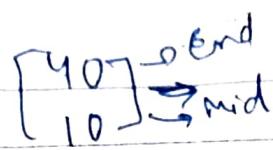
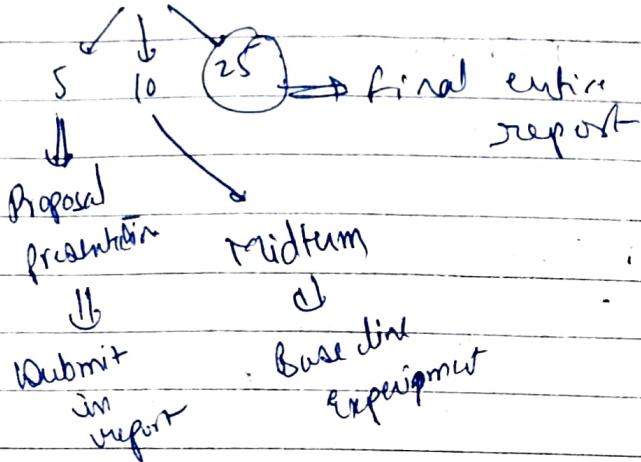


No relative grading



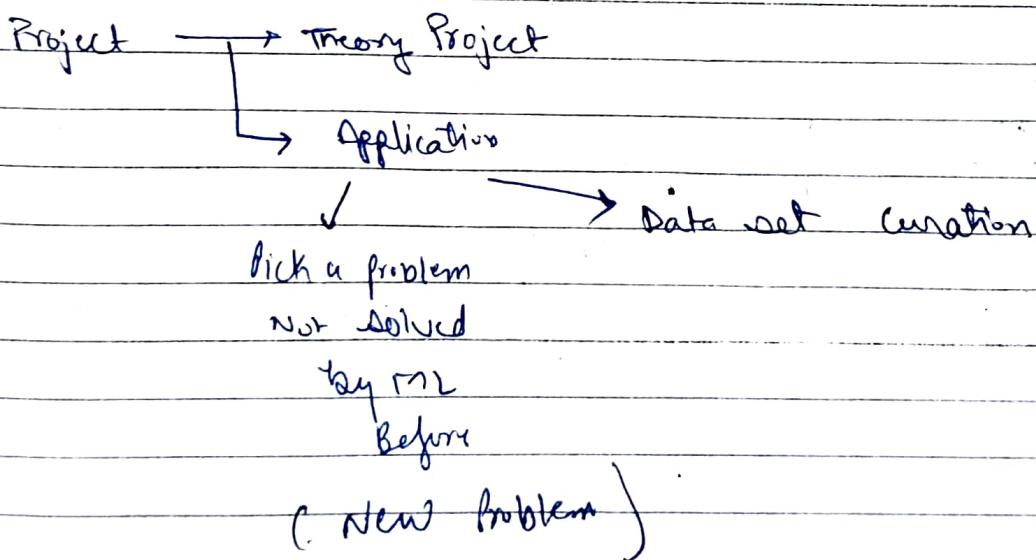
10 → Assignment → Poster (5th April)
40 → Project (A2)



[5, 10] → Bonus for publication

Midterm before Saturday
Endterm before Saturday

Timelines



Eduardo, Camila & Igoutube

[AI] → somebody as good as.

What ever humans
can do

[AI] → Rule-Based

e.g.
ATM

→ Data-Based
(Machine learning)

Next
target

[AI] → Artificial General Intelligence

Human-like

Next
target

[SI] → Superior intelligence

Table

Data

+

→ Easy

↳ systematic

↳ fast processing

O(1)

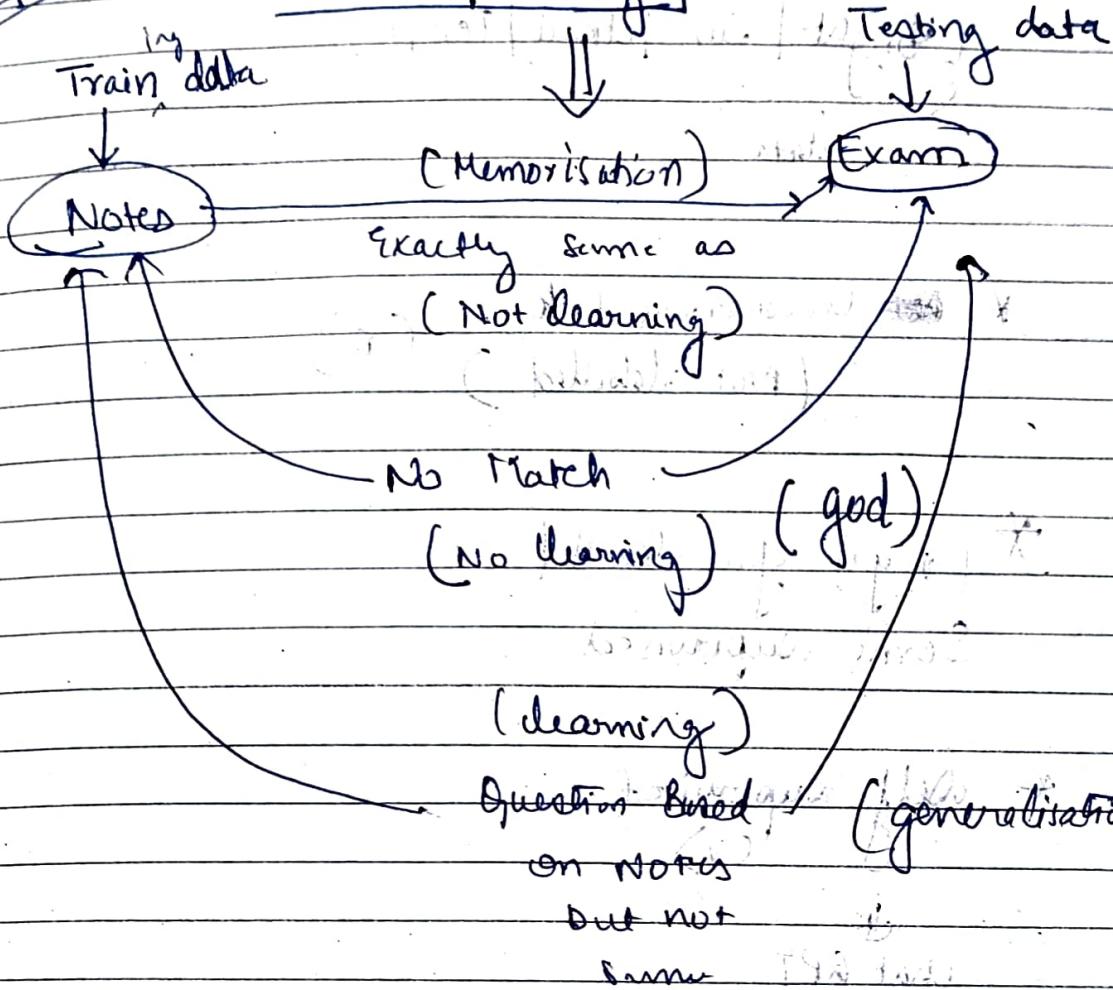
Data → Table

Page No. _____
Date / / 20

use ML where data is not defined

9/1/2023

Machine Learning



Not do correct model based on test set

Training data

Trained → Validation → Testing

These data should be in similar space

Machine L based on types of Notes it has

* Supervised (x, y) Label

bright | sun | cloud | Temp

data

* ~~Unsupervised~~ (x)
(Not labelled)

(x, y') few labelled

Semi Supervised

* Self Supervised
 (x)

that APT

I love Indian food

I love India Spain

I love _____ food

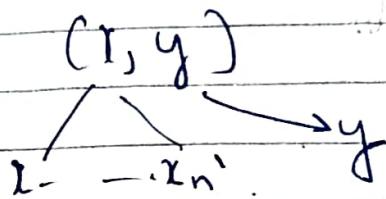
I love Indian _____

* Reinforcement Learning

(x, y)

Bunch of steps, I could lose some, but
(game of chess)

Overall I have to win the game.



