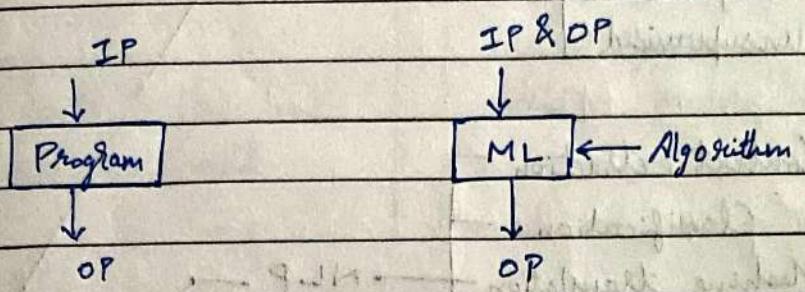


5th August, 2024

Machine Learning



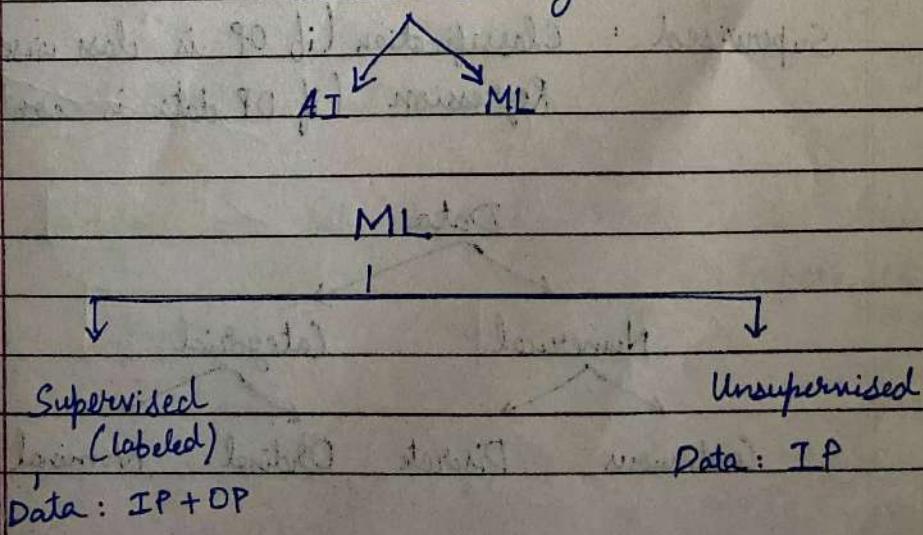
* ML is a subset of AI that focuses on the development of algs & statistical models that enables computer to make predictions based on data.

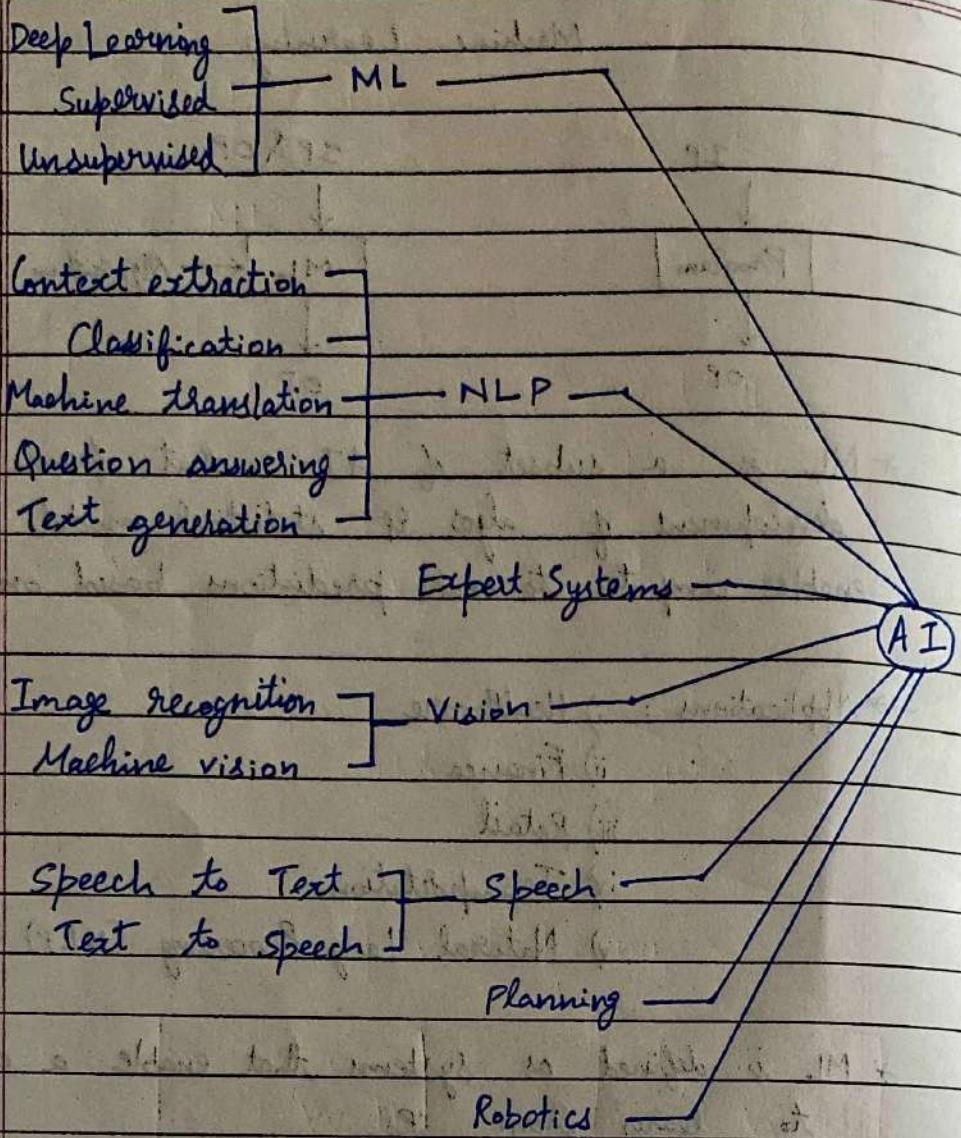
- Applications :
- Healthcare
 - Finance
 - Retail
 - Transportation
 - Natural Lang Processing (NLP)

* ML is defined as systems that enable a computer to learn from IPs.

* AI is composed of systems that allow computers to imitate human cognitive processes.

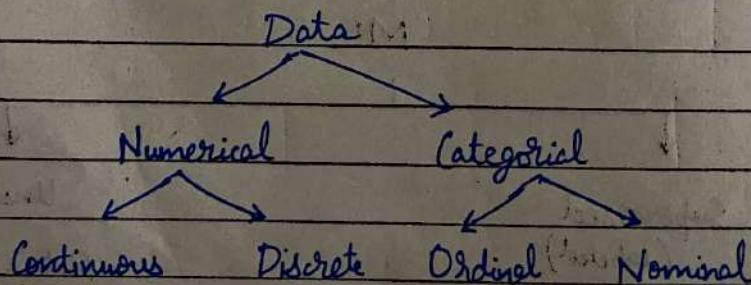
Machine Intelligence





Unsupervised : Clustering
Associativity

Supervised : Classification (if OP is class wise)/category
Regression (if OP data is continuous)



Training

- * using data
to create & fine tune
a predictive model

Prediction

- * answer questions
predictive model is used
to derive up predictions
based on unseen data

→ 7 steps to ML

- i) Gathering data
- ii) Data preparation
- iii) Choosing a model
- iv) Training
- v) Evaluation
- vi) Parameter tuning
- vii) Prediction

→ PTE V ML Model Definition

ML is the study of algs that:

- improve their PERFORMANCE (P)
- at some TASK (T)
- with EXPERIENCE (E)

