

MACHINE LEARNING in systems

TM MITCHELL - A computer learns a task

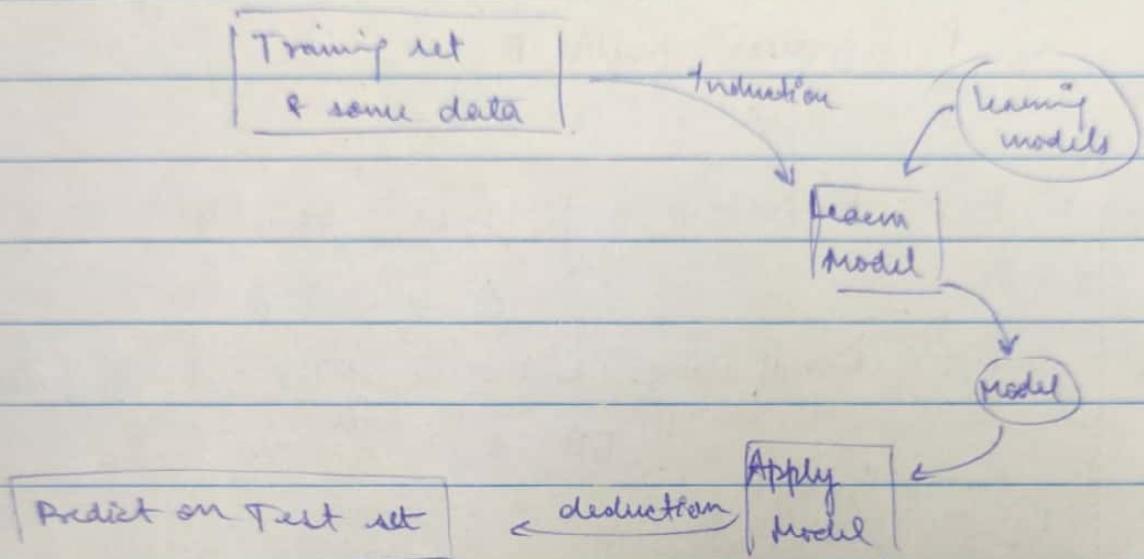
T from experience E, if its performance
P improves with E.

E: Experience of past results.

T: Predicting outcome of next or future
events

P: Performance → Many metrics,
accuracy, F₁ score, R₂ score, etc.,
probability.

Interaction of Data & Learning Algos.



It is primarily divided into

3 types of learning

Supervised

with ~~any kind~~
of target to
predict

Unsupervised

without
any kind of
target of
predt.

Re-inforcement

- If model
does something
right, it is
rewarded.

- If model
does something
wrong,
it incurs
a penalty.

- DEPENDS
on training env.

SUPERVISED LEARNING → output is present.

↳ Regression → predict the value of a target feature which may have Eg. monthly continuous values. price.

→ Classification → classify the data entry into one or more sets of classes.

Eg. classifying whether a person has breast cancer or not.

CROSS-VALIDATION of data

There may be a structure present in our data such that train data DOES NOT CAPTURE the actual distribution of the full data.

What do we do?

Cross validation - why? To access the whole feature space.

• simply take different

training sets & different

test sets to eliminate bias.

$$\{C_1 | C_2 | C_3 | \dots | C_{10}\}$$

first $C_1 \rightarrow$ train
next $C_2 \rightarrow$ train, next $C_3 \rightarrow$ train,
----- ?

CLASSIFICATION: when the output var.
is a category i.e. 'red' or 'blue'

REGRESSION: A regression problem
is when the output is a [real] value.
Eg 'dollars' or 'weights'

UNSUPERVISED ALGORITHMS

No target features to predict,

CLUSTERING: whenever we want to
discover [groupings] in data

ASSOCIATION: An association rule
learning is where you want to
discover rules that describe
large portions of data.

Find an Associate relationships b/w some
features to categorise.

SEMI - SUPERVISED & TRANSFER LEARNING

↓
some target values
are known but
not for
others.

Learning of one
domain helped
in learning of
another.

Eg movie recommends
through novel
likings.

Data, Info & Knowledge

)
raw
facts &
figures.
)
ordered
data

↓
combination of
past experiences &
current info.

PROPERTIES of Data

- ① VOLUME - quantity of the data you have.
- ② VARIETY - The different forms of data.
Eg. health care, images, video,
audio.
- ③ VELOCITY - freq / Rate of data generation.
- ④ VALUE - How much meaning the data adds
to our model.
- ⑤ VERACITY - Certainty & correctness in data.

TRAIN, VALIDATE & TEST

↓
the data
on which
we want
to build
our model

↓
the data
which gives
review while
the training
phase of
model

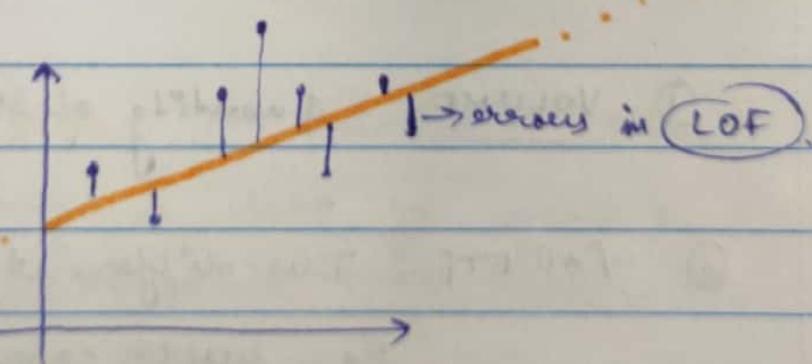
→ final set of
data on
which we
evaluate our
model.

REGRESSION

We want to predict a continuous
RESPONSE VARIABLE given some
PREDICTOR VARIABLES.

→ Linear Regression

- ① A Predictive modelling technique
to predict a CONTINUOUS response var.
- ② The independent vars. are predictors.



The relationship is thought of to be
linear.

THE LEAST SQUARES METHOD

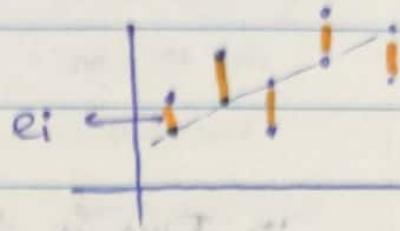
$$\Rightarrow m = \frac{\sum (x - \bar{x})(y - \bar{y})}{\sum (x - \bar{x})^2}$$

\Rightarrow This is the slope of the line which minimizes the SQUARED ERRORS.

Draw

and draw a straight line through
the two given values and
find the area

MEAN SQUARED ERROR



Actual values $\rightarrow y_i$

Predicted values \rightarrow green line $\rightarrow Y$

$$E = \frac{[\sum(Y - y_i)^2]}{n}$$

this is a metric to evaluate how
bad our values are w.r.t the
ACTUAL VALUES.

R-squared

→ It is a statistical measure of
how good our model is predicting
at the response variable.

or say $x, y \rightarrow 0.85 R^2$

↓
85% of variability of
 y could be explained
by x .

VARIABILITY: It tells us that the behaviour of y could be explained with x .

85% ~~variance~~ R^2 means \rightarrow 85% of the time, this model can predict the response behaviour values. ~~values~~ correctly.

$$R^2 \rightarrow \text{formula} = 1 - \frac{\sum (y_p - y)^2}{\sum (\bar{y} - y)^2}$$

- 1 - RSS \rightarrow mean
 $\text{RSS} \rightarrow \text{avg. sum of squares.}$

RSS: Residual sum of squares \rightarrow the sum of squares of the residues (errors) of the points wrt predicted values.

$$\rightarrow \sum [(y_p - y)^2]$$

ASS: On avg., how does our model reduce the error in performance of an average model.

different methods to do
LReg.

Linear Regression - gradient descent
+ regularization.

Linear regression has 2 parameters \rightarrow

$$Y = mX + c$$

m & c are called parameters.

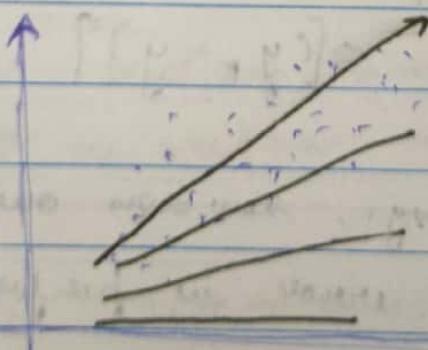
\hookrightarrow we need to find m & c such
that the Error of our model is
minimized. \hookrightarrow RESIDUAL

\Rightarrow Gradient Descent

Slowly move in some direction.

which direction?

The direction in which the
error reduces the MAX.



The values of
 m & c are
updated

\rightarrow let us say residual error \rightarrow residual error
with the first step is E_1

\rightarrow If the next step gives us an error E_2 ,
then we update our line

- We keep on moving our line to new parameters till we are getting a SMALLER ERROR.

Loss f" in gradient descent

→ The loss is the error in our predicted value of m & c.

→ we need to minimize the value of THIS LOSS FN to get most accurate (m) & (c) values.

We will need - mean squared error.