Exam

Testing

Regression

Continuous
Value

eg Temp (continuous value)
 $22, 63^{\circ}\text{C}$

Classification

grouped

eg high temp
or low temp

into multiple
classes

multi labelled

Binary classification

Dog, cat, no animal, human

Page No.

Date: / /

most inf

INDUCTIVE BIAS

Inherently Bias regardless of data

On various models

What is

structured data

which is in format
on which we apply statistical
principle

eg Table

tidy → ① Every column has to have
one information

② Every row has to have
one observation

③ Cell has to have one

information

23!

Support Vector Machine → e.g. of supervised learning

label
 Email → spam (1) → supervised, classified
 → Not-spam (0)

SVM prefers 1 8 - 1

Convert non numeric into numeric value

letter to word mapping
 ASCII → numeric

but we need word convert to numeric

[bag of words]

It takes entire English dictionary

(May be only popular

words they which

Eng Dict

a apple

ball boy is in the room, 2 boys

have (corresponding huge vectors)

Biggest vector problem → computationally complex

One word \Rightarrow not able to accommodate all words.

	x_1	x_2	x_3	x_4	\dots	x_n
Email	1	1	0	0	= 0	0 1

Vector will have 1 when word present in it.

↓
Limitation

↳ like "frequency of words" will be lost

↳ Order of words doesn't matter, which is disadvantage

* One hot encoding *

Email converted into vector, now what to do?

Inductive Bias

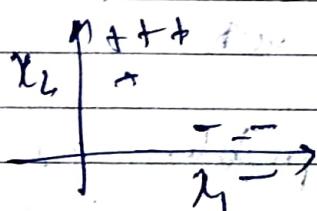
Also expect away
specific my people
context

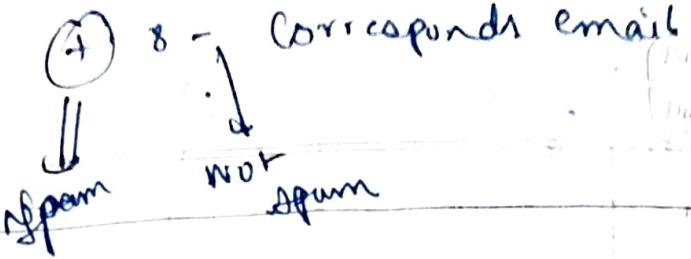
Can be represent

in ~~space~~ space

Vector consist of only 2.

$x_1 \quad x_2$





Goal is find line that cuts space into half



* Inductive Bias is in SVM

they have bias that either of
it's fine if it has to
be left or right
side

$y \rightarrow \hat{y}$ an predicting

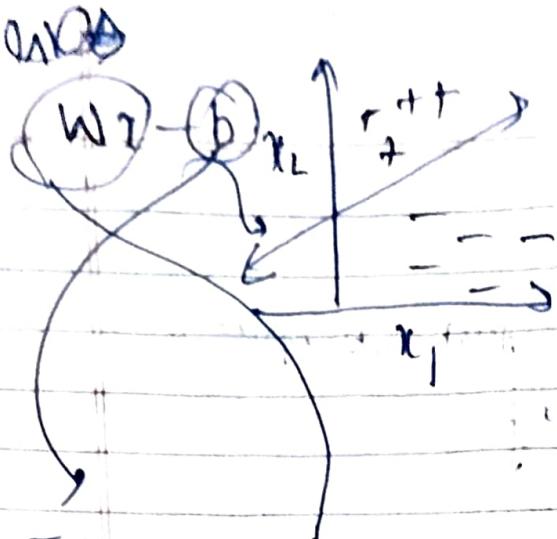
$$y = f(x)$$

= sign(x)
function

$$+ w \Rightarrow +1$$

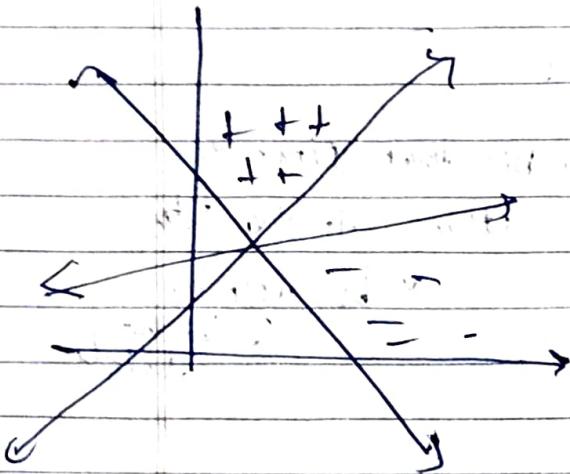
$$y = \text{sign}(w^T x + b)$$

↓ ↓ ↗
vector vector bias



It is a no.

It will
make
a no.



multiple such lines
many

constraints

$$w x - b \geq 1, y = +1$$

$$w x - b \leq -1, y = -1$$

$$\text{sign} = -1$$

that is the problem

if constraints
 ≥ 0 or
 ≤ 0

and in to find best
 w & b

try to find best possible

Combination of w^* & b^*

$$y = \text{Sign}(w^* x - b^*)$$

in email

SVM

marginal maximum
classified

(Page 16)

Date _____

$$w^T b = 1$$

$$w^T b = 0$$

$$w^T b = -1$$

find two points which
are closest to
each other

Support
vectors

Bad Centroids

By eqn there could be chance that there doesn't
exist a line but in reality there exist a line

So take the closest points so nothing can ignore

distance between the two lines = $\frac{2}{\|w\|}$ correct?
 $\frac{2}{\|w\|}$ in min?

Maximize this so that achieve classify
here.