6th August, 2024

* if we have labeled data along with labels
⇒ Supervised learning

Unsupervised learning

Semi-supervised learning

small amount: labeled data
large amount: unlabeled data

Semi supervised - only a portion of data is labeled
Reinforcement learning

→ Types of ML

i) Supervised (inductive) learning -

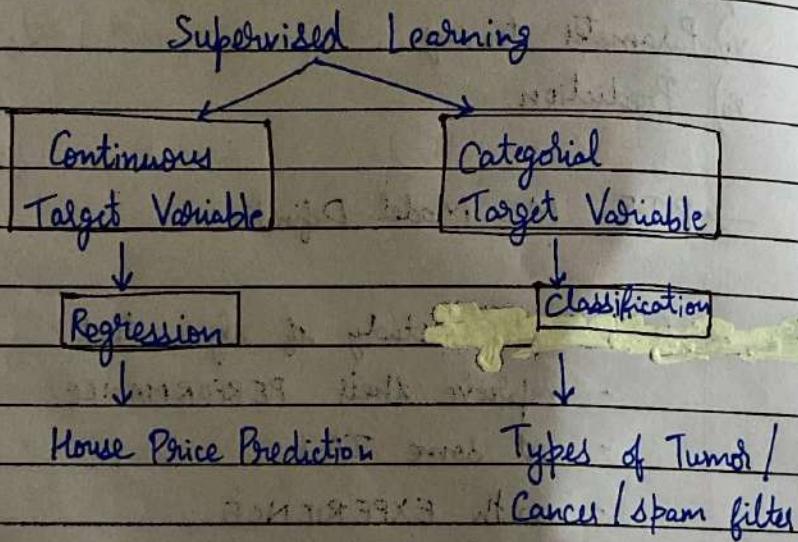
training data + desired outputs (labels)

ii) unsupervised learning -

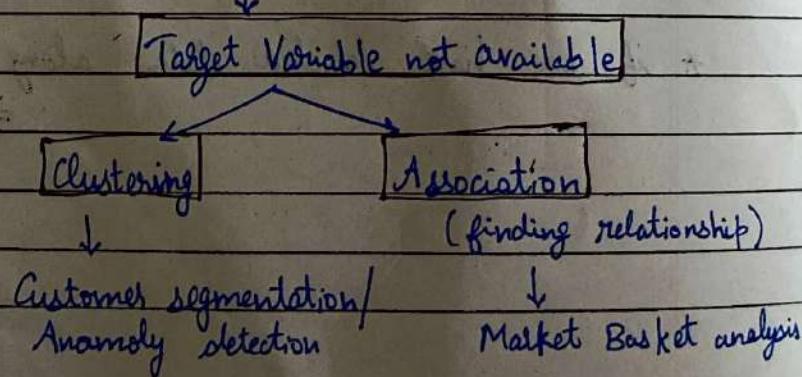
training data (without desired OP)

iii) Reinforced learning -

rewards from sequence of actions



Unsupervised learning / Exploratory learning



- * initially to choose classification / categorical we check OP label

→ Supervised Learning

$$Y = f(X)$$

OP variable IP variable

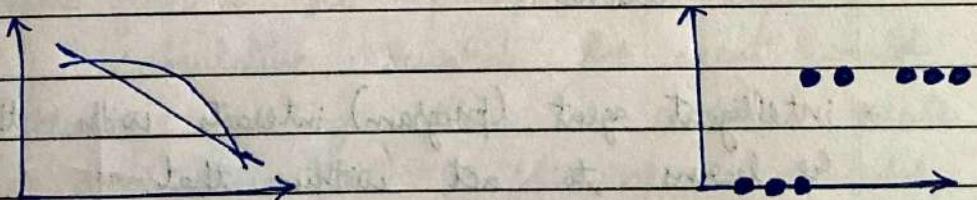
- * algorithm to learn the mapping func. from IP to OP.

- * func. approximation is also called hypothesis func. (an "/ claim")

- * goal : to approximate the mapping func. so well that when you have new IP data (X) you can predict the OP var. (Y) for that data

— Regression

— Categorical - Classification



→ Unsupervised Learning

given - training data (w/o desired OP)

- * unlabeled IP data is fed to ML model to train it.
- * Interpret raw data → hidden patterns → clustering

Association:

An association rule is an unsupervised learning method which is used for finding relationship b/w variables in a large data set.

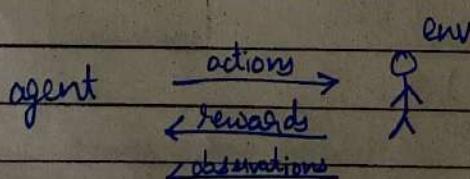
→ Reinforced Learning

- * RL involves lot of computing
- * No training data set is there for training any model.

RL: what to do and how to map situations to actions.

* intelligent agent (program) interacts with the env
if learns to act within that.

* learner is not told which action to take but instead it must discover which action yields max reward



Basic Terms

1. Agent : Learner / decision maker that interacts with the env.
 2. Environment : External system with which the agent interacts.
 3. State : A representation of current situation of the env.
 4. Action : The choices available to the agent at each state.
 5. Policy : A strategy / rule that the agent follows to select actions.
 6. Reward : A scalar feedback signal from env indicating how good / bad the action taken by user was.
 7. Value func : It predicts the expected cumulative reward the agent would receive starting from a certain state (or state-action pair) if following a particular policy.
- * The goal of the agent is to learn the optimal policy that maximizes cumulative rewards.

- Q-table in RL

* look up table - max expected future rewards
for an action

tells us which action is best at each state

* Q func: takes 2 ips - state & action
it uses Bellman's eqn:

$$Q(s_t, a_t) = E[R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} + \dots | (s_t, a_t)]$$

→ Q learning algo

initialize the Q-table



choose an action ←



Perform an action

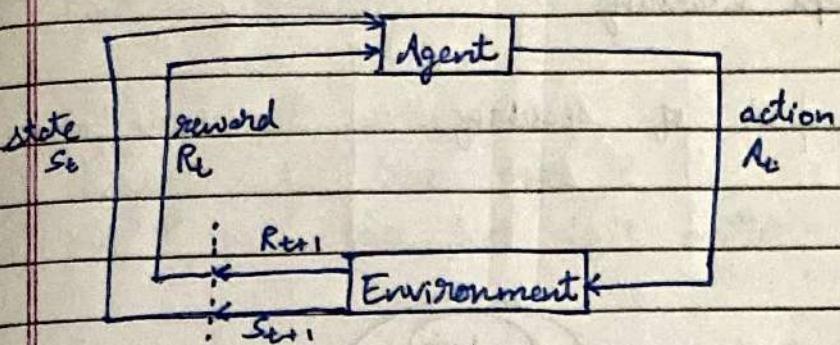


compute the reward

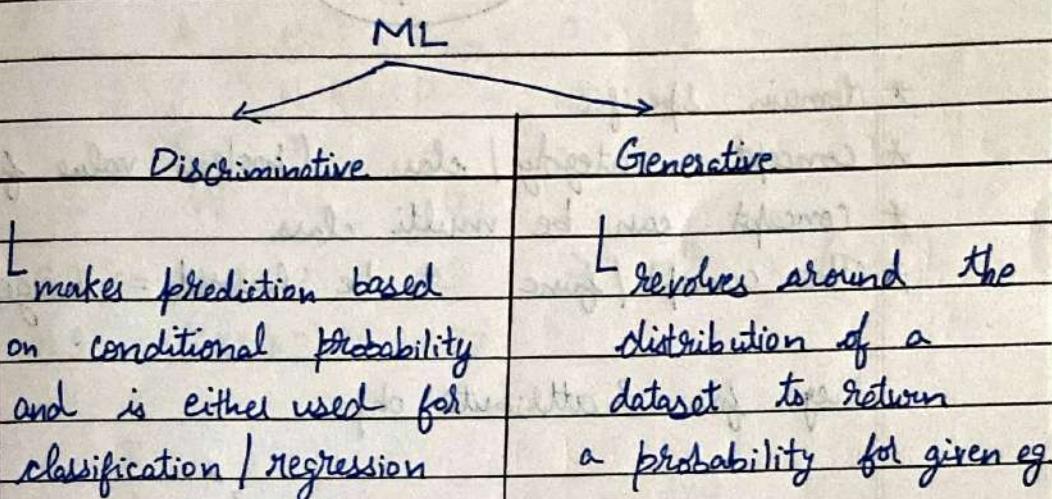


update the reward

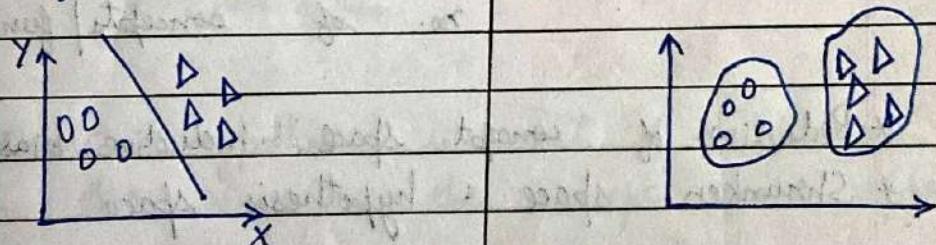
many iterations
to get good
Q-table



Generative and Discriminative models



* Logistic regression = binary classification



- | | |
|---|--|
| <ul style="list-style-type: none"> * draws boundaries in data space * simple data more data | <ul style="list-style-type: none"> * models how data is placed throughout the space * complex data less data |
|---|--|

Types :