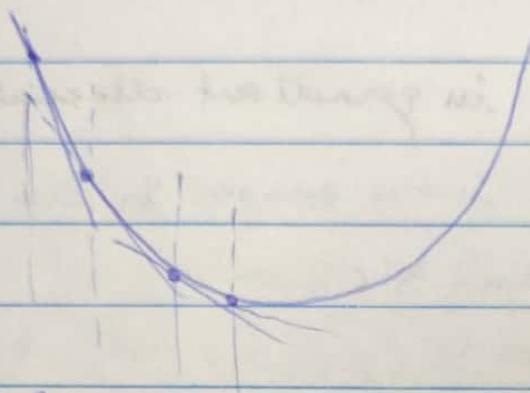
$$\text{MSE / RMSE} \rightarrow \frac{\sum_i (y_i - \hat{y}_{p,i})^2}{N}$$

↓
this is our
loss f"

$$E = \frac{\sum_i (y_i - mx_i - c)^2}{N}$$

VISUALIZE

- Steep slope: take large steps.
- less slope & take small steps.



→ GD is an iterative algo. to find min. of a f^n .

→ It's like you get the slope at a point & if the slope is very high \rightarrow [more faster]

→ If it is smaller slope, you [more slower].

$$E \rightarrow f(m, c).$$

$$0 = \frac{\partial f(m, c)}{\partial m}, \quad \frac{\partial f(m, c)}{\partial c} = 0$$

[at the minima.]

Applying gradient descent \rightarrow

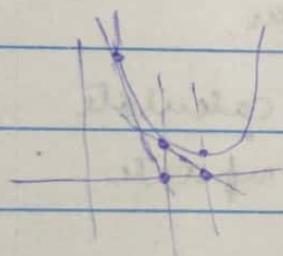
form $m = m$ & $c = c$,

$$\frac{\partial E}{\partial m} = \frac{1}{n} \sum_{i=0}^n 2(y_i - (m x_i + c)) \cdot x_i$$

$$D_m = \frac{2}{n} \sum_{i=0}^n (y_p - y_i) \cdot x_i$$

This gradient descent also has

a LEARNING RATE, L - how big steps we take.

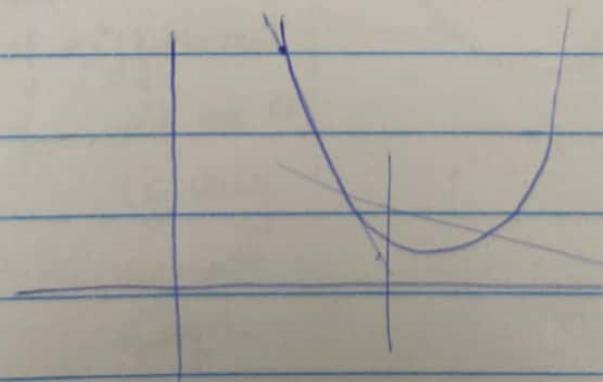


$$m_{\text{new}} = m_{\text{old}} - L \times D_m$$

$$c_{\text{new}} = c_{\text{old}} - L \times D_c$$

→ simultaneous update

We go in the reverse direction of the slope of E w.r.t m & c .



→ large LRs mean that we may overshoot the minimum of our cost f.

There are 3 types.

* BATCH GD

Take the whole dataset at each iteration.

→ E_i is evaluated for all the points at each iteration.

* STOCHASTIC GD

→ Here we are selecting ONE ENTRY at an ~~and~~ iteration & then do not consider that for future iterations.

→ TAKE only 1 point, calculate error w.r.t that & update.

* MINI BATCH GD

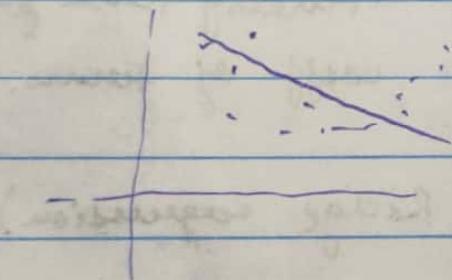
→ Here, it uses n samples instead of 1, at each iteration.

REGULARIZATION -

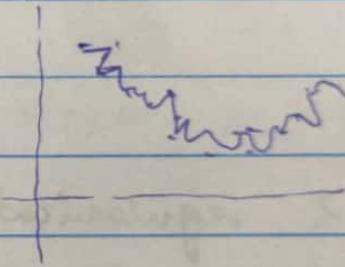
It is a technique for reducing model complexity & prevent over-fitting which may result from simple linear regression.

underfit & overfit are problems which regularization deals with.

UNDERFIT



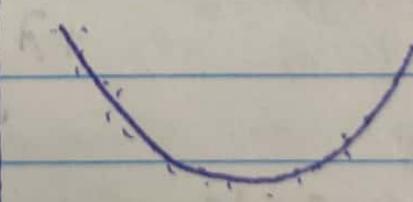
OVERFIT



tries to

make the contribution
of different
powers of 0,

if they
are not
good.



Ex → Let us say we have

$$h_{\theta}(x_i) = \theta_0 x_0 + \theta_1 x_1^2 + \theta_2 x_2^3$$

⇒ Then, the mse becomes -

$$\text{mse} = \frac{1}{N} \sum_{i=1}^N (y_i - h_{\theta}(x_i))^2 + \lambda \underbrace{\sum_{j=1}^n \theta_j^2}_{\text{Penalty term for coeff. of terms.}}$$

L2 regularization (Ridge regression).

$$\text{mse} = \frac{1}{N} \sum_{i=1}^N (y_i - h_{\theta}(x_i))^2 + \lambda \sum_{j=1}^n \theta_j^2$$

L1 regularization (Lasso Regression)

$$\text{mse} = \frac{1}{N} \sum_{i=1}^N (y_i - h_{\theta}(x_i))^2 + \lambda \sum_{j=1}^n |\theta_j|$$

Okay, we will have →
partial derivatives wrt.
each parameter & then
we start updating.

Ex. $h(x) = c + m_1 x_1 + m_2 x_2 + m_3 x_3$

$N = 300$, $\text{Iterations} = 10$

1. Initialize c, m_1, m_2, m_3 as 0.

2. SGD for $i = 10$,

3. Dm & $Dc \rightarrow m_1, m_2, m_3$ & c .

$$\frac{\partial \text{mse}}{\partial m_1}, \frac{\partial \text{mse}}{\partial m_2}, \frac{\partial \text{mse}}{\partial m_3}, \frac{\partial \text{mse}}{\partial c}$$

Put values of x_1, x_2, x_3 for
 $i = 10$ in h & then

update m_1, m_2, m_3 & c

$$m_1 = m_1 - L \times Dm_1$$

$$m_2 = m_2 - L \times Dm_2$$

$$c = c - L \times Dc$$

LOGISTIC REGRESSION

Ways to express probability:

$$p(0_1) = p$$

$$\bar{p}(0_1) = 1 - p = q.$$

standard prob.

notation

values

$$p$$

$$p=0 \quad p=0.5 \quad p=1$$

$$q=1 \quad q=0.5 \quad q=0.$$

odds

$$p/q$$

$$0 \quad 1 \quad \infty$$

log. prob.

$$\log(p/q)$$

$$-\infty \quad 0 \quad \infty$$

LOG ODDS

→ if the prob. of selection is equal,

$$\text{logit}(p/q) = 0, \text{ log odds} = 0$$

⇒ if the odds in favour of one event

$$\text{is } +\gamma, \quad x = \log(p/q)$$

then odds in

$$\text{disfavour of other event } q \text{ are } n = -\log(q/p)$$

event q are

$$-n \rightarrow x = \log(q/p).$$

From prob. to log odds.

$$z = \ln \frac{p}{1-p} \rightarrow \text{logit of } p.$$

$$c = \frac{p}{1-p} \quad z \text{ - is a hypothesis}$$

$$(1-p)e^z = p$$
$$e^z = \frac{p}{1-p}$$

$$p = \frac{e^z}{1+e^z} = \frac{1}{1+e^{-z}}$$

$$P = \frac{1}{1+e^{-z}} \rightarrow \text{logistic fn}$$

How to use as classification model?

→ if for a value $z > z_1$, $p = 0.7$.

Now, our basis

→ $p > 0.7 \rightarrow 1$

→ $p < 0.7 \rightarrow 0$.

→ $z_1 < z \rightarrow p > 0.7, \rightarrow 1$.

