Observe) \downarrow In this way, we capture maximum variance.

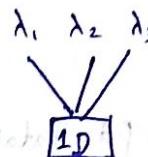
Capture maximum Variance.

Suppose,

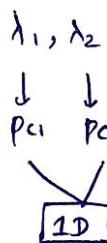
Change $3D \rightarrow 2D$



Change $3D \rightarrow 1D$



Change $2D \rightarrow 1D$



Likewise, we change the dimension on PCA.

PRACTICAL IMPLEMENTATION ON PCA:

```
import matplotlib.pyplot as plt
```

```
import Seaborn as sns
```

```
import numpy as np
```

```
import Pandas as pd
```

```
%matplotlib inline
```

```
# Load the dataset
```

```
from sklearn.datasets import load_breast_cancer
```

```
Cancer_dataset = load_breast_cancer()
```

```
Cancer_dataset.keys()
```

```
print(Cancer_dataset.DESCR)
```

```
df = pd.DataFrame(Cancer_dataset['data'], columns =
```

```
Cancer_dataset['feature_names'])
```

```
df.head()
```

```
# Feature Extraction (PCA)
```

```
# Standardization
```

```
from sklearn.preprocessing import StandardScaler
```

```
Scalar = StandardScaler()
```

```
Scalar.fit(df)
```

```
Scaled_data = scalar.transform(df)
```

```
Scaled_data
```

Applying PCA algorithms

```
from sklearn.decomposition import PCA
```

`Pca = PCA(n_components=2)` # n-components are the change dimension need.

data_pca = pca.fit_transform(scaled_data)

data_pca

Output: array ([[9.128 , 1.9458]])
Change to two
Columns

PCA - Explained Variance

Output: array ([13.3049907, 5.7013746])

\downarrow \downarrow
 feature 1 feature 2
 Captured Variance of data Captured Variance of data

fit PCA without n_components, it gives more feature and explained variance for each and every feature.

Finally, plot the figure

Plt. figure (figsize = (8,6))

```
plt.scatter(data_pca[:,0], data_pca[:,1], c=cancer_dataset['target'], cmap='plasma')
```

plt.xlabel ('First Principal Component')

plt.ylabel ('Second Principal Component')

Practical implementation of k-means on kerish naik Day-06

[This is very long code].

Definition of Bias and Variance:

Suppose, Training Dataset = 90% Acc }
Test Dataset = 70% Acc } \Rightarrow Overfitting.

!!

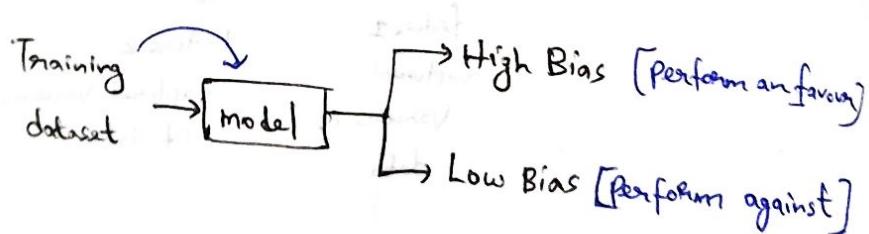
Bias:

Low Bias

High Variance

It is a phenomenon that skews the result of an algorithm in favour or against an idea.

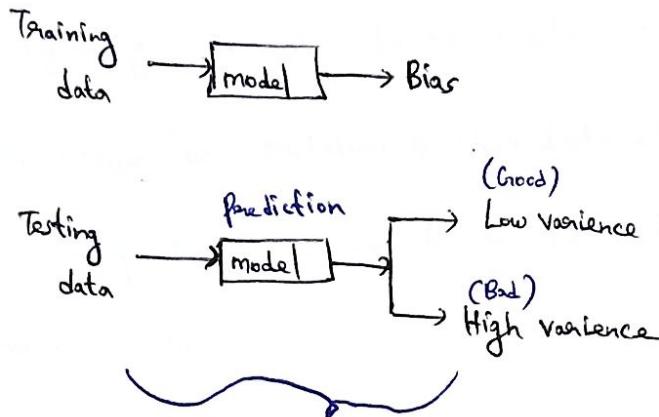
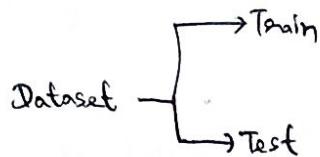
↳ Training dataset.



Variance:

Variance refers to the changes in the model when

Using different portions of the training or test data.



This says about the Perfect definition.

Model 1

Training Acc = 90%

Test Acc = 75%



{
Low Bias
High Variance}

Overfitting

Model 2

Train Acc = 60%

Test Acc = 55%



{
High Bias
High Variance}

Underfitting

Model 3

Train Acc = 90%

Test Acc = 92%



{
Low Bias
Low Variance}

Generalised
Model