



Decision Tree and Support Vector Machine HandWritten Notes

- Darshan R M

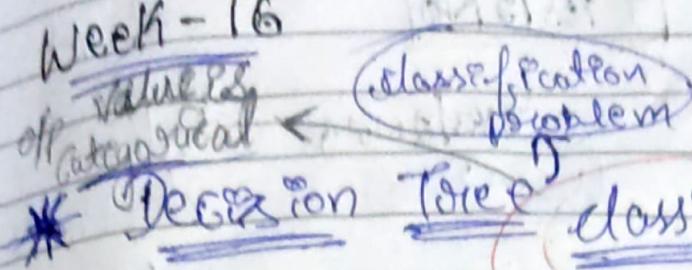
"Learn, Share,
and
Collaborate"

Week - 16

Decision Trees

and

Support Vector Machines



* Consider, Problems Perfect poison play Tennis (or) not.

Dataset

Dependent	Outlook			Temperature / Humidity / Wind		
	Day	Sunny	Overcast	Hot	Weak	No
1	Sunny			Hot	Weak	No
2	Overcast			Mild	Strong	Yes
3	Rain					

Independent

Play Tennis

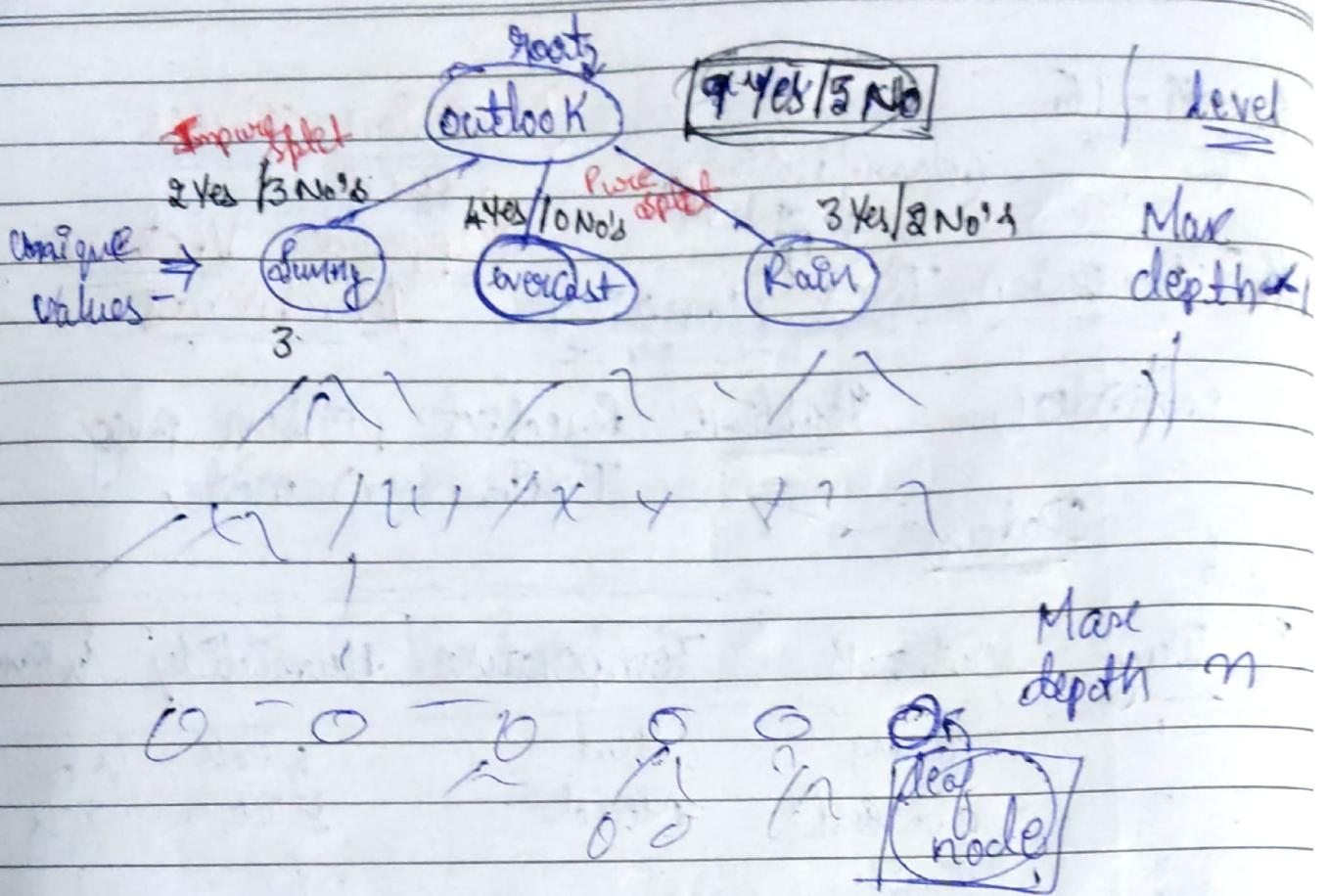
No
Yes

Categorical
Output

Decision Tree

- ① Selects root node.

↓
from available independent features.



After first ~~Pure~~ Max depth level:

⑪ Purity split check - we check

Pure Split (01)
Impure Split

No need to
shift further

we have to
split further

Ex :-
Number of
Yes/No = 0

e.g.
Both
Yes & No ≠ 0

then it will
become
leaf node.

* But in maths we use some techniques

B. Pure split / Impure

check

(i) Entropy

(ii) Gini Impurity

? Measure
of purity

Q(3) what feature we need to select to start the split?

→ Information gain

* Split happens till we reach leaf node.

① Purity check

In chemistry, entropy means

measure of randomness of system.

(i) Entropy

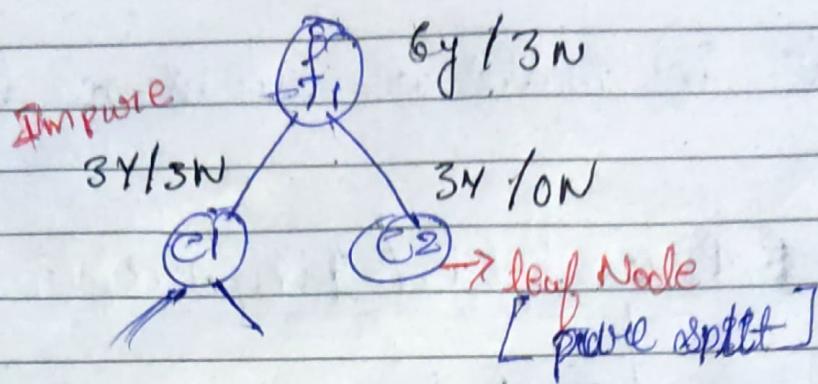
for binary

classification

$$H(S) = -(P_+) (\log_2 P_+) - (P_-) (\log_2 P_-)$$

P_+ → Probability of positive category
 P_- → Probability of negative category

Consider,



$$* H(C_1) = -(P_+)(\log_2 P_+) - (P_-)(\log_2 P_-)$$

* $P_+ = \frac{3}{6} = \frac{1}{2}$ ^{+ve values}

$$P_- = \frac{3}{6} = \frac{1}{2}$$

~~Value range below~~ $\Rightarrow H(C_1) = -(\frac{1}{2})(\log_2(\frac{1}{2})) - (\frac{1}{2})(\log_2(\frac{1}{2}))$

$H(C_1) = 1$ \leftarrow entropy

Value 1 indicates impure split.
as $P_+ = P_-$

$$* H(C_2) = -(P_+)(\log_2 P_+) - (P_-)(\log_2 P_-)$$

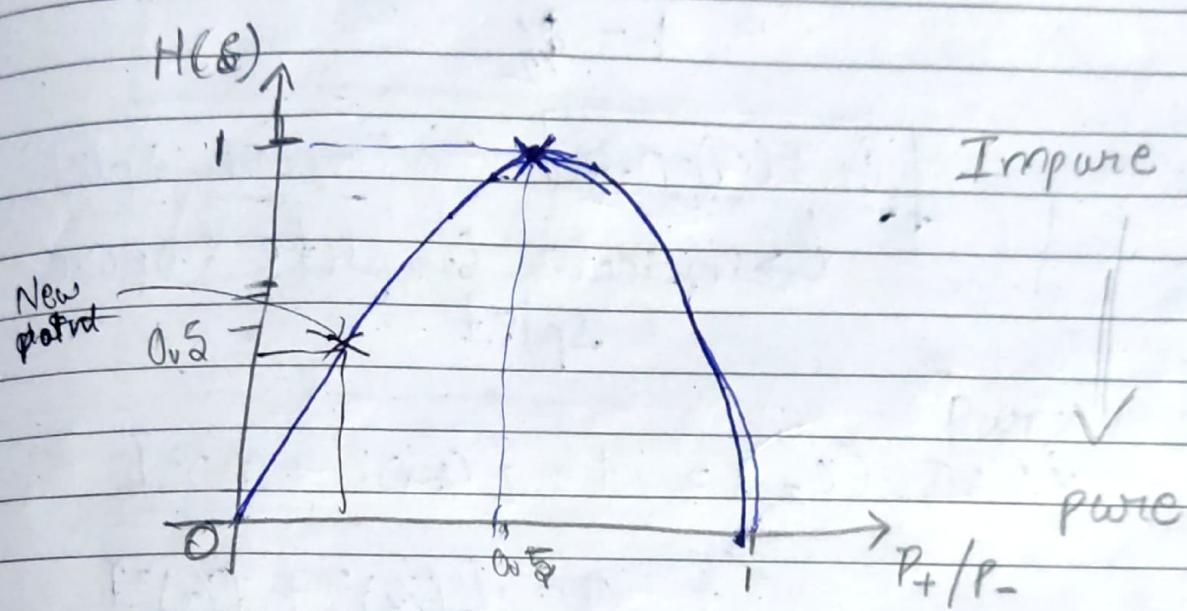
$$* P_+ = \frac{3}{8} = \frac{1}{2}$$

$$P_- = \frac{5}{8} = \frac{1}{2}$$

$$\Rightarrow H(c_2) = -(1)(\log_2 1) - (0)(\log_2 0)$$

$$[H(c_2) = 0]$$

↳ indicates pure split
 ↳ Completely pure split



(ii) Binary Impurity [B.I]

