



Page NO.

1. 1 to 18 – Mid Term
2. 19 to 45- Final Term

MD IFTAKHAR KABIR SAKUR

25th BATCH

COMPUTER AND COMMUNICATION ENGINEERING

International Islamic University Chittagong

COURSE CODE: CCE-4829

COURSE TITLE: Machine Learning

COURSE TEACHER:

Saiful Islam

Adjunct Faculty
Computer and Communication Engineering

Machine Learning

MidTerm

24.02.24

CEE-4829

Differences between ML Algorithm & Traditional Rule Based Algorithms:-

Rule Based

- ⇒ ① Human crafted
- ② Interpretability
- ③ Limited Adaptability
- ④ Limited Scalability
- ⑤ Domain specificity

Machine Learning:-

- Data Driven Approach
- Automated Feature
- Adaptability
- Scalability
- Complexity
- Black Box Nature

Main Problems Solved by Machine Learning:-

(i) Classification:-

In ML it is about training a model to predict the category or class of an input based on its features.

Ex:- Email is spam or not spam

(ii) Regression:-

→ Drawing a line or curve through scatterplot of points to understand relationship between two or more variables.

Regression models predict continuous numerical values based on input features.

Ex:- Ice cream sales based on temperature.

(iii) Clustering:-

Finding natural groupings within a set of objects based on their characteristics & similarities.

Ex:- A collection of ^{shapes} boxes of different colors.

Where similar shapes will be grouped then based on their similarities inside the group

Machine Learning Classification

① Supervised Learning:- (Classification/Questions)

means, having a teacher who teaches about which one is what.

The algorithm learns from labeled data. Where each example is associated with outcome. It allows the algorithm to make predictions on new, unseen data.

② Unsupervised Learning:- (Clustering)

The algorithm learns patterns & structure from unlabeled data. The goal is to find similarities.

Ex:- 6 Fruits are given. Now label them based on similarities.

③ Semi-Supervised:-

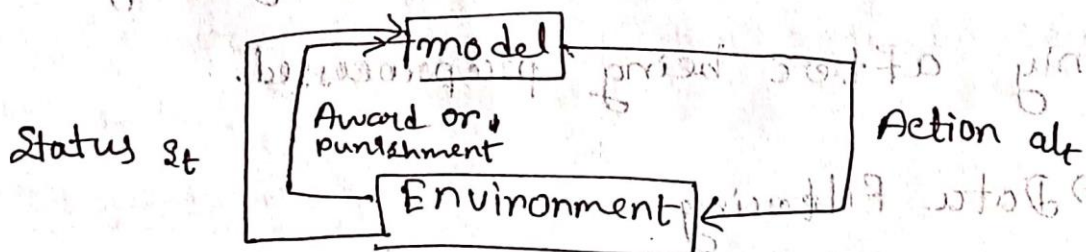
The algorithm learns from a combination of labeled & unlabeled data.

Use both the labeled & unlabeled images of animals. While you have some examples

OF dogs & cats labeled, many others are unlabeled. You ~~can~~ use both labeled & unlabeled data.

(iv) Reinforcement Learning:

Learns to make decisions by interacting with an environment. Receives Feedback in the form of rewards or penalties based on actions & its goal is to learn a policy.



Basic Machine Learning Concept

(i) Data Set: A collection of data used in machine learning ~~training~~ tasks. Each data record is called a sample.

(ii) Training Set: - Dataset used in training process. Where each sample is referred to as a

training sample. The process of creating model from data.

(ii) Test set - process of using the model obtained after learning for prediction.

Data Cleansing

Most machine learning models process features which are numeric.

Collected data can be used by algorithms only after being preprocessed.

