

Unit - 1

Introduction to machine learning

Enable the machine to automatically learn from data improve performance from experience and predict things without being explicitly program.

Features of machine learning

- Machine can learn itself from past data & automatically improve.
- Machine learning is data driven technology. Large amount of data generated by organizations on daily basis. So by notable relationship in data organization makes better decisions.
- From the given dataset it detects various patterns on data.
- It is similar to data mining because it is also deals with the huge amount of data.

Machine Learning life Cycle involve 7 major steps

1. Gathering Data
2. Data Preparation
3. Data Wrangling
4. Analyze data
5. Train the model
6. Test the model
7. Deployment

(Q) Intense based learning or Memory learning or Lazy learning
Model based learning.

Bias and Variance

Bias is simply defined as the inability of model because of that there is some difference or error occurring between the model predicted values.

These difference between actual or expected values & predicted values are known as error or bias. Bias or error due to its biased.

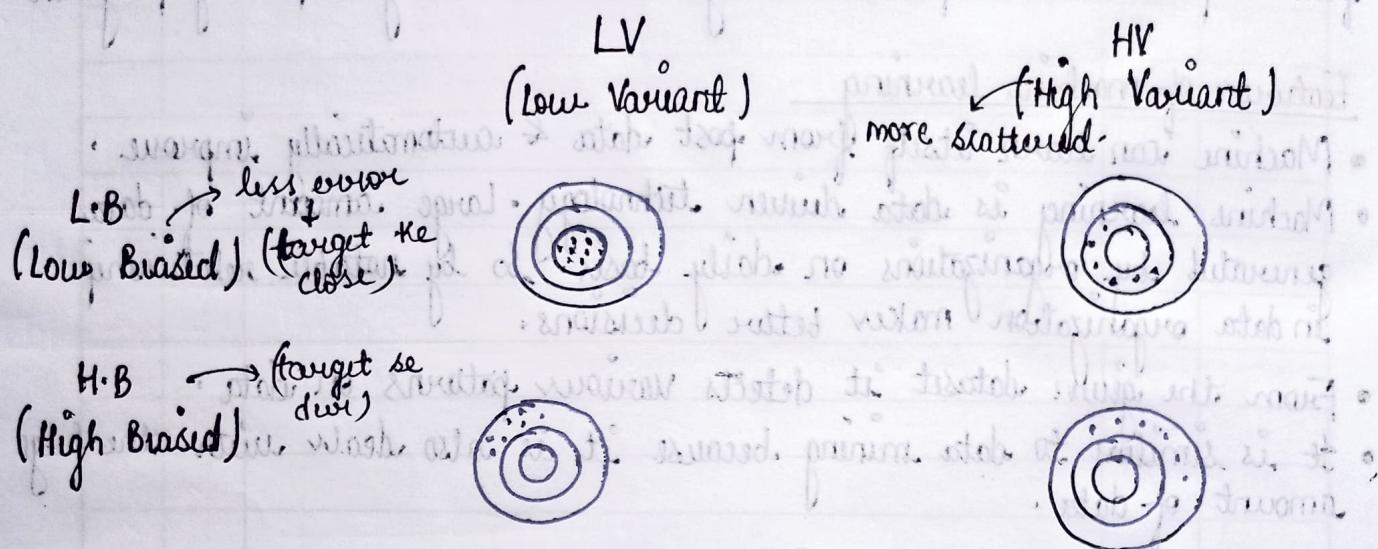
Low Bias

Low bias values means assumption are taken to build the target function. In these case the model will closely match the training dataset.

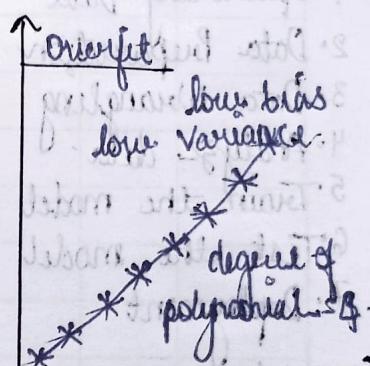
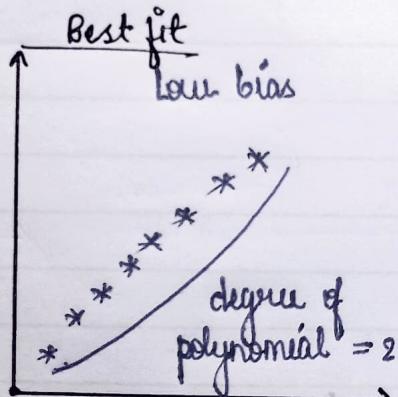
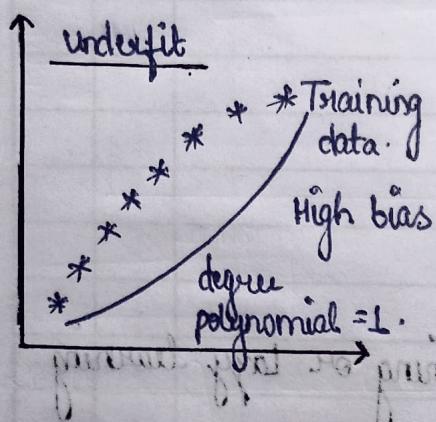
High Bias

I think.

High Bias values mean assumption are taken to build the target function. In these cases the model will not closely match the training dataset. margin width will decrease as relation with data mapping function will be less. High bias may.