Minimize $\|w\|$

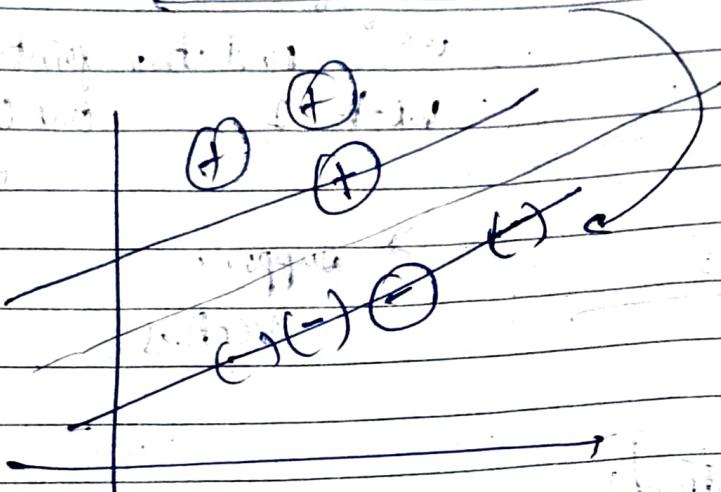
$$\text{s.t. } w^T b \geq 1, y=1$$

$$w^T b \leq -1, y=-1$$

At least two support vectors do
not have to be there

Page No. _____
Date: _____

More than 2 can exists



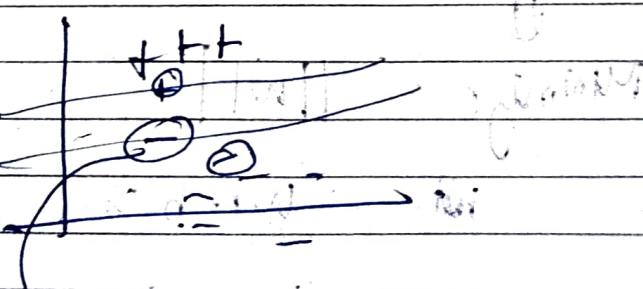
Inductive Bias of the algorithm

Guidelines
for
algorithms

Advantage of algo \Rightarrow little complex

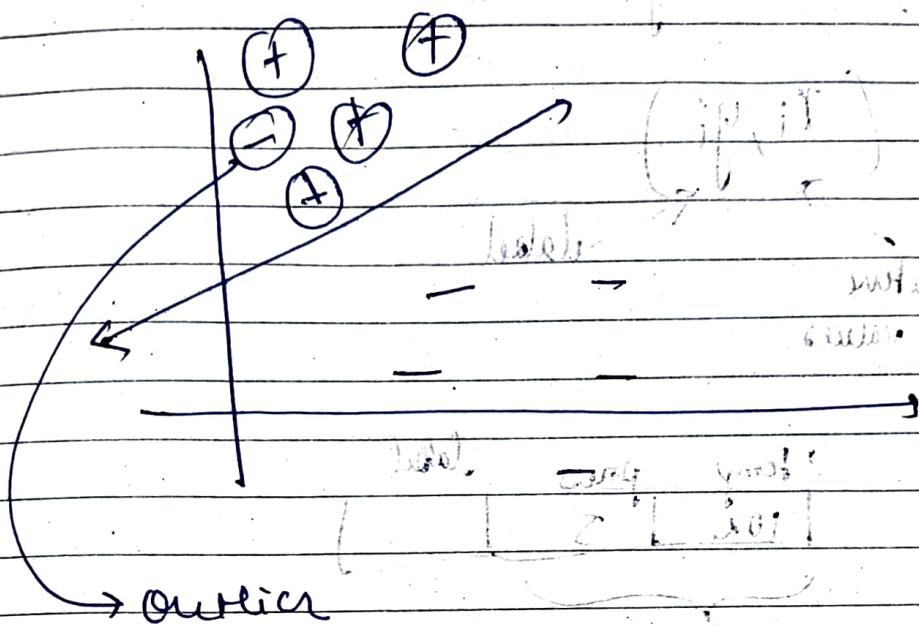
Model based learning \Rightarrow w₁, b, only depends upon the support vector.

bad support vectors are there then day

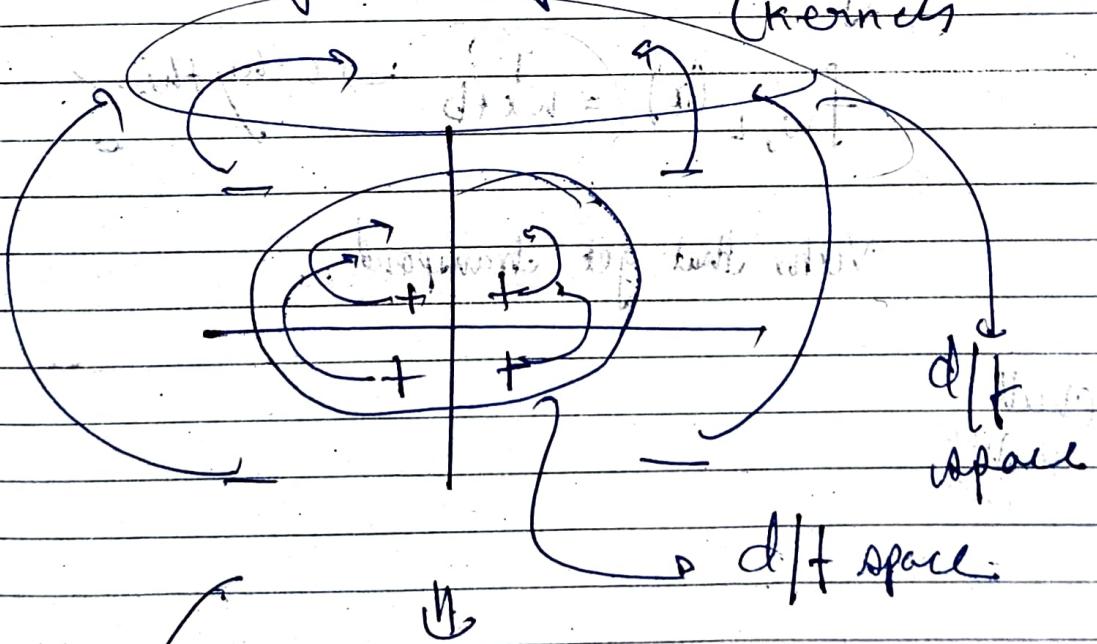


bad support
vector

How to Cure Outliers



① Transformation of non-linear to linear
functions



No line can separate them

so square things up

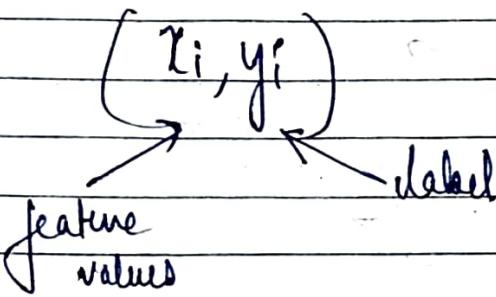
(2) Penalty / Error bar

Page No.

Date: / / 20

9/1/23

Linear Regression



temp	pour	label
102	5	

r_1

A table labeled r_1 with three columns: 'temp', 'pour', and 'label'. The 'temp' column contains the value 102, and the 'pour' column contains the value 5. The 'label' column is empty.

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

$$f_{w,b}(x) = w^T x + b \quad \text{why this?}$$

Vector that get transpose

Check
Final
notes

Linear Regression

Page No. _____

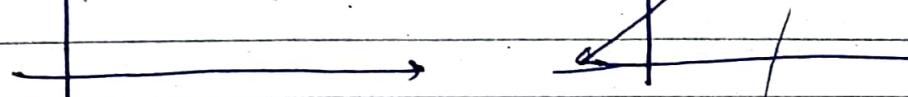
Date: / / 20

used for Regression Values

→ find w & b s.t. most of the points fall on the line

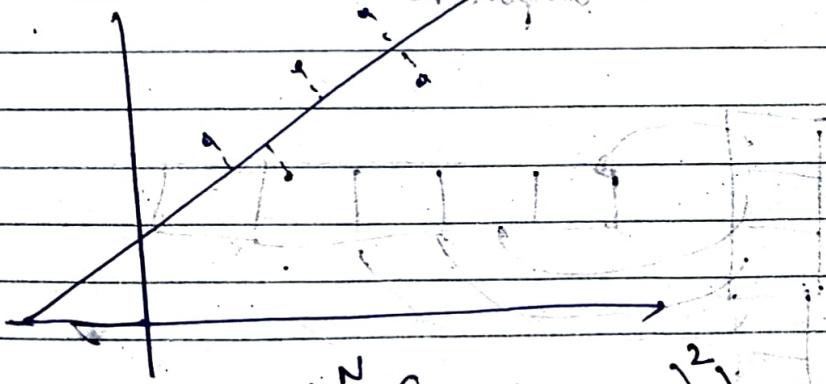
I have the points, we have to fit a straight line with minimum error.

Subtract normal equation. \rightarrow cost function



Cost function \rightarrow sum of squared errors

\rightarrow minimize the cost function \rightarrow least square fit



$$\min \left(\sum_{i=1}^N (f(x_i) - y_i)^2 \right)$$

Cost function

$$\min \left(\frac{1}{N} \sum_{i=1}^N (f_{w,b}(x_i) - y_i)^2 \right)$$

Weight error

Empirical loss function

Loss function

Analytical solution \rightarrow the solved form of
least square error.