→ $z_1 > z \rightarrow p < 0.7, \rightarrow 0$.

$$\Rightarrow \begin{cases} z > z_1, & 1 \\ z < z_1, & 0 \end{cases}$$

Now, this model uses a multi-dim space \rightarrow

$$z = \alpha + \beta \cdot x$$

$$z = \alpha + [\beta_1 x_1 + \beta_2 x_2 + \dots + \beta_k x_k]$$

\hookrightarrow Now, map z to the range 0 to 1 using logistic fn -

$$p = \frac{1}{1 + e^{-z}}$$

Now, let us say we have

A point $z_1 \rightarrow (x_1, x_2, \dots, x_3)$
 \downarrow
 $\begin{pmatrix} 3 & 4 & 10 \end{pmatrix}$

Here, we will have a decision boundary

or better, $x_1 > 3$

any combo. $x_2 > 4$

which gives $x_3 > 10$

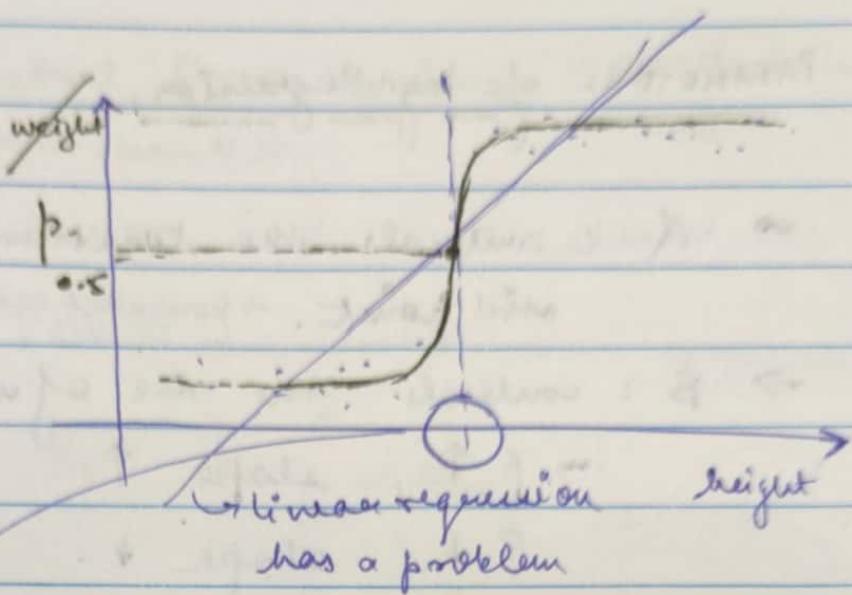
$$f(z) > 0.7$$

Gas 1

&

$$f(z) < 0.7 \text{ as } 0.$$

Eg.



→ In logistic regression,
we just input the height.

→ When $f(H) > 0.5$

→ OBESSE

else

this

provides a

decision boundary

in feature space. → means

$$z_0 = \alpha + \beta_1 x_1 + \beta_2 x_2 + \dots$$

- z_0 is boundary value.
- $p(z_0)$ is threshold.

Any combo. of input vals such that
 z value is $\geq z_0$, classified
as 1, else 0.

PARAMETERS of log Regression.

- α : controls the location of mid point. → how much shift has to be there.
- β : controls the size of model.
 - $\beta \uparrow$, slope \uparrow .
 - $\beta \downarrow$, slope \downarrow .

Adjusted R-squared

- ↳ This is used to take into account the dimensionality of model too.
- ↳ Not necessary that the higher parameter model is better.

	R^2	adj. R^2
1 param	0.67	0.63
2 params	0.83	(0.82) → maybe the best
3 params	0.88	0.7

④ ASSUMPTIONS for LINEAR MODELS to work.

④ There

↓
errors must hold true for
linear model to WORK correctly

Constant Error Variance - Needs to hold
true

↳ Homoscedasticity.

↳ Variance in the errors should be

similar.

constant variance

varying variance

Eg.

✓ P1



↳ same

variance.

The model behaves in

fit correctly in \rightarrow first case

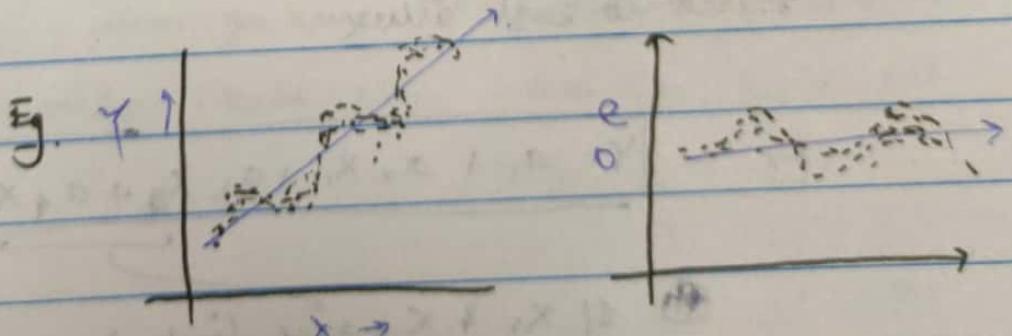
But NOT in [second]

Independent Error Terms. - Needs to hold
true.

Each point must have an independent
error (not dependent on any other points).

↳ Errors must be UN-CORRELATED.

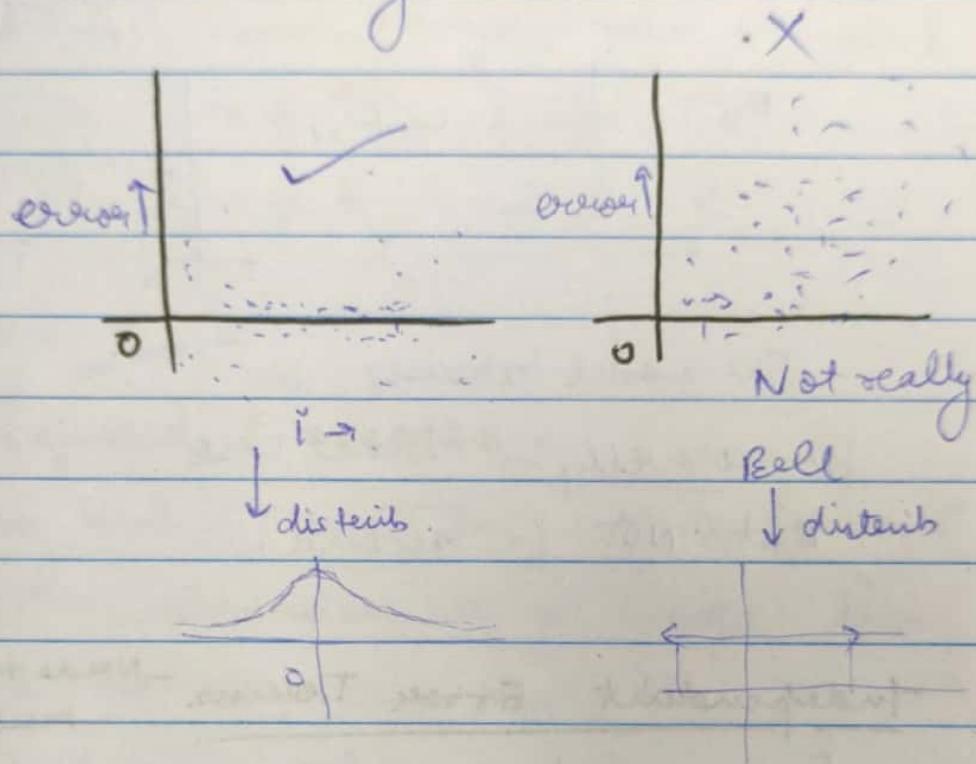
Eg. $y = 1$



Linear model does not work very
well.

Normality of errors:

most of the points must lie around the best fit line to for the model to work nicely.



No multi-collinearity

Multi-collinearity occurs when X 's are related to each other.

$$\text{Ex} \quad Y = a_1 + a_2 X_1 + a_3 X_2 + a_4 X_3$$

- If X_2 & X_3 are linked or correlated, then if we include only 1, then it is going to be redundant.

→ We try to avoid collinearity amongst vars,
* only keep one of 2 vars if they are
highly related.

$$\text{Eg. Motor acc} = \beta_0 + \beta_1(\text{caes}) + \beta_2(\text{no.of residents}) + \varepsilon$$

* If caes & no.of residents , then we may
only predict Motor accidents ~~as~~ only
by caes or no.of residents.

→ we can ~~only~~ also try to remove
the variable & then try to see how
model behaves.

(i) Variable Inflation Factors
associated with removed features.

NEW THOUGHT : A feature inclusion would matter
when you have 100 features & a billion
records. Data size changes by a lot.

EXOGENITY : Dependence on unknown features.

There may be other features or variables
affecting the dependent variable.

When these features are omitted, model
doesn't behave very well.

$$\text{Eg} \quad \text{salary}_i = c_0 + c_1(\text{years of edn}) + \epsilon_i$$

↳ This is not a very good inference.

↳ We may know that the years of
edu affects the salary but it
may not be a full CAUSE for the salary.

ADJUSTED R² → detail.

→ It penalizes ns for adding features
[which DO NOT HELP] in PREDICTING
the dependent variable.

→ The formula for adjusted R² is \rightarrow

$$R^2 = 1 - \frac{\frac{SS_{res}}{(d-1)}}{\frac{SS_{tot}}{(n-(p)-1)}} \rightarrow$$

deg. of freedom
 $(n-(p)-1)$
 \rightarrow p removed features.
 $\rightarrow (n-1) \rightarrow$ total deg.

Eg. Features = 10 df freedom.

Removed 1, p = 4

$$df = 10 - 4 - 1$$

$$\boxed{df = 5}$$

$$\boxed{dt = 9}$$

using R^2 itself for $R^2_{adjusted} \rightarrow$

$$R^2_{adjusted} = 1 - \frac{(R^2 - 1)}{...}$$

write...

Challenges in ML

→ Avoid OVER-LEARNING.

① Hold Out: keep some test data & don't even show that data to model.

② Cross validation:

③ Data Augmentation: to handle class imbalance

④ Cause of Dimensionality

⑤ L₁, L₂ regularization.

⑥ Feature Selection & Reduction:

we just select out from the set

$$\left(A_1, A_2, A_3 \right) A_4$$



take these.

Here, we COMBINE variables

& transform values.