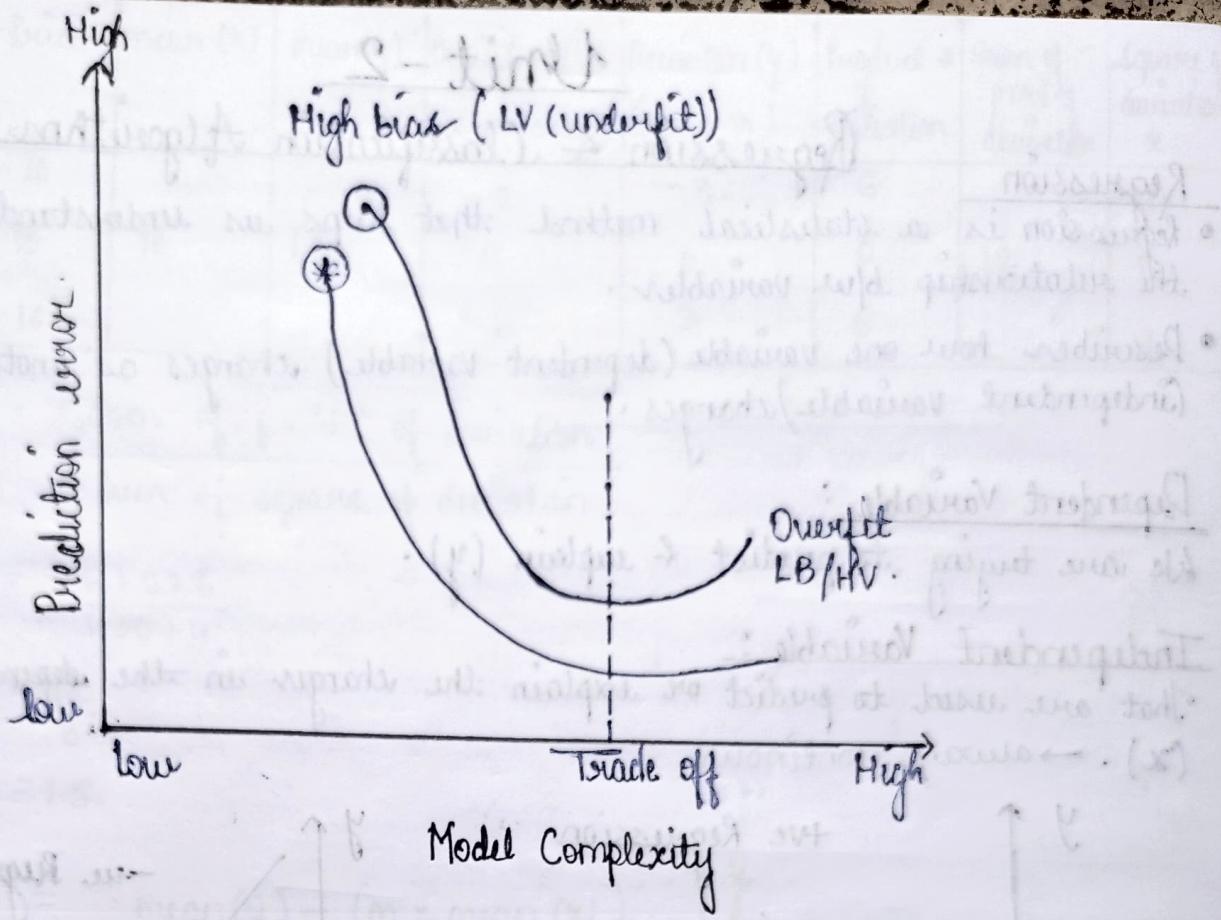


Regression



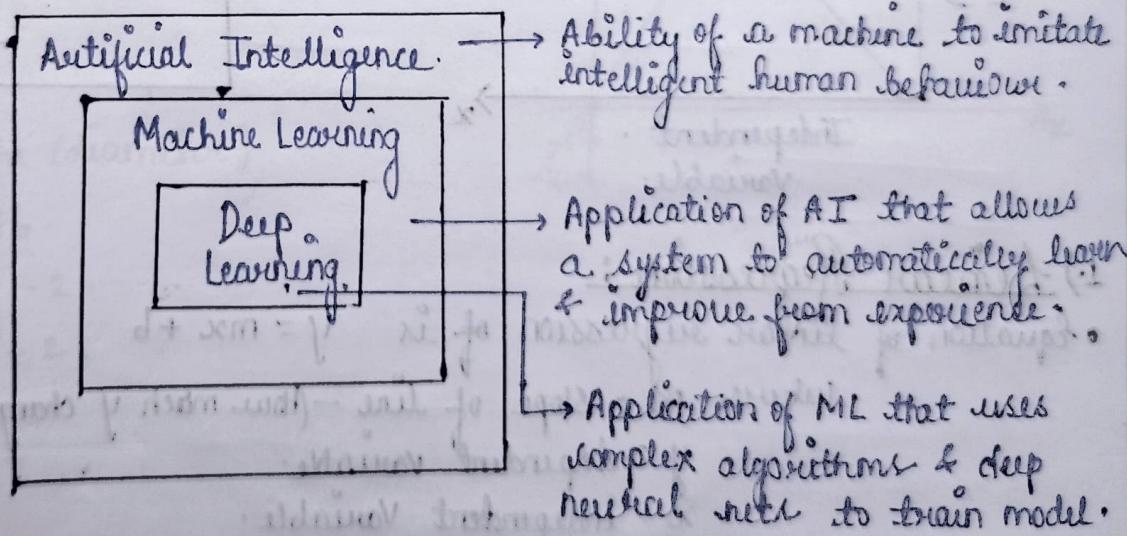
Classification problem specification

Joint Model	Model 1	Model 2	Model 3
Training error: 1%	Training error: 25%	Training error: <10%	Training error: <10%
Testing error: 20%	Testing error: 26%	Testing error: <10%	Testing error: <10%
Underfitting - unable to learn from data	High Bias - High Variance	High Variance	Low Variance - Low Bias
Overfitting - fit to noise	Underfitting - fit to noise	Optimal fitting	Optimal fitting



13/02/2024

Difference between AI / ML / Deep learning.



Unit - 2

Regression & Classification Algorithm

Regression

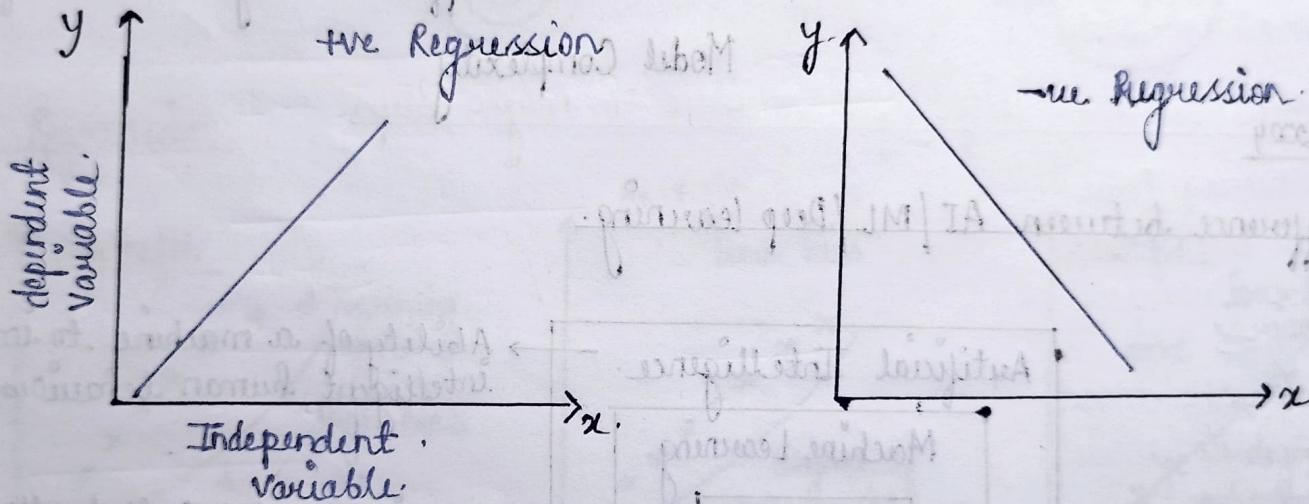
- Regression is a statistical method that helps us understand & predict the relationship b/w variables.
- Describes how one variable (dependent variable) changes as another variable (independent variable) changes.

Dependent Variable :

We are trying to predict & explain (y).

Independent Variable :

That are used to predict or explain the changes in the dependent variable (x). → always continuous.



1.) Linear Regression :

- Equation of linear regression is $y = mx + b$

where m = Slope of line - How much y changes for unit change in x .

y = dependent variable.

x = independent variable.

b = intercept - (when $x=0$ then what is the value of y)

(Q.) Predicting Pizza prices

Diameter (x)	Price (y)
8	10
10	13
20	16

→ dependent : continuous means linear regression.

Sel	Diameter	Price	mean(x)	mean(y)	Deviation(x) $(x - \text{mean}(x))$	Deviation(y) $(y - \text{mean})$	Product of deviation	Sum of prod of deviation	Square of deviation x^2
8	10				-2	-3	6		4
10	13	10	13		0	0	0	12	0
12	16				+2	3	6		4

i) Calculate $m = \frac{\text{sum of product of deviation}}{\text{sum of square of deviation}}$

$$= \frac{6+0+6}{4+0+4}$$

$$= \frac{12}{8}$$

$$m = 1.5$$

ii) Calculate $b = \text{mean}(y) - (m \times \text{mean}(x))$

$$= 13 - (1.5 \times 10)$$

$$= 13 - 15$$

$$= -2$$

iii) Let $x = 20$ inch (diameter)

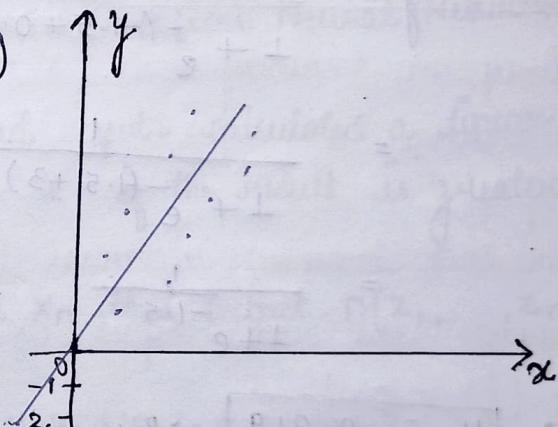
$$y = mx + b$$

$$y = 1.5 \times 20 - 2$$

$$= \frac{15}{10} \times 20 - 2$$

$$= 30 - 2$$

$$\boxed{y = 28}$$



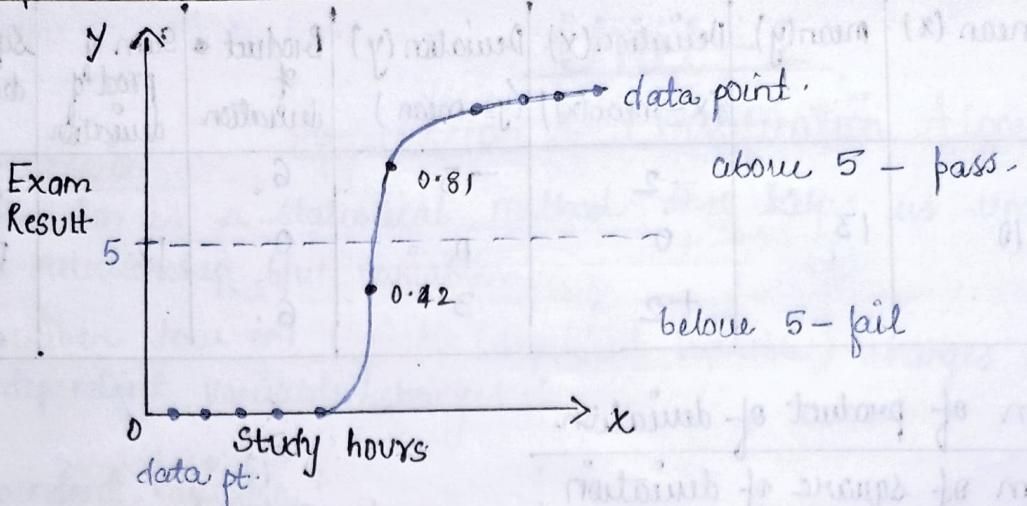
2) Logistic Regression :

	Study Hours	Exam Result
2		0
3		0
4		1
5		1
6		1
7		1
8		1

Given $a_0 = -1.5$

$a_1 = 0.6$

dependent variable : categorical means
in logistic regression



$$y = \frac{1}{1 + e^{-(a_0 + a_1 x)}}$$

sigmoid function.

where x = independent variable
 y = dependent variable
 a_0 = intercept
 a_1 = slope