Page No. _____

Date: / / 20

\rightarrow gradient descent

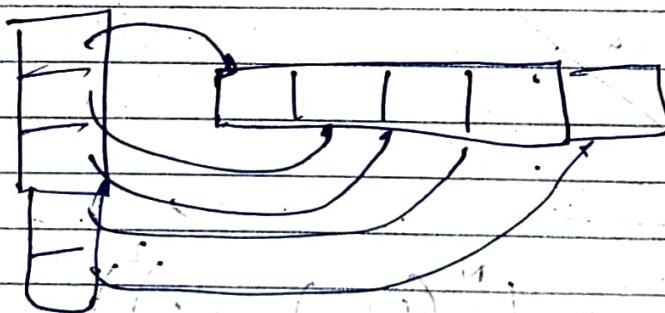
~~Why square, why not mod?~~

To penalize the point which are doing a bigger mistake

Inductive bias :- Assume linear relation w.r.t. data points

Why to choose linear regression?

tells which feature is important



100	1
flump	Press
102	5

More weight \rightarrow More weight \rightarrow Explainability

so we can remember this feature

form of
where

Page No.

Date: / / 20

What is the need of other algorithm

① less error

(1)

Empirical

Theoretical

Page No.

Date: / / 20

bounds
on the
error

Meaning

standard
dataset

How to prove that my algo is better?

Perform algo on standard dataset

How to make your own dataset standard?

Feature engineering

① Dataset should not contain features which are not useful

② Computer only understand numbers

Convert to Number

Assign numbers

problem with this → 2 & 3

Continuous flow of information

Computer thinks 1 is closer to 2

& 3 is closer to 2

but 1 to 3 is far which is wrong

→ Explaination

can

on this

One hot encoding

Page No.

Date: 1/1/20

Instead of giving no. represent them
in vector of 3

$$C [0, 0, 1]$$

$$A [0, 1, 0]$$

$$B [1, 0, 0]$$

disadvantage is the data become bigger

③

Binning

$$\begin{bmatrix} 22.5 \\ 20.1 \\ 19.2 \end{bmatrix}$$

→ normal

→ convert

$$\begin{bmatrix} 32.2 \\ 33.4 \\ 35.5 \end{bmatrix}$$

→ bin

one-hot
encoding

④

Normalisation

$$\begin{bmatrix} 0, 1 \end{bmatrix}$$

Making max to the ② [1]

$$\text{Normalisation} = \frac{x_i - \min}{\max - \min}$$

⑤ Standardization

means.

$$\frac{x - \bar{x}}{\sigma} \Rightarrow \text{z-score} \cdot \text{Normalization}$$

↑ deviation

~~$$\text{if } \mu = 10 \text{ and } \sigma = 5$$~~

- * Why normalization is done? \rightarrow To bring down them in same scale
- Why Standardization?

Question is important

⑥ Data imputation

There can be missing values

If you miss some value then what to do?

11	11	11	11
11	11	11	11
11	11	11	11

Missing

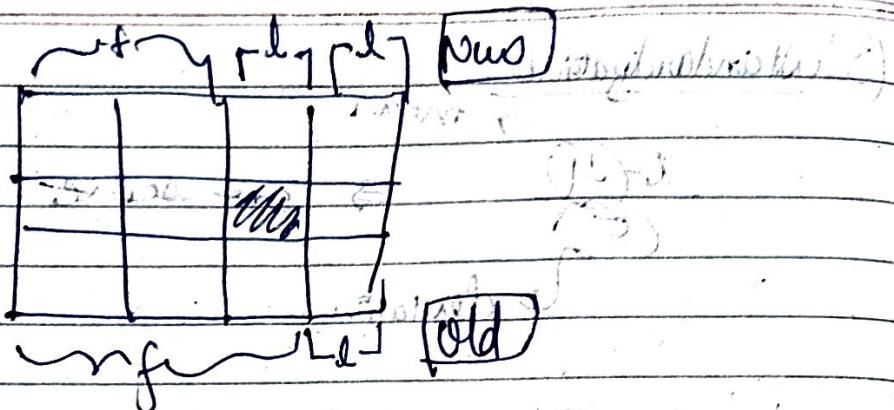
\rightarrow fix it with average of this

Or just predict the missing data , put it as

label

Page No.

Date: / / 20



10/1/23

K-NN

$k = \text{Nearest Neighbour}$

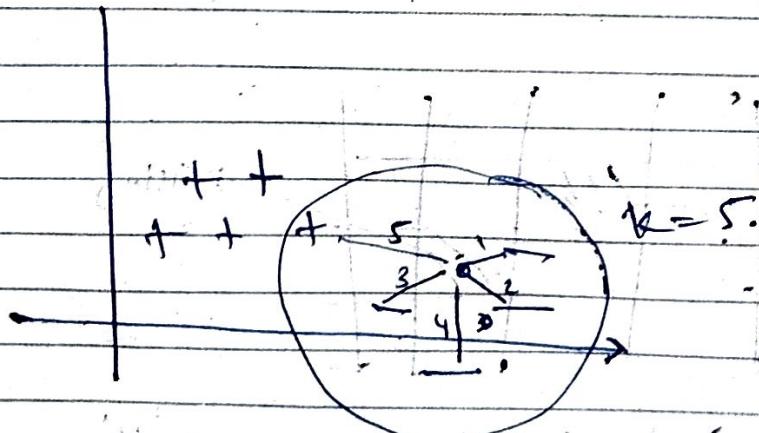
Homophilic & Heterophilic

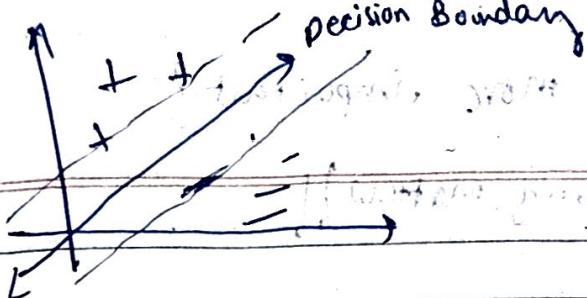
I behave the same way like my neighbour belongs

* find K NN

+ take Majority

+ Label with majority

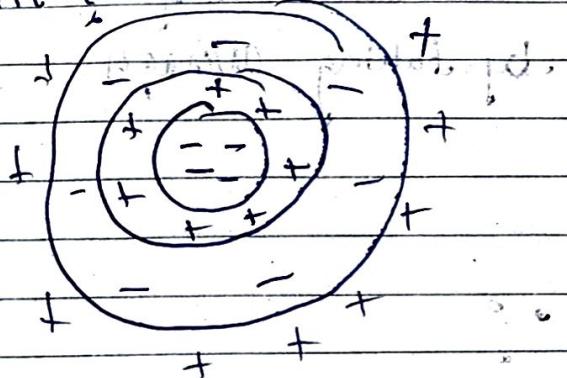




Page No.

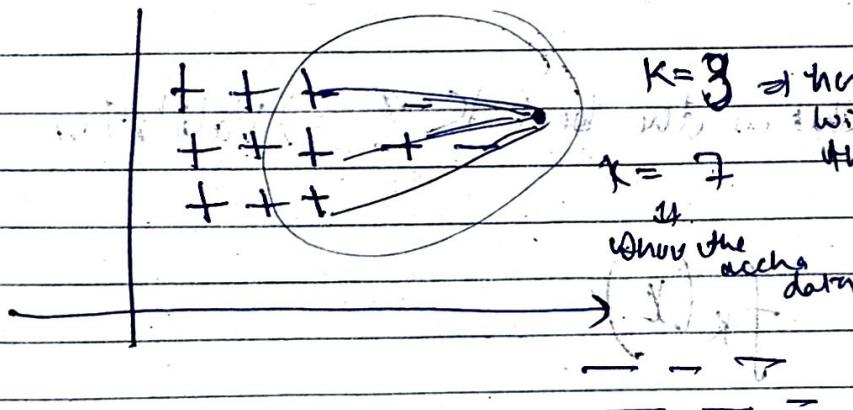
Date: / / 20

Can we have decision boundary which is
discrete?