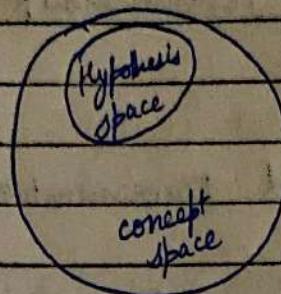
1. Logistic Regression
2. Support Vector Machine
3. Decision Tree
4. Random Forest

Types :

1. Bayesian Network
2. Hidden Markov Model
3. Autoregressive Model
4. Generative adversarial NW

Concept Learning

The prob. of searching



- * domain specific
- * concept: category / class (boolean value func.)
- * concept can be multi class
- * The concept / func. to be learned = target concept

e.g. for 3 attribute obj

for binary feature concept space has 2^d
no. of concepts / functions

- * Reduction of concept space : inductive bias
- * Shrunken space : hypothesis space

→ Inductive Bias / Learning Bias

- * set of assumptions that learner uses to predict or
- * concepts are of a particular kind called conjecture concepts because of which the hypothesis space is restricted. Hence also called as Restriction Bias

Find S algorithm

- initialize all with \emptyset
- choose only true ones
- put ? with ones that don't matter

* Version Space : Set of all consistent hypothesis

8th August, 2024

(A, B) (M, N, O) (C, D) (O, I)

e.g.: X Y Z o/P

A (M C O) (N, D)

B (N D) (M, O)

C (O) (M, N)

(2) (3) (2)

$$CS = 2 \times 3 \times 2 = 12 \Rightarrow 2^{12} (2^2^d)$$

$$HS = (3 \times 4 \times 3) + 1 = 37$$

A: w.r.t. 3rd feature - total boundary

B

A

? → either or

\emptyset → none of them

* Consistent hypothesis is set of assumption that classifies the concept space appropriately

* Accuracy of a model



CONFUSION METRICS

		Actual class		
		Disease (1)	No Disease (2)	
Model Prediction class	+	8	4	Type I error
	-	2	16	
		FN	TN	

Type II error

$$\text{accuracy} = \frac{\text{TP} + \text{TN}}{\text{TP} + \text{FP} + \text{FN} + \text{TN}} \stackrel{\text{(correct Predict)}}{=} \frac{24}{30}$$

Precision

8	4
2	16

$$\text{Recall (Total +ve classes)} = \frac{\text{TP}}{\text{TP} + \text{FN}}$$

$$(\text{True Positive Rate}) = \frac{\text{TP}}{\text{Total Positive}}$$

sensitivity

True by actual

$$\text{Precision (model predicted +ve)} = \frac{\text{TP}}{\text{TP} + \text{FP}}$$

True by model

Balanced data - Equally +ve & -ve in data op



Accuracy ✓

Imbalanced data → Recall, Precision ✓

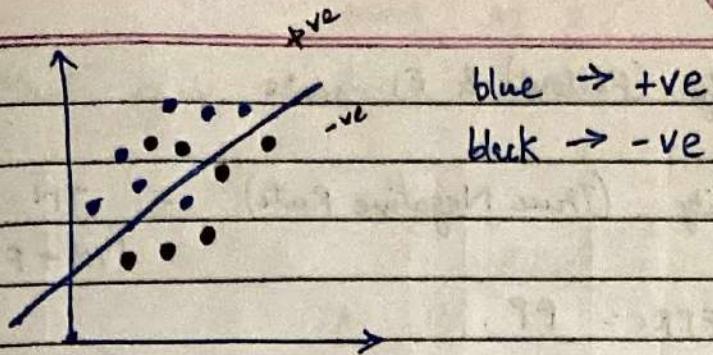
$$\text{specificity} = \frac{\text{TN}}{\text{FP} + \text{TN}}$$

$$\text{False Negative Rate} = \frac{\text{FN}}{\text{TP} + \text{FN}}$$

Type-I error

reduce this

slides 9:



actual	+	-
prediction	(+)	(-)
(+)	6 TP	2 FP
(-)	1 FN	5 TN

Prediction

actual	(+)	(-)
data	(+)	(-)
(+)	6 TP	1 FN
(-)	2 FP	5 TN

→ Precision Recall trade off

↓ scenario when Recall ↓ Precision ↑

eg:	-ve	+ve
(0VA)	X X	O O
	X X	X
		O O

→ Specificity (Fallout) & F1 score

$$\text{Specificity (True Negative Rate)} = \frac{TN}{TN + FP} = 1 - FPR$$

$$FPR = \frac{FP}{TN + FP}$$

$$F1 \text{ score} = \frac{2 * (\text{recall} * \text{precision})}{(\text{recall} + \text{precision})}$$

Harmonic mean

$$= \frac{2TP}{(2TP + FP + FN)}$$

* Higher TNR & a lower FPR is desired since we want to correctly classify -ve class

→ ROC - Received Operating Characteristics

X → FRR (specificity)

Y → TPR (sensitivity)

Plotted with
multiple TPR, FPR

used when threshold

comes into picture: up set by user manually, since model should work correctly for multiple thresholds we need ROC graph.

* Maximum area under curve (AUC)

(more area \Rightarrow better model)

e.g: converting probabilities to Yes/No

set threshold

False Positive \Rightarrow Actual : Actually X
Model : But model ✓

False Negative \Rightarrow Actual : Real ✓
Model : AI X



Multi-class confusion matrix

		Actual			
		A	B	C	D
Predicted	A	100	0	0	0
	B	80	9	1	1
	C	10	0	8	0
	D	10	1	1	9

$$\text{Class A : } TP = 100$$

$$TN = 9 + 1 + 1, + 0 + 8 + 0 + 1 + 1 + 9$$

$$FP = 0 + 0 + 0$$

$$FN = 80 + 10 + 10$$

$$\text{Class B : } TP = 9$$

$$TN = 100 + 10 + 8 + 10 + 1 + 1 + 9$$

$$FP = 80 + 1 + 1$$

$$FN = 0 + 0 + 1$$

$$\text{Class C : } TP = 98$$

$$TN = 100 + 0 + 0 + 80 + 9 + 1 + 10 + 10 + 1 + 1 + 9$$

$$FP = 10 + 0$$

$$FN = 0 + 1 + 1$$

$$\text{Class D : } TP = 9$$

$$TN = 100 + 80 + 9 + 1 + 10 + 8$$

$$FP = 10 + 1 + 1$$

$$FN = 0 + 1 + 0$$

11th August, 2024

Decision Tree - ID3 algorithm

Graphical representation of all possible solns. to a decision

e.g.: IP

Gender	Age	OP
--------	-----	----

F	old	snaf
---	-----	------

M	middle	insta
---	--------	-------

F	young	PUBG
---	-------	------

M	young	PUBG
---	-------	------

F	young	PUBG
---	-------	------

F	middle	snaf
---	--------	------

* Classification model \Rightarrow Decision Tree (works best)

Supervised learning

* Decision Tree creates a tree structure with decision node as root node,

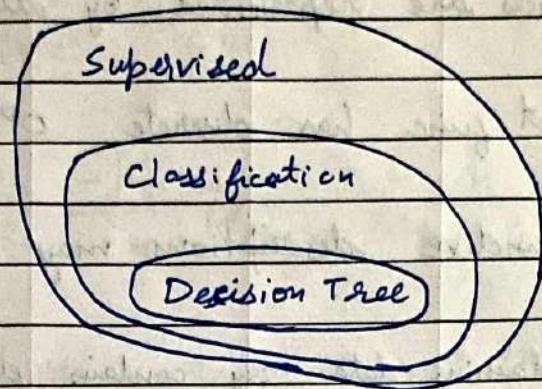
* pruning : makes decision tree less complex

\rightarrow Entropy : tells impurity factor (checks purity of node alone) is not enough

\rightarrow Information Gain :

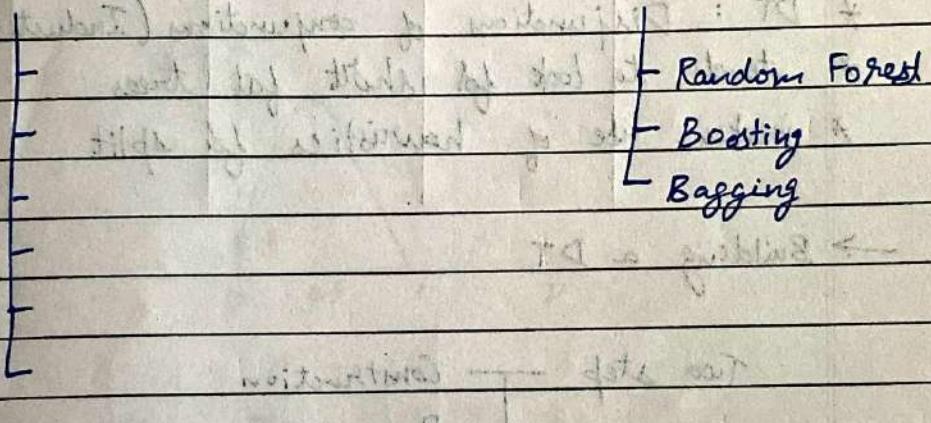
* with numbers DT take more time , CART algo.
(regression)