$$G_{BI} = 1 - \sum_{i=1}^n (p_i)^2$$

~~$$G_{BI} = 1 - \sum_{i=1}^n (p_i)^2$$~~

For Binary class

$$G_{BI} = 1 - [(p_+)^2 + (p_-)^2]$$

* Consider previously example

$$= \frac{8}{3} \cdot I(\varepsilon) = 1 - \left[(P_+)^2 + (P_-)^2 \right]$$

$$= 1 - \left[\left(\frac{3}{6} \right)^2 + \left(\frac{3}{6} \right)^2 \right]$$

$$= 1 - \frac{1}{2}$$

en förtiden

~~Value
range b/w
P to 0.5~~

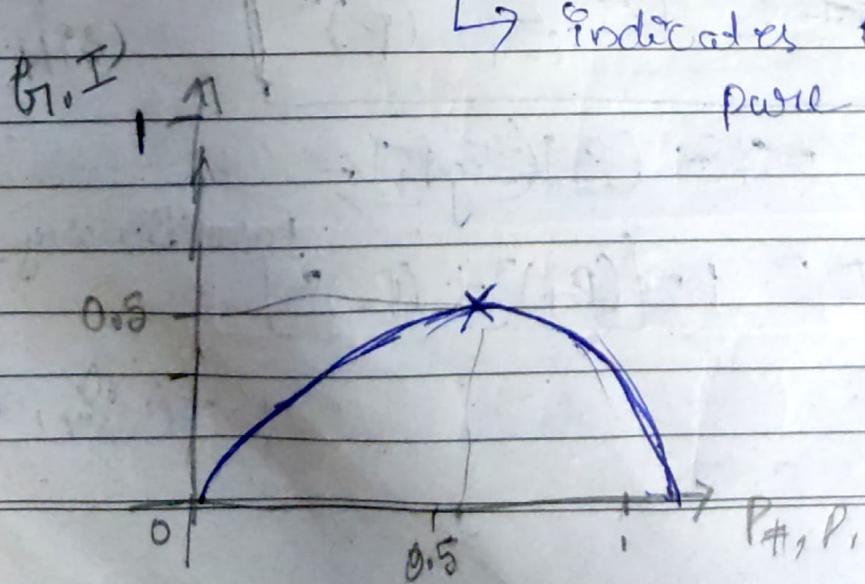
$G_1, I(G) = 0.5$ \Rightarrow impure split
↳ indicates complete impure split

$$\text{G.I. } (C_2) = 1 - [(P+)^2 + (P-)^2]$$

$$= 1 - [(53)^2 + (05)^2]$$

$$= 1 - 1$$

BIT(C2)=0 \Rightarrow pure split
↳ indicates complete pure split



Multiclass classification problem

consider, w3 categories in o/p

$$H(S) = -(P_{C_1})(\log_2 P_{C_1}) - (P_{C_2})(\log_2 P_{C_2}) - (P_{C_3})(\log_2 P_{C_3})$$

$$G.I = 1 - [(P_{C_1})^2 + (P_{C_2})^2 + (P_{C_3})^2]$$

③ Information gain :

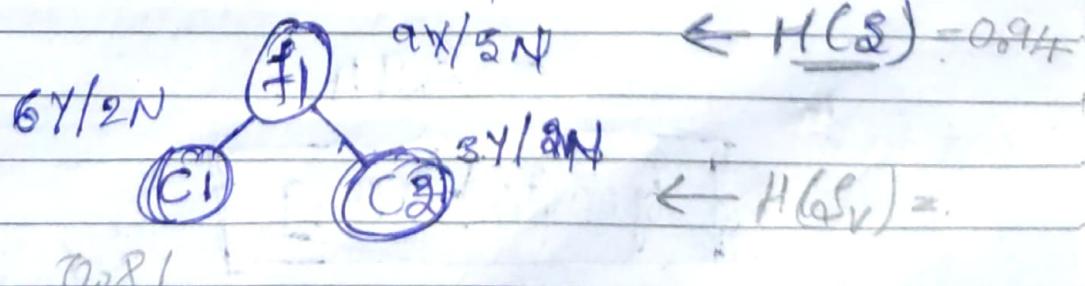
→ helps to select which feature to select to start the split.

$$\text{Gain}(S, f_i) = H(S) - \sum_{V \in \text{val}} \frac{|S_V|}{|S|} H(S_V)$$

Entropy
parent
node

f_1	f_2	f_3	O/P
$6Y/2N$			

entropy of
child node



~~for Root Node~~

$$* \text{Root Node, } H(S) = -(P_+)(\log_2 P_+) - (P_-)(\log_2 P_-)$$

$$= -\frac{9}{14} (\log_2 \frac{9}{14}) - \frac{5}{14} \log_2 \left(\frac{5}{14}\right)$$

$$\approx \underline{\underline{0.94}}$$

$$* H(C_1) = -\frac{6}{8} (\log_2 \frac{6}{8}) - \frac{2}{8} (\log_2 \frac{2}{8})$$

$$\approx \underline{\underline{0.81}}$$

$$* H(C_2) = -\frac{3}{6} (\log_2 \frac{3}{6}) - \frac{3}{6} (\log_2 \frac{3}{6})$$

$$H(C_2) \approx 1$$

for first feature

$$\text{gain}(S, f_1) = H(S) - \sum_{v \in V(f_1)} \frac{|S_v|}{|S|} H(S_v)$$

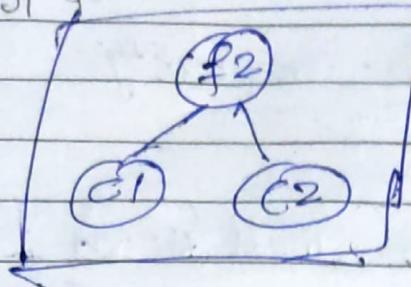
$$= 0.94 - \left[\frac{8}{14} \cdot (0.81) + \frac{6}{14} \cdot (1) \right]$$

$$= 0.94 - [0.606 + 0.428] = 0.94 - 1.034 = -0.094$$

for first split $\Rightarrow \boxed{\text{gain}(S, f_1) = 0.049}$

first split

\$0.51

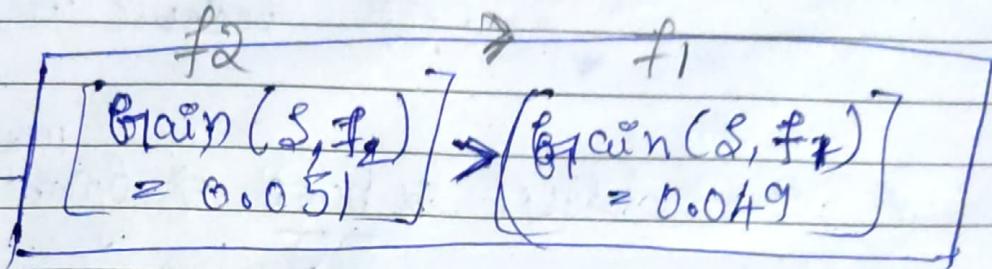


Considering

$$G(f_2)$$

Information gain = 0.051

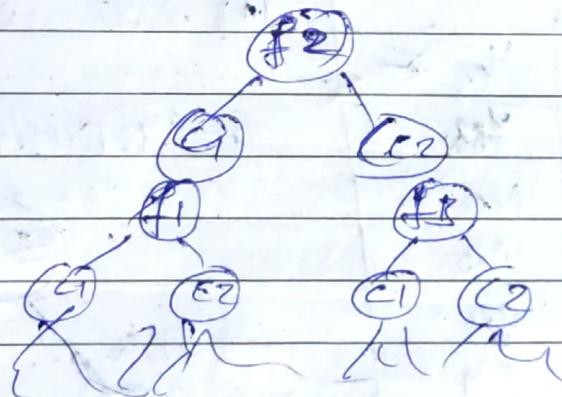
* Here,



Conclusion:

we need to start splitting by using f_2 feature.