What if K is larger

When larger K will help \Rightarrow outliers



$K=8$ so here it

will show the false alarm

when the actual data

is mixed with outliers

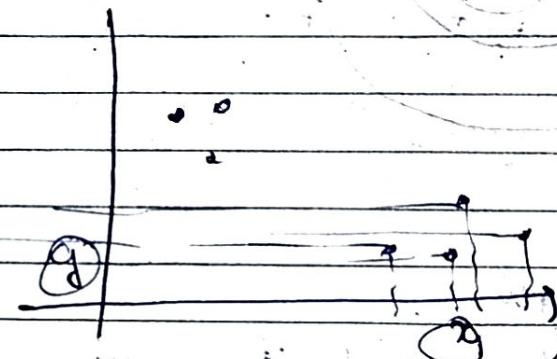
KNN is powerful because of non-linear
boundary

Where small now is more important!
Minority matters!!

Page No. _____

Date: / / 20

In linear regression we can also find
by taking average



$$y = \bar{y}_1 + \bar{y}_2 + \bar{y}_3$$

What is the best K ? \Rightarrow Validation

$$f_K(x)$$

Model v/s instance based learning

Inductive Bias & regardless of data it only considers k nearest neighbours

Page No.

Date: / / 20

Weighted k nearest neighbours

Radial Nearest Neighbour

Problem with KNN

Data structure which can store closest points together

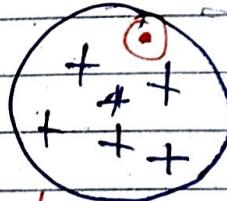
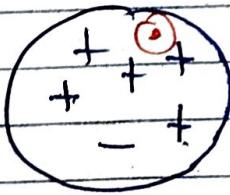
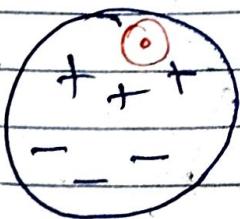
To know the closest distance you have to go through every point.

~~Is KNN explainable ??~~

~~Even less function of KNN ?~~

Decision Tree

Non linear model



More confident
to decide
the class

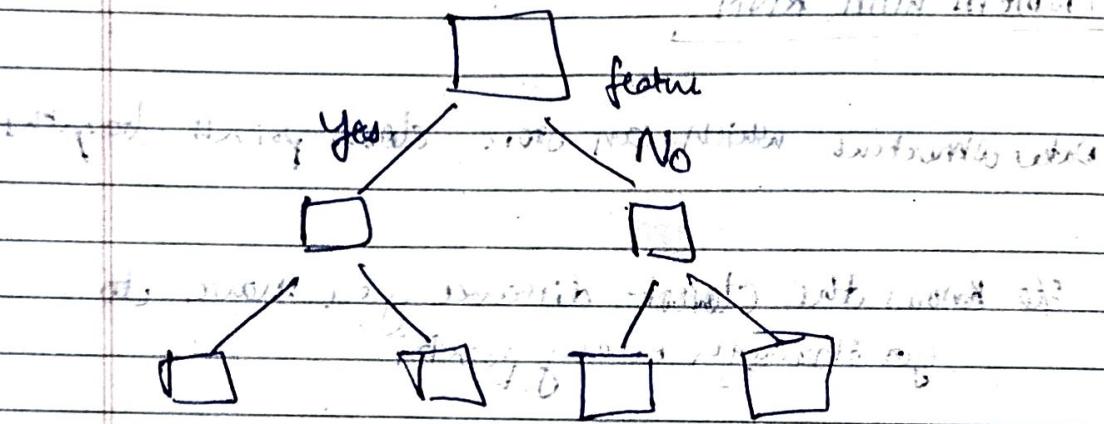
Q What features to look first?

Page No.

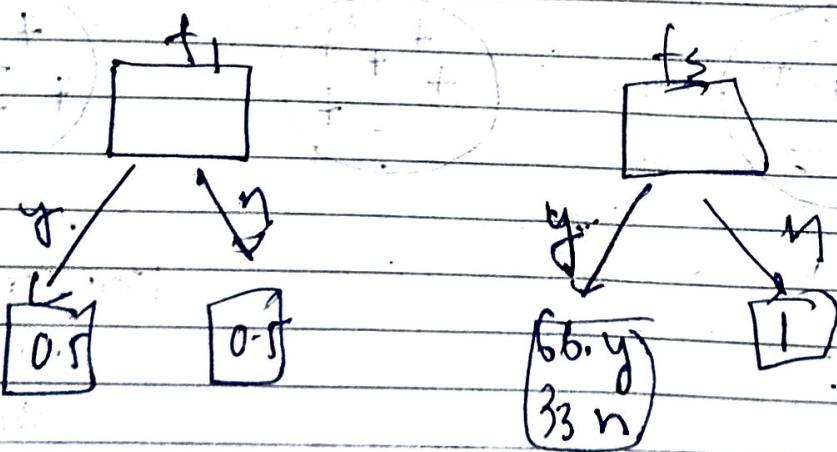
Date: / /20

1st step to do

To in like dialogue



y	y	y
n	y	y
y	y	n
n	n	n



Entropy

Gini's Impurity

Page No. _____

Date: / /20

Advantage

- ① We don't have to ask all the features

less questions good

Very high question \rightarrow Memorize

Normal ^{no. of} question \rightarrow learning

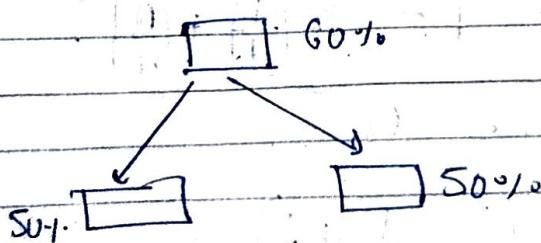
(Inductive Bias) of Decision tree

- Answer a question by asking a lot of question.

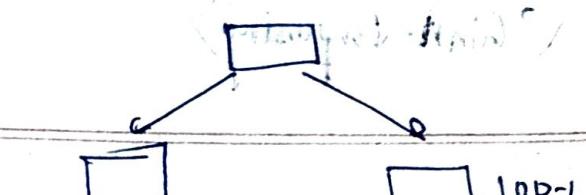
* Causality \rightarrow your Bio people

* hyperparameters \rightarrow check notes

IS/1/29



Certainty should not reduce.



Ex. (a) example

(1) example

Gini Impurity for each node task 3/6

$$G(\cdot) = 1 - \sum p_i^2 \rightarrow$$

0
1
Very
few
separations

Money

yes

No

Poly

100

Yes

100

Should sign

0

Eg 1

Eg 2

50

50

$$G(1) = 1 - \sum \left(\frac{1}{2} \right)^2 + \left(\frac{1}{2} \right)^2 = 0$$

$$G(2) = 1 - \left(\frac{1}{4} + \frac{1}{4} \right)$$

$$= 1/2$$

No. of branches depend on feature value

Page No
a

Date / / 20

$$g(\text{Money}) = \text{Weight avg}$$

$$g(\text{Money} = \text{yes})$$

$$g(\text{Money} = \text{no})$$

$$g(\text{Money} = \text{maybe})$$

$$= \frac{100}{150} \times g(\text{Money} = \text{yes})$$

$$+ \frac{50}{150} \times g(\text{Money} = \text{no})$$

$$1 - \varepsilon(p_{\text{yes}}^2 + p_{\text{no}}^2)$$

↳ yes for party

→ If Money is high what possibility that it will go for party or not?