Gini impurity
Gini index



Base Classifiers
(1 model)

Ensemble classifier
(more than 1 model)

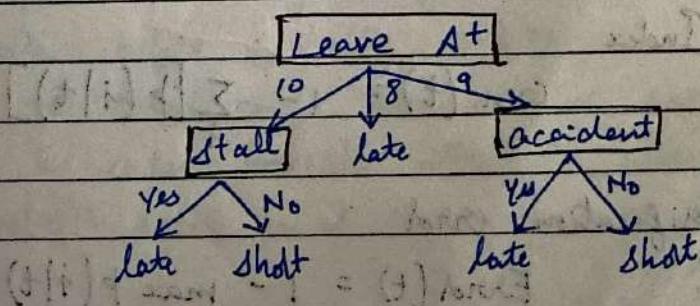


* DT works nicely with text data (NLP)
(decision)

Root, intermediate : IP

Leaf nodes : OP

categories of IP / classes / attributes : edges / branches



→ CART : Classification And Regression Tree

→ Appropriate Problem for DT Learning

- i) instances are represented by attribute-value pairs
- ii) target func has discrete OP values
- iii) disjunctive descriptions may be required
- iv) the training data may contain errors
- v) the training data may contain missing attribute values

- * DT : Disjunction of conjunctions (Inductive Bias)
- * tends to look for short fat trees
- * makes use of heuristics for split

→ Building a DT

Two step — Construction
[Pruning

→ Impurity Measures

1. Entropy

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

2. Gini Index

$$\text{Gini}(t) = 1 - \sum_i [p(i|t)]^2$$

3. Misclassification error

$$\text{Error}(t) = 1 - \max_i p(i|t)$$

Entropy

given collection S

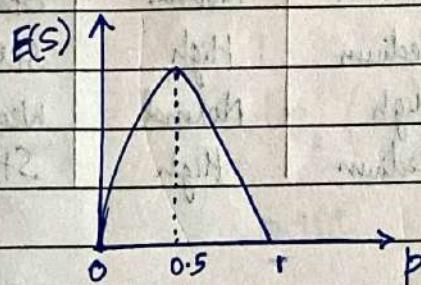
$$\begin{array}{c} n \\ \swarrow \quad \searrow \\ p \quad n-p \\ (+ve) \quad (-ve) \end{array}$$

$$\text{Entropy } (S) = -\frac{p}{p+n} \log_2 \frac{p}{p+n} - \frac{n-p}{p+n} \log_2 \frac{n-p}{p+n}$$

All -ve classes \Rightarrow Entropy = 0

All +ve classes

equal +ve, -ve \Rightarrow Entropy = maximum (0.5)



Information Gain



used to test goodness of a split

$$\Delta = I(\text{parent node}) - \sum_{i=1}^k \frac{N(v_i)}{N} I(v_i)$$

* compare the degree of impurity of parent node before & after splitting

N: total no. of instances at parent node

k: total possible attribute values

$N(v_i)$: total no. of instances associated with child node v_i

Slides

ID3 algo problem

Outlook	Temp	Humidity	Windy	Play Tennis
Sunny	High	High	Weak	No
Sunny	High	High	Strong	No
Overcast	High	High	Weak	Yes
Rainy	Medium	High	Weak	Yes
Rainy	Cool	Normal	Weak	Yes
Rainy	Cool	Normal	Strong	No
Overcast	Cool	Normal	Strong	Yes
Sunny	Medium	High	Weak	No
Sunny	Cool	Normal	Weak	Yes
Rainy	Medium	Normal	Weak	Yes
Sunny	Medium	Normal	Strong	Yes
Overcast	Medium	High	Strong	Yes
Overcast	High	Normal	Weak	Yes
Rainy	Medium	High	Strong	No

Step 1. $\sum(S) = -\frac{9}{14} \log_2 \left(\frac{9}{14}\right) - \frac{5}{14} \log_2 \left(\frac{5}{14}\right)$
 $= 0.94$

Entropy of whole data set

Step 2. $\sum_{\text{outlook (sunny)}} = -\frac{2}{5} \log_2 \left(\frac{2}{5}\right) - \frac{3}{5} \log_2 \left(\frac{3}{5}\right)$
 $= 0.91$

$\sum_{\text{outlook (overcast)}} = -\frac{4}{4} \log_2 \left(\frac{4}{4}\right) - \frac{0}{4} \log_2 \left(\frac{0}{4}\right)$
 $= 0$

$$\sum_{\substack{\text{outlook (having)} \\ 5}} = -\frac{3}{5} \log_2 \left(\frac{3}{5} \right) - \frac{2}{5} \log_2 \left(\frac{2}{5} \right)$$

$$= 0.971$$

$$\text{avg Entropy} = \frac{5}{14} (0.971) + \frac{4}{14} (0) + \frac{5}{14} (0.971)$$

$$= 0.693$$

$$\text{Information Gain} = \text{Parent E} - \text{Avg}$$

$$= 0.94 - 0.693$$

$$= 0.24$$

$$\sum_{\substack{\text{temp (High)} \\ 4}} = -\frac{2}{4} \log_2 \left(\frac{2}{4} \right) - \frac{2}{4} \log_2 \left(\frac{2}{4} \right)$$

$$= 0.1$$

$$\sum_{\substack{\text{temp (Medium)} \\ 6}} = -\frac{2}{6} \log_2 \left(\frac{2}{6} \right) - \frac{4}{6} \log_2 \left(\frac{4}{6} \right)$$

$$= 0.918$$

$$\sum_{\substack{\text{temp (cool)} \\ 4}} = -\frac{1}{4} \log_2 \left(\frac{1}{4} \right) - \frac{3}{4} \log_2 \left(\frac{3}{4} \right)$$

$$= 0.811$$

$$\text{avg Entropy} = \frac{4}{14} (1) + \frac{6}{14} (0.918) + \frac{4}{14} (0.811)$$

$$= 0.910$$

$$\text{Information Gain} = 0.94 - 0.91 = 0.03$$

$$\sum_{\text{Hum}(H)} = -\frac{4}{7} \log_2 \left(\frac{4}{7} \right) - \frac{3}{7} \log_2 \left(\frac{3}{7} \right)$$

$$= 0.985$$

$$\sum_{\text{Hum}(N)} = -\frac{1}{7} \log_2 \left(\frac{1}{7} \right) - \frac{6}{7} \log_2 \left(\frac{6}{7} \right) = 0.591$$

$$\text{avg Entropy} = \frac{7}{14} (0.985) + \frac{7}{14} (0.591)$$

$$= 0.788$$

$$\text{Information Gain} = 0.94 - 0.788 = 0.152$$

$$\sum_{\text{windy}(W)} = -\frac{2}{8} \log_2 \left(\frac{2}{8} \right) - \frac{6}{8} \log_2 \left(\frac{6}{8} \right) = 0.811$$

$$\sum_{\text{windy}(S)} = -\frac{3}{6} \log_2 \left(\frac{3}{6} \right) - \frac{3}{6} \log_2 \left(\frac{3}{6} \right) = 1$$

$$\text{avg Entropy} = \frac{8}{14} (0.811) + \frac{6}{14} (1) = 0.892$$

$$\text{Info Gain} = 0.94 - 0.892 = 0.048$$

12th August, 2024

q.	Salary	Location	Job acceptance
	T1	MUM	Yes
	T2	BLR	Yes
	T1	BLR	No
	T1	HYD	No
	T2	MUM	Yes
	T1	HYD	No
	T1	HYD	No

$$\textcircled{1} E(S) = -\frac{3}{7} \log_2 \left(\frac{3}{7} \right) - \frac{4}{7} \log_2 \left(\frac{4}{7} \right)$$

$$= 0.985$$

salary: $\textcircled{2} \sum_{\text{salary}(T1)} = -\frac{1}{5} \log_2 \left(\frac{1}{5} \right) - \frac{4}{5} \log_2 \left(\frac{4}{5} \right)$