verando



d. classification measure

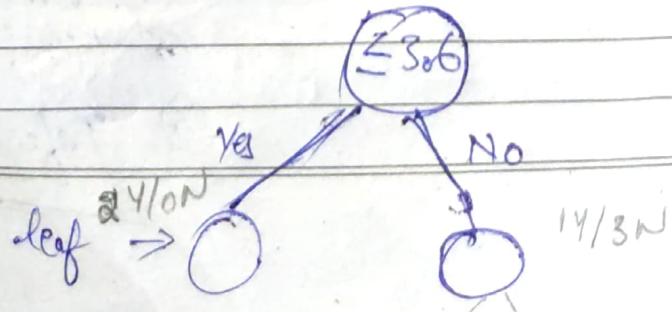
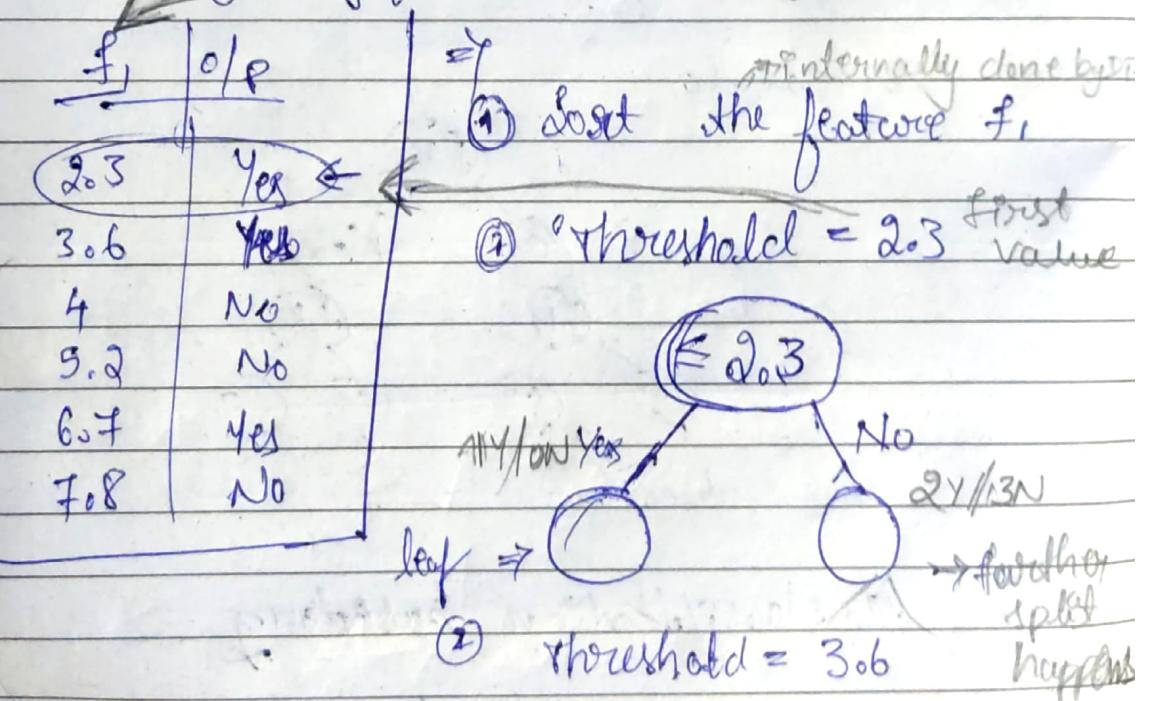
- ① Entropy
- ② Gini Impurity
- ③ Information gain

Q) Entropy v/s gini - impurity

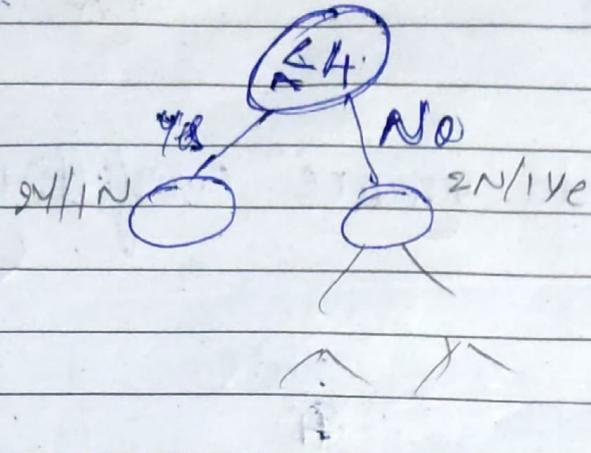
* when dataset is small (10,000) \Rightarrow Entropy
 as we use "log" because Time \downarrow . \uparrow use

* when dataset is huge \Rightarrow Gini impurity
 as we use "Simple maths" \downarrow use

Q). what if my features are in. Continuous?



③ threshold = 4



and Information gain is calculated which ever value ~~feature~~ has the highest gain that will be selected then further splitting takes place.

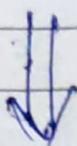
Disadvantage:

① Time Complexity ↑↑
when dataset is huge.

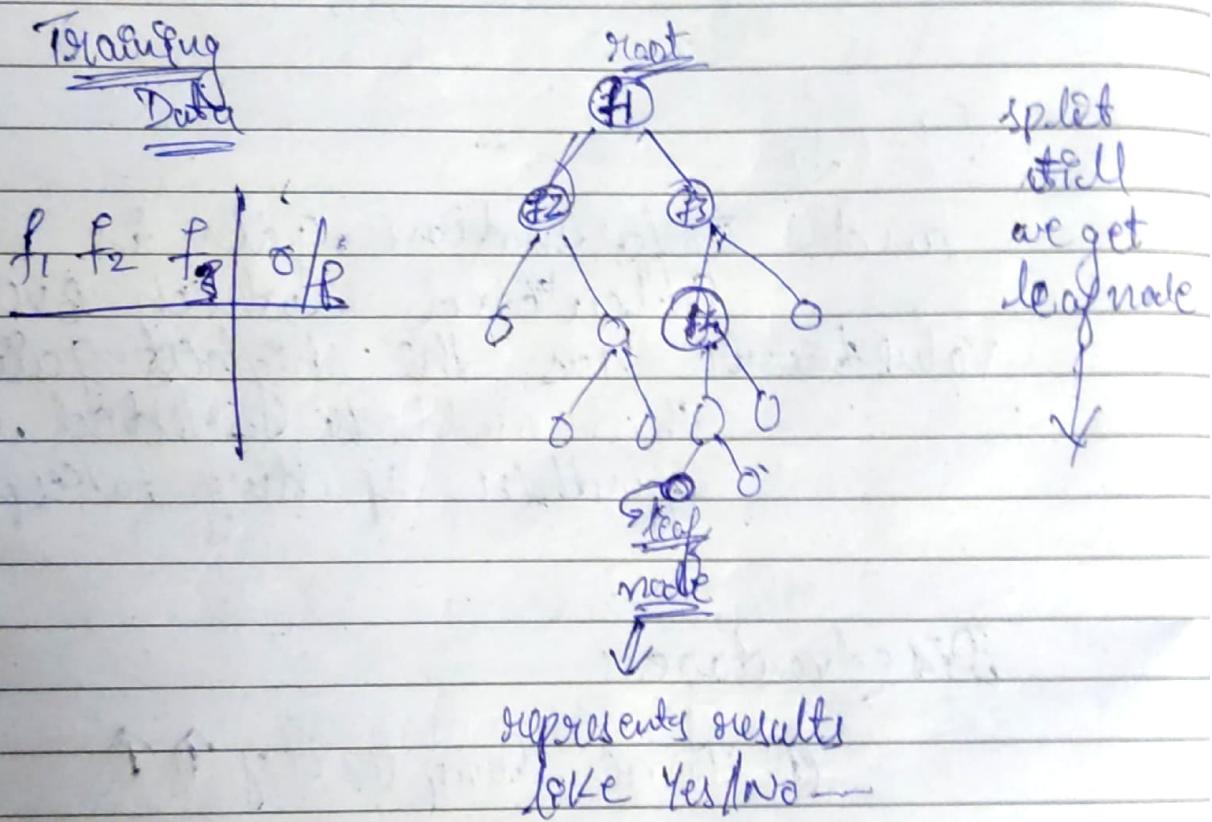
solution

we have some parameters
in SKLearn we can
use them .

* Post-pruning and pre-pruning in Decision Tree



Used to "reduce overfitting".

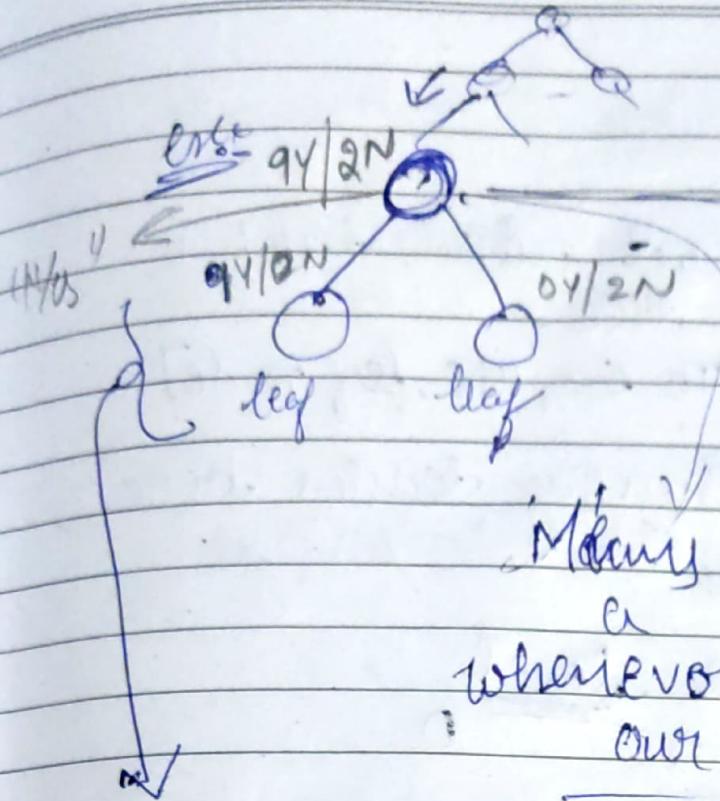


~~overfitting~~

whenever, we split decision tree till depth-1 leafnode it leads to "overfitting"

in case of splitting till depth.

low bias & high variance
e.g. Training data \rightarrow Acc $\uparrow\uparrow$
Test data \rightarrow Acc $\downarrow\downarrow$



"Yes"
ratio also
high

Here, Max value
of o/p is Yes
and Min value
of o/p is No

Meaning we can create
a generic model say
whenever we reach this node
our o/p is "Yes".