Currently we have 150 rows with Money & party column

(mini) improving in your feature

$$g(f_1)$$

$$g(f_2)$$

$$\vdots \quad g(f_n)$$

} weighted average

to
which choose which has less gini impurity

Page No. _____

Date: / / 120

PT	OI	TV	Return
+VC	low	high	up -
-VC	high	low	Down
+VC	down	high	up -
+VC	high	down	up -
-VC	low	high	Down
+VC	low	low	Down \ominus
-VC	high	high	Down
+VI	low	low	Down \ominus
+VC	high	high	up -
-VC	low	high	down

$$G(f_{value}) = 1 - \sum p_i^2$$

Value of data

i f value of feature

$$h(f_{value}) = \sum w_i G(f_{value})$$

$$G(PT) = ?$$

$$G(OI) = ?$$

$$G(TV) = ?$$

$$G(PT) = w_i \text{ } f(\text{value}_i)$$

$$= w_i G(PT_+) + w_{i+} G(PT_-)$$

$$= \frac{\text{No. of trc}}{10} G(PT_+) + \frac{\text{No. of wr}}{10} G(PT_-)$$

$$= \frac{6}{10} G(PT_+) + \frac{4}{10} G(PT_-)$$

$$G(PT_+) = 1 - (p_{up}^2 + p_{down}^2)$$

$$= 1 - \left(\frac{16}{36} + \frac{4}{36} \right)$$

$$= 1 - \left(\frac{20}{36} \right) = 1 - \frac{5}{9} = \frac{4}{9}$$

$$G(PT_-) = 1 - (0^2 + 1^2)$$

$$= 0$$

$$\frac{6}{10} \times \frac{4}{9} + \frac{4}{10} \times 0$$

$$G(PT) = \frac{24}{90} = \frac{8}{30} = \frac{4}{15}$$

$$G(OI) = \sum w_i G(OI_i)$$

$$G(OI_i) = 1 - \sum p_j^2$$

Calculate
IT

$$G(0.5)$$

$$G(0.5) = \frac{1}{10} + \left(\frac{1}{8} + \frac{1}{4}\right) \cdot \frac{1}{10}$$

$$G(0.5) = \frac{24}{70} = \frac{12}{35}$$

$$G(0.5) = \frac{1}{10} \times \left(\frac{1}{8} + \frac{1}{4}\right) + \frac{6}{10}$$

$$\left(\frac{1}{8} + \frac{1}{4}\right) - 1 = \frac{1}{4}$$

$$1 - \sum \frac{1}{8} = \frac{1}{4} - 1$$

$$1 - \sum \frac{1}{8} = \frac{1}{4} - 1$$

$$\left(\frac{1}{4} - 1\right) - 1 = \frac{1}{4} - 1$$

$$0.25 + 0.25 - 1$$

$$0.25 + 0.25 - 1 = \frac{1}{4} - 1$$

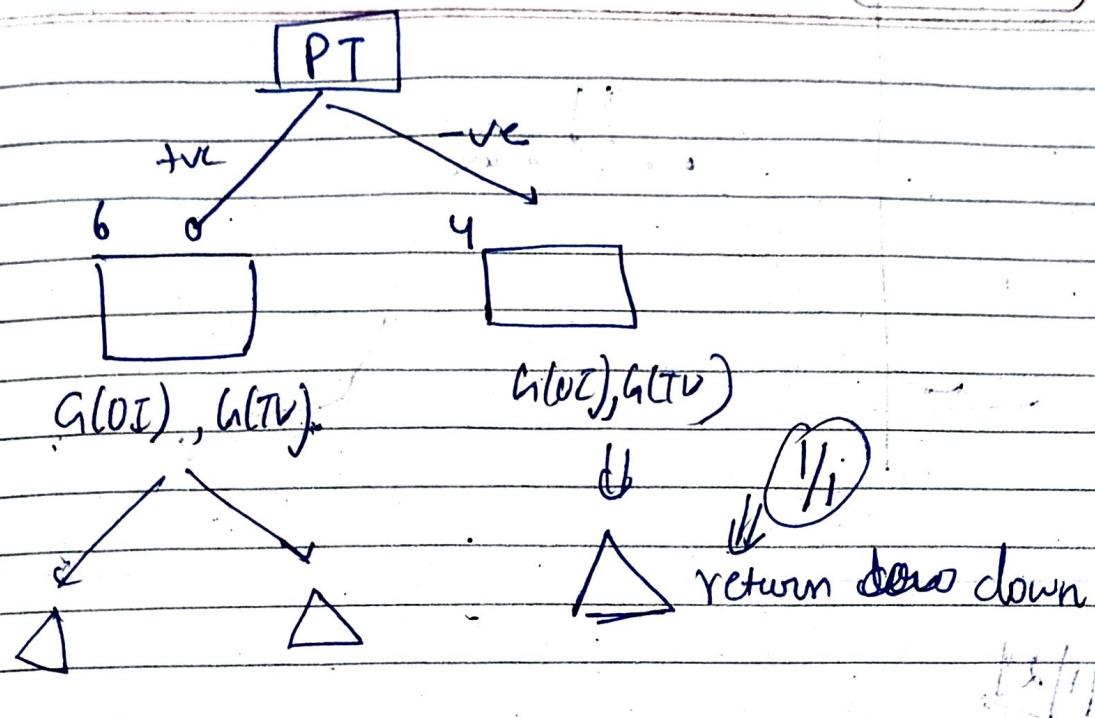
$$0.25 + 0.25 - 1 = \frac{1}{4} - 1$$

$$0.25 + 0.25 - 1 = \frac{1}{4} - 1$$

choose the lowest one

Page No. _____

Date: / / 20



Only split when purity is 0

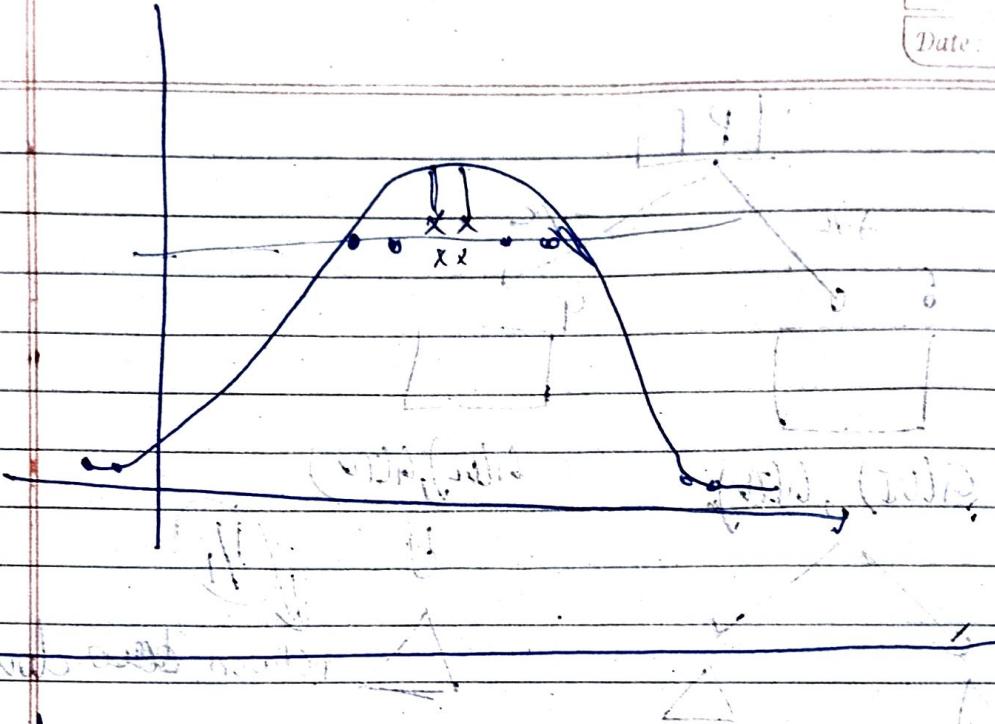
if it matches hyperparameters stop splitting it