→ Data Filtering

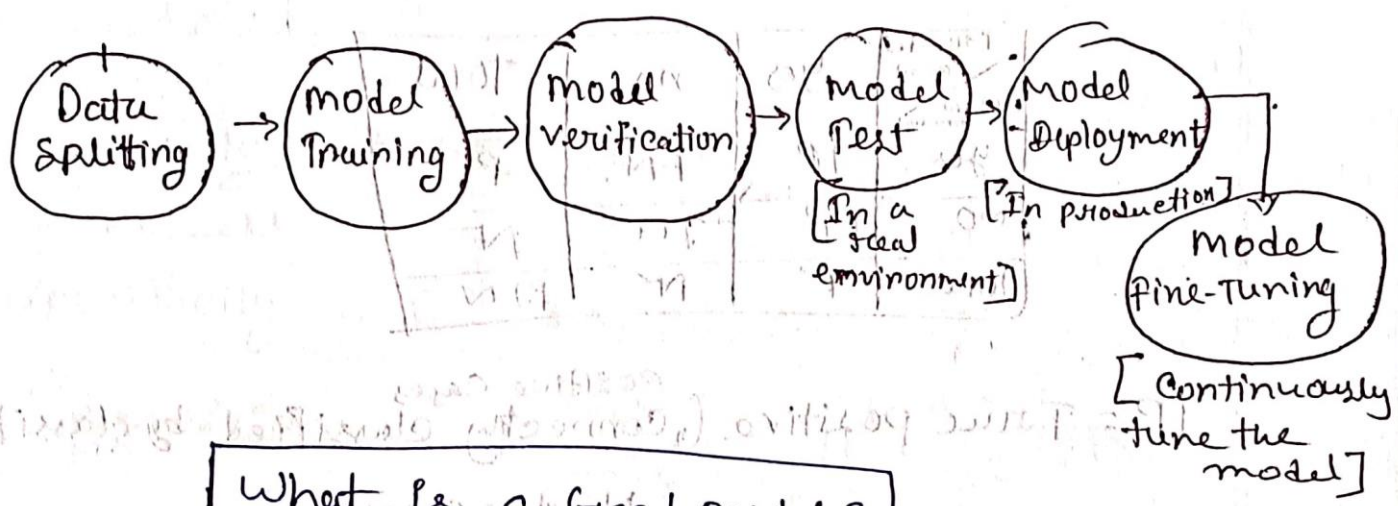
→ processing of lost data

→ processing of possible exceptions, errors

→ Combination of data from multiple data sources.

→ Data Consolidation

Overall procedure of Building a model



What is a Good Model?

- Generalization Capabilities:
 - ↳ Accurately predict the actual service data.
- Interpretability:
 - The prediction result easy to interpret.
- Prediction Speed:
 - ↳ How long does it take
- Practicability:
 - ↳ Is the prediction result acceptable or not.

Machine Learning Performance Evaluation

Estimated Amount Actual Amount	yes	no	Total
yes	TP	FN	P
no	FP	TN	N
Total	P'	N'	P+N

TP = True positive (Positive cases correctly classified by classifier)

FN = False Negative (Negative cases incorrectly classified)

P = Positive (The number of real positive cases)

FP = False positive (positive cases incorrectly classified by classifier)

TN = True Negative (neg. cases correctly classified by classifier)

N = Negative (Neg. cases in data)

(a) Accuracy & recognition rate: $\frac{TP+TN}{P+FN}$

(b) Error rate & missclassification rate: $\frac{FP+FN}{P+FN}$

(c) Sensitivity, true positive rate & recall: $\frac{TP}{P}$

(d) Specificity & True Neg. rate: $\frac{TN}{N}$

(e) precision = $\frac{TP}{TP+FP}$ | Recall $\rightarrow \frac{TP}{TP+FN}$

(f) Harmonic mean of the recall state = $\frac{2 \times \text{precision} \times \text{recall}}{\text{precision} + \text{recall}}$

(g) F_β ($\beta = \text{Non negative real Number}$) = $\frac{(1+\beta) \times \text{precision} \times \text{recall}}{\beta \times \text{precision} + \text{recall}}$

Hyperparameters of a Model

→ Learning rate, Number of layers, neurons in each layer, epoch these are hyper parameters.

→ Manually specified

→ Code:- `model.coef_`
`model.intercept_`

→ Often used in model parameter estimation process.

→ Often specified by the practitioner

→ Can often be set using heuristics.

Hyperparameter Search procedure & Method

- ① Dividing dataset into a training set, validation set & test set.
- ② Optimizing the model parameters using training set.
- ③ Searching for the model's hyper-params. using the validation.
- ④ perform step-2 & step-3 alternately.

Search Algorithm (Step-3):-

- ① Grid Search & Exploring a grid of possible combinations of hyperparameters.
- ② Random Search: - Randomly exploring different regions of the hyperparameters space.
- ③ Heuristic Intelligent Search: - To find approximate solutions to complex problems efficiently.
- ④ Bayesian Search: It is like using past observations to guide the search to get better hyperparameters.

Cross validation

Cross validation is a powerful technique for estimating a machine learning model's performance & assessing its ability to generalize to new data.