$$= 0.722$$

$$\sum_{\text{salary}(T2)} = 0$$

$$\text{avg entropy} = \frac{5}{7} (0.722) + \frac{2}{7} (0) = 0.515$$

$$IG_{\text{salary}} = 0.985 - 0.515 = 0.47$$

location: $\sum_{\text{location}(MUM)} = 0$

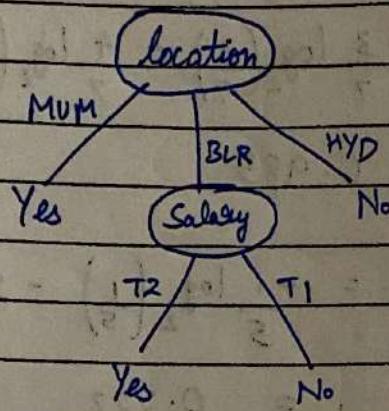
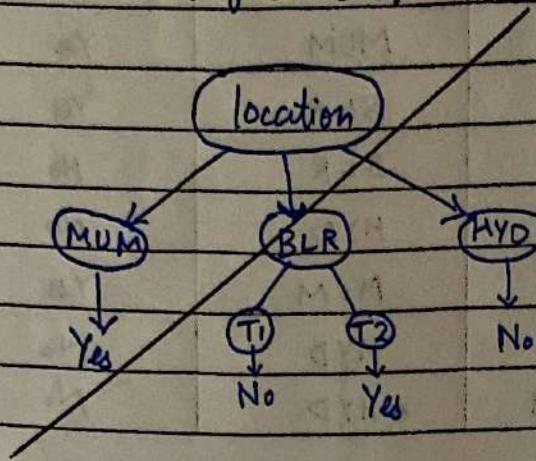
$$\sum_{\text{location}(BLR)} = 1.0$$

$$\sum_{\text{location}(HYD)} = 0$$

$$\text{avg entropy} = 0 + \frac{2}{7} (1) + 0 = 0.286$$

$$IG_{\text{location}} = 0.985 - 0.286 = 0.699$$

③ Consider location as root
(highest IG)



q: numerical IP attribute

Temp	Playing
40	No
48	No
60	Yes
72	Yes
80	Yes
90	No

40	No
48	No
60	Yes
72	Yes
80	Yes
90	No

option 1: take decision for all values 40, 48 ...

Step 1. Sort the numerical attribute

Step 2. Whenever there is a change in OP consider threshold

eg:	22	No	-
	23	Yes	←
	26	No	←
	28	No	-
	29	Yes	←
	30	Yes	-
	32	No	←

Temp	Playing
40	No
48	No
50	Yes
72	Yes
80	Yes
90	No

$$E(S) = 1$$

Step 3. ≤ 54 , > 54 ≤ 85 , > 85

to find which is a better split find IG

$$E(\leq 54) = 0$$

$$E_{\frac{4}{4}}(> 54) = -\frac{1}{4} \log_2 \left(\frac{1}{4} \right) - \frac{3}{4} \log_2 \left(\frac{3}{4} \right) = 0.811$$

$$\text{avg } E = \frac{2}{6}(0) + \frac{4}{6}(0.811) = 0.54$$

$$IG_{\text{at } 54} = 1 - 0.54 = 0.46$$

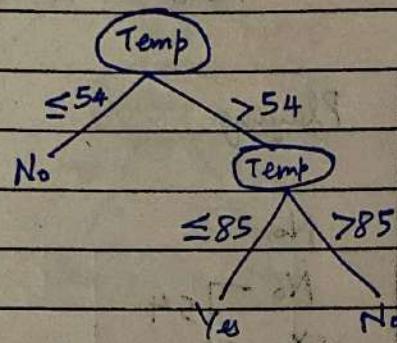
$$E(\leq 85) = \frac{2}{5} \log_2\left(\frac{2}{5}\right) - \frac{3}{5} \log_2\left(\frac{3}{5}\right) = 0.971$$

$$E(>85) = 0$$

$$\text{avg } E = \frac{5}{6}(0.971) + 0 = 0.809$$

$$IG_{\text{wt } 85} = 1 - 0.809 = 0.19$$

Take 54 as root node



- * overfitting
- * underfitting

CART algo.

uses Gini Index instead of Entropy

$$\text{Gini Impurity} = 1 - \sum [p(i|t)]^2$$

Step 1: $G \cdot I(s)$ Gini Impurity

Step 2: Attribute $G \cdot I$

Step 3: Find Gini Index

Step 4: Gain = Overall $G \cdot I$ - Gini Index

q.	Income	Student	Credit	Buys
1	High	No	Fair	No
2	High	No	Excellent	No
3	High	Yes	Fair	Yes
4	High	Yes	Excellent	Yes
5	Medium	Yes	Fair	Yes
6	Low	Yes	Fair	Yes
7	Low	No	Excellent	No
8	Low	Yes	Excellent	Yes
9	Medium	No	Fair	No
10	Low	No	Fair	No
11	Medium	Yes	Excellent	Yes
12	Medium	No	Excellent	Yes

$$\text{① GI (S)} = 1 - \left[\left(\frac{7}{12} \right)^2 + \left(\frac{5}{12} \right)^2 \right] = 0.486$$

income

$$\text{② GI income} = 1 - \left[\left(\frac{2}{4} \right)^2 + \left(\frac{2}{4} \right)^2 \right] = 0.5$$

$$\text{GI income} = 1 - \left[\left(\frac{3}{4} \right)^2 + \left(\frac{1}{4} \right)^2 \right] = 0.375$$

$$\text{GI income} = 1 - \left[\left(\frac{2}{4} \right)^2 + \left(\frac{2}{4} \right)^2 \right] = 0.5$$

$$\text{avg GI} = \frac{4}{12} (0.5) + \frac{4}{12} (0.375) + \frac{4}{12} (0.5) = 0.485$$

$$\text{Grain} = 0.486 - 0.485 = 0.001$$

student

$$\text{GI student} = 1 - \left[\left(\frac{5}{6} \right)^2 + \left(\frac{1}{6} \right)^2 \right] = 0.277$$

$$\text{GI student} = 1 - \left[\left(\frac{6}{6} \right)^2 + \left(\frac{0}{6} \right)^2 \right] = 0$$

$$\text{avg GI} = \frac{6}{12} (0.277) = 0.1385$$

$$\text{Grain student} = 0.486 - 0.1385 = 0.3475$$

credit

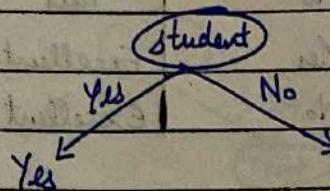
$$GI \text{ credit} = 1 - \left[\left(\frac{3}{6} \right)^2 + \left(\frac{3}{6} \right)^2 \right] = 0.5$$

$$GI \text{ credit} = 1 - \left[\left(\frac{2}{6} \right)^2 + \left(\frac{4}{6} \right)^2 \right] = 0.444$$

$$\text{avg GI} = \frac{6(0.5)}{12} + \frac{6(0.444)}{12} = 0.472$$

$$\text{Gain credit} = 0.486 - 0.472 = 0.014$$

Take student as root



income	student	credit	buys
High	No	Fair	No
High	No	Exc	No
Low	No	Exc	No
Medium	No	Fair	No
Low	No	Fair	No
Medium	No	Exc	Yes

$$GI(S) = 1 - \left[\left(\frac{1}{6} \right)^2 + \left(\frac{5}{6} \right)^2 \right] = 0.277$$

income

$$GI \text{ income} = 0$$

(H) 2

$$GI \text{ income} = 1/2$$

(M)

$$GI \text{ income} = 0$$

(L)

$$\text{avg GII} = \frac{2}{6} \left(\frac{1}{2} \right) = 0.166$$

$$\text{Grain} = 0.277 - 0.166 = 0.111$$

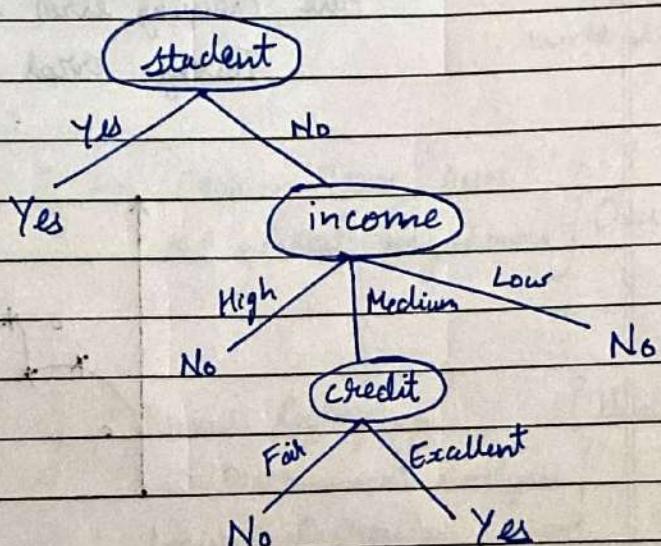
(take)

credit GII credit = 0
(F) 3

$$\text{GII credit} = 1 - \left[\left(\frac{1}{3} \right)^2 + \left(\frac{2}{3} \right)^2 \right] = 0.444$$

$$\text{avg} = \frac{3}{6} (0.444) = 0.222$$

$$\text{Grain} = 0.05$$



- * Hypothesis which understands data completely
 - Consistent hyp
- * Hypothesis which understands some/most data
 - Satisfied hyp

Syntactic

Syntactically diff hyp = no along with ? \oplus

$$3 \times 2 \times 2$$

$$\begin{array}{c} ? \\ \oplus \\ \oplus \end{array}$$

$$\overline{\Rightarrow 5 \times 4 \times \dots \dots}$$

Semantically diff = $\oplus \oplus \oplus \Rightarrow \Rightarrow (4 \times 3 \times 3) + 1$

22nd August, 2024

* Decision Tree creates problem of over fitting which can be overcome by pruning.