Here we are further split storing
data
so Model is "by" heaving storing
data

Aim: To build "generic model"

low bias & low variance
[Train Acc] [Test Acc]

Solution

now to reduce overfitting

- ① Post pruning.
- ② Pre pruning..

Cutting

① Post pruning :

step

step ① constructs entire decision tree.

[To complete leaf node]

step ② Post Pruning decision tree.

② Pruning

level

0

Can cut,
prune

(*)

full
max depth = 3

[One of the
parameters
for cutting]

* Max depth = 4

* Suitable for smaller dataset.

parameters in AB learner

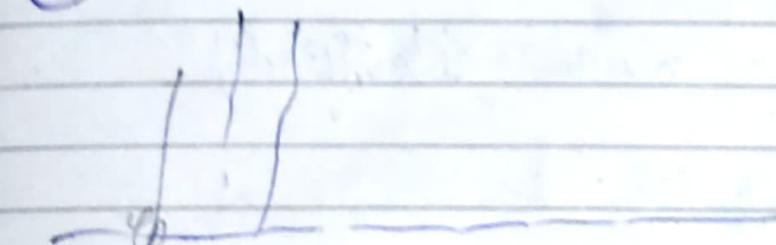
② Criterion : { "gini", "entropy", "log-loss" }

Cost functions of
logistic regression

③ Splitter : { 'best', 'random' }

Take max gain split
and split

④ max_depth



⑤ Preprocessing : Using GridSearchCV
and RandomizedSearchCV.

Step ① : Hyperparameter Tuning,
to select best parameters.

② Train model using best
parameters.

* Implementation using Python

① Segregate independent and dependent features

② Train Test split

Do not do standard scaling

③ Decision Tree

> from ~~sk~~learn.tree import
DecisionTreeClassifier

> classifier = DecisionTreeClassifier()

> classifier.fit()

→ can display
constructed tree

→ tree.plot_tree(classifier)

> Post-pruning!
[based on depth]

> classifier = DecisionTreeClassifier(max_depth=2)
from observation

④ Prediction

⑤ accuracy score, classification report

Preprocessing + Hyperparameter Tuning

① → parameters

> from sklearn.model_selection import GridSearchCV

> clf = GridSearchCV(classifier,
param_grid=parameters,
scoring='accuracy')

> clf.fit(Xtrain, ytrain)

> clf.best_params_

> clf.best_score_

> ypre =

>

↳ values are
continuous.

↳ values are
categorical.

* Decision Tree regression

- ~~Classification~~ DT classification
- ① Entropy. ↳ Measure of impurity
 - ② Gini Impurity ↳ impurity
 - ③ Information gain.

DT Regression

- ① Variance reduction.
- ② Variance.

* Dataset

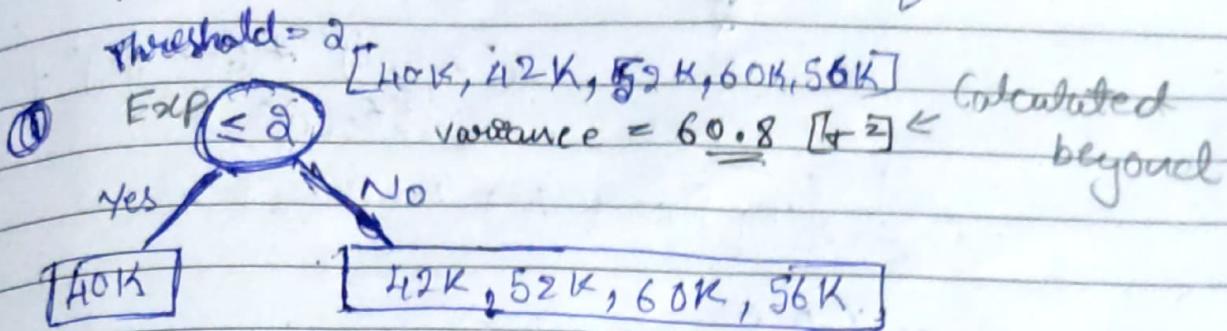
Exp	Career Gap	Regression	
		Salary	Continuous value
2	Yes	40K	
2.5	Yes	42K	
3	No	50K	
4	No	60K	
4.5	Yes	56K	

$\bar{y} = 50K$

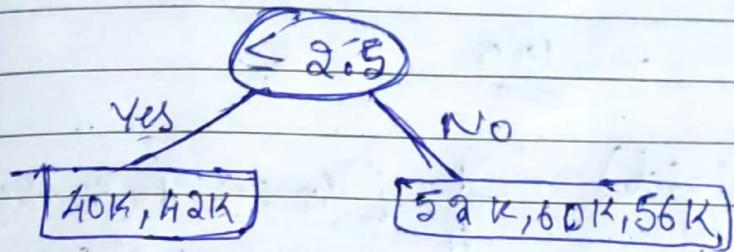
(average mean)

Now, we have to select the feature to split that feature can be done using 'Information gain'

* Consider, 'Exp' G.S. stat. feature continuous value



② Threshold = 2.5



③ Threshold = 3

reduction should be more

* Final aim: To calculate "Variance reduction".

$$\text{Variance (O.M)} = \frac{\sum_{i=1}^n [x_i - \bar{x}]^2}{n}$$

Now, calculate Variance for root, child

① ~~first threshold~~
2

$$(y_i - \bar{y}_P)^2 \quad (\sum y_i = 50)$$

$$\text{Variance (root)} = \frac{1}{5} \cdot [(40-50)^2 + (42-50)^2 + (52-50)^2 + (60-50)^2 + (56-50)^2]$$

$$= \frac{1}{5} [100 + 64 + 4 + 100 + 36]$$

$$\approx \underline{\underline{60.8}}$$

$$\text{Variance (left)} = \frac{1}{4} [(40-50)^2]$$

$$= \underline{\underline{100}}$$

$$\text{Variance (right)} = \frac{1}{4} [(42-50)^2 + (52-50)^2 + (60-50)^2 + (56-50)^2]$$

$$= \underline{\underline{51}}$$

imp

$$\boxed{\text{Variance reduction} = \text{Var}(\text{root}) - \sum w_j \text{Var}(\text{child})}$$

$$= 60.8 - \left[\frac{1}{5} \times 100 + \frac{4}{5} \times 51 \right]$$

$$\boxed{\text{Variance reduction} = 0}$$

② for threshold
2.5

$$\text{Variance}_{\text{cost}} = \underline{\underline{60.8}}$$

$$\text{variance}_{\text{left}} = \underline{\underline{82}}$$

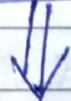
$$\text{variance}_{\text{right}} = \underline{\underline{16.66}}$$

$$\text{Variance reduction} = 60.8 - \left[\frac{2}{5} \times 82 + \frac{3}{5} \times 16.66 \right]$$

$$\underline{\underline{\text{Variance reduction}}} = 0.004$$

* Q) which split to ~~use~~ left

$$\nabla \quad \boxed{VR_{T=2} < VR_{T=2.5}}$$



∴ we will use

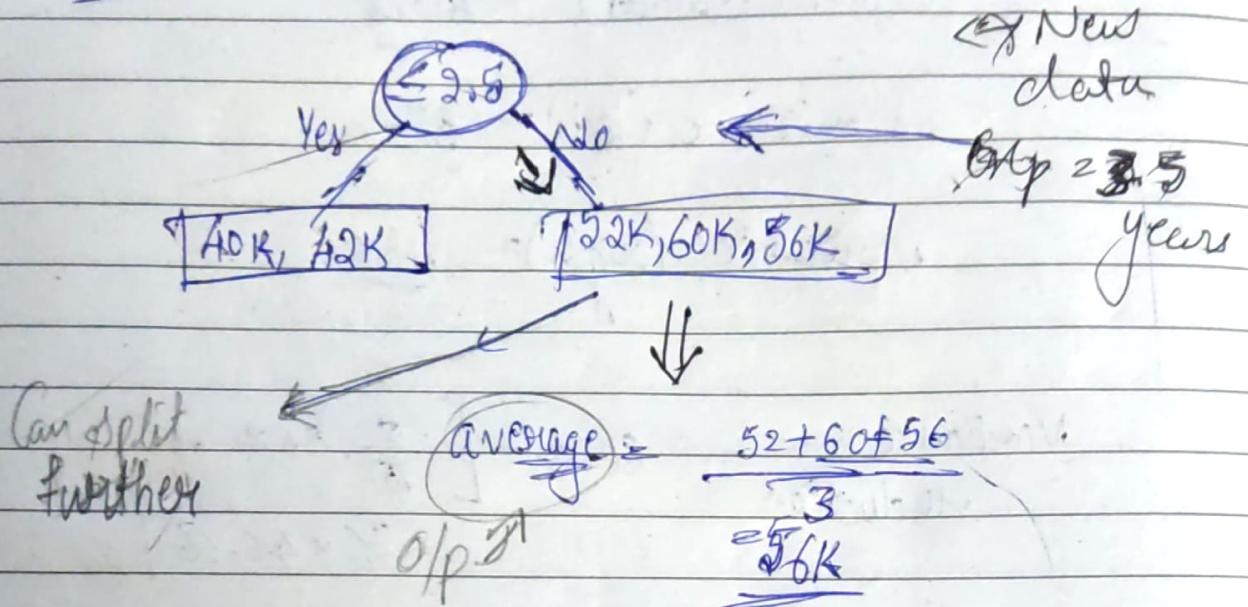
threshold = 2.5

split cost feature

cost

More
variance
reduction
happening

→ selected split:



* Implementation with Python

* only ~~'criterion'~~ parameter will change compared to classifier.