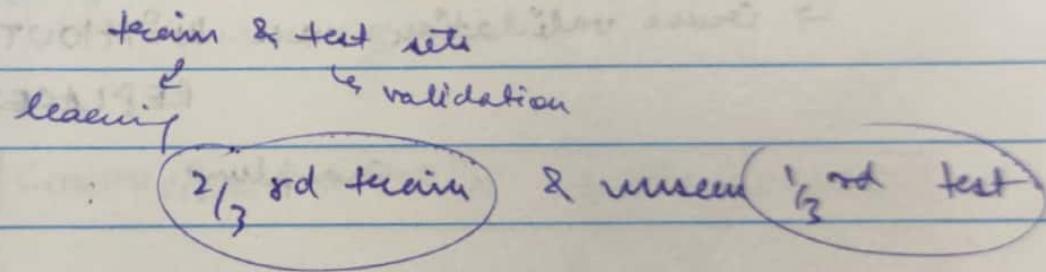
$$A_1, A_2, A_3, A_4$$

$$\left. \begin{array}{l} A_1 \times A_2 \\ A_3 \times A_4 \end{array} \right\}$$

→ MODEL EVALUATION.

→ leave one out : Train on all but one & keep on doing this.

→ Hold out : Divide dataset into 2 sets:



: Repeated Holdout.

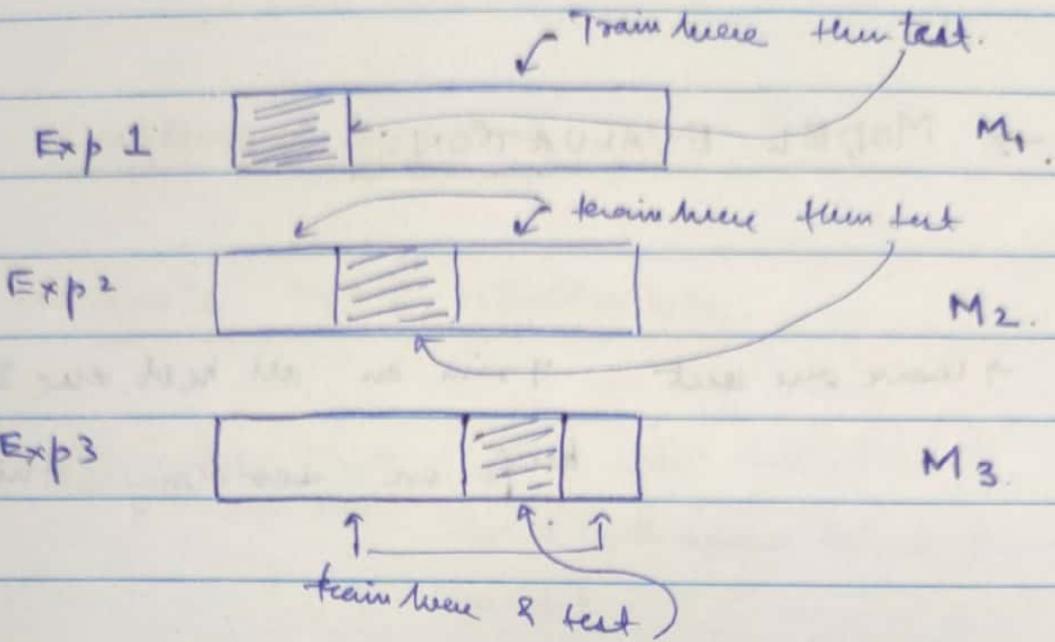
→ k-fold Cross validation:

→ Divide data into ~~equal~~ parts.

→ ONE for test

→ OTHER for training.

→ Process is repeated for all folds & accuracy is averaged.



BOOTSTRAP SAMPLING

→ Cross validation uses WITHOUT REPLACEMENT sampling.

→ Bootstrap sampling is WITH REPLACEMENT to form the training set.

→ n size : TRAIN set.

→ Take out n instances with replacement

→ Use this data as training.

→ USE instances from original

which are NOT IN NEW TRAIN SET

Evaluation Metrics

→ **ACCURACY** → $\frac{\text{correct classified}}{\text{total no. of instances}} \times 100$.

↳ assumed equal cost for each of the classes.

↳ what does that mean?

→ 1000 total → 990 healthy, 10 unhealthy

↳ 99% accurate if we just predict healthy.

↳ This is wrong...

→ **CONFUSION MATRIX & F1-SCORE**

		predicted	
		T	F
Actuals	T	True pos. Falses	False Negative → Overlooked Dangerous.
	F	False Positive	True Negative

↓

false alarms

True Positive Rate : $\frac{TP}{TP+FN}$ = TPR.

why? because
 $FN \uparrow$,
TPR \downarrow

False Positive Rate : $\frac{FP}{FP+TN}$

why? because

$TN \uparrow$, false
negative
mean False
Positive

SENSITIVITY

$\downarrow \rightarrow \frac{TP}{TP+FN}$ = How correctly
we detect positive
classes.

SPECIFICITY

$\downarrow \rightarrow 1 - \text{False Positive rate}$
True Negative
detection = $\frac{TN}{TN+FP}$

		Pos	Neg. & Predicted
Actuals	Pos	22	0
	Neg	17	0.

$$\text{True Pos} = 22$$

$$\text{True Neg} = 0$$

$$\text{False Pos} = 17$$

$$\text{False Neg} = 0$$

$$\rightarrow \text{Accuracy} = \frac{22}{22+0+17+0} = \frac{22}{39}$$

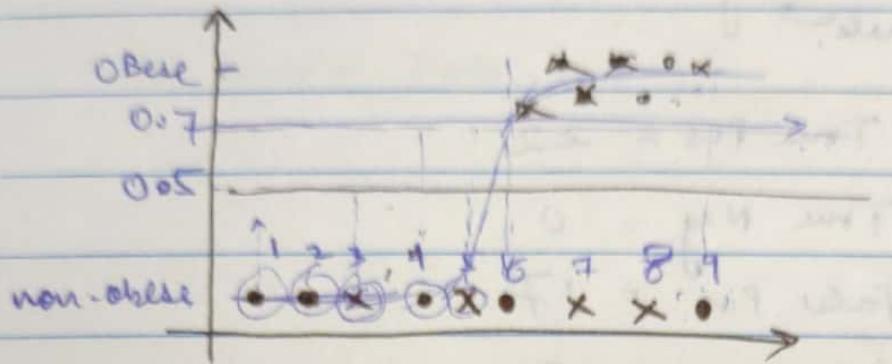
or sensitivity \rightarrow True Positive Rate

$$\Rightarrow \frac{\text{TP}}{\text{TP} + \text{FN}} = \frac{22}{22} = 1.$$

or Specificity \rightarrow True Negative Rate

$$\Rightarrow \frac{\text{TN}}{\text{TN} + \text{FP}} = \frac{0}{0 + 17} = 0$$

Logistic Regression : cont.



when we have a threshold of 0.5 this is the particular graph.

Case 1 $\rightarrow 0.5$ threshold - how does it effect?

		Predicted	
		0	1
Actual	0	3 (1, 2, 1)	2 (6, 9)
	1	2 (3, 5)	2 (7, 8)

where non-observe classified as observe.

0 \rightarrow Non-Observe

1 - Observe

Case 2 $\rightarrow 0.7$ threshold - how does it affect.

		Pred	
		0	1
Actual	0	4 (1, 2, 1)	1 (9)
	1	2 (3, 5)	2 (7, 8)

→ CONCLUSION

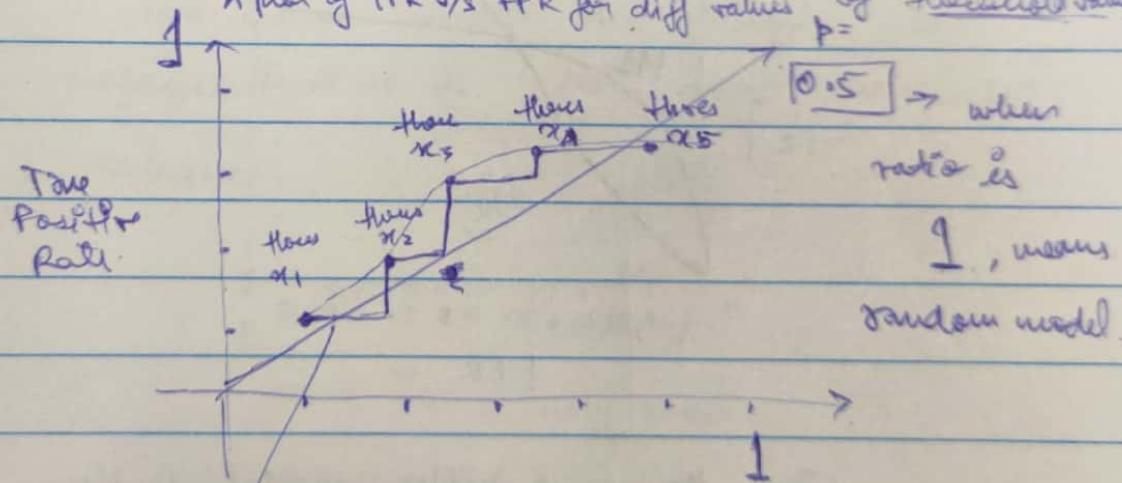
For the same model, if I CHANGE THE THRESHOLDS, the CONFUSION MATRIX changes.

→ This is nice !!

- ④ How to do this effectively, instead of changing threshold & calculating confusion matrix again & again,
PLOT !