The basic idea is to divide the original dataset into two parts: Training set & validation set.

Entire set

Training set

Test set

Training set

validation set

Test set

K-Fold Cross validation is a hyperparameters

K-Fold Cross validation:-

used to evaluate the performance of a ML model more reliably than a simple train-test split.

~~K-Fold Cross~~

→ Divide the raw data into k groups

→ Each subset as validation test.

$k-1$ subsets as the training set.

→ k models as the performance indicator of the k -ev classifier.

Reasons for underfitting:-

- High bias & low variance
- Training dataset not enough
- Model is too simple
- Training data is not cleaned & having noise in it.

Techniques to reduce underfitting:-

- Increase model complexity
- Increase number of features, performing Feature Engineering
- Remove noise from the data
- Increase the number of epochs or increase the duration of training to get better results.

Reasons for overfitting:-

- High variance low bias
- The model is too complex
- Size of the training data

Technique to reduce overfitting:-

- Increase training data
- Reduce complexity.

Avoid overfitting:-

- Cross validation
- Training with more data
- Removing features
- Early stopping the training

K-means Clustering

→ Similar उत्सार्क कमार्क सार्क क Clustering
कर 24।

Make 2 Clustering based on given values:-

<u>Serial no</u>	<u>Age</u>	<u>Amount</u>
C ₁	20	500
C ₂	40	1000
C ₃	30	800
C ₄	18	300
C ₅	28	1200
C ₆	35	1400
C ₇	45	1800

⇒ Centralized values:

(C₁)
(20, 500)

(C₂)
(40, 1000)

For, c_3

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

(centroid) (centroid)

So, for $c_3 = \sqrt{(30 - 20)^2 + (800 - 500)^2}$ [This calculation based on c_1]

= $\sqrt{(10)^2 + (300)^2}$

= 300...

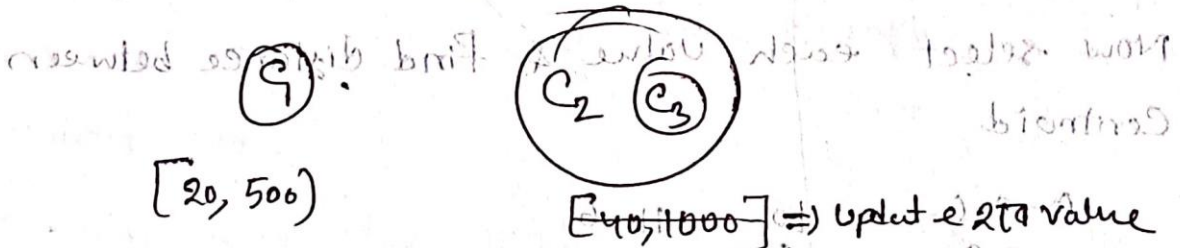
c_3 , value based on c_2

$c_3 = \sqrt{(30 - 40)^2 + (800 - 1000)^2}$

= 200

So, $200 < 300$

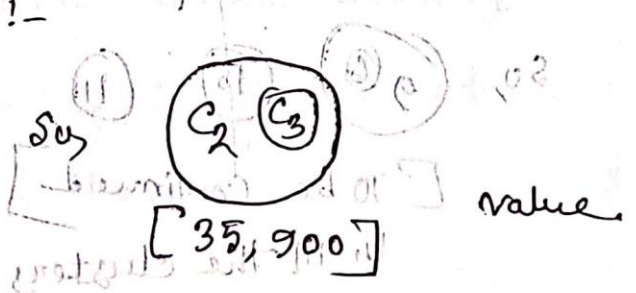
So, c_3 will be in c_2



Centroid value updation!-

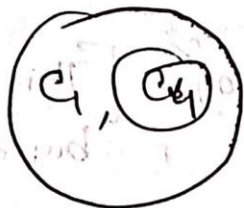
$\frac{40 + 30}{2} = 35$

$\frac{800 + 1000}{2} = 900$



$$\text{For, } C_1 = \sqrt{(18-20)^2 + (300-500)^2} =$$

$$C_1 = \sqrt{(18-35)^2 + (300-900)^2} =$$



→ Update:-

$$\frac{18+20}{2} = 19.5$$

$$\frac{500+300}{2} = 400$$

$$[19.5, 400]$$

Example-2

[6, 7, 8, 9, 10, 11, 12, 13, 14]



C_1

9



C_2

10



C_3

11

[Randomly Chosen]

Now select each value & find distance between Centroid

$$9-6 = 3 \quad 10-6 = 4 \quad 11-6 = 5$$



[To be Continued]

until the clusters don't change anymore.