Sol

$$y = \frac{1}{1 + e^{-(1.5 + 0.6 \times 5)}} \quad (i) \quad [x = 5]$$

$$= \frac{1}{1 + e^{-(1.5 + 3)}}$$

$$= \frac{1}{1 + e^{-1.5}} = \frac{1}{1 + 0.22}$$

$$\therefore \boxed{y = 0.819} \rightarrow \text{Pass} \\ = 81.9\%$$

$$(ii) x = 2.$$

$$y = \frac{1}{1 + e^{-(1.5 + 0.6 \times 2)}} = \frac{1}{1 + e^{-(1.5 + 1.2)}} = \frac{1}{1 + e^{-0.3}}$$

$$= \frac{1}{1 + 1.349} = \frac{1}{2.349}$$

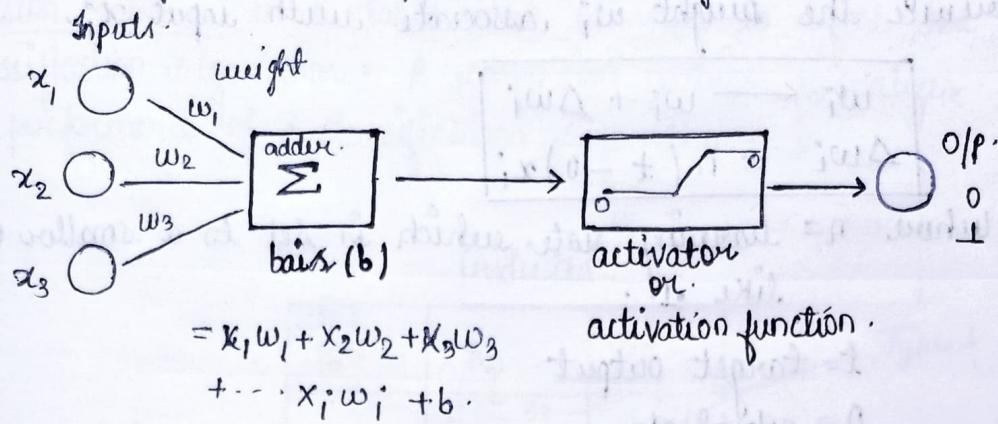
$$\therefore \boxed{y = 0.42} \rightarrow \text{fail}$$

15/02/2024

Perception : In machine learning the perceptron or McCulloch-Pitts neuron is an algorithm for supervised learning of binary classifiers.

A binary classifier is a function which can decides whether or not an input

represented by a vector of numbers belongs to some specific class.



$$a = \sum_{i=1}^n x_i w_i + b$$

or

$$a = X \cdot w + b$$

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

- A perception unit is used to build the ANNS (Artificial Neural Network system).

- A perception takes a vector of real valued inputs calculates a linear combination of these inputs then outputs a '1' if the result is greater than some threshold & -1 otherwise.

- Most precisely given if from x_1 through x_n then output $O(x_1, \dots, x_n)$ computed by the perception is

$$O(x_1, \dots, x_n) = \begin{cases} 1 & \text{if } w_1 x_1 + w_2 x_2 + \dots + w_n x_n > 0 \\ -1 & \text{otherwise} \end{cases}$$

- where each w is a real valued constant or weight that determines the contribution of input x to the perception output.

Perception Training Rule

- The way to learn an acceptable weight vector is to begin with random weights, then iteratively apply the perception to each training example, modifying the perception weights whenever it mis classifies an example.

- This process is repeated iteratively through the training examples as many times as needed until the perceptions classifies all the training examples correctly.

3.) Weights are modified at each step according to the perceptron training rule which revise the weight w_i associate with input x_i :

$$w_i \leftarrow w_i + \Delta w_i$$
$$\Delta w_i = n(t - o)x_i$$

where n = learning rate which is set to a smaller value like ± 1 .

t = target output

o = actual o/p.

x_i = I/P associated with w_i .

* Using this equation we will modify a weight until in this iteration until all the examples are classified correctly.

Perception Algorithm

Perception training rule (X, n)

initialize w ($w_i \leftarrow$ an initial (small) random value).

repeat

for each training instance $(x, t^x) \in X$ do

Compute the real output $o_x = \text{Activation}(\text{Summation}(w \cdot x))$

if $(t^x \neq o_x)$

for each w_i

$w_i \leftarrow w_i + \Delta w_i$

$\Delta w_i \leftarrow n(t^x - o_x)x_i$

end for

end if

end for

Until all the training instances in X are correctly classified

return w

Once all the training instances are classified correctly final weight are returned as the learned parameter by the perceptron network.

16/02/2024

Confusion Matrix: is a table that is used to define the performance of a classification algorithm. A confusion matrix visualizes & summarizes the performance of a classification algorithm.

		Predicted		
		No	Yes	
Actual	No	TN 50	FP 10	60
	Yes	FN 5	TP 100	105
		55	110	165

Type 1 error
if someone is not having any disease but on diagnosis it shows person is having disease.

Type 2 error
False Negative. True Positive

- The matrix displays the no. of instances produced by the model on the test data.

True positive (TP): Occurs when the model accurately predicts a true data-point.

True -ve (TN): Occurs when model predicts a -ve datapoint.

False +ve (FP): Occurs when model predicts a +ve data point incorrectly.

False -ve (FN): Occurs when a model miss predicts a -ve data point.

Accuracy: used to measure the performance of the model. It is the ratio of total correct instance to the total instances.

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

$$= \frac{100 + 50}{165}$$

$$\therefore \text{Accuracy} = \frac{150}{165} = 0.91 = 91\%$$

Precision :- is a measure of how accurate a model's predictions are. It is defined as the ratio of true true predictions to the total no. of the predictions made by the model.

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

$$= \frac{100}{10+100} = \frac{100}{110}$$

$$\therefore \text{Precision} = 0.90 = 90\%$$

Recall :- measures the effectiveness of a classification model in identifying all relevant instances. It is the ratio of the no. of true +ve instances to the sum of true +ve & false -ve instances.

$$\text{Recall} = \frac{TP}{TP + FN}$$

$$= \frac{100}{100+5} = \frac{100}{105}$$

$$\therefore \text{Recall} = 0.95 = 95\%$$

F₁ score :- it used to evaluate the overall performance of a classification model. It is the harmonic mean of precision & recall.

$$F_1 \text{ score} = \frac{2 \times \text{precision} \times \text{recall}}{\text{precision} + \text{recall}}$$

$$= \frac{2 \times 0.90 \times 0.95}{0.90 + 0.95}$$

$$= \frac{1.71}{1.85}$$

$$\frac{100 + 5}{105}$$

$$\therefore F_1 \text{ score} = 0.92 = 92\%$$

Error :-

$$\text{error} = \frac{\text{FP} + \text{FN}}{\text{Total}}$$

Type 1. Type 2
or

$$\text{error} = 1 - \text{accuracy}$$

$$\text{error} = \frac{10 + 5}{150} = 0.09 \quad \text{or} \quad \text{error} = 1 - 0.91 = 0.09.$$

17/02/2024

(Q)

	Predicted yes	Predicted No
Actual No	45 FP	5 TN
Actual yes	5 TP	95 FN

50
100
- 150

Sol Accuracy = $\frac{TP + TN}{\text{Total}} = \frac{5 + 5}{150} = \frac{10}{150} = 0.06$

$$\text{Precision} = \frac{TP}{FP + TP} = \frac{5}{45 + 5} = \frac{5}{50} = 0.1$$

$$\text{Recall} = \frac{TP}{TP + FN} = \frac{5}{5 + 95} = \frac{5}{100} = 0.05$$

$$F_1 \text{ score} = \frac{2 \times \text{precision} \times \text{recall}}{\text{precision} + \text{recall}} = \frac{2 \times 0.1 \times 0.05}{0.1 + 0.05} = \frac{0.01}{0.15} = 0.06$$

Error is $1 - \text{accuracy} = 1 - 0.06 = 0.94$

(Q)

	Predicted yes	Predicted No
Actual yes	45 TN	5 FP
Actual no	5 FN	95 TP

MV3 session

Sol Accuracy = $\frac{TP + TN}{\text{Total}} = \frac{45 + 95}{150} = \frac{140}{150} = 0.90 = 90\%$

$$\text{Precision} = \frac{\text{TP}}{\text{TP} + \text{FP}} = \frac{95}{95+5} = \frac{95}{100} = 0.95$$