ROC curves. → Receiver Op^{er} Characteristic

A plot of TPR v/s FPR for diff values of threshold value $p =$

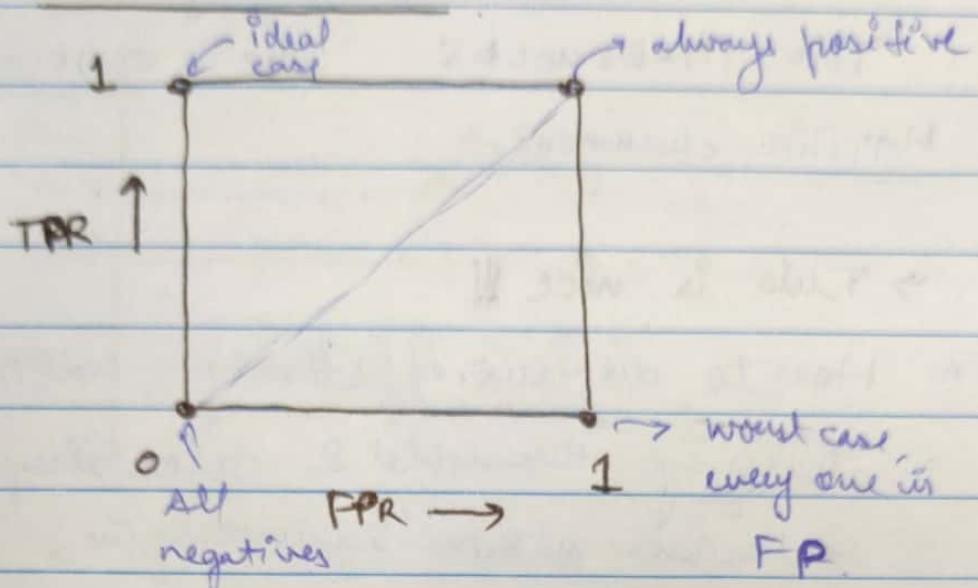


④ thresholds need not be related.

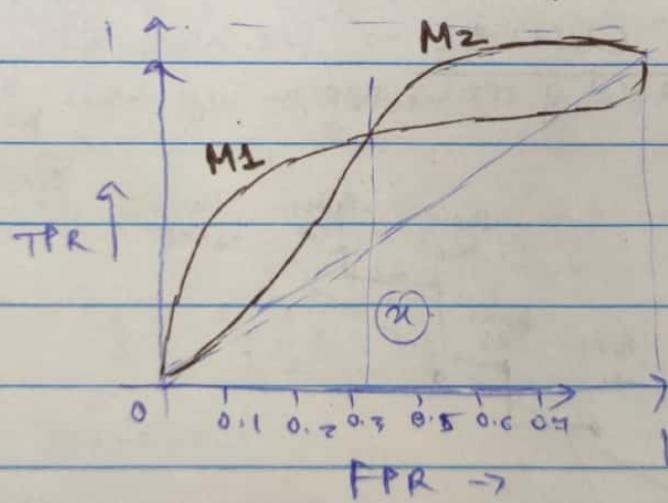
True Positive Rate

False Positive Rate
Each point is a model evaluation against a prob threshold p .

Roc curve characteristics.



Roc curve evaluation.



④ M₂ is a better model in the sense that it has a LESSER FPR after a point for & greater TPR.

④ If we want a very high TPR & low FPR, then M_1 would be better as M_1 would although not give a very high TPR (BUT FPR is very low).

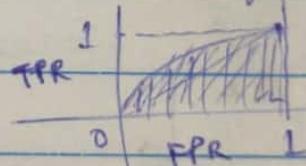
less False Alarms $\rightarrow M_1$.

Better TPR $\rightarrow M_2$.

(regardless of what is FPR).

AREA under the ROC curve

The area under the graph is representative of how good the model behaves.



0.9 - 1 - Excellent

0.8 - 0.9 - very good.

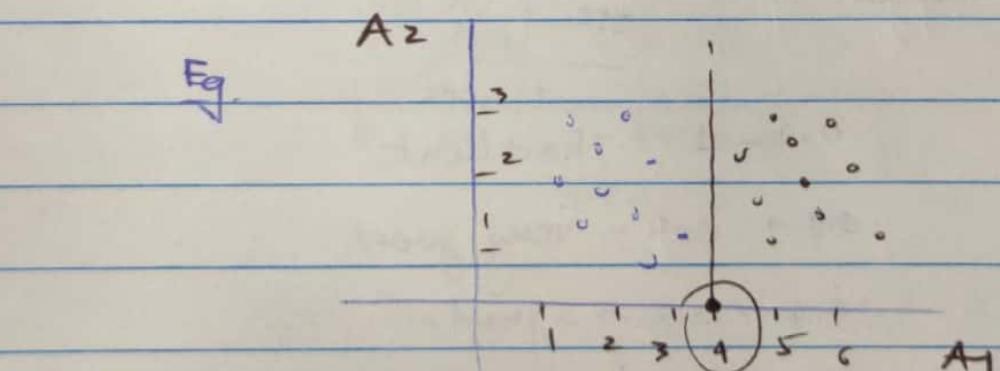
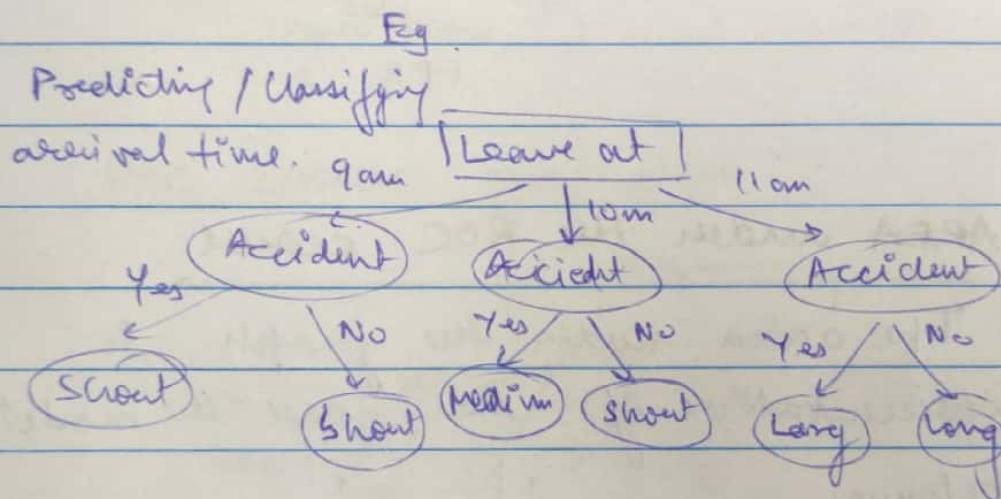
0.7 - 0.8 - good

< 0.7 - okayish ..

DECISION TREES

- ① It is an Inductive learning task
↳ we use a dataset to find a
SPLIT POINT.

- ② Each ~~non~~ non-leaf node represents
a DECISION.

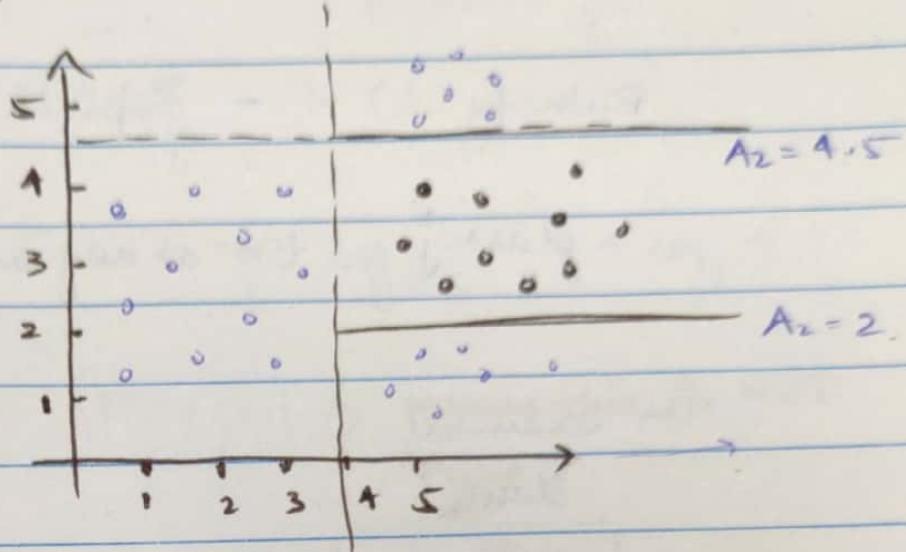


Decision trees try to
find the SPLIT POINTS

$A_1 = 4$ is a DECISION
BOUNDARY.

Examples

$$A = 4$$



Decision Tree

Decision Trees try to find regions of Homogeneity of data.

How is homogeneity defined?

→ A measure is ENTROPY of data.

ENTROPY is the metric which is used for finding out split points.

→ The attrib. value which gives the min. entropy is taken as a DECISION BOUNDARY.

ENTROPY

$$\text{Entropy}(t) = - \sum_j p(j|t) \log_2(p(j|t))$$

class j , $t \rightarrow$ at any node t .

For example.

classes.