* { "squared_error", "friedman_mse",
* "absolute_error", "poisson" }

> df = df.sample(frac=0.25)

① segregate independent and dependent feature.

② train test split

③ Model training

> from sklearn.tree import

DecisionTreeRegressor

① Prediction

② Accuracy

> from sklearn.metrics import

RidgeScore

fit > RScore(y_true, y_pred)

(~~def~~)

Hyper parameter tuning

* include ignore warnings

① Best parameters selection.

> from sklearn.model_selection import

GridSearchCV

* In regression problem - we have
"scoring parameter different"

regressor = GridSearchCV(~~regressor~~,

param_grid = parameters,

scoring = "neg_mean_squared_error"

regressor

> fit(-, -)

tree.plot_tree(regressor, filled=True)

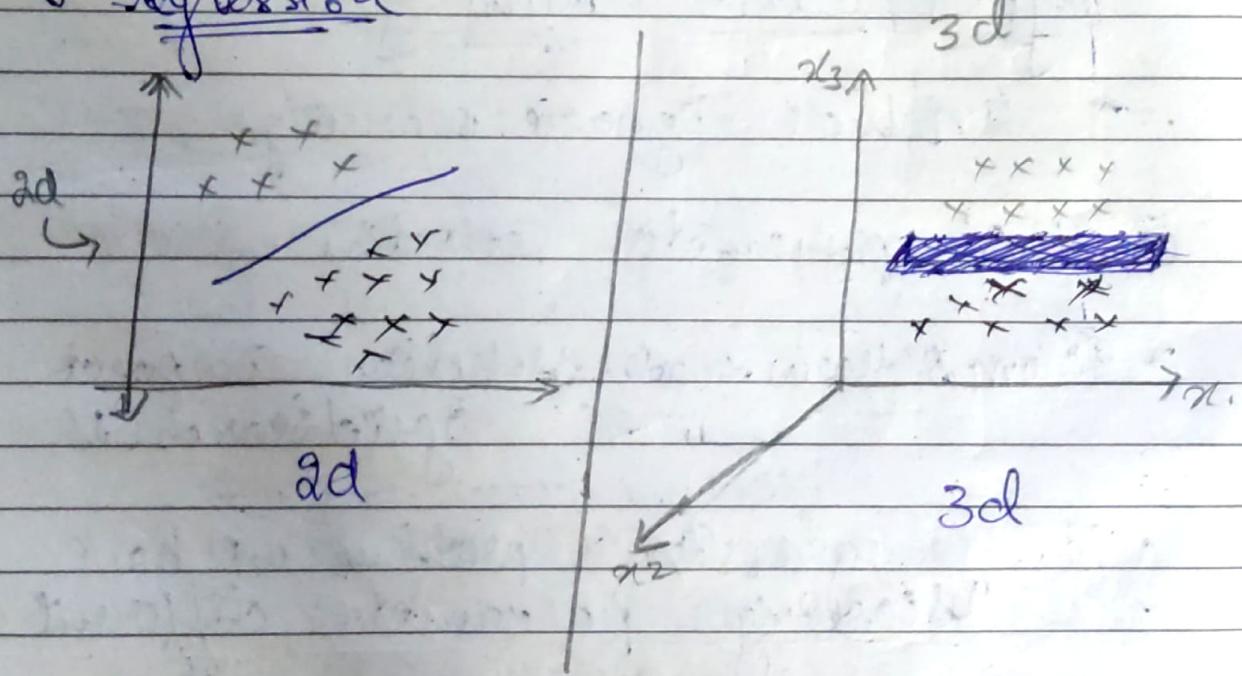
* Support vector classifier Machines

* Support vector Machines

2 Types:

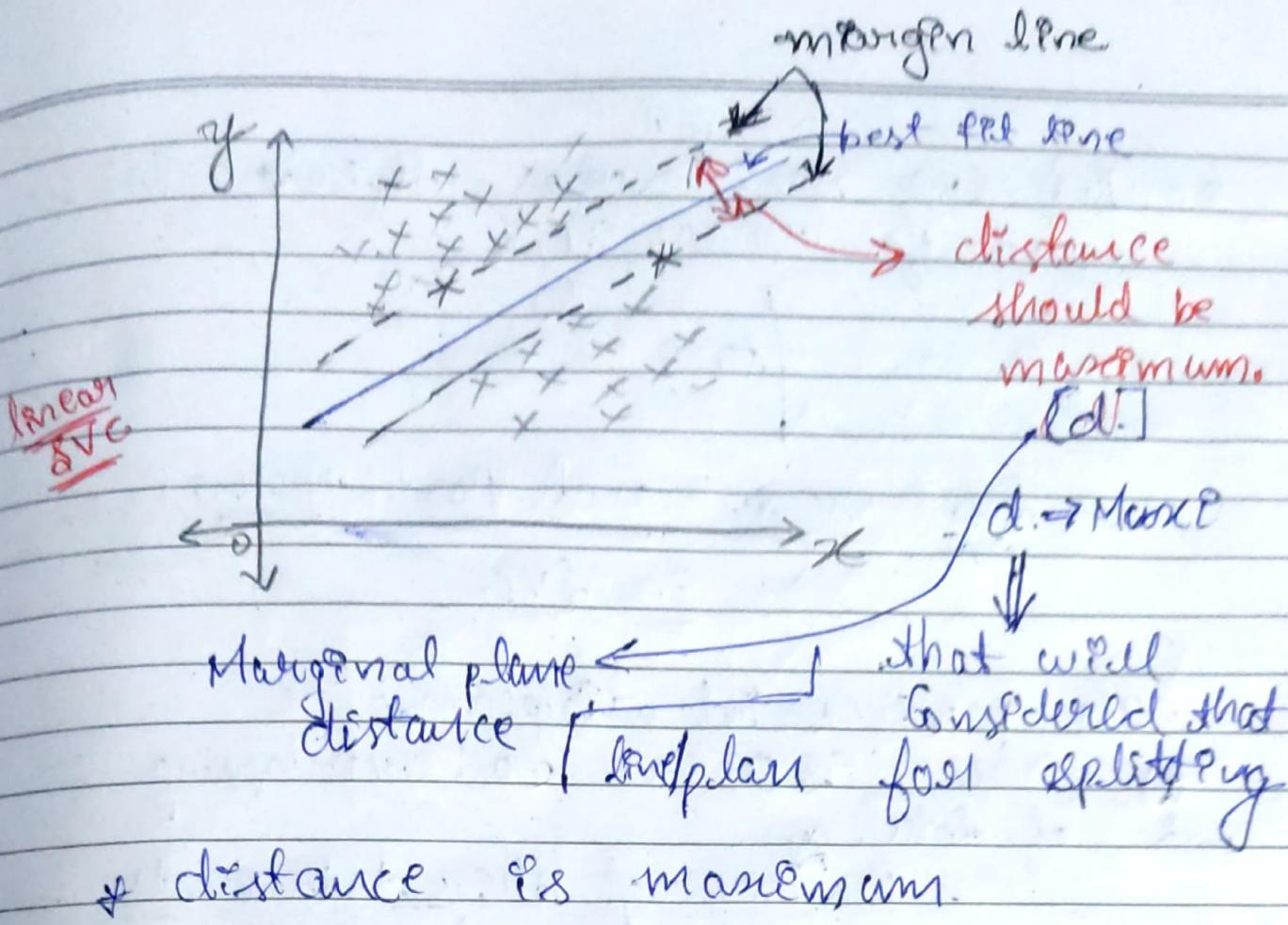
- (SVC) ① support vector classifier \rightarrow classification
(SVR) ② support vector regression \rightarrow Regression