* : Training data

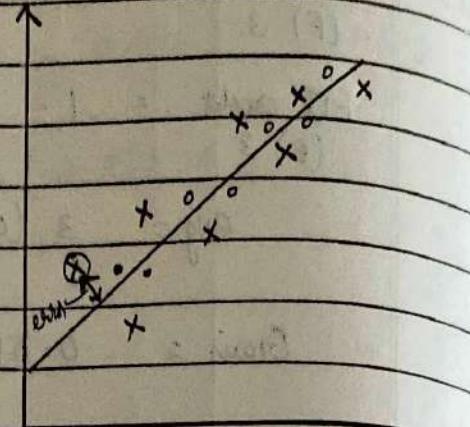
o : Testing data

/ : model

Training Error



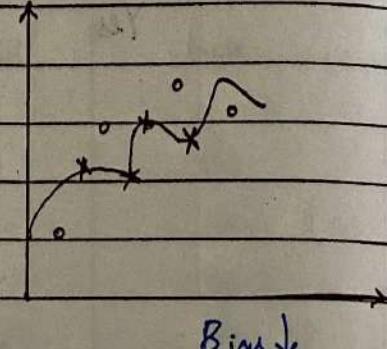
Bias



Here Training error is more \Rightarrow Bias \uparrow

Testing error is less \downarrow

Underfitting



More Testing error \Rightarrow Variance \uparrow



Overfitting \Rightarrow model highly complex

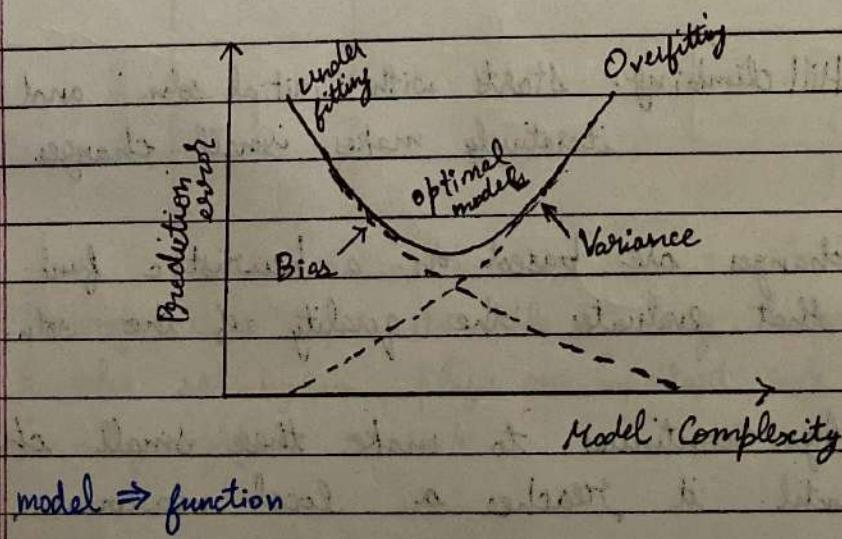
* Model we prefer : Low Bias (less training error)
Low Variance (less testing error)

$$E_{\text{err}}(s) = \text{Bias}^2 + \text{Variance} + \text{Irreducible error}$$

noise

$$Err(s) =$$

- * Bias = diff b/w avg prediction & correct values
 - * Variance = how are prediction made on same observation different from each other.
 - * Decision Tree \rightarrow complex \rightarrow High Variance \rightarrow Overfitting
it tries to understand training data more & more \downarrow
make simple



Hypothesis search and Inductive bias - ID₃

- * ID₃ can be characterized as searching a hypotheses space for one that fits the training examples.
 - * The hypothesis space (searched by ID₃) is the set of all possible DTs.
 - * ID₃ performs a simple-to-complex, hill climbing search.
 - ↓
way of identifying func.
searching & trying to reach last leaf node
height of tree ↑ complexity ↑
 - * Hill climbing → simple optimization algo to find best possible solution.
 - ↓
 - * It belongs to local search algorithm
- Goal : Find best soln. from set of all possible soln.
- * Hill climbing: starts with initial soln. and then iteratively makes small changes
 - * changes are based on a heuristic func. that evaluates the quality of the soln.
 - * Algo continues to make these small changes until it reaches a local maximum.

- * The evaluation func. that guides hill climbing search is information gain measure

→ Advantages of Favoring Shallow Trees

- i) Simplicity and Interpretability
- ii) Generalization
- iii) Reduced Overfitting

→ Limitations

- i) Underfitting - Did not understand training data
- ii) Loss of Information
- iii) Bias against complex pattern

26th August, 2024

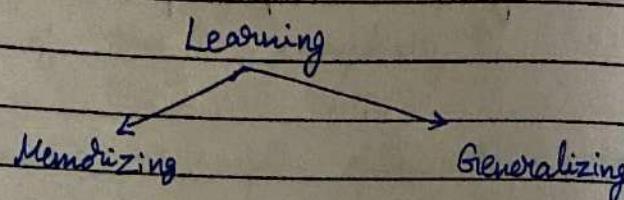
KNN : K Nearest Neighbour

Height	Weight	Target

- * Training time it sits idle - lazy learner
- * k-value set by us ; plays an important role.
how many neighbours to check

→ Instance based Learning

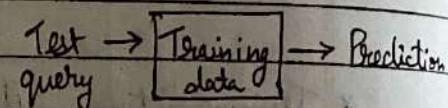
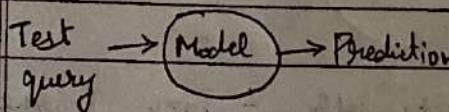
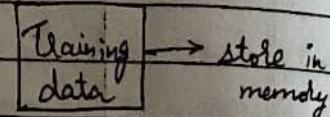
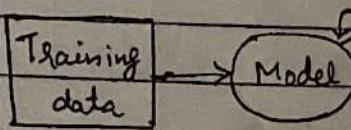
- * It is an ML technique in which the model learns the training examples by heart and then generalizes to new instances based on some similarity measure
- * number as IP attribute
- * Memory-based learning / lazy-learning
- * Worst case $O(n)$ → size of training data



- instance based learning

- Eager Learning

- Lazy Learning



→ KNN

* instance based, lazy learner ML algo,
supervised learning

KNN Inductive Bias

requires 3 things : i) set of stored records
 ii) distance metric : compute dist b/w records
 iii) value of k : no. of nearest neighbours to determine

To classify an unknown record :

- i) compute distance to other training records
- ii) identify k nearest neighbors
- iii) use class labels of nearest neighbours to determine class of unknown record

* for odd classes take even k-value
 for even classes take odd k-value

- K-NN working algo

Step 1: Select k of the neighbours

Step 2: Calculate Euclidean dist. of k no. of neighbours

Step 3: Take k nearest neighbours as per calculated euc dist.

Step 4: Among k neighbours, count no. of data points in each category

Step 5: Assign the new data point to that category for which no. of neighbour is maximum

Step 6: Model is ready.

- * All instances correspond to points in n-D space
- * nearest neighbors : dist measure - $\text{dist}(x_i, x_j)$
- * Target func : discrete or real valued
- * for discrete : k NN returns most common value (classes)
- * for real valued : k NN returns mean values of (numbers) k nearest neighbours
- Choosing the value of k
 - * if k too small \rightarrow overfitting
 - * if k too large \rightarrow underfitting
- Consider learning discrete-valued target functions of the form $f: \mathbb{R}^n \rightarrow V$ where V is finite set
- * Training Algo : For each training algo example $(x_i, f(x_i))$ add the eg. to the list of training egs.
- * Classification Algo : Given a query instance x_q to be classified

1. Let x_1, \dots, x_k denote k instances from training egs. that are nearest to x_q

2. Return

$$\hat{f}(x_q) \leftarrow \operatorname{argmax}_{v \in V} \sum_{i=1}^k \delta(v, f(x_i))$$

where $\delta(a, b) = 1$ if $a = b$ &
 $\delta(a, b) = 0$ otherwise

* Argmax - operation that finds the arg that gives the max from a target func.