Inductive V/s transductive learning

do look at
which space

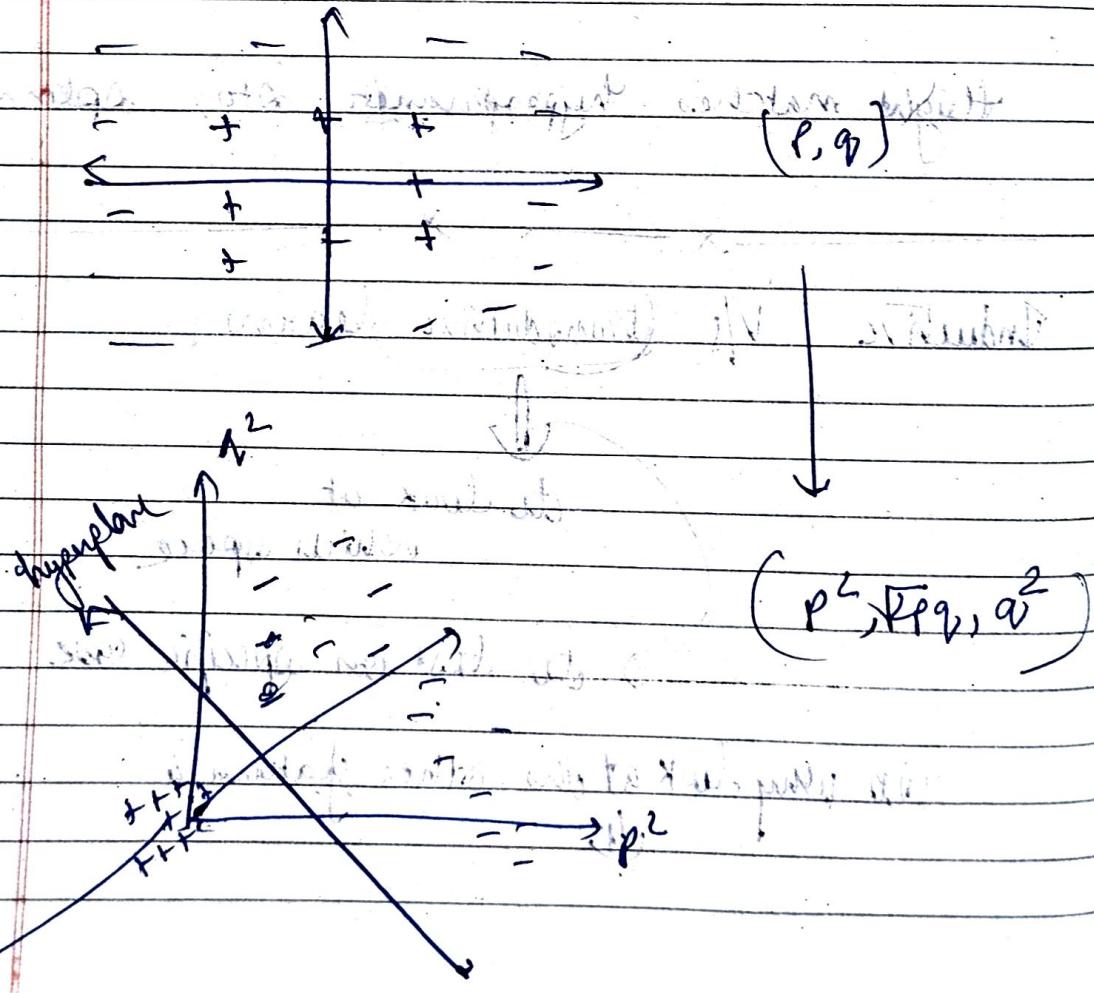
do test on specific use

do why look at the other features?



~~16/1/21~~

* Non Linearly Separable data



If something is not linear in some dimension then

We can make it linear in higher dimension. /20

(Page No.)

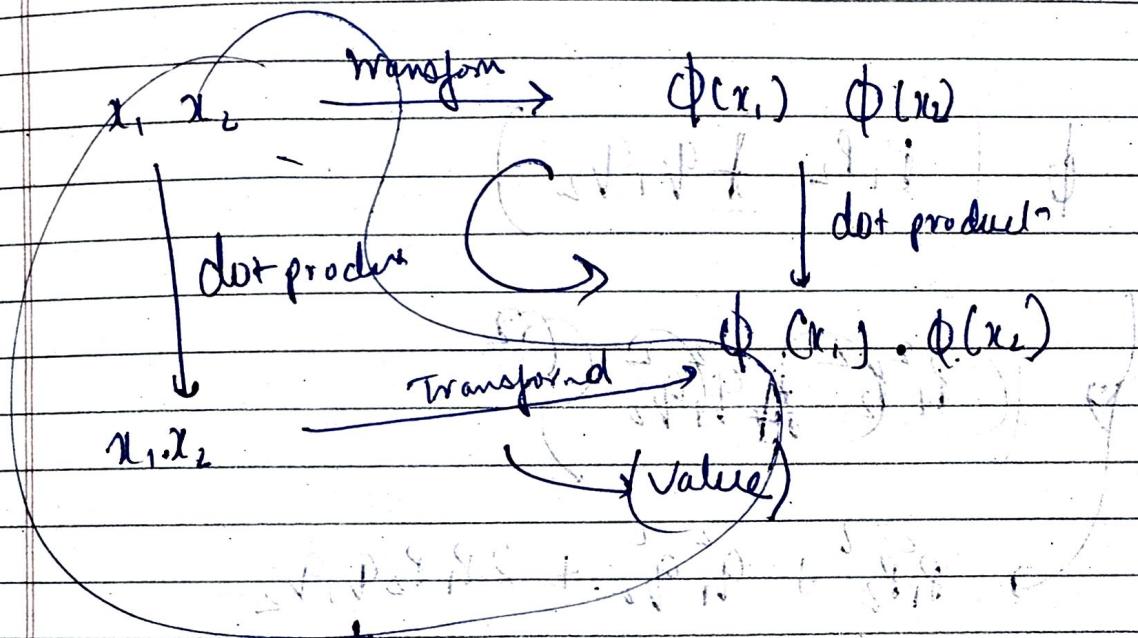
~~exan
Dr~~

If something is non linear can it become linear in higher dimension? ?

Kernel trick

$$x_1, x_2 \rightarrow (\phi(x_1), \phi(x_2))$$

↳ Transform



In case this function satisfies the
dot product

$$\phi(x_1, x_2) = \phi(x_1) \cdot \phi(x_2)$$

This is kernel

$$P_1^2 + \sqrt{2} P_1 Q_1, Q_1^2$$

Date: / /

$$x_1 = (P_1, Q_1)$$

quadratic

transformer.

$$(P_1^2, 2P_1Q_1, Q_1^2)$$

$$x_2 = (P_2, Q_2)$$

$$(P_2^2, 2P_2Q_2, Q_2^2)$$

$$(P_1^2 P_2^2 + 4P_1 Q_1 P_2 Q_2 + Q_1^2 Q_2^2)$$

$$\therefore (P_1 P_2 \pm Q_1 Q_2)$$

$$(P_1^2 P_2^2 + Q_1^2 Q_2^2)$$

$$P_1^2 P_2^2 + Q_1^2 Q_2^2 + 2P_1 P_2 Q_1 Q_2$$

$$x_1(p_1, q_1) \rightarrow p_1^2, \sqrt{2}p_1q_2, q_2^2$$

$$x_2(p_2, q_2) \rightarrow p_2^2, \sqrt{2}p_2q_2, q_2^2$$

P.I

~~Affine Function~~

Why do we need other kernels

Linear SVM

$$\min_{w,b} \frac{1}{2} \|w\|^2$$

$$\text{s.t. } y_i(wx_i + b) \geq 1$$

$\Rightarrow \min_{w,b}$
margin
(a.k.a. ||w||)

$$\|w\|^2 = w^T w = w \cdot w$$

$$L(w, b; a) = \frac{1}{2} \|w\|^2 - \sum_i a_i [y_i(wx_i + b) - 1]$$

$$\Rightarrow \sum a_i y_i = 0$$

$$x_i \geq 0.$$

Page No. _____
Date _____
Page _____

Tricky
way on
right as dot produce
129

$$\max \sum \alpha_i \rightarrow \frac{1}{2} \sum_n \sum_i y_i \alpha_i (x_i \cdot x_i) + \sum_k \alpha_k$$