* ~~logistic~~
~~regression~~



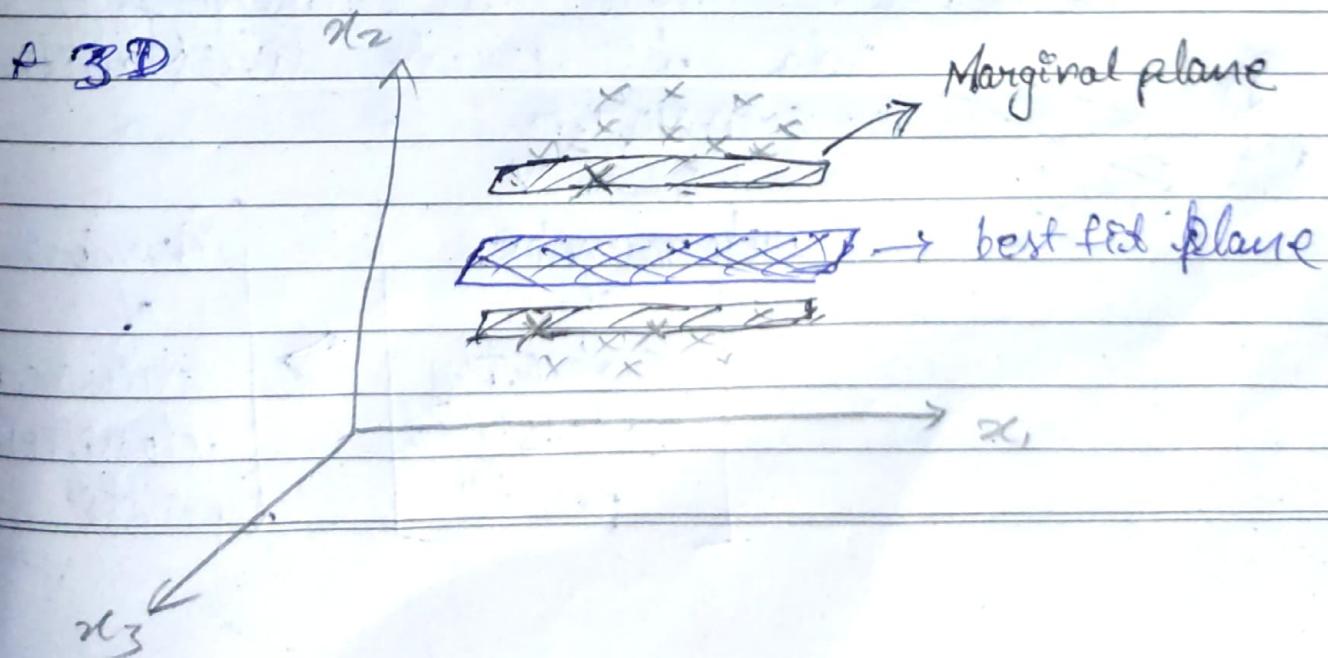
① Support vector classifier (SVC):

classification problem.

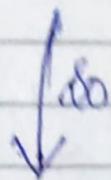


→ Aim: To create best fit line among which we create margin line. which passes through nearest points.

Called "support vectors"



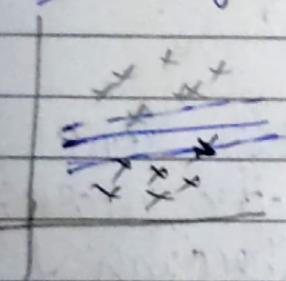
→ But in real world there will be lot of overlapping.



→ Soft margin and Hard margin in SVC

→ In previous Ad diagram, the ① occur if & no. overlapping of data.

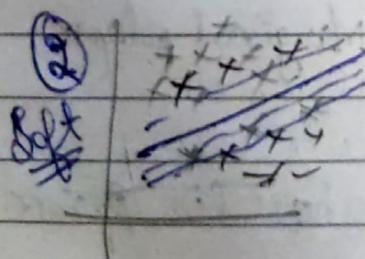
② None of the data points are misclassified!



we say at "Hard margin".

almost impossible in our data,

→ when there is Overlapping of data we create lines with some sloping

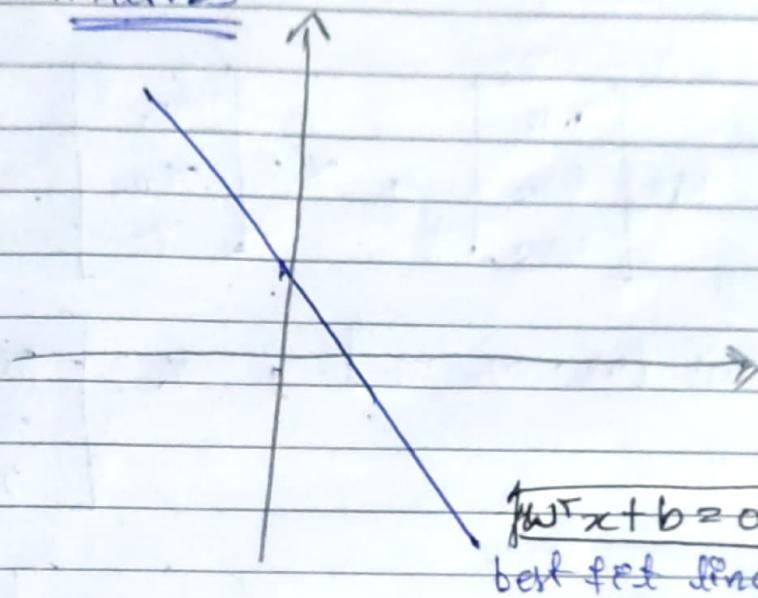


but number of correctly classified points	>	Number of incorrectly classified points
-------------------------------------------	---	-----------------------------------------

Soft margin \Leftrightarrow some data point are misclassified

we take [Error] as hyperparameter \rightarrow

SVC maths



$$w^T x + b = 0$$

best fit line

Equation of a straight line

$$y = mx + c$$

(or)

$$h_0(x) = \theta_0 + \theta_1 x$$

$$\Leftrightarrow ax + by + c = 0$$

$$by = -ax - c$$

$$y = \frac{-a}{b}x - \frac{c}{b}$$

[slope]

[intercept]

Multilinear regression

↓ 3 ip feature

$$y = \theta_0 + \theta_1 x_1 + \theta_2 x_2 + \theta_3 x_3$$

search for θ

$$y = b + [w_1 x_1 + w_2 x_2 + w_3 x_3]$$

$$\Rightarrow \text{vec} \begin{bmatrix} w_1 \\ w_2 \\ w_3 \end{bmatrix}, \quad x = \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix}$$

$$w^T = [w_1 \cdot w_2 \cdot w_3], \quad x = \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix}$$

Dot product, $w^T \cdot x$

$$= [w_1 \cdot w_2 \cdot w_3] \cdot \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix}$$

$$\underline{w^T \cdot x = w_1 x_1 + w_2 x_2 + w_3 x_3}$$

$$y = b^T x + b \quad \Leftrightarrow \quad y = mx + c$$

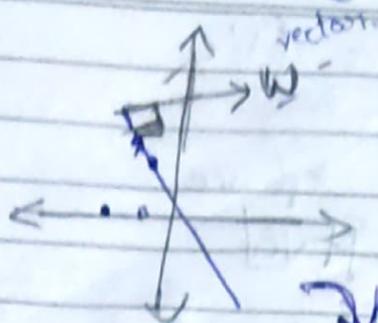
↓

$$ax + by + c = 0$$

$$\boxed{w^T x + b = 0}$$

- [\square] scalar \rightarrow only magnitude
- * [\checkmark] vector \rightarrow magnitude and direction \Rightarrow represents

* slope

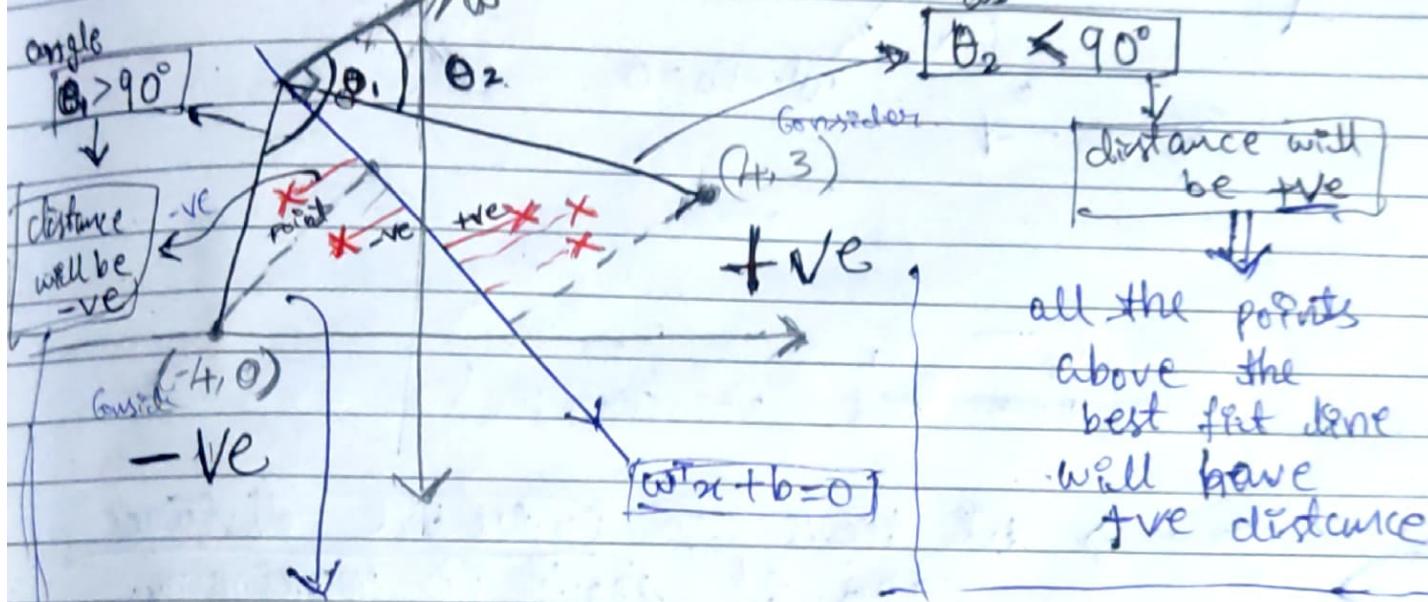


$$\text{slope} = \frac{\Delta y}{\Delta x}$$

which had both magnitude and direct $\log_{10} z$

-ve slope \Rightarrow time
~~decreasing~~ (Pending
downwards)