Model validity

Generalization Capabilities:-

The goal of ML is the model should be working well on new samples not only on the training sample. And this capability of applying a model to new samples is called generalization or robustness.

Error:-

The difference between the sample that are predicted by the model obtained after learning & actual sample ~~data~~ result.

→ Training error:- Error during training

→ Generalization error:- Error when the model is run on new sample.

Underfitting:-

occurs when the model or the algorithm does not fit the data well enough.

Overfitting:- When learning is small but the generalization error is large.

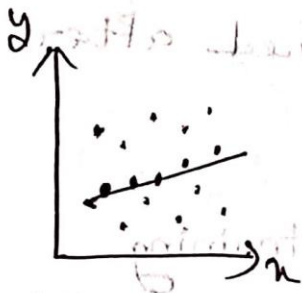
Model Capacity - [Bias Variance]

Model's capacity of fitting functions which is called as model complexity.

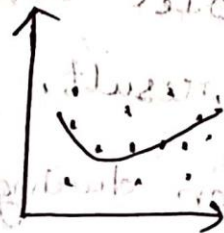
→ If all model suits the task the algorithm is optimal

→ Model with insufficient capacity can't solve complex tasks.

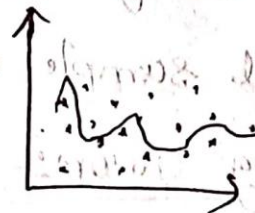
→ High capacity can solve but solve complex tasks.



Underfitting



Good Fitting



Overfitting

Underfitting: Occurs when the model is too simple and cannot capture the underlying pattern of the data. It results in high bias and high error.

Overfitting: Occurs when the model is too complex and captures noise in the training data. It results in low bias and low error on training data, but high error on new data.

1 (b) working principle OF ML:-

- Data Collection
- Data preprocessing
- Model selection
- Model Training
- Model Evaluation
- Model deployment

1 (c) Image recognition most Common Application

3 (a)

Clustering play in unsupervised learning:-

- Grouping similar data
- NO prior labels
- Clustering Algorithm
- ~~App~~ Evaluation
- parameter tuning.



**KEEP
CALM
ITS TIME FOR THE
FINAL
EXAM**

[Machine Learning]

Q) What is Hidden Markov Model?

⇒ A statistical model that represents a system containing hidden states where the system evolves over time.

It is hidden cause the state of the system is not directly visible to the observer. But the observer can see some outputs that depend on the state.

→ Markov models are characterized by the Markov property which states that the future state of a process only depends on its current state, not on the sequence of it.

It is used in: Speech, Handwriting, gesture recognition, part-of-speech tagging, musical score following.

There we represent the probabilities of the system being in a particular state at the beginning of the process.

[Date]

[Page No.]

[Machine Learning]

Components of Hidden Markov Model:-

(1) Hidden States:- Not directly observable. ex:- Speech recognition system.

(2) Observations:- Data point or output that can be observed. ex:- In speech recognition, the observations could be the audio signal feature.

(3) Transition probabilities:-

Transition from one hidden state to another.

These are usually represented in a matrix form.

(4) Emission probabilities:-

These describe the probability of an observed output, given a hidden state.

(5) Initial state probabilities

These represent the probabilities of the system being in a particular state at the beginning of the process.

Q2] How Hidden Markov Model works?

① Evaluation problem - Given the model parameter

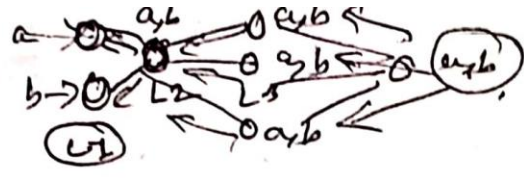
& an observation sequence, determine the probability of the observation sequence. This is solved using forward-backward algorithm.

② Decoding problem - The decoding problem is an

observed sequence of data; this decoding problem is to determine the most likely sequence of hidden state.

③ Learning - The learning problem is to

estimate the model parameters. The Baum-welch algorithm is often used for this.