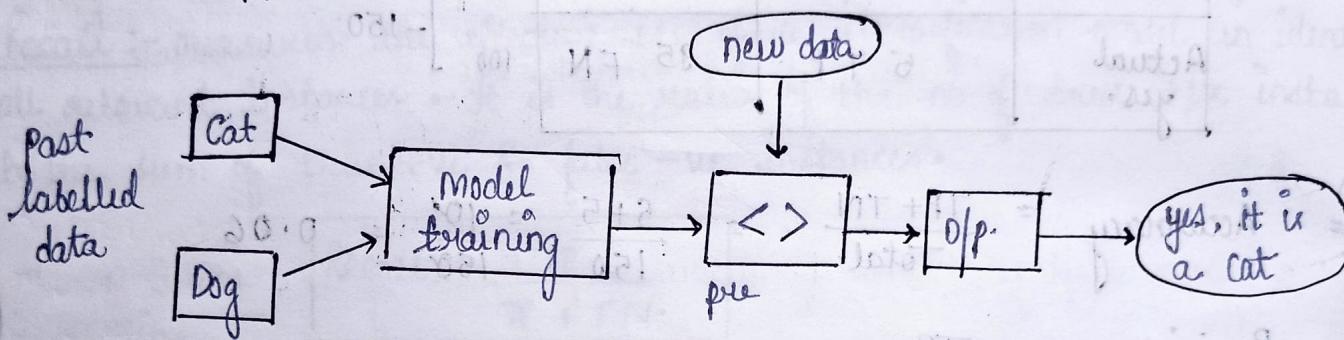
$$\text{Recall} = \frac{\text{TP}}{\text{TP} + \text{FN}} = \frac{95}{95+5} = \frac{95}{100} = 0.95$$

$$\text{F}_1 \text{ score} = \frac{2 \times \text{precision} \times \text{recall}}{\text{precision} + \text{recall}} = \frac{2 \times 0.95 \times 0.95}{0.95 + 0.95} = \frac{1.805}{1.90} = 0.95$$

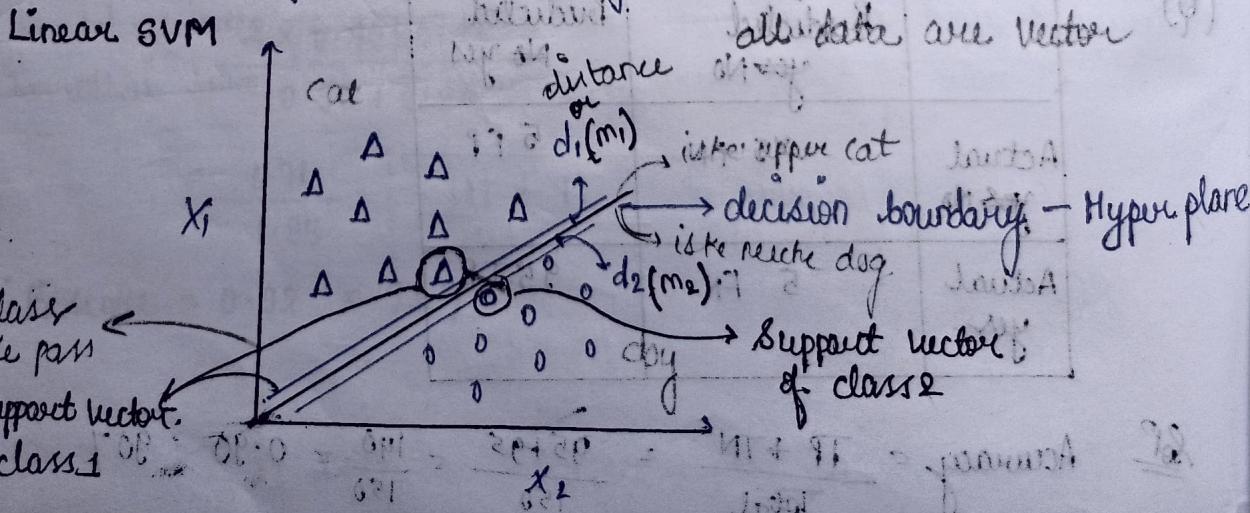
$$\text{Error} = 1 - \text{accuracy} = 1 - 0.9 = 0.1 = 10\%$$

20/02/2024

SVM (Support Vector Machine)



- A support vector m/c is a type of supervised m/c algo. used for in m/c learning to solve a classification problem.
- SVM are particularly good at solving binary classification problem which require classifying the elements of a dataset into two groups.
- The goal of the SVM algorithm is to create the best line or decision boundary that can segregate n-dimension space into classes so that we can easily put the new data point in the correct category in future.
- This best decision boundary is called a hyperplane.

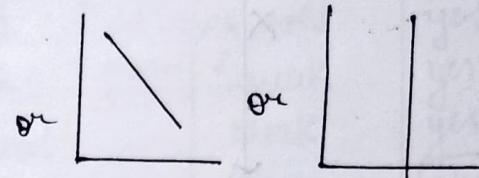


- Compare m_1 or d_1 & m_2 or d_2 which will having a minimum distance or margin from support vector will consider best line(hyperplane) for these data.

- Hyperplane can be shown mathematically as —

$$\vec{w} \cdot \vec{x} + b = \pm 1$$

intercept.



$$\frac{2}{(\text{max distance}) |\vec{w}|}$$

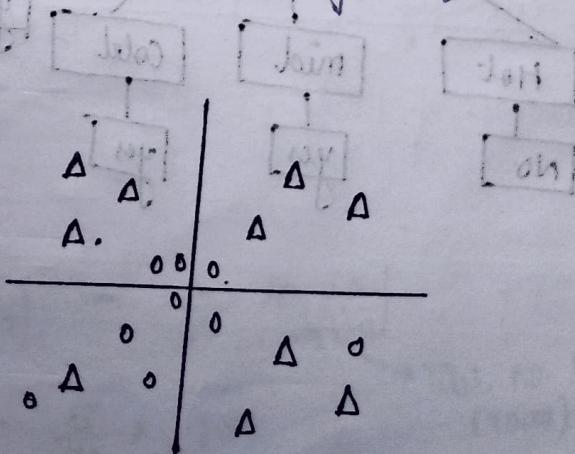
or $(\vec{w} \cdot \vec{x} + b) \geq 1 \quad \forall \vec{x} \text{ of class 1}$

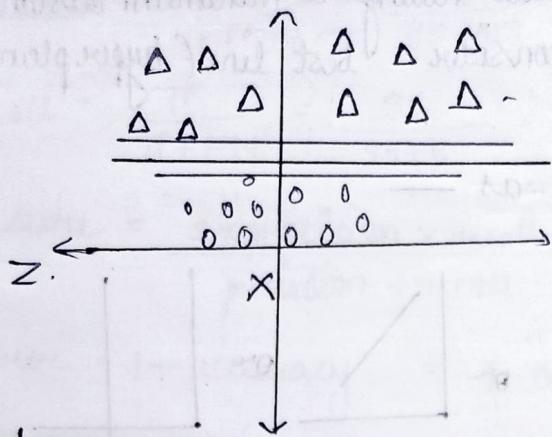
$(\vec{w} \cdot \vec{x} + b) \leq -1 \quad \forall \vec{x} \text{ of class 2}$

Non Linear SVM

- If data is linearly separable then we can separate it by using straight line but for non-linear data we can't draw a single straight line.
- So to separate those data points we need to add one more dimension.
- For linear data we have used two dimension x & y so for non-linear data we will add 3rd dimⁿ z .
- It can be calculated as

$$z = x^2 + y^2$$



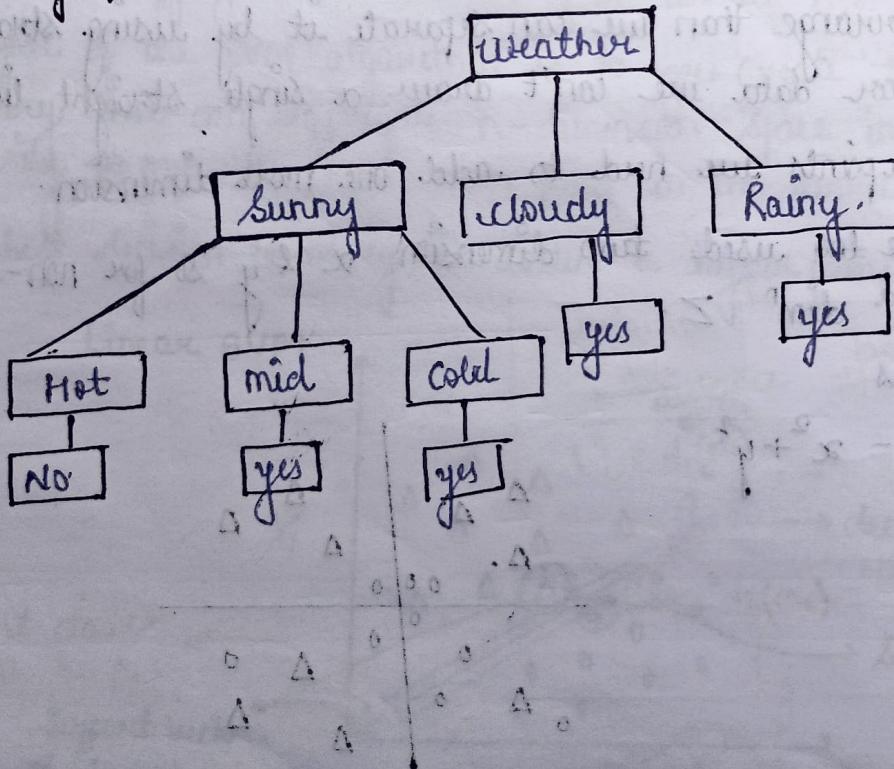


21/02/2024

Decision Tree

- Decision Tree is a machine learning algorithm used for both classification & regression.
- It is a tree structure
 - nodes
 - edges
 - root node
 - leaf node
- Decision nodes.
- Leaf nodes.
- Splitting.
- Entropy & information gain
- Pruning.

e.g.: - Play football outside



ID-3 Algorithm for Decision Tree

Day	Weather	Temperature	Humidity	Wind	Play
Day 1	Sunny	Hot	High	weak	No
Day 2	Sunny	Hot	High	Strong	No
Day 3	Cloudy	Hot	High	weak	yes
Day 4	Rain	Mild	High	weak	yes
Day 5	Rain	Cold	Normal	weak	No
Day 6	Rain	Cold	Normal	Strong	Strong
Day 7	Cloudy	Cold	Normal	weak	No
Day 8	Sunny	Mild	High	weak	yes
Day 9	Sunny	Cold	Normal	weak	yes
Day 10	Rain	Mild	Normal	Strong	yes
Day 11	Sunny	Mild	Normal	Strong	yes
Day 12	Cloudy	Mild	High	Strong	yes
Day 13	Cloudy	Hot	Normal	weak	yes
Day 14	Rain	Mild	High	Strong	No

Entropy :

Entropy is a metric to measure the impurity in a given attribute. It specifies randomness in data. Entropy can be represented as :

$$\text{Entropy } (S) = - P(\text{yes}) \log_2 P(\text{yes}) - P(\text{No}) \log_2 P(\text{No})$$

where S = Total no. of samples.