- * vector will be perpendicular
to ~~perp~~ line. as



all the points
above the
best fit line
will have
the same

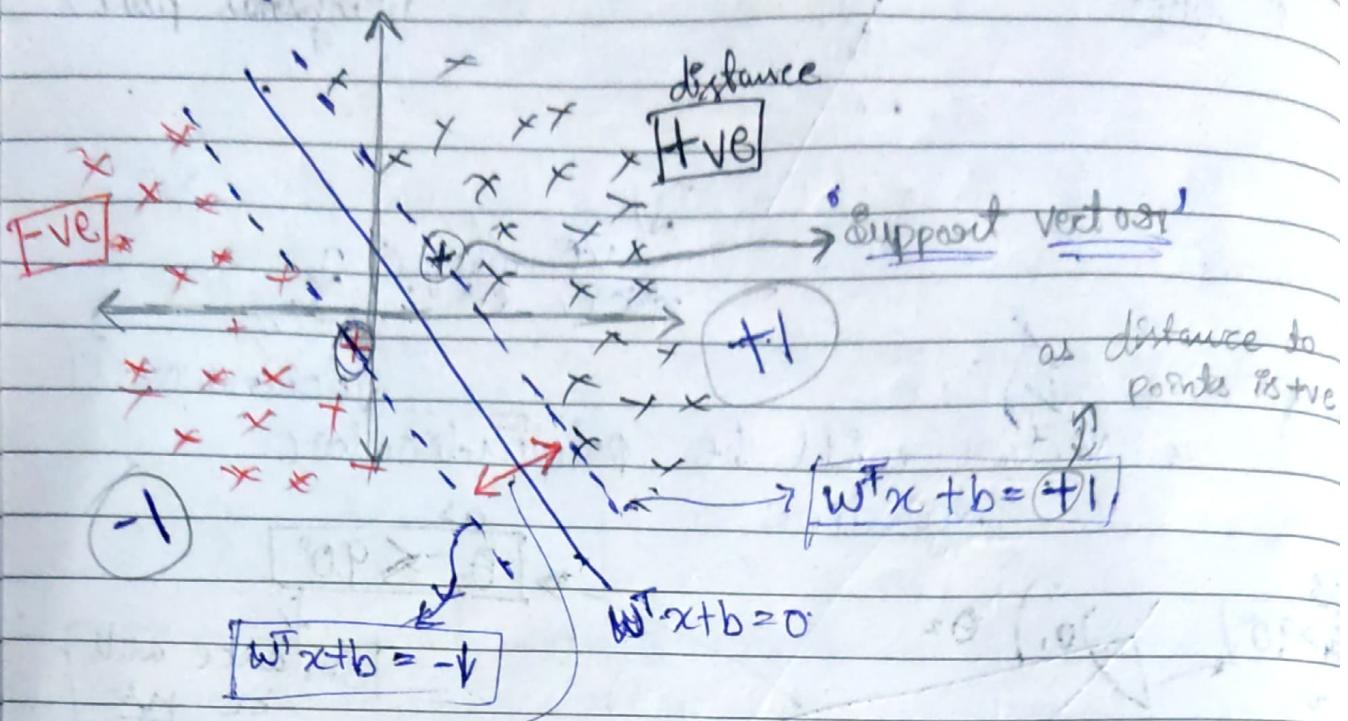
* How to calculate the distance?
↳ Using vectors

magnitude will be negative.

all the points below the best fit line will have -ve distance

represents direction.

* Marginal plane in SVC



point, $P_1 \leftarrow x_1 \Rightarrow (x_1, x_2)$
 $P_2 \leftarrow x_2 \Rightarrow (x_1, x_2)$

→ we have to calculate distance
 and it should be maximum.

$$\Rightarrow W^T x_1 + b = +1$$

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

(+) (-) (+)

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

← divide by $\|w\|$ magnitude

$$\Rightarrow \left(\frac{w^T}{\|w\|} \right) (x_1 - x_2) = \left(\frac{2}{\|w\|} \right) \text{distance b/w } \underset{\text{marginal planes.}}{\cancel{\text{marginal planes.}}}$$

* work out

unit vector,

$$\hat{w} = \frac{w}{\|w\|}$$

\Rightarrow aim \rightarrow To get maximize the distance b/w marginal plane.

* Cost function

Maximize
by changing
 w, b
slope, intercept

$$\frac{2}{\|w\|}$$

\Rightarrow distance b/w marginal plane

+ constraint such that

$$\begin{aligned} \text{Hard margin} \Rightarrow y_i &= \begin{cases} +1, & \text{if } w^T x_i + b \geq 1 \\ -1, & \text{if } w^T x_i + b < -1 \end{cases} \end{aligned}$$

* For all correctly classified data point

$$y_i * (w^T x_i + b) \geq 1$$

$$(-ve) * (w^T x_i + b) \leq -1$$

incorrectly classified

$y_i \Rightarrow -ve$

$$(+ve) * (w^T x_i + b) \leq +1$$

$+ve \Leftarrow y_i$

Modified Cost function of SVC

Can be written

$$\text{Maximize}_{w, b} \frac{1}{2} \|w\|$$

$$\Leftrightarrow \text{Minize}_{w, b} \frac{1}{2}$$

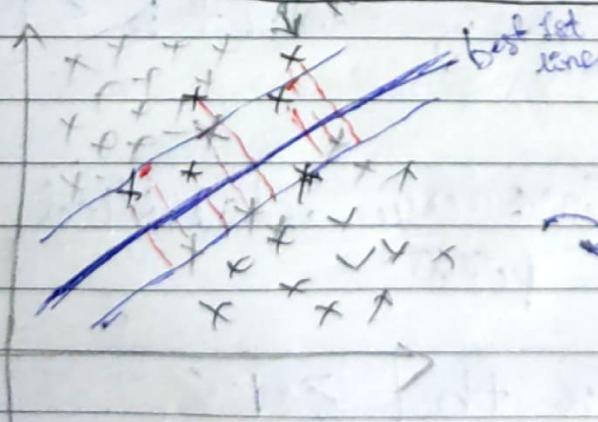
we write in this form because in other also we write cost function.

- * what happens when we have overlapped datapoint / 2 How many misclassified datapoint

~~Ex~~

Cost function of Soft margin SVC

$$\text{Cost fun.} = \text{Max}_{w, b} \frac{1}{2} \|w\| + C \sum_{i=1}^n \text{misclassified points}$$



Hinge loss

we follow here as we have the correct classification

Some option also depends
 as it can classify
 correctly
 remove wrong
 points \rightarrow if $\epsilon = 0$ then it says
 incorrect data points
 can remain if exceeds
 than we tune
 model

$\zeta_i \rightarrow$ by per parameter
 If cross threshold
 then we change
 hyperparameters

which say how many
 points we can consider
 for classification

$\text{heta } \sum \zeta_i \rightarrow$ summation of the distance
 of incorrect data points
 from the marginal plane.

If the threshold more then
 we will change the best fit line

$$\left[C = \frac{1}{\lambda} \right] \rightarrow \text{in logistic regression}$$

also

*Implementation of SVC with Python

① Data in. / output → import
make classification

② -train test split

③ ^{Model} Training

> from sklearn.svm import SVC

> svc = SVC(kernel='linear')

> svc.fit(x_train, y -)