$c_1 \rightarrow 0$
$c_2 \rightarrow 1$

How, feature set does not matter.

$$① \rightarrow E(S) = - \sum_j p(j|t) \log_2(p(j|t))$$

$$p(c_1) = 0 \text{ for all,}$$

$$p(c_2) = 1 \text{ for all.}$$

$$E(S) = - (0 \times \log_2(0) + 1 \times \log_2(1))$$

$$\boxed{E(S) = 0}$$

$$\textcircled{2} \quad \begin{aligned} C_1 &= 2 & p(C_1) &= \frac{2}{6} \\ C_2 &= 4 & p(C_2) &= \frac{4}{6} \end{aligned}$$

Entropy - ?

$$E(S) = - \left[2 \log_{2/6} \frac{2}{6} + 4 \log_{2/6} \frac{4}{6} \right]$$

$|E(S) \neq 0| \rightarrow$ Means this is NOT Pure.

EXAMPLE

Let us say we have a data set →

Day	outlook	Temp	Humid	wind	Play
Sunny,	Hot, mild,	High	weak		
overcast,	cool	Normal	steamy		Yes/No
Rain					

How do we determine split point?

For each attribute, we can calculate the entropy value & $GAIN$

$GAIN(S, A) = Entropy(S) -$
 At an attribute \downarrow Attribute
 how much data \downarrow
 do we gain.

$$\sum \left[\frac{|S_V|}{|S|} \times \text{entropy}_{\text{value } v \text{ of } A} \right]$$

Eg. Gain (S, out)

= Entropy (S) - sum of all entropies for the attrib. out.

$$G(S, \text{out}) = \text{Entropy}(S) - [\text{Entropy}(\text{sunny}) + \text{Entropy}(\text{overcast}) +]$$

$E(S)$ is just the ENTROPY of data set (class values)

① $E(S) = - \left[\frac{9}{14} \log \left(\frac{9}{14} \right) + \frac{5}{14} \log \left(\frac{5}{14} \right) \right]$

$E(S) = 0.939$

② Entropies for each attribute value.

OUTLOOK $E(\text{sunny}) + E(\text{overcast}) + E(\text{rain})$

How to calculate ↑ ?

For each attribute value calc. the no. of yeses & nos.

	Yes	No	Total
sunny	2	3	5
overcast	4	0	4
Rain	3	2	5

how many data rows have sunny outlook & play as Yes.

Now, entropy for sunny \rightarrow

$$E(\text{sunny}) = \frac{5}{14} \left[\frac{2}{5} \log \frac{2}{5} + \frac{3}{5} \log \frac{3}{5} \right]$$

5 because we are considering the case in which the outlook is sunny.

$$E(\text{overcast}) = \frac{4}{14} \left[\frac{4}{4} \log \frac{4}{4} + 0 \right]$$

$$E(\text{Rain}) = \frac{5}{14} \left[\frac{3}{5} \log \frac{3}{5} + \frac{2}{5} \log \frac{2}{5} \right]$$

TEMPERATURE

$$\text{Gain}(S, \text{temp}) = E(s) - [E(\text{hot}) + E(\text{mild}) + E(\text{cold})]$$

	Yes	No	Total
hot	2	2	4
mild	4	2	6
cold	3	1	4

$$E(\text{hot}) = \frac{4}{14} \left[\frac{2}{4} \log \left(\frac{2}{4} \right) + \frac{2}{4} \log \left(\frac{2}{4} \right) \right]$$

$$E(\text{mild}) = \frac{6}{14} \left[\frac{4}{6} \log \frac{4}{6} + \frac{2}{6} \log \frac{2}{6} \right]$$

$$E(\text{cold}) = \frac{4}{14} \left[\frac{3}{4} \log \frac{3}{4} + \frac{1}{4} \log \frac{1}{4} \right]$$

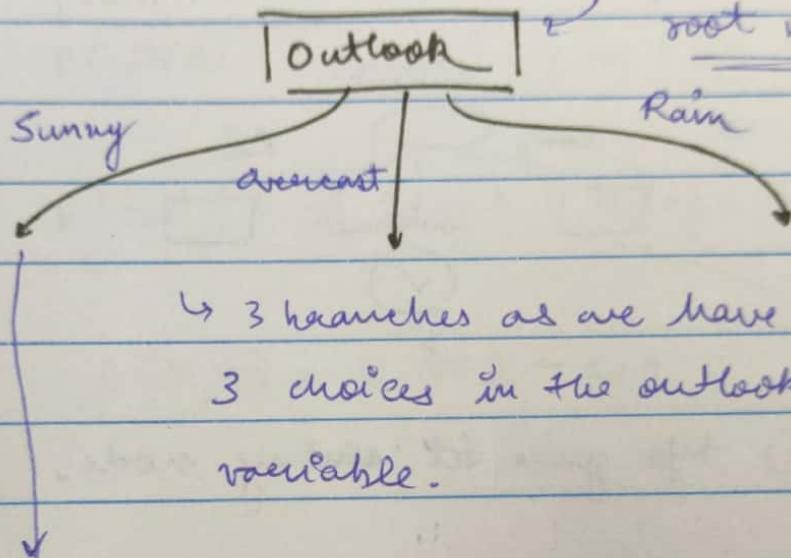
* All of the above calculation is to find a split point at the root level!

$$IG(S, \text{outlook}) = 0.247 \rightarrow \text{highest } IG.$$

$$IG(S, \text{temp}) = 0.029$$

$$IG(S, \text{hum}) = 0.152$$

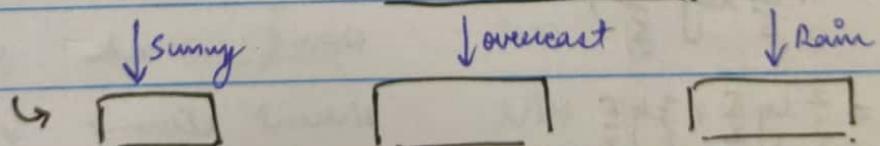
$$IG(S, \text{wind}) = 0.018$$



this branch means we have all the elements having 'SUNNY' outlook.

Now, dataset is reduced. How?

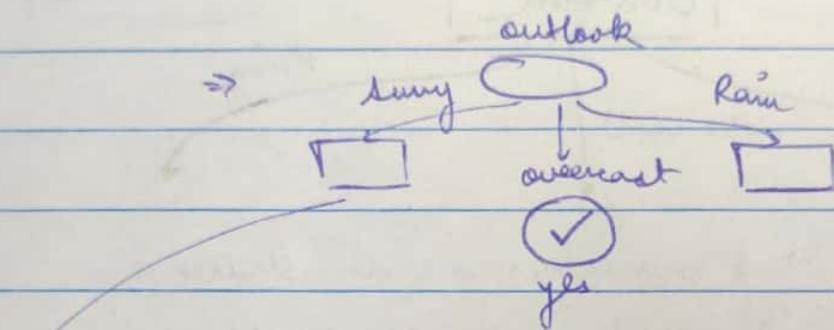
We take SUNNY outlook rows!



↳ these boxes are info.gain nodes FOR EACH dataset pertaining to that node.

Some Inferences →

Since each & every row in the
OVERCAST rows have target as
YES, there is NO ENTROPY.



→ ① Info gain at sunny node.

Dataset S'

Temp	Humidity	Wind	Play
------	----------	------	------

Hot	High	Weak	No
-----	------	------	----

$$E(S') = -\left(\frac{2}{5} \log \frac{2}{5}\right) \text{ Hot } \text{ High } \text{ Strong } \text{ No}$$

$$+ \left(\frac{3}{5} \log \frac{3}{5}\right) \text{ Mild } \text{ High } \text{ Weak } \text{ No}$$

$$\text{ cool } \text{ Normal } \text{ Weak } \text{ Yes}$$

$$= \frac{2}{5} \log \frac{2}{2} + \frac{3}{5} \log \frac{3}{3} \text{ Mild } \text{ Normal } \text{ Strong } \text{ Yes}$$

$$= 0.2 \times 1.322 + 0.6 \times 0.737 \text{ for this.} \quad \text{Yes} \quad \text{No} \quad \text{Total}$$

① $IG(S', \text{Temp}) \rightarrow$	Hot	0	2	2
	Mild	1	1	2
	cool	1	0	1

$$IG(s', T) = E(s') - [\text{Entropy (Hot)} + \\ \text{Entropy (Mild)} + \\ \text{Entropy (cold)}]$$

$$E(\text{Hot}) = - \left[\frac{0}{2} \log \frac{0}{2} + \frac{2}{2} \log \frac{2}{2} \right] \times \frac{2}{5} = 0$$

$$E(\text{Mild}) = - \left[\frac{1}{2} \log \frac{1}{2} + \frac{1}{2} \log \frac{1}{2} \right] \times \frac{2}{5} + \left[\frac{2}{5} \right] = 0.4$$

$$E(\text{cold}) = - \left[\frac{1}{1} \log 1 + 0 \right] \times \frac{1}{5} = 0$$

$$E(s', T) = 0.971 - 0.4 \\ = 0.571$$

② $IG(s', \text{Humidity})$

s'	Temp Yes	Humid No	Wind Yes	
High	0	3	3	
Normal	2	0	2	

$$E(\text{Humidity}) = - \left[\frac{3}{5} \times \left[0 + \frac{3}{3} \log 1 \right] \right] \\ - \left[\frac{2}{5} \left[\frac{2}{2} \log 1 + 0 \right] \right]$$

$$E(\text{Humid}) = 0$$

② E(wind)