$P(\text{yes})$ = prob. of yes.

$P(\text{No})$ = prob. of No.

Step 1 : Entropy of entire dataset.

$$S = 9, - 5.4 = - \frac{9}{14} \log_2 \left(\frac{9}{14} \right) - \frac{5}{14} \log_2 \left(\frac{5}{14} \right)$$

Total no. Sample
(rows)

$$= - \frac{9}{14} \times -0.63 - \frac{5}{14} \times -1.48$$

= 0.405 + 0.528 = 0.933 \approx 0.94

Step 2 Entropy of all attributes

(i) Entropy of weather

- Entropy of sunny = $\{+2, -3\}$

$$= \frac{-2}{5} \log_2 \left(\frac{2}{5} \right) - \frac{3}{5} \log_2 \left(\frac{3}{5} \right)$$

Total no.
of sunny = 0.97.

- Entropy of cloudy $\{+4, -0\}$ = $-\frac{4}{4} \log_2 \left(\frac{4}{4} \right) - 0$

$= 0$

- Entropy of rain $\{+3, -2\}$ = $-\frac{3}{5} \log_2 \left(\frac{3}{5} \right) - \frac{2}{5} \log_2 \left(\frac{2}{5} \right)$

$= -\frac{3}{5} \times 0.73 - \frac{2}{5} \times -1.32$

$= 0.438 + 0.528$

$= 0.966$

$= 0.97$

Information Gain: weather

total no. sunny total no. of cloudy # rain

$$= \text{Entropy (whole dataset)} - \frac{5}{14} \text{Ent}(S) - \frac{4}{14} \text{Ent}(C) - \frac{5}{14} \text{Ent}(R)$$

Total no. of dataset

$= 0.94 - \frac{5}{14} \times 0.97 - \frac{4}{14} \times 0 - \frac{5}{14} \times 0.97$

$= 0.246$

Step 3 Entropy of Temperature

- Entropy of hot $\{+2, -2\}$ = $-\frac{2}{4} \log_2 \left(\frac{2}{4} \right) - \frac{2}{4} \log_2 \left(\frac{2}{4} \right)$

$= 1$

- Entropy of cold $\{+3, 1\}$ = $-\frac{3}{4} \log_2 \left(\frac{3}{4} \right) - \frac{1}{4} \log_2 \left(\frac{1}{4} \right)$

$=$

- Entropy of mild $\{+4, -2\}$ = $-\frac{4}{6} \log_2 \left(\frac{4}{6} \right) - \frac{2}{6} \log_2 \left(\frac{2}{6} \right) = 0.91$

(ii) Entropy of humidity

Information gain of Temp

$$\text{Entropy (WD)} = \frac{\text{Total \# Hot}}{\text{Total \# dataset}} \times \text{Ent}(H) - \frac{\text{Total \# cold}}{\text{Total \# dataset}} \times \text{Ent}(C)$$

$$= \frac{\text{Total \# Mild}}{\text{Total \# dataset}} \times \text{Ent}(M).$$

$$= 0.94 - \frac{4}{14} \text{Ent}(H) - \frac{6}{14} \text{Ent}(M) - \frac{4}{14} \text{Ent}(C)$$

$$= 0.029.$$

(iii) Entropy of humidity

$$(i) \text{Entropy of high } \{+3, -4\} = -\frac{3}{7} \log_2 \left(\frac{3}{7}\right) - \frac{4}{7} \log_2 \left(\frac{4}{7}\right)$$

$$= 0.5 + 0.46 \approx 0.96.$$

$$\text{Entropy of normal } \{+5, -2\} = -\frac{5}{7} \log_2 \left(\frac{5}{7}\right) - \frac{2}{7} \log_2 \left(\frac{2}{7}\right)$$

$$= 1$$

$$\text{Information gain} = 0.15$$

(iv) Entropy of wind

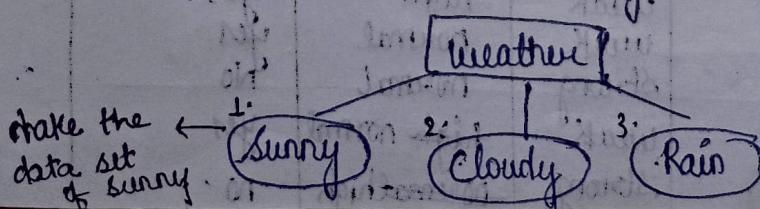
$$\text{Entropy of weak } \{-6, -2\} = -\frac{6}{8} \log_2 \left(\frac{6}{8}\right) - \frac{2}{8} \log_2 \left(\frac{2}{8}\right) \approx 0.81$$

$$\text{Entropy of strong } \{-3, -3\} = -\frac{3}{6} \log_2 \left(\frac{3}{6}\right) - \frac{3}{6} \log_2 \left(\frac{3}{6}\right) = 1$$

$$\text{Information gain} = 0.94 - \frac{8}{14} \times \text{Ent}(W) - \frac{6}{14} \times \text{Ent}(S)$$

$$= 0.0478.$$

\therefore Information gain of weather is higher so it will become root node



<u>Now, for sunny</u>	Day	Weather	Temp	Wind	Play.	Humidity
	Day 1	sunny	hot	weak	No	High
	Day 2	sunny	Hot.	Strong	No	high
	Day 8	sunny	mild	weak	No	High
	Day 9	sunny	cold.	weak.	yes	normal.
	Day 11	sunny	mild	Strong	yes.	normal.

Step 1: Entropy of sunny $f(2, -3) = \frac{2}{5} \log_2 \frac{2}{5} + \frac{3}{5} \log_2 \left(\frac{3}{5}\right) = 0.97$.

Step 2: Entropy of all attributes.

Entropy of hot $f(0, -2) = 0$

Entropy of mild $f(1, -1) = 1$

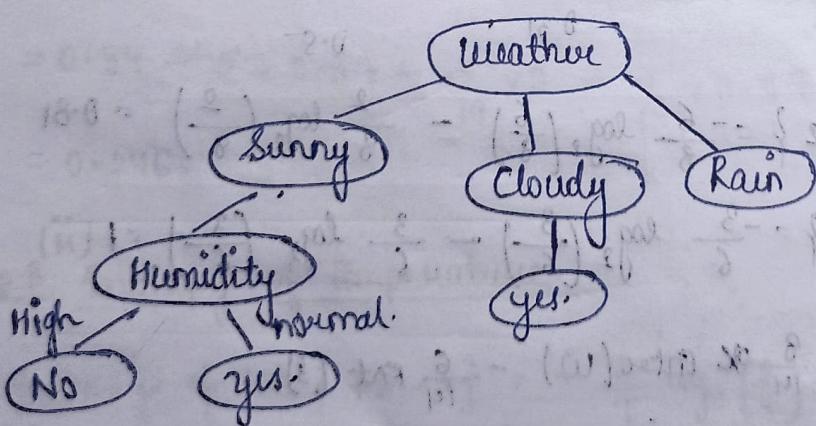
Entropy of cold $f(1, 0) = 0$

Information Gain

$$\text{Gain}(S_{\text{sunny}}, \text{Temp}) = 0.57$$

$$\text{Gain}(S_{\text{sunny}}, \text{Humidity}) = 0.97$$

$$\text{Gain}(S_{\text{sunny}}, \text{wind}) = 0.019$$



Now for Rain

Day	Weather	Temp	Wind	Humidity	Play
4	Rain	mild	weak	high	yes
5	Rain	cold	weak	normal	yes
6	Rain	cold	Strong	normal	No
10	Rain	mild	weak	high normal	yes
14	Rain	mild	Strong	normal high	No

Step 1: Entropy of rain = 0.97.

Step 2: (i) Entropy of temp.

$$\text{Entropy of cold} = -\frac{1}{2} \log_2\left(\frac{1}{2}\right) - \frac{1}{2} \log_2\left(\frac{1}{2}\right) = 0.5 + 0.5 = 1$$

$$\text{Entropy of mild} = -\frac{2}{3} \log_2\left(\frac{2}{3}\right) - \frac{1}{3} \log_2\left(\frac{1}{3}\right) = 0.38 + 0.38 + 0.52 = 0.90$$

(ii) Entropy of wind.

$$\text{Entropy of weak } \{+3, 0\} = 0 \quad \text{Entropy of strong } \{0, -2\} = 0$$

(iii) Entropy of humidity.