> svc.coef_

④ prediction

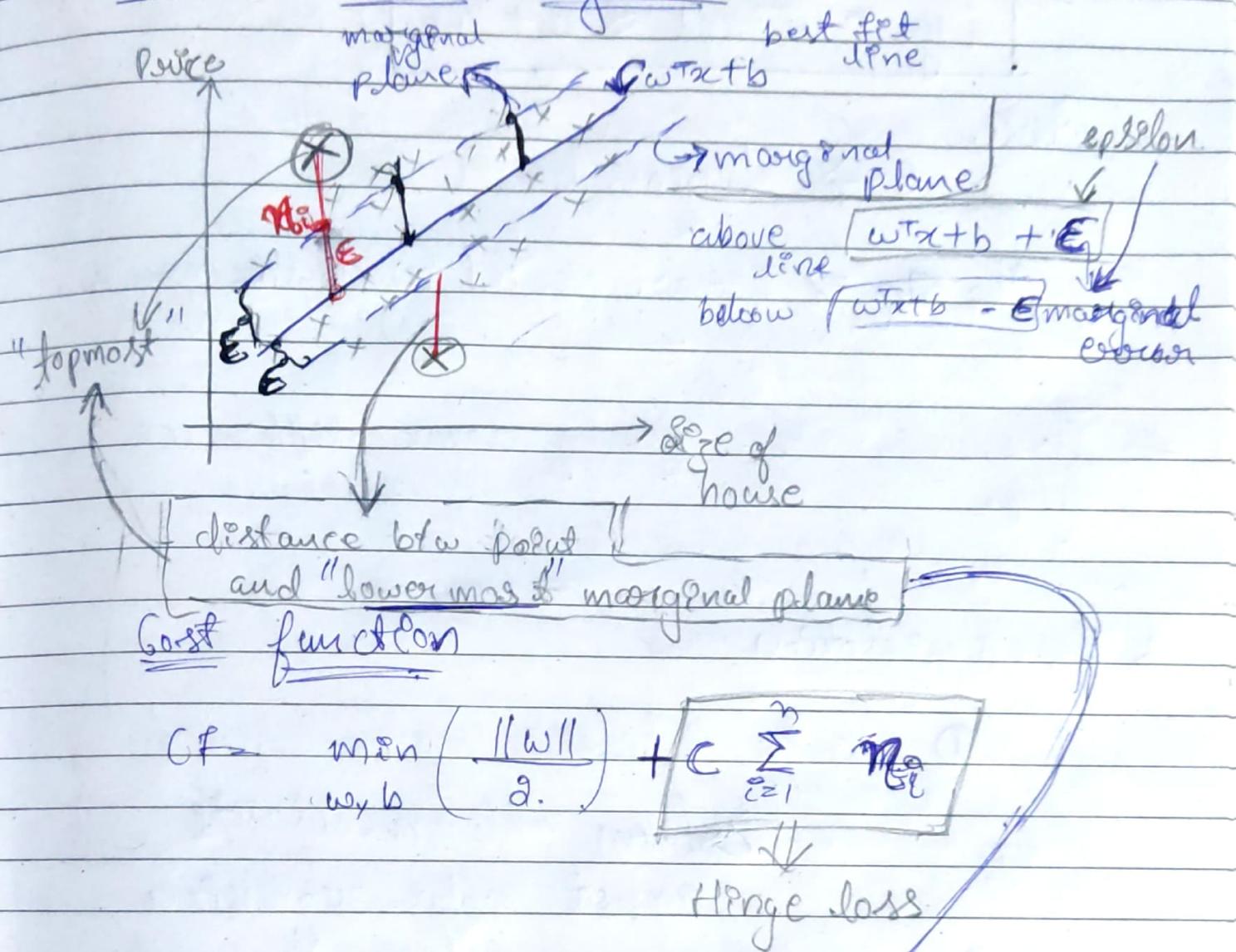
⑤ Performance. ~~metrices~~ metrices

⑥ Hyperparameter tuning

grid_cv = GridSearchCV(SVC(), param_grid=param_grid,
refit=True, verbose=3)

all info
of tuning
(each step)

* Suppose vector regression [SVR]



Constraint

actual - predicted

$$|y_i - w^T x_i| \leq E + n_i$$

We created a best fit line and marginal line, and make sure that all points in marginal line. So make ~~large~~ distance between marginal line more.

$$|y_i - \hat{y}_i| \leq E + \gamma_0$$

where,

$E \Rightarrow$ Marginal Error

$\gamma_0, \gamma_1 \Rightarrow$ Error above/below the margin.

Here, we are giving some difference gap.

* Implementation

① Create dataset

↳ from sklearn.datasets

import make_regression

② $x, y = \text{make_regression} (n_samples=1000, n_features=2, n_targets=1, noise=3.0)$

③ Train test split

④ Model training

> from sklearn.svm import SVR

④ Prediction.

⑤ Performance metrics.

→ from sklearn.metrics import f1Score

⑥ Hyperparameter tuning

> from sklearn.model_selection import GridSearchCV

> parameters = {
 "epsilon": [0.1, 0.2, 0.3],
 "gamma": [1]}

> grid_cv = GridSearchCV(SVR(), param_grid=
 {"gamma": [1]}, refit=True,
 scoring="neg_mean_squared", verbose=3)

> grid_cv.fit(x_train, y_train)

> grid_cv.best_params_

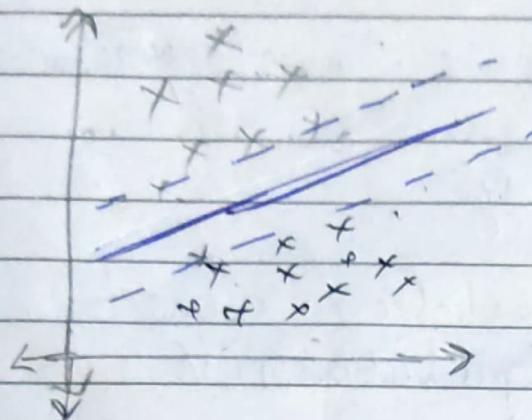
> y_pred = grid_cv.predict(x_test)

> print(f1Score(y_test, y_pred))

(intuition)

QMR

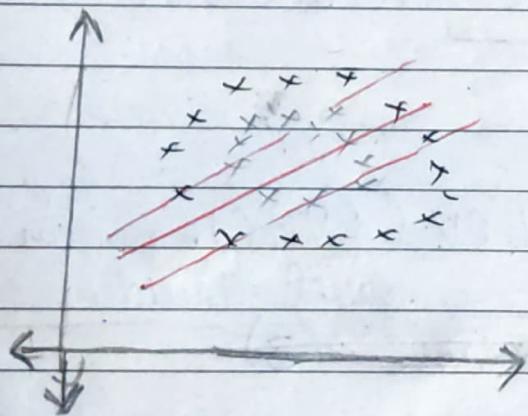
* SVM Kernels



If there were any outliers
then we handled those
through hyperparameter

Linear SVC

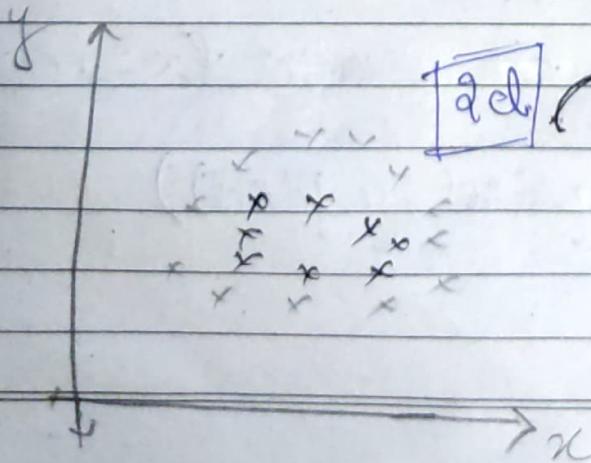
↓
Straight
Straight best fit + Marginal
line



→ Cannot use Linear
SVC.

→ more overlapped

Solve use SVM Kernels
to solve it



$$[2 \text{ cl}]$$

⇒ Transformation ⇒
using
Mathematical
formulas

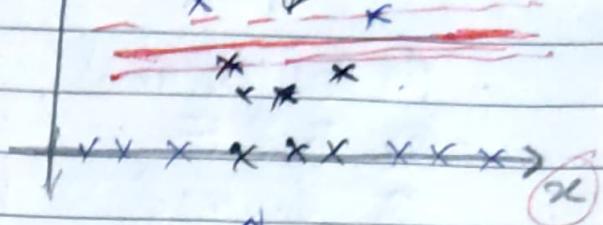
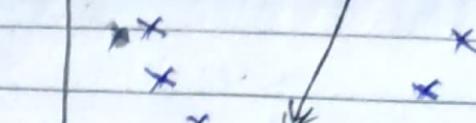
Applied formula
on x to make
 $1d$ as $2d$

Consider Dataset

x	y	$I = x^2$
2	Yes	4
3	No	9
4	Yes	16
-	yes	-
-	-	-
-	-	-

$$I = x^2$$

Linear SVC



$$I = x^2$$

1d points \rightarrow SVM Kernel

Data transformation

Technique

means

applying mathematical
formula

cannot
separate
by a
line
hence

formulas on x and y

$$I = f(x, y)$$

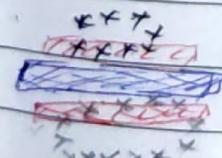
I is obtained

by applying
the formula

Now we

can easily

classify datapoints
after transforming
into 3d.

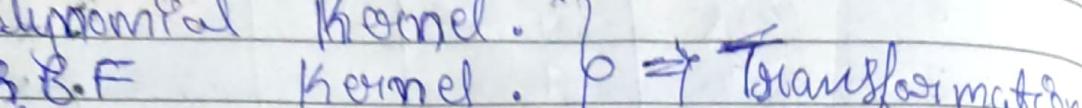


3d

then we apply
"Linear SVC"

we may get nearly
100% accuracy.

There are 3 types of Torsion formation given

- ① Polynomial Kernel .
② R.B.F Kernel .
③ Sigmoid Kernel . } \Rightarrow Transformation


* This basically convert lower dimension to higher dimension and easily separate points using linear SVC.

~~* Interaktion~~

① Polynomial Kernel

$$f(x, y) = (x^T y + c)^d$$

where, $C \neq$ Constant

\Rightarrow degree of polynomial?

depends on no of

o/p ~~feature~~ Category
if & then \rightarrow $d=2$

* Considerably
 $x = [x_1, x_2]$

$y = [x_1, x_2]$

then

$$\Rightarrow x^T \cdot y = \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} \cdot [x_1 \ x_2]$$

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

Created 3 new features
before we having (x_1, x_2)

2d

$x_1, x_2 | y$

Transformation

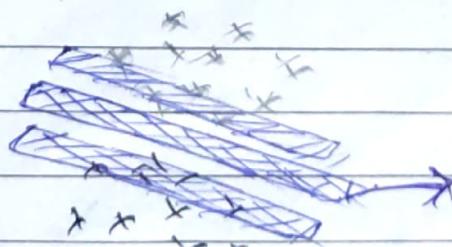
↓ Transforming [polynomial kernel
Transformed]

$x_1^2, x_1 x_2, x_2^2$

y

3d

x_1, x_2



best fit plane

x_2^2

x_1^2

IX

② RBF Kernel:

$$K(x_1, x_2) = \exp\left(\frac{|x_1 - x_2|^2}{2\sigma^2}\right)$$

③ Digmoid kernel:

$$\delta(x) = \frac{1}{1 + e^{-x}}$$