	Y	N	T
weak	1	2	3
steady	1	1	2

$$E(\text{wind}) = - \frac{3}{5} \left[\frac{1}{3} \lg \frac{1}{3} + \frac{2}{3} \lg \frac{2}{3} \right]$$

$$= - \frac{2}{5} \left[\frac{2}{3} \lg \frac{1}{2} \right]$$

$$= + \frac{3}{5} \left[\frac{1}{3} \lg 3 + \frac{2}{3} \lg \frac{3}{2} \right]$$

$$+ \frac{2}{5} \lg 2$$

$$= 0.4 + 0.6 \times \left[\frac{1}{3} \lg 3 + \frac{2}{3} \lg \frac{3}{2} \right] - \frac{2}{5} \lg 2$$

$$= 0.4 + 0.6 \left[\lg 3 - \frac{2}{3} \lg 2 \right]$$

$$= 0.4 + 0.6 [1.5849 - 0.6777]$$

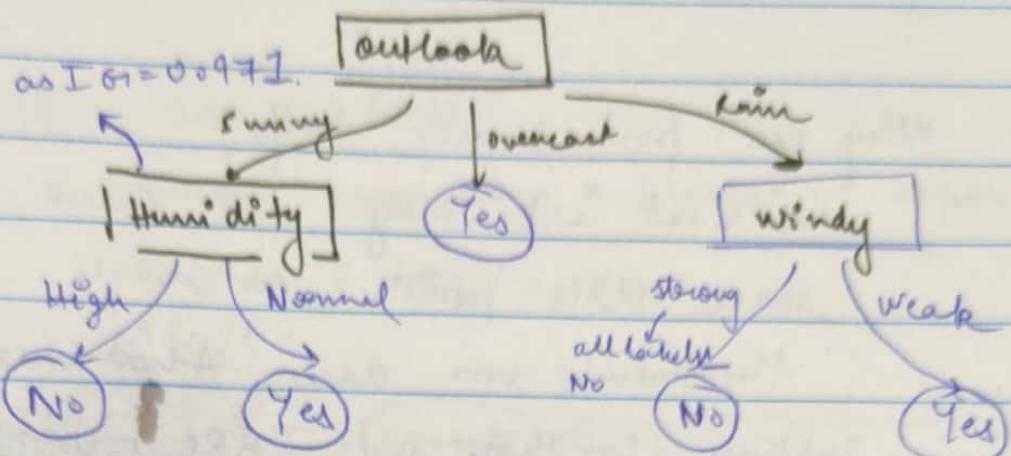
$$= 0.4 + 0.6 [-]$$

$$E(\text{wind}) = 0.94432$$

$$IG(S', \text{Temp}) = 0.571$$

$$IG(S', \text{Humidity}) = 0.971 \quad \checkmark$$

$$IG(S', \text{Wind}) = 0.019$$



You do not need to calculate
 IG if you have a **PURE** branch
 No noise.

What if we have very large no. of features?
 ⇒ It would be very computationally expensive.

PRUNING → If a branch or a root-to-leaf path or a subtree doesn't give any performance increase, we can remove.

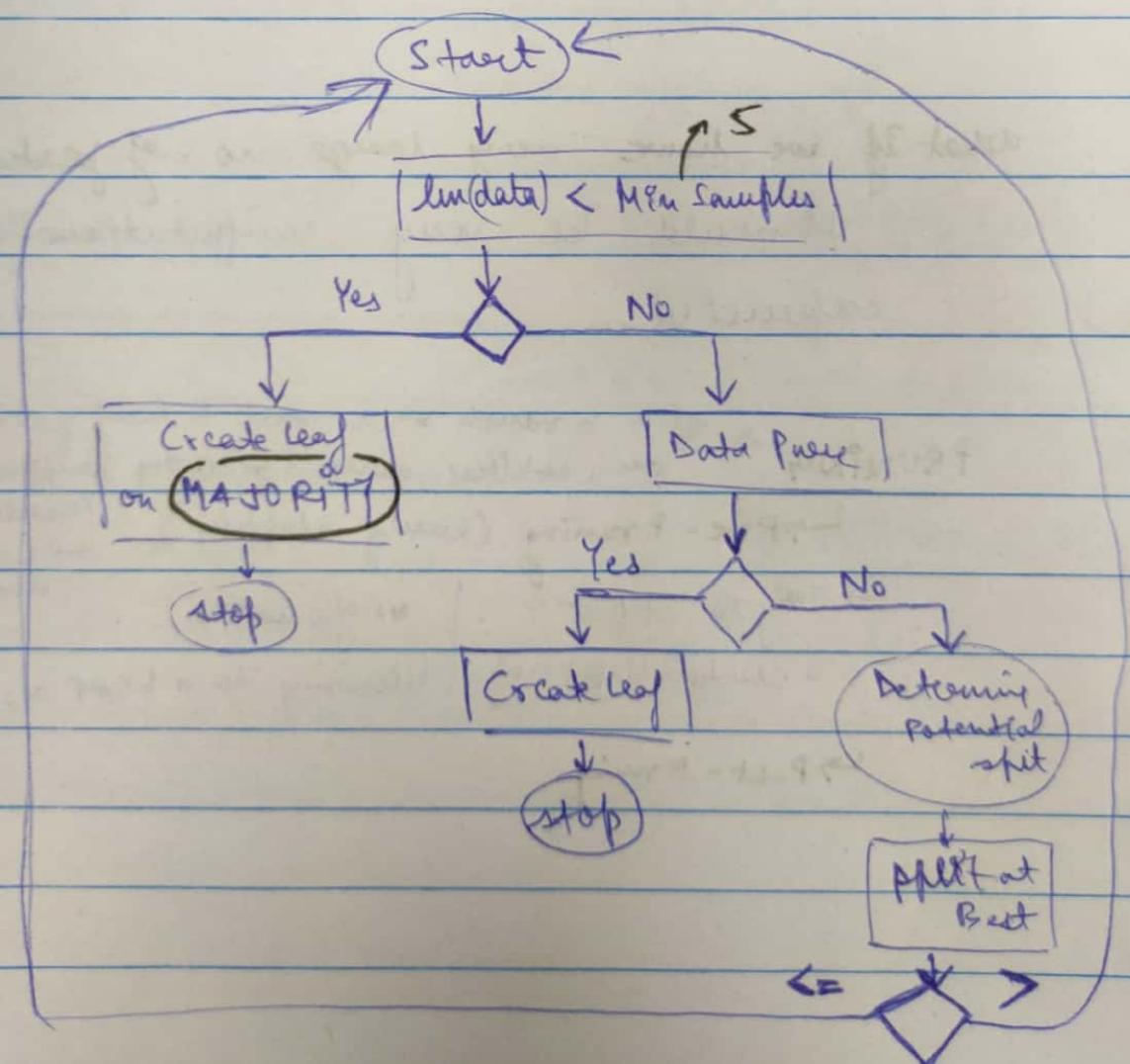
- ↳ Pre-Pruning (early stopping)
- ↳ Try to stop at a certain DEPTH | No. of samples belonging to a LEAF NODE
- ↳ Post-Pruning

Why can post-pruning help?

→ Test set may not have some data points on which the training was done, that means subtrees for those are IRRELEVANT.

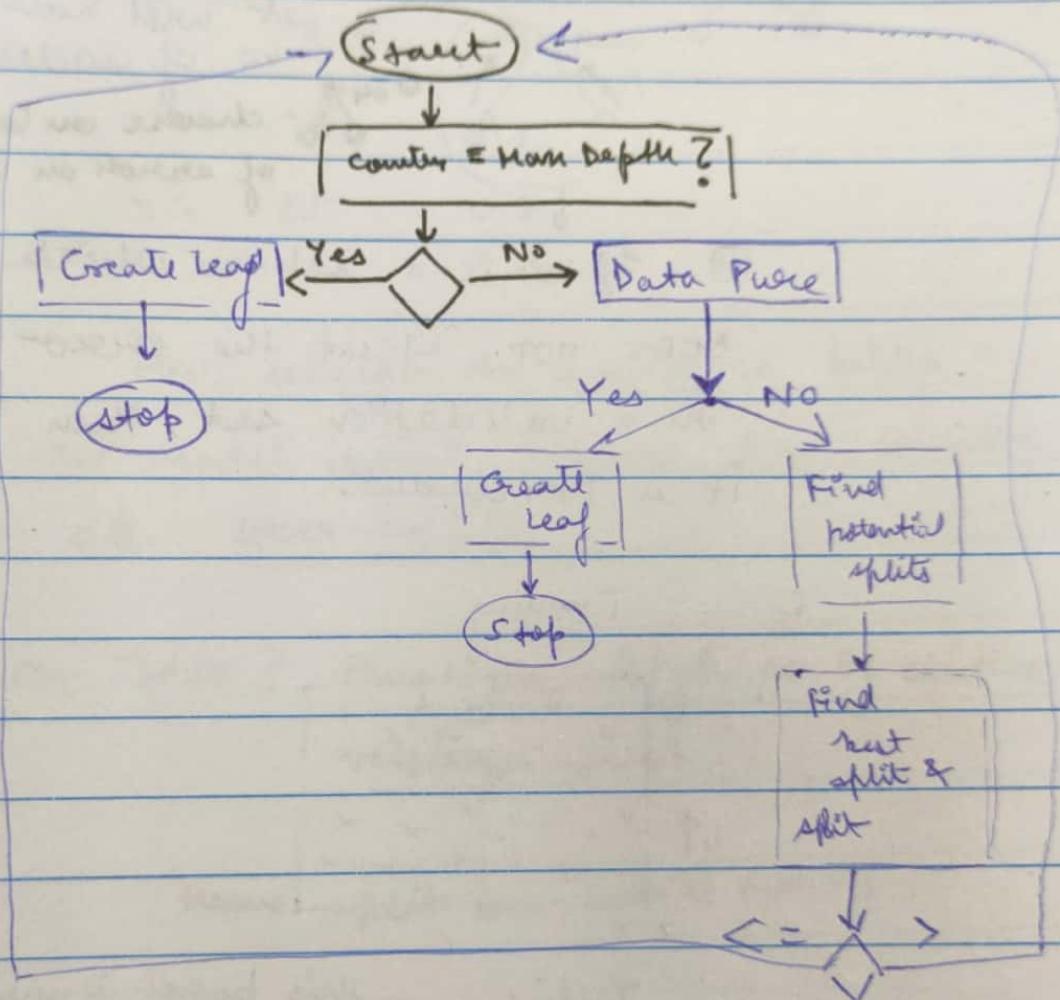
some pre-pruning methodologies.