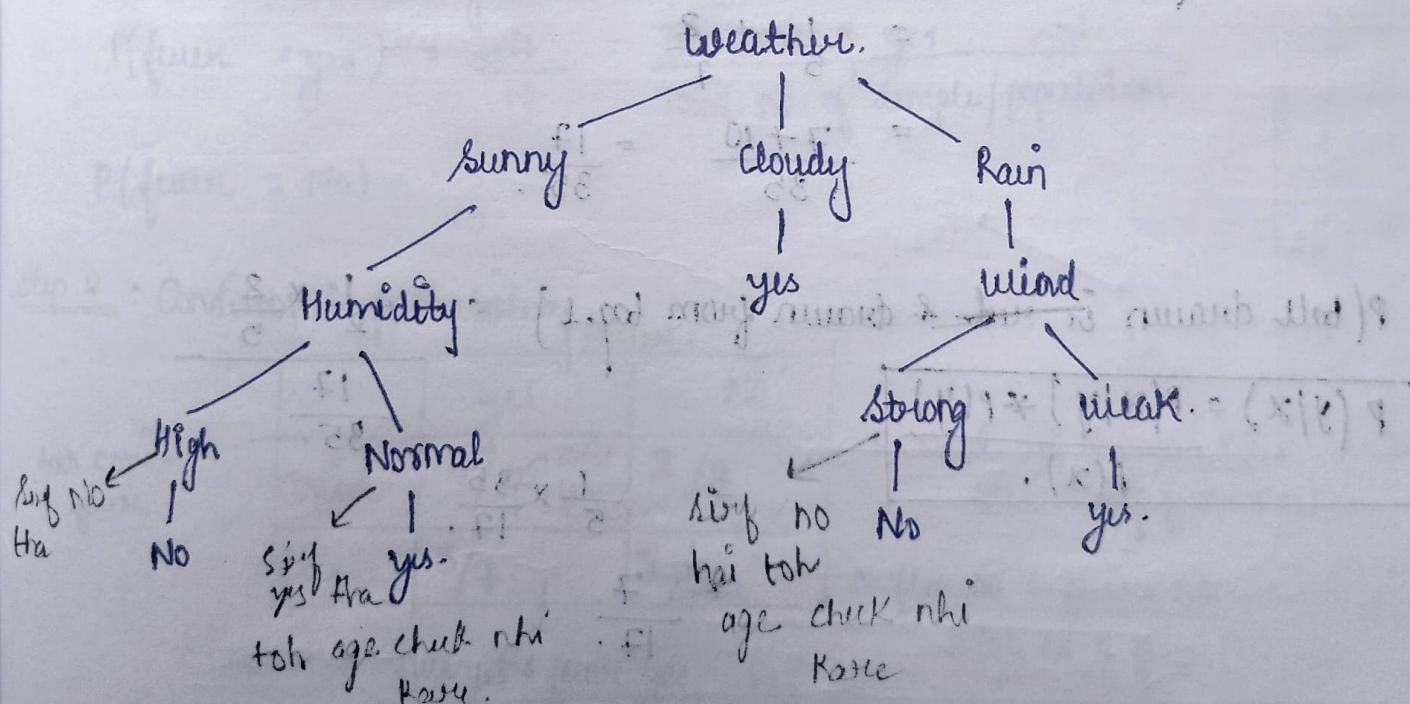
$$\text{Entropy of high } \{+1, -1\} = 1$$

$$\text{Entropy of normal } \{+2, -1\} = -\frac{2}{3} \log_2\left(\frac{2}{3}\right) - \frac{1}{3} \log_2\left(\frac{1}{3}\right) = 0.38 + 0.52 \\ = 0.90$$

$$\text{Information Gain (temp)} = 0.97 - \frac{2}{5} \times 1 - \frac{3}{5} \times 0.90 = 0.97 - 0.36 - 0.4 \\ = 0.03$$

$$\text{Information Gain (Wind)} = 0.97 - \frac{3}{5} \times 0 = 0 \quad \boxed{\approx 0.97}$$

$$\text{Information Gain (humidity)} = 0.97 - \frac{2}{5} \times 1 - \frac{3}{5} \times 0.90 \\ = 0.97 - 0.4 - 0.54 = 0.03$$



26/02/2024

Naïve - Bayes - Classification

- It is a supervised machine learning.
- It is based on applying Bayes theorem.

Bayes Theorem

Q1.) Bag 1 contains 2 red & 3 black balls. Bag 2 contain 4 red & 3 black balls. One ball is drawn at random from one of these bags & it's red. Find the probability that it is drawn from bag 1.

Sol'

Bag 1

2 red
3 black

Bag 2

4 red
3 black

 $P(B_1/R)$.

$$E_1 = \frac{1}{2}$$

$$E_2 = \frac{1}{2}$$

$$P(\text{red ball}) = \frac{2}{2+3} = \frac{2}{5}$$

$$P(\text{Red ball}) = \frac{4}{4+3} = \frac{4}{7}$$

→ Selecting one bag.

$$\text{Total probability} = \frac{1}{2} \times \frac{2}{5} + \frac{1}{2} \times \frac{4}{7}$$

$$= \frac{1}{5} + \frac{2}{7}$$

$$= \frac{7+10}{35} = \frac{17}{35}$$

$$P(\text{ball drawn is red & drawn from bag 1}) = \frac{\frac{1}{2} \times \frac{2}{5}}{\frac{17}{35}}$$

$$P(y/x) = \frac{P(x/y) * P(y)}{P(x)}$$

$$= \frac{1}{5} \times \frac{35}{17}$$

$$= \frac{7}{17}$$

If more than one variable :

$$P(y/x_1, x_2, \dots, x_n) = \frac{P(x_1/y) * P(x_2/y) \dots * P(x_n/y) * P(y)}{P(x_1) * P(x_2) * \dots * P(x_n)}$$

$$P(N/x) = \frac{P(x/N) * P(N)}{P(x)} \quad [\text{Probability of no. of } x]$$

$$P(N/x) = \frac{P(x_1/N) * P(x_2/N) * \dots * P(x_n/N) * P(N)}{P(x_1) * P(x_2) * \dots * P(x_n)}.$$

$$P(x_1) * P(x_2) * \dots * P(x_n).$$

(Q.)

Person	Covid yes/NO	Flu yes/NO	Fever yes/NO
1	yes	No	yes
2	No	yes	yes
3	yes	yes	yes
4	No	No	No
5	yes	No	yes
6	No	No	yes
7	yes	No	yes
8	yes	No	No
9	No	yes	yes
10	No	yes	No
11	yes	yes	yes

if new person is added
it will be yes, bcz yes
has higher probability

Step 1 : Prior Probability

$$P(\text{fever} = \text{yes}) = \frac{7}{10} = \frac{\text{Total no. of yes}}{\text{Total no. of sample/population}}$$

$$P(\text{fever} = \text{No}) = \frac{3}{10}.$$

Step 2 : Conditional probability

both covid & fever.

	yes	No	
covid.	$\frac{4}{7}$	$\frac{3}{7}$	$\frac{\text{covid. no} + \text{fever no}}{\text{Total no. of fever (no)}}$
flu	$\frac{3}{7}$	$\frac{2}{3}$	$\frac{\text{flu no} + \text{fever no}}{\text{Total no. of fever (no)}}$

flu yes & fever yes
Total no. of fever (yes)

$$P(\text{yes} / \text{flu, covid}) = P(\text{flu/yes}) + P(\text{covid/yes}) * P(\text{yes})$$

$$= \frac{3}{7} * \frac{4}{7} * \frac{7}{10} = \frac{12}{70} = 0.17$$

$$\begin{aligned}
 P(\text{NO flu, covid}) &= P(\text{flu}|\text{NO}) * P(\text{covid}|\text{NO}) * P(\text{NO}) \\
 &= \frac{2}{3} * \frac{2}{3} * \frac{5}{10.5} \\
 &= \frac{2}{15}
 \end{aligned}$$

$$P(\text{NO flu, covid}) = 0.13.$$

27/02/2024

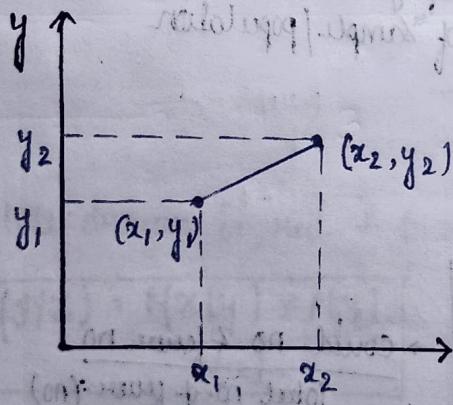
KNN - K-nearest neighbor algorithm (Classification)

Data Set

K always given in question.