↓
used to find class with largest predicted probability.

→ KNN Applications

* In KNN, every data instance is assumed to be a point in d-dimensional space (R^d) where d is the no. of features.

* KNN can be used for both classification & regression

- KNN classification

* target func. is discrete valued (m classes)

* each of the data instances may belong to one of the m classes

* We define set V as set of all classes

* target func. $f: R^d \rightarrow V$

* any data instance x_i will have a label $f(x_i)$
(one of the class from set V)

- KNN regression

* target func. is real valued (continuous value)

* target func. $f: R^d \rightarrow R$

* query instance is represented as x_q

* we are interested in estimating the $f(x_q)$ value, denoted by $\hat{f}(x_q)$

→ Minkowski Distance

$$= \left(\sum_{i=1}^k (|x_i - y_i|)^q \right)^{1/q}$$

$q=1 \Rightarrow$ Manhattan
 $q=2 \Rightarrow$ Euclidean

k : no. of features/attributes

x_i, y_i : i^{th} attributes

→ How to choose no. of neighbours

- * k : odd if classes-even
- * k : even if classes-odd

BEST PRACTICE:

- * start with $k = \text{no. of classes} + 1$
- * in case of tie decrease k by 1
- * $k = \sqrt{n}$ performs well too

Elbow method : used to determine best k value

- Calc. error rate for diff. k -value
- choose optimal k -value as elbow value of curve
- retain with new k value
- retain model with best k value

→ KNN algo for real valued target

- * store each training example $\langle x, f(x) \rangle$
- * find k nearest neighbors x_1, \dots, x_k of g

$$f(x_g) \leftarrow \frac{\sum_{i=1}^k f(x_i)}{k}$$

mean value of k nearest
neighbours is returned as
appropriate value of x_g

$$MSE = \frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2$$

→ Weighted KNN

* More closer \Rightarrow more importance \Rightarrow more weight

$$w_i = \frac{1}{d(x_q, x_i)^2}$$

Distance wtd. Voting

$$\hat{f}(x_q) \leftarrow \operatorname{argmax}_{v \in V} \sum_{i=1}^k w_i \delta(v, f(x_i))$$

$$w_i = \frac{1}{d(x_q, x_i)^2}$$

→ Issues with KNN

i) * KNN slow algo : as data set \uparrow KNN struggles

ii) * Curse of dimensionality : no. of IP features \uparrow KNN struggles

→ Overcoming curse of dimensionality

* Assigning weights

* Leave-one-out

→ Issues with KNN

iii) Optimal no. of neighbours

iv) Imbalanced data causes problems

v) Outlier sensitivity

vi) Missing value treatment

vii) K-NN needs homogeneous features

→ KNN computational complexity

- * Basic KNN algo stores all examples
- * n examples of dim d

compute dist to one eg: $O(d)$

find 1 nearest neighbour: $O(nd)$

find k closest egs: $O(knd)$

↓

probably expensive for large no. of samples

20th August, 2024

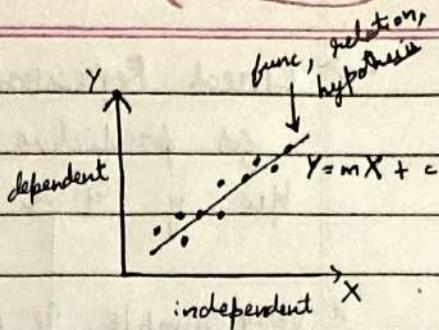
Date _____
Page _____

Linear Regression

Simple Linear Regression



1 IP, 1 OP



it assumes linear relation

Optimization - improving the model \rightarrow decreasing error
 \hookrightarrow Gradient Descent also ~~except~~

Regression : method / statistical measure for understanding relationship b/w independent variable (X) of dependent variable (Y)
(features) (outcome)

- Types of Regression
- Linear
 - Polynomial
 - Support vector
 - Decision Tree
 - Random Forest
 - Ridge
 - Lasso
 - Logistic

\rightarrow Scatter Plot

\downarrow ((x + 2y) - 1) = 0 is best line

- * Visualization tool to choose 'type of regression'
- * Shows individual points as dots in 2D plane

→ Linear Regression : statistical regression method
for predictive analysis, shows linear relationship
b/w y & x (one/more)

- * very simple & easy
- * used when continuous numeric data

Types of Linear Regression : i) Simple 1 IP 10p
ii) Multiple many IP 10p

→ Simple Linear Regression

only 1 IP variable

$$y = mx + b$$

dependent slope \ independent
y-intercept

→ Best Fitted Line

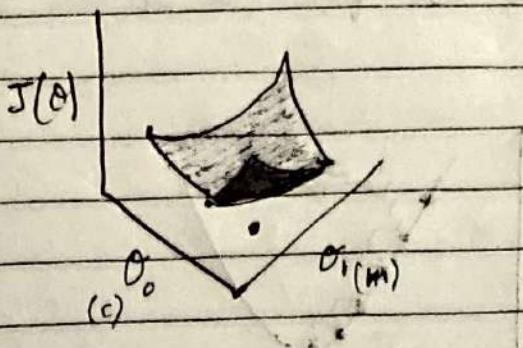
- * Straight line that minimizes the vertical distance b/w the data points & the regression line.

$$\begin{aligned} \text{Error} &= \text{Actual value} - \text{Predicted value} \\ &= y_i - (mx_i + b) \end{aligned}$$

$$\text{Squared Error} = (y_i - (mx_i + b))^2$$

$$\text{Sum of Squared Error (SSE)} = \sum_{i=1}^n (y_i - (mx_i + b))^2$$

$$\text{Mean Squared Error (MSE)} = \frac{1}{n} \sum_{i=1}^n (y_i - (mx_i + b))^2$$



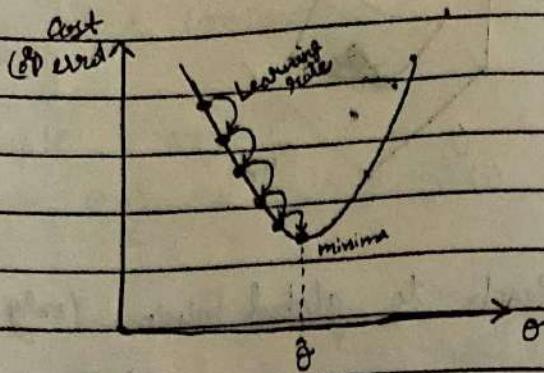
Moto: Reach to global minima (only 1 global minima)

rate of improving \rightarrow learning rate

\rightarrow Gradient Descent

- i) initialize
- ii) derivative of loss func.
- iii) find min of loss func \Leftrightarrow best-fit value of func.

→ Visualization of Cost func. in 2D



Learning rate alpha : +ve constant determining how much bigger the steps are gonna be.

Note : If α is very large then it is gonna overshoot the min val of cost func (J) & might diverge.

If α is smaller value then smaller steps are taken to reach global min.

→ Steps to find optimum values of $m \& b$

1. Initialize $m \& b$ with random values.

2. Define the loss func :

$$MSE = \frac{1}{n} \sum_{i=1}^n (y_i - (mx_i + b))^2$$

n : no. of data points

y_i : actual y value

3. Calculate the gradient : partial derivatives of loss func with respect to $m \& b$

$$\frac{\partial MSE}{\partial m} = -2 \sum_{i=1}^n x_i (y_i - (mx_i + b))$$

$$\frac{\partial MSE}{\partial b} = -2 \sum_{i=1}^n (y_i - (mx_i + b))$$

4. Update parameters.

$$m \leftarrow m - \alpha \frac{\partial \text{MSE}}{\partial m}$$

$$b \leftarrow b - \alpha \frac{\partial \text{MSE}}{\partial b}$$

5. Iterate : Repeat ③, ④ until parameters converge

g. $y = mx + b$, $m = 10$, $b = 300$, $\text{LR}(\alpha) = 0.0001$

$y = 10x + 300$	Age ^(x)	Salary ^(y)
600	30	800
670	37	950
550	25	600
730	43	1050
800	50	1200
590	29	740
760	46	1100

$$\text{MSE} = \frac{1}{n} \sum_{i=0}^n (y_i - \hat{y}_i)^2$$

$$= 74485.71$$

$$\text{GD} \rightarrow \frac{\partial \text{MSE}}{\partial m} = \frac{-2}{n} \sum_{i=1}^n x_i (y_i - (mx_i + b)) \\ = -20388.57$$

$$\frac{\partial \text{MSE}}{\partial b} = \frac{-2}{n} \sum_{i=1}^n (y_i - (mx_i + b)) \\ = -497.14$$

$$m_1 = m - \alpha \left(\frac{\partial \text{MSE}}{\partial m} \right)$$

$$= 12.038$$

$$b_1 = b - \alpha \left(\frac{\partial \text{MSE}}{\partial b} \right)$$

$$= 300.0497$$

x	$mx + b$	$m_1x + b_1$
30	600	661.189
37	670	745.455
25	550	600.999
43	730	817.6837
50	800	901.949
29	590	649.152
46	760	853.798

$$\text{MSE} = \frac{1}{7} \sum_{i=1}^7 (y_i - \hat{y})^2$$

$$= 38957.2$$

$$\frac{\partial \text{MSE}}{\partial m} = -\frac{2}{n} \sum x_i (y_i - (mx_i + b)) = -1443.2$$

$$\frac{\partial \text{MSE}}{\partial b} = -\frac{2}{n} \sum (y_i - (mx_i + b)) = -345.59$$

$$m_2 = m_1 - \alpha \left(\frac{\partial \text{MSE}}{\partial m_1} \right) = 13.4831$$

$$b_2 = b_1 - \alpha \left(\frac{\partial \text{MSE}}{\partial b_1} \right) = 300.084$$