① Min samples - To limit depth / decision boundaries

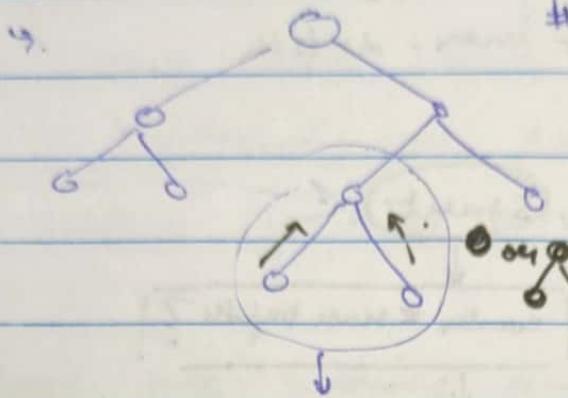


② MAX DEPTH

→ Stop processing a branch if depth reaches $>$ max-depth.



Some past pruning methodologies.



will start from

Last LEVEL as

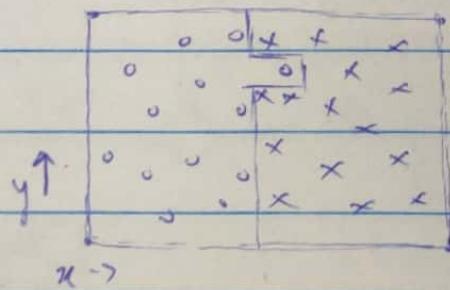
it will have least no. of nodes.

- choose on basis
of scores on test set.

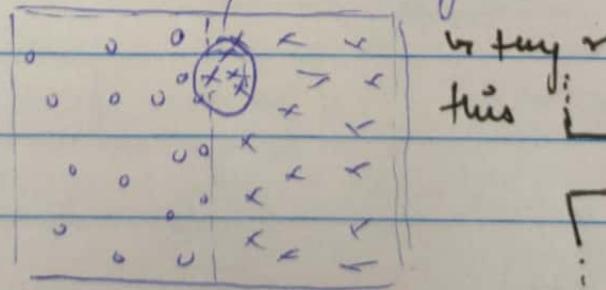
④ If this is a subset which
DOES NOT affect the error
on a validation set, then
it is irrelevant.

Ey

Training



Testing → this bunch is misclassified.



SOME CAVEATS & DIFF. METHODS FOR
SELECTING ATTRIBUTES & SPLITS.

① What $IG_{A_i} = IG_{A_j}$?

If we have same Info Gain for eg.
3 values \Rightarrow values

	A_1	A_2
IG	0.7	0.7

Here, selecting A_2 would be better as
the model would become less complex
at lower splits.

Better Idea? Penalize High no. of distinct
values in attrib.

Means split on GAIN RATIO

$$\text{Gain Ratio}_{A_i} = \frac{\text{Info Gain}_{A_i}}{\text{Entropy}_{A_i}}$$

② GINI INDEX - calculate for any attribute t ,

NOTE: A minimum split gini index is a preferable way for split.

Algos like CART, SPRINT...

Lazy learning / Instance based learning.

KNN Algorithm (K-Nearest Neighbors).

→ Assigning the target features value/
class on the basis of the K-nearest points
in the data set.

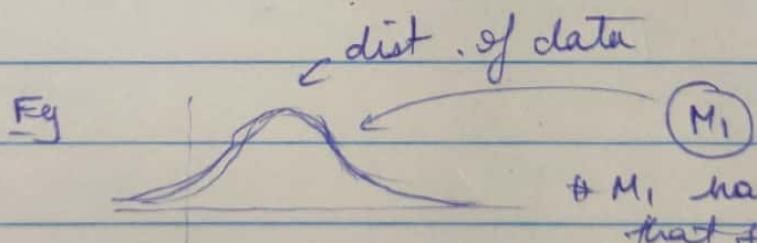
* Supervised learning.

→ Basis is a similarity measure (distance)
→ Stores all available cases & classifies
new cases based on a similarity measure.

There are 2 types of models -

① Parametric - When there is an underlying
assumption for data to come from some distn.

② Non-parametric - No assumptions.

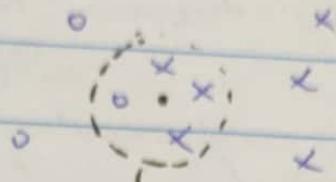


→ M₁ has been told
that the data has
a gaussian distribution

M₂ → no assumption.

\Rightarrow KNN is a NON-PARAMETRIC model.

Eg We have taken $[K=4]$,



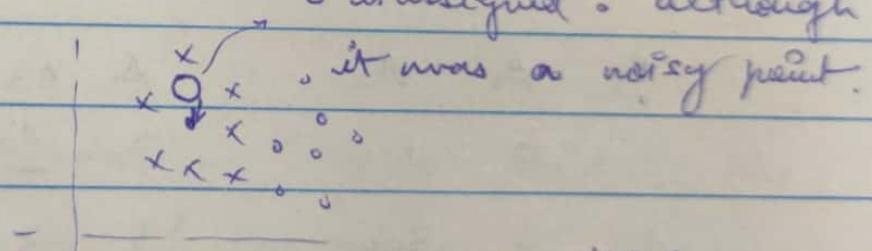
If $K=4$, x is assigned to

How to choose K ?

i) K should not be very small.

why? To avoid classification due to noise.

0 is assigned. although



something

K should be selected as \sqrt{n} .

some distance measures \rightarrow

Euclidean, Manhattan & Minkowski

(3) - KNN : Ex!

use the 3 closest people for classification.

customer	Age	income	No. of credit cards	Class.
	X			Y

→ One thing to note is that

NORMALIZATION of variables is important, why?

(SCALE affects distances).

$$\text{Ex. } \begin{matrix} A_1 & A_2 \\ 2 & 1000 \\ 3 & 3000 \\ 1 & 4000 \\ 5 & 6000 \end{matrix} \quad \begin{matrix} A_1 & A_2 \\ 2 & 1 \\ 3 & 3 \\ 1 & 4 \\ 5 & 6 \end{matrix} \quad \text{scaled down}$$

$$\begin{matrix} & \\ & \\ 7 & x \end{matrix}$$

One way to standardize \rightarrow

$$X_s = \frac{X_i - \min}{\max - \min} \rightarrow \text{Min-Max scale.}$$

Another \rightarrow

$$X_s = \frac{X_i - \text{mean}}{\sqrt{\text{var}}} \rightarrow \text{Std. scale.}$$

CONS - computational time

\downarrow for euclidean distance.

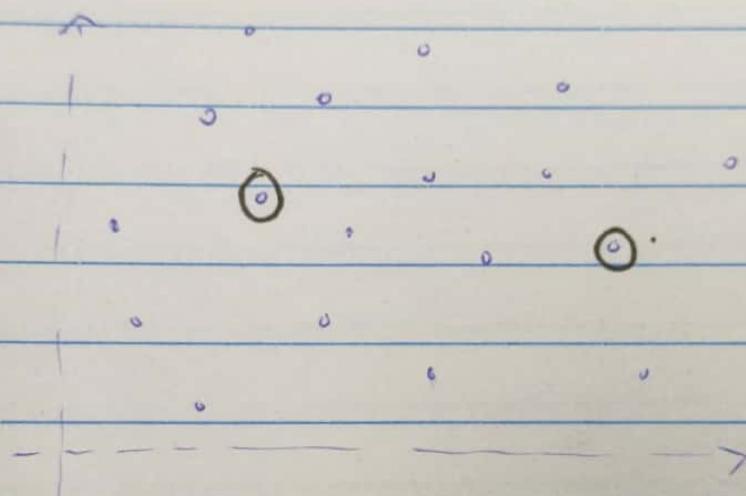
$O(N \log k \times k)$ for finding the first k ~~the~~ minimums out of the training set.

WORK AROUND - Try to reduce the no. of features of the dataset to reduce computational time.

* Also, note that distance measures do not work with categorical vals. Need some data transformation steps.

K-Means Clustering

Selecting the first K-centroids &
iteratively improving upon the cluster
centres.



Init. - Initialize K-data points