$$x_i \neq 0 \quad (y_i \neq 0)$$

Non linear SVM

Linear SVM

Non-linearly
separable
data

Kernel
trick

SVM

$$w = \sum \alpha_i y_i x_i \rightarrow \text{this is a support vector}$$

y should be -1 or 1

x_i should exist for Support vector

x_i should only be zero if it is not SVM

RBF \Rightarrow other Kernels

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

sigma

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

variance across the axis

This particular function
is convex to all dimensions

hard margin classifiers

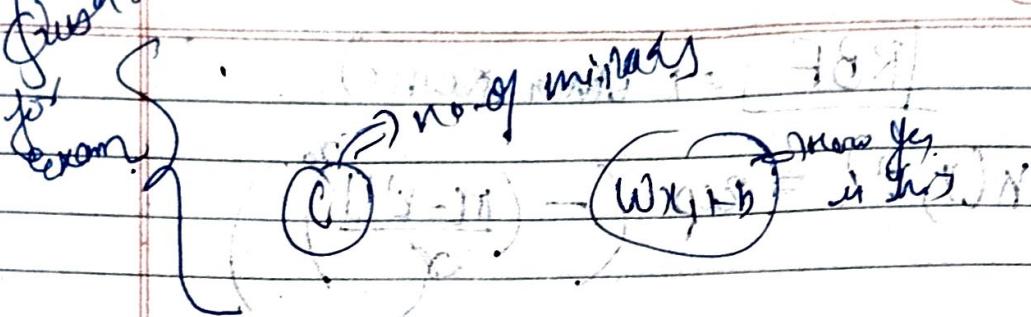
soft margins classifiers

Hinge loss

$$\min_w \frac{1}{2} \|w\|^2 + C \sum_{i=1}^n \max(0, 1 - y_i(w^T x_i + b))$$

constant but scalar

Question
from exam



Trying to reduce the mistake

Now decision boundary looks -

17/1/24

Performance

perform well independent of data

Regression

Classification

→ MSE or RMSE

→ Bayes idea ①

	+ve	-ve
+ve	TP	FP
-ve	FN	TN

TP \Rightarrow True positive

FP \Rightarrow False positive

Predicted

Confusion Matrix

	+ve	-ve
+ve	10	15
-ve	6	5

Classification

Performance Matrix

Accuracy

$$\text{Accuracy} = \frac{TP + TN}{\text{All}}$$

Problem with this
Data unbalanced

Accuracy is still high but doing not good in other cases

No biology people think Accuracy

$$\text{Precision} = \frac{TP}{TP + FP} = \frac{1}{1 + \frac{FP}{TP}}$$

Total positive predicted

mark + true
rest -ve

$$\text{Recall} = \frac{TP}{TP + FN} \Rightarrow \text{mark everything}$$

Actual positive

Covid Test

How can I cheat Precision

→ find one which has ~~no~~ full chance of positive
and rest negative

recall ↑ → precision ↓

what if precision \rightarrow recall?

Page No.:

Date: / / 20

Make both right !!

$$F1 \text{ score} = \frac{2 \cdot \text{Precision} \cdot \text{Recall}}{\text{Precision} + \text{Recall}}$$

Better than f1 score (Poster idea 2)

Other matrix

* Imbalanced data

SMOTE & Synthetic minority oversampling
techniques

very simple like KNN

, Poster Idea (3)

$$1 + \frac{1}{r^2} + \frac{1}{r^2} + \frac{1}{r^2}$$

$$\frac{2 \times 1 \times 1}{1 \times 1} = \frac{0}{0}$$

Page No. _____
Date: / /

In what is the issue
with F1 score?

x^+
 x^+
 x^+
 x^+
 x^-
 x^-
 x^-

$\text{random}(0,1)$

p (q, new point)

— Pick sample in pair

— $n_p = p + q(0,1) (q-p)$

→ Pick sample

→ pick 1 in the KNN

problem of random outlier

undersampling

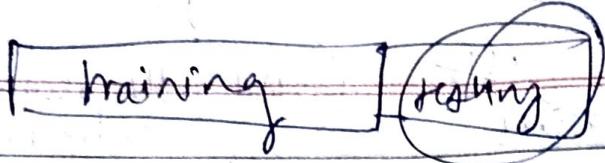
* Validation

→ parameter

→ hyperparameter → height of decision tree

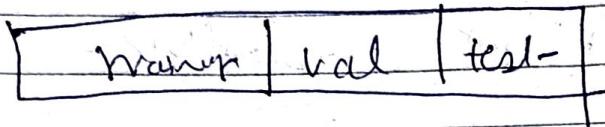
↳ C in non-linear SVM kernel, or

dataset



Page No. 1 / 120
Date: 1/1/2023

↳ hold out -



How to choose C?

* grid search

Table where C value written

2 kernel value is written

written



K-cross fold validation

train | test

train | test



K-fold

5-fold

fold 0



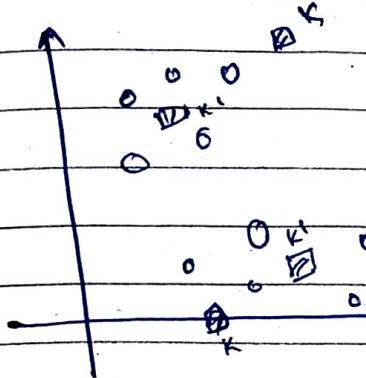
18/1/24

JK bold

Page No.

Date: / /20

Unsupervised



① Clustering

K Means

DBSCAN
(ϵ, n)

Randomly pick k centroids

Ensemble Learn

Bias

Variance

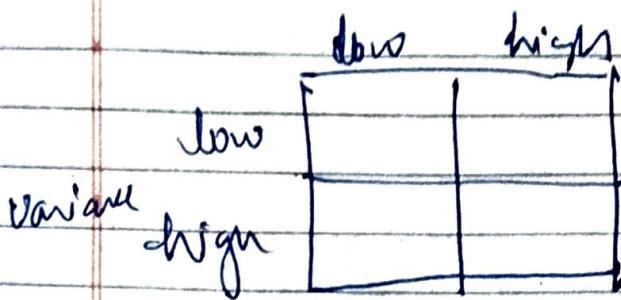
How much change
in the data can
change the
model

↓ ↗
SOT DT

Assumptions
that we
make lead
to an
error

1) Bagging → $D \xrightarrow{v_1} D_1 \xrightarrow{v_2} D_2 \xrightarrow{v_3} D_3 \xrightarrow{v_n} D_n$ } bootstrapping

Bias



2) Boosting $\Rightarrow y = y_P + \text{residual}$

\uparrow
Prediction

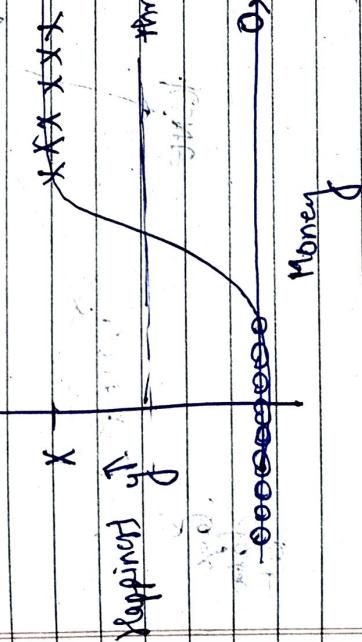
3)
 Mad Boost
 Ada Boost
 Kh Boost

Classification Function

Logistic Regression

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

Binary classification



$$f(x) = \frac{1}{1 + e^{-(wx+b)}}$$

How to find error, find the least MSE

$$\text{Loss} = \prod_i f_{w,b}(x_i, y_i) (1 - f_{w,b}(x_i))$$

Maximum likelihood
log (max likelihood)