→ Measuring model Performance

- * The goodness of fit determines how the line of regression fits the set of observations. Process of finding best model out of various models: optimization

R-squared method

L determines goodness of fit

- * lies b/w 0 and 1
- * R-squared = 1 model perfectly fits
- * R-squared = 0 model does not predict any variability
- * measures strength of relationship b/w dependent & independent variable. (0 - 100%)

$$SST_{\text{tot}} = \sum_{i=1}^n (y_i - \bar{y})^2$$

$$R^2 = 1 - \frac{\text{sum squared regression (SSR)}}{\text{total sum of squares (SST)}}$$

- R squared method

- i. calculate mean of the target/dependent variable
 y & denote it with \bar{y}

$$SS_{\text{tot}} = \sum_{i=1}^n (y_i - \bar{y})^2$$

$$3. R^2 = 1 - \frac{SSR}{SST}$$

$$= 1 - \frac{\sum (y_i - \hat{y}_i)^2}{\sum (y_i - \bar{y}_i)^2}$$

- * High $R^2 \Rightarrow$ less diff b/w predicted & actual values
 \Downarrow
 good model

- * for multiple regression

$$R^2 = \frac{\text{Explained variation}}{\text{Total variation}}$$

↓
 coefficient of determination or
 multiple determination

Linear Regression

Advantages

Disadvantages

- * Simple to implement & easier to interpret OP efficiencies

Outliers can have a huge effect on the regression

- * We know the relationship b/w ind & dep variable (linear) algo is the best to use because of its less complexity

LR assumes linear rel.
 \Downarrow
 it assumes straight line relationship, assumes independence b/w attributes

- * LR is susceptible to over fitting but it can be avoided using some dimensionality reduction tech, regularization

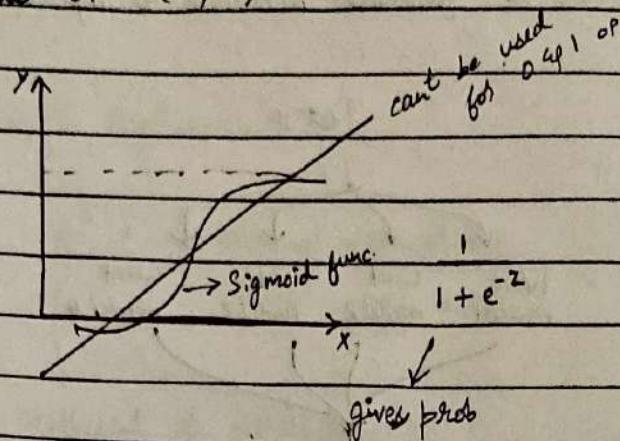
It also looks at the relationship b/w the mean of dep & ind vars, mean is not a complete desc of relationship b/w vars

2nd September, 2024

KNN \rightarrow classification

Logistic Regression - Not a regression model
↓ but a classification model

discrete OP (0, 1)



$$\text{odd} = \frac{\text{Prob. of event occurring}}{\text{Prob. of event not occurring}} = \frac{p}{1-p}$$

$Z \rightarrow$ linear eqn: for 1 parameter $Z = mx + c$

for multiple: $Z = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \beta_3 x_3$

q: Predict likelihood of admission

Patient ID	Age	BMI	BP	Glucose level	days in hospital	no. of prev admission
6	45	27	138	102	4	1
7	62	25	132	106	5	2

Q:

Restaurant	Location	Food Quality	High Rating
1	3	2	0
2	4	3	1
3	2	1	0
4	5	2	1
5	4	2	1

$$w_0 = -3, w_1 = 0.4, w_2 = 2$$

Predict whether restaurant with location = 2 & food quality = 2 will receive high rating or not

$$z = -3 + (0.4)(2) + 2(2) = -1.8$$

$$y = \frac{1}{1 + e^{-1.8}} = 0.858$$

If threshold = 0.5 $\Rightarrow y$ is high rating

Logistic Regression: supervised ML algo mainly used for binary classification tasks

where goal is to predict the prob. that an instance belongs to one of the 2 given classes

- * IP: independent vars can be categorical / numerical
- * Logistic Regression: * Binary classification problems
 - * handling imbalanced data
 - * capturing non-linear relationships
- * Linear Regression: * continuous numeric data
 - * not well suited for classification tasks

* Logistic Regression more robust to outliers
L due to sigmoid func.

3 types of Logistic Regression

- i) Binomial : 2 possible types of dependent variables
- ii) Multinomial : 3 or more possible unordered types of dependent variable
- iii) ordinal : 3 or more possible ordered types of dependent variable

$$p = \frac{\text{outcome of interest}}{\text{all possible outcomes}}$$

$$\text{Odds} = \frac{p \text{ of event occurring}}{p \text{ of event not occurring}} = \frac{p}{1-p}$$

Odds ratio : ratio of 2 odds

* dependent variable in logistic regression follows Bernoulli distribution

Logit func: maps the linear combination of IP natural log of odds to variables to OP variable

goal of logistic regression: estimate p

$$\text{Logit}(p) = \ln(\text{odds}) = \ln\left(\frac{p}{1-p}\right) = \ln(p) - \ln(1-p)$$

prob value b/w 0 & 1

- Relationship bw Logit & Sigmoid func.

$$\text{Logit}^{-1}(z) = \frac{1}{1 + e^{-z}}, \quad z: \text{linear combination of IP variables}$$

$$z = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \dots + \beta_n x_n \\ (x_1, x_2, \dots, x_n : \text{IP features})$$

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

sigmoid func

$\sigma(z)$ tends to 1 as $z \rightarrow \infty$: $P(y=1) = \sigma(z)$
 $\sigma(z)$ tends to 0 as $z \rightarrow -\infty$: $P(y=0) = 1 - \sigma(z)$

- Loss Function in Logistic Regression

Binary cross-entropy loss in logistic regression :

$$\text{Jobs} \cdot L(y, \hat{y}) = -[y \log(\hat{y}) + (1-y) \log(1-\hat{y})]$$

$$\frac{\partial L}{\partial w} = -x(y - \hat{y}) \quad \frac{\partial L}{\partial b} = -(y - \hat{y})$$

- Finding coefficients using Gradient Descent

1. Initialize coefficients

2. Calculate Predicted Probabilities

$$P(Y=1|X) = \frac{1}{1 + e^{-(\beta_0 + \beta_1 x)}}$$

3. Calculate log-likelihood

$$L(y, \hat{y}) = -[y \log(\hat{y}) + (1-y) \log(1-\hat{y})]$$

4. Calculate Gradient $\frac{\partial L}{\partial m} = -x(y - \hat{y})$

$$\frac{\partial L}{\partial b} = -(y - \hat{y})$$

5. Update coefficients

$$m = m - \alpha \left(\frac{\partial L}{\partial m} \right)$$

$$b = b - \alpha \left(\frac{\partial L}{\partial b} \right)$$

repeat ② to ⑤ until coefficients converge

6. Convergence & Termination

7. Final coefficient values

- Assumption of Logistic Regression

- i) Independent observations
- ii) Binary dependent variables
- iii) Linearity b/w ind vars & log odds
- iv) No outliers
- v) Large sample size

- ### - Advantages :
- i) performs better when data is linearly separable
 - ii) doesn't req. too many computational resources
 - iii) doesn't req. tuning
 - iv) easy to implement & train
 - v) gives measure of how relevant a predictor is

- ### - Disadvantages :
- i) fails to predict continuous outcome
 - ii) assumes linearity b/w dep & ind vars
 - iii) may not be accurate if sample size is too small.