If not given generally K is taken odd value, by default $K=5$, mostly $K=5$ is used.

x	y	Height (cm.)	Weight (Kg.)	Class	Dist.	Rank.
167	51			Underweight	6.7	5.
182	62			Normal	13	2
176	67			Normal	11.6	8
173	64			Normal	7.6	6
172	65			Normal	8.2	7
174	56			Underweight	4.1	4
169	58			Normal	1.4	1
173	57			Normal	3	3
170	55			Normal	2	2
170	57			Normal		



Step 1 : Calculate the distance & sort it

$$d_1 = \sqrt{(x_2 - x_1)^2 + (y_2 - y_1)^2} \quad x_1, y_1 - \text{variable}$$

$$d_1 = \sqrt{(170 - 167)^2 + (57 - 51)^2} = 6.7 \quad x_2, y_2 - \text{Value for which we have to calculate}$$

$$d_2 = \sqrt{(170 - 182)^2 + (57 - 62)^2} = 6.7$$

$$d_3 = \sqrt{(170 - 176)^2 + (57 - 67)^2} = 13.4$$

$$d_4 = \sqrt{(170 - 173)^2 + (57 - 64)^2} = 7.6$$

$$d_5 = \sqrt{(170 - 172)^2 + (57 - 65)^2} = 8.2$$

$$d_6 = \sqrt{(170 - 174)^2 + (57 - 56)^2} = 4.1$$

$$d_7 = \sqrt{(170 - 169)^2 + (57 - 58)^2} = 1.4$$

$$d_8 = \sqrt{(170 - 173)^2 + (57 - 57)^2} = 3$$

$$d_9 = \sqrt{(170 - 170)^2 + (57 - 55)^2} = 2$$

\rightarrow rank the nearest
Step 2 If $K=1$ = Normal = new value can be classify as normal.

If $K=3$ = Normal = new value can be classify as normal

If $K=5$ = 3 normal & 2 underweight = consider the majority

\therefore new value can be classify as normal

Step 3 : Majority Voting.

$K=5$: 3 normal & 2 underweight. If 3 normal, then new value is normal.
 & consider majority.

Q2 Predicting movie genre.

IMDb Rating	Duration	Genre	Distance	Rank
8.0 (Mission Impossible)	160	Action	46.0	2
6.2 (Gadar 2)	170	Action	56.01	4
7.2 (Rocky & Rani)	168	Comedy	54.0	3
8.2 (OMG 2)	155	Comedy	41.0	1

Now predict the genre of barbie movie with IMDb rating = 7.4 & duration 114

Q3 Step 1 : Calculate the distance.

$$d_1 = \sqrt{(7.4 - 8.0)^2 + (114 - 160)^2} = 46.0$$

$$d_2 = \sqrt{(7.4 - 6.2)^2 + (114 - 170)^2} = 56.0$$

$$d_3 = \sqrt{(7.4 - 7.2)^2 + (114 - 168)^2} = 54.0$$

$$d_4 = \sqrt{(7.4 - 8.2)^2 + (114 - 155)^2} = 41$$

Step 2 If $K=3$

2 comedy & 1 action.

If $K=1$. = comedy.

Step 3 A Barbie movie comes under Comedy generic (majority voting).

29/02/2024

(KNN) K-means Clustering

$A_1(2, 10), A_2(2, 5), A_3(8, 4)$

$B_1(5, 8), B_2(7, 5), B_3(6, 4)$

$C_1(1, 2), C_2(4, 9)$.

Sol Suppose initially we assign A_1, B_1, C_1 as the center of each cluster respectively.

Datapoints		Distance to						cluster	New cluster
x_2	y_2	2	10	5	8	1	2	smallest	
2	10	0	3.61	8.06				1	
2	5	5	4.24	3.16				3	
8	4	8.49	5.00	7.28				2	
5	8	3.61	0	7.21				1	
7	5	7.07	3.61	6.71				2	
6	4	7.21	4.12	5.39				2	
1	2	8.06	7.21	0				3	
4	9	2.24	1.41	7.62				2	

Step 1 $c_1 \quad c_2 \quad c_3$

Calculate the distance (c_1)

$$d = \sqrt{(x_2 - x_1)^2 + (y_2 - y_1)^2}$$

$$d = \sqrt{(2-2)^2 + (10-10)^2} = 0$$

$$d = \sqrt{(2-2)^2 + (5-10)^2} = \sqrt{25} = 5$$

$$C_2(\text{cluster 2}) : d = \sqrt{(x_2 - 5)^2 + (y_2 - 8)^2}$$

Step 2 : Under cluster column, find smallest value among 3 clusters & write the cluster no. having smallest value.

Step 3 : Calculate new centroid

$$\underline{1^{\text{st}} \text{ centroid}} : x = \frac{2}{1} = 2$$

$$y = \frac{10}{1} = 10$$

$$\underline{2^{\text{nd}} \text{ centroid}} : y = \frac{4+8+5+4+9}{5} = 6.$$

$$xy = \frac{8+5+7+6+4}{5} = 6.4$$

$$\underline{3^{\text{rd}} \text{ centroid}} : x = \frac{2+1}{2} = 1.5$$

$$y = \frac{5+2}{2} = 3.5$$

Initial Centroid

$$A_1 : (2, 10)$$

$$B_1 : (5, 8)$$

$$C_1 : (1, 2)$$

New Centroid

$$A_1 : (2, 10)$$

$$B_1 : (6, 6)$$

$$C_1 : (1.5, 3.5) \text{ from prev. table.}$$

Data points	Distance to						Cluster	New Cluster
	2	10	6	4	1.5	3.5		
2, 10	0	5.68	6.52	3.4	8	1, 3, 4		
2, 5	5	4.12	1.58	3	8	3		
8, 4	6.49	2.83	6.52	2	8	2		
5, 8	3.61	2.24	5.40	2	8	2		
7, 5	7.07	2.41	5.70	2, 3, 4, 8	2	2		
6, 4	7.21	2.80	4.53	2	8	2		
1, 2	8.06	6.41	1.58	3	8	3		
4, 9	2.24	3.61	6.04	2	8	1	C ₁ k value small	cat has

The datapoint is moved to one cluster to another cluster here so this is not considered here so we need to calculate new centroid.

New Centroid

$$A_1 : (3, 9.5)$$

$$B_1 : (6.5, 5.25)$$

$$C_1 : (1.5, 3.5)$$

Datapoints	Distance to						Cluster (value of new new cluster)	New Cluster
	3	9.5	6.5	5.25	1.5	3.5		
A ₁	2	10	1.12	6.54	6.52		1	1
A ₂	2	5	4.61	4.51	1.58		3	3
A ₃	8	4	7.43	1.95	6.52		2	2
B ₁	5	8	2.50	3.13	5.70		2	1
B ₂	7	5	6.02	0.56	5.70		2	2
B ₃	6	4	6.26	1.35	4.53		2	2
C ₁	1	2	7.76	6.39	1.58		3	3
C ₂	4	9	1.12	4.51	6.04			

New Centroid

$$A_1 : x = \frac{2+4+5}{3} = 3.67, y = \frac{8+5+4}{3} = 5.67$$

$$y = \frac{10+8+9}{3} = 9$$

$$B_1 : x = \frac{8+7+6}{3} = 7, y = \frac{4+5+4}{3} = 4.33$$

$$C_1 : (1.5, 3.5)$$

Data point	Distance to						(Old) cluster	(New) cluster
	1	2	3	4	5	6		
A ₁	2	10	1.94	7.56	6.52	1	1	1
A ₂	2	5	4.83	5.04	1.53	3	3	3
A ₃	8	4	6.82	1.05	6.52	2	2	2
B ₁	5	8	1.67	4.18	5.70	1	1	1
B ₂	7	5	5.21	0.67	5.70	2	2	2
B ₃	6	4	5.52	1.05	4.53	2	2	2
C ₁	1	2	7.49	6.44	1.58	3	3	3
C ₂	4	9	0.33	5.55	6.04	1	1	1

1st cluster : Data point A₁, B₁, C₂

2nd cluster : Data point A₃, B₂, B₃

3rd cluster : Data point A₂, C₁

jume 1 : 1st cluster.
jume 2 : 2nd

01/03/2024

(Q.)	Data point	Coordinates
A ₁		(2, 10)
A ₂		(2, 6)
A ₃		(11, 11)
A ₄		(6, 9)
A ₅		(6, 4)
A ₆		(1, 2)
A ₇		(5, 10)
A ₈		(4, 9)
A ₉		(10, 12)
A ₁₀		(7, 5)
A ₁₁		(9, 11)
A ₁₂		(4, 6)
A ₁₃		(3, 10)
A ₁₄		(3, 8)
A ₁₅		(6, 11)

Centroids

Centroid 1 : (2, 6)

Centroid 2 : (5, 10)

Centroid 3 : (6, 11)

Sol

Datapoints		Distance to							Cluster
x_1	y_1	2	6	5	10	6	11		
2	10	8.4		3		4.12		2	
2	6	0		5		6.40		1	
11	11	10.2		6.08		5		3	
6	9	5		1.41		2		2	
6	4	4.47		6.08		7		1	
1	2	4.12		8.94		10.2		1	
5	10	5		0		1.41		2	
4	9	3.6		1.41		2.82		2	
10	12	10		5.38		4.12		3	
7	5	-5.09		5.38		6.08		1	
9	11	8.60		4.12		3		3	
4	6	2		4.12		5.38		1	
3	10	4.12		2		3.16		2	
3	8	2.23		2.82		4.24		1	
6	11	6.40		1.41		0		3	

$$\therefore A_1 = \sqrt{(2-2)^2 + (6-10)^2} = \sqrt{0 + 16} = 4 \text{ (miles)}$$