Hidden Markov Model Algorithm

Step-1:- Define state space & observation space (1)
 (All possible hidden state & observation space)

Step-2:- Define the initial state (The probability distribution)

Step-3:- Define the state transition
 (Transition from one state to another)

Step-4:- Define the ~~the~~ observation likelihoods.
 (Probabilities of generating each observation from each state)

Step-5:- Train The model.
 (Baum Welch Algorithm, Forward-backward algorithm)

Step-6:- Evaluate model (Metrics Accuracy, precision, F1 score)

Limitation of Hidden Markov Model

① Limited Modeling Capabilities:-

→ designed to model sequences of data.

→ Ex:- In speech recognition, the complex relationship between the speech sounds & the

corresponding acoustic signals may not fully

captured by the simple structure of an HMM.

② Overfitting:- when the model fits too well

for the training data & can't generalize

to new data. As a result, high error rates

can result. So, needs to be careful to choose

the number of hidden states.

③ Robustness:- ^{lack of} limited in their robustness to

noise & variability in the data.

Ex:- In speech recognition the acoustic signal

there the data can have noise in it, which

can result in poor performance.

Computational Complexity

→ limited in its ability to handle large amounts of data.

→ Specially when dealing with large amounts of data.

of data.

→ And to address this limitation, it is

often necessary to use parallel computing techniques or to use

approximate that reduce the computational complexity of the model.

approximate that reduce the computational complexity of the model.

complexity of the model.

For the training data & test data.

to new data. As a result high error rates

can result. So, needs to be careful to choose

the number of hidden states.

of hidden states is limited in their robustness to

noise & variability in the data.

For the speech recognition the acoustic signal

from the data can have noise in it. which

can result in poor performance.

Topic:

Topic: Decision Tree (Theory)

⇒ A supervised learning algorithm which looks like an inverted tree, each node represents a predictor variable, the link between the nodes represents decision tree & each leaf node represent an outcome.

☐ How does the decision tree algorithm works?

⇒ Step-1 Select the feature that best classifies the data set into desired classes.

Step-2 Traverse down from root node.

Step-3 Route back to step-1.

☐ Decision Tree Structure

→ Root Node: The starting point of a tree.

→ Internal node: Represents a decision point that eventually leads to the prediction of the outcome.

→ Leaf Node: Represents the final class of the outcome.

Branches : Are connections between nodes.

Expressiveness of decision trees :-

Can represent any boolean function of the input attributes.

A	B	$A \& B$
F	F	F
F	T	F
T	F	F
T	T	T

And

A	B	$A \text{ OR } B$
F	F	F
F	T	T
T	F	T
T	T	T

OR

A	B	$A \oplus B$
F	F	F
F	T	T
T	F	T
T	T	F

X-OR

$$A \cdot \bar{B} + \bar{A} \cdot B$$

$$F + 0 \cdot F$$

$$F \cdot F + T$$

Q] How a decision tree is created:-

→ Using ID3 Algorithm.

→ ID3 = Iterative Dichotomiser 3 Algorithm

→ Uses the concept of Entropy & information gain to generate a decision tree for a given set of data.

ID3 Algorithm:

① Best Attribute selection. (Separate the dataset into different classes. Can get

↳ Information Gain
↳ Entropy

② Assign A as a decision variable for the root

③ For each A, build a descendant of the

④ Assign classification levels to the leaf node.

⑤ If data is correctly classified :- stop

⑥ else! iterate over the tree.

□ Entropy: Measures the impurity or uncertainty of present in the data.

$$\text{Entropy} = - \sum P(x) \cdot \log P(x)$$

□ Info. gain: How much information a particular feature/variable gives us about the final outcome.

$$\text{Info gain}(c) = \text{entropy}(\text{parent}) - [\text{weighted avg}]^* \text{entropy}(\text{children})$$

$$\text{entropy}(S) = \text{Entropy}(S, T)$$

□ Entropy(children) with weighted averages:

$$[\text{weighted avg}] \text{Entropy}(\text{children}) = \left(\frac{\text{no. of outcomes in left child node}}{\text{total no. of outcome in parent node}} \right) \cdot (\text{entropy of the left node})$$

$$+ \left(\frac{\text{no. of outcome in right child node}}{\text{total no. of outcomes in parent node}} \right) \cdot (\text{entropy of right node})$$

So

□ Entropy for two variable:

$$\text{Entropy}(S, T) = \sum P(c) \cdot E(c)$$

Topic: - Classification & Analysis

=> Classification algorithms are used to predict/ classify the discrete values.

Types of ML classification:

1) Logistic algorithm

2) K-Nearest Neighbours

3