New Centroid

$$C_1: x = \frac{2+6+1+7+4+3}{6} = 3.85$$

$$y = \frac{6+4+2+5+6+8}{6} = 5.316$$

$$C_2: x = \frac{2+5+3+4}{4} = 3.333\bar{3}$$

$$y = \frac{10+10+10+9}{4} = 9.75$$

$$C_3: x = \frac{11+6+10+9+6}{5} = 8.4$$

$$y = \frac{9+12+11+11+11}{5} = 10.8$$

$$x = \frac{2+6+5+4+3}{5} = 4$$

$$y = \frac{10+9+9+10+10}{5} = 9.6$$

$$x = \frac{11+10+9+6}{4} = 9$$

$$y = \frac{11+12+11+11}{4} = 11.25$$

Data points	Distance to					Cluster	New cluster
	3.8	5.16	4.96	9.11.25	10.12.13		
2 10	5.06	2.03	7.11	2	CP-D	2	2,3
2 6	1.79	4.11	8.75	1	NH-A	1	2,4
11 11	9.5	7.13	2.01	3	SP-B	3	11,12
6 9	4.5	2.08	3.75	2	BP-C	2	11,12
6 4	2.75	5.94	7.84	1	SP-C	1	3,4
1 2	4.02	8.17	12.22	1	BP-C	1	3,4
5 10	5.06	3.07	4.19	2	BP-B	2	1,2
4 9	3.87	0.6	5.48	2	BP-B	2	1,2
10 12	9.43	6.46	1.25	3	P-A	3	1,2
7 5	3.5	5.49	6.56	1	SP-A	3	1,2
9 11	8.02	5.19	0.25	3	SP-B	3	1,2
4 6	0.97	3.6	7.25	1	SP-B	1	3,4
3 10	4.86	1.07	6.12	2	BP-C	2	1,2
3 8	2.86	1.88	6.82	1	P-A	2	3,4
6 11	6.35	2.44	3.01	3	BP-C	2	1,2

New centroid

$$C_1: x = \frac{2+6+1+7+4}{5} = 4$$

$$y = \frac{6+4+2+5+6}{5} = 4.0$$

$$C_2: x = \frac{2+6+5+4+3+3+6}{7} = 4.19$$

$$y = \frac{10+9+10+9+10+8+11}{7} = 9.57$$

$$C_3: x = \frac{11+10+9}{3} = 10$$

$$y = \frac{11+12+11}{3} = 11.33$$

3.8	5.16	4.96	9.11.25	10.12.13	3.8	5.16	4.96	9.11.25	10.12.13
F	8	F	F	A	3	P	E	E	S
S	C	H	S	V	S	F	S	P	S

Data points	Distance to					Cluster	New Cluster
	4	4.6	4.14	9.57	10		
A ₁ 2 10	5.75	2.18	8.10			2	2
A ₂ 2 6	2.44	4.16	9.61			1	1
A ₃ 11 11	9.48	7.00	1.05			3	3
A ₄ 6 9	4.83	1.94	8.35	4.62		2	2
A ₅ 6 4	2.08	5.87	8.35			1	1
A ₆ 1 2	3.96	8.19	12.96			1	1
A ₇ 5 10	5.49	0.96	5.17			2	2
A ₈ 4 9	4.4	0.58	6.43			2	2
A ₉ 10 12	9.52	6.34	0.67			3	3
A ₁₀ 7 5	3.02	5.39	7.00			1	1
A ₁₁ 9 11	8.12	5.06	1.05			3	3
A ₁₂ 4 6	1.4	3.57	8.02			1	1
A ₁₃ 3 10	5.49	1.21	7.12			2	2
A ₁₄ 3 8	3.54	1.94	7.75			2	2
A ₁₅ 6 11	6.70	2.34	4.01		2	2	

$\therefore 1^{\text{st}}$ cluster data point A₂, A₅, A₆, A₁₀, A₁₂

$$\frac{x}{z} \rightarrow \frac{1+1+1+2+2}{5} = 1.6$$

2^{nd} cluster data point A₁, A₄, A₇, A₈, A₁₃, A₁₄, A₁₅

3^{rd} cluster data point A₃, A₉, A₁₁

04/03/2024

KNN - K-Means Clustering

(Q) Apply K-means clustering algorithm to form two clusters

i	x ₁	x ₂	x ₃	x ₄	x ₅	x ₆	x ₇	x ₈	x ₉	x ₁₀
x ₁	2	3	3	4	6	6	7	7	8	7
y ₁	6	4	8	7	2	4	3	4	5	6

- Step 1: Select two medoids (Random).
- $c_1 = (3, 4)$.
 - $c_2 = (7, 4)$

Step 2:

- use the Manhattan distance to find distance between datapoint & medoid
- (x_1, y_1) and (x_2, y_2) are datapoint.
- Manhattan distance = $|x_1 - x_2| + |y_1 - y_2|$.
- $Mdist_{c_1}[(2, 6), (3, 4)] = |2-3| + |6-4| = 3$.
- $Mdist_{c_1}[(3, 4), (3, 4)] = |3-3| + |4-4| = 0$.
- $Mdist_{c_1}[(3, 8), (3, 4)] = |3-3| + |8-4| = 4$

Similarly $Mdist_{c_2}[(2, 6), (7, 4)] = |2-7| + |6-4| = 7$.

$$Mdist_{c_2}[(3, 4), (7, 4)] = |3-7| + |4-4| = 4$$

i	x	y	$c_1(3, 4)$	$c_2(7, 4)$	Smaller cluster
x_1	2	6	3	7	C_1
x_2	3	4	0	4	C_1
x_3	3	8	4	8	C_1
x_4	4	7	4	6	C_1
x_5	6	2	5	3	C_2
x_6	6	4	3	1	C_2
x_7	7	3	5	1	C_2
x_8	7	4	4	0	C_2
x_9	8	5	6	2	C_2
x_{10}	7	6	6	2	C_2

$$C_1 : \{(2, 6), (3, 4), (3, 8),$$

$$(4, 7), (6, 2)\}$$

$$C_2 : \{(6, 4), (7, 3), (7, 4),$$

$$(8, 5), (7, 6)\}$$

Step 3: calculate the total cost.

$$\text{cost} = \sum_i |c_i - x_i|$$

Total cost = Cost $((3,4), (2,6)) + \text{Cost}((3,4), (3,8)) + \text{Cost}((3,4), (3,4))$
 $+ \text{Cost}((3,4), (4,7)) + \text{Cost}((4,7), (6,2)) + \text{Cost}((7,4), (6,4)) +$
 $\text{Cost}((7,4), (7,3)) + \text{Cost}((7,4), (7,4)) + \text{Cost}((7,4), (8,5)) +$
 $\text{Cost}((7,4), (7,6)).$

$$= 3+4+4+2+3+1+1+2.$$

$$\text{Cost} = 20.$$

Step 4: Randomly select one non medoid point & recalculate the cost.

$$C_1 = (3,4) \quad \& \quad C_2 = (7,4)$$

$$\text{Now } C_2 = 0 = (7,3).$$

Now, calculate the distance.

$$M \text{ dist}_0 = [(2,6), (7,4)] = |(2-7)| + |(2-7)| + |6-4| = 9.$$

i	x	y	$C_1 (3,4)$	$C_2 (7,4)$	cluster
x_1	2	6	3	8	C_1
x_2	3	4	0	5	C_1
x_3	3	8	4	9	C_1
x_4	4	7	4	7	C_1
x_5	6	2	5	2	C_2
x_6	6	4	3	2	C_2
x_7	7	3	5	0	C_2
x_8	7	4	4	1	C_2
x_9	8	5	6	3	C_2
x_{10}	7	6	6	3	C_2

New cluster are :

$$C_1 (3,4) = \{(2,6), (3,4), (3,8), (4,7)\}$$

$$C_2 (7,3) = \{(6,2), (6,4), (7,3), (7,4), (8,5), (7,6)\}$$

Calculate the ^{new} cost = Cost $((3,4), (2,6)) + \text{cost}((3,4), (3,4)) + \text{cost}((3,4), (3,4)) +$
 $\text{cost}((3,4), (4,7)) + \text{cost}((7,3), (6,2)) + \text{cost}((7,3), (6,4)) +$
 $\text{cost}((7,3), (7,3)) + \text{cost}((7,3), (7,4)) + \text{cost}((7,3), (8,5)) +$
 $\text{cost}((7,3), (7,6))$.

$$\text{new cost} = 22$$

Based on this new & prev. total cost we have to decide.

$$S = \text{Current total cost} - \text{previous total cost}$$

$$= 22 - 20 = 2 > 0$$

Cost S is greater than 0, so C_2 replacing with 0 is not a good idea
 bcoz prev. cost was less than new cost.

So the final medoids are $C_1 (3,4)$ & $C_2 (7,4)$ i.e the prev. medoid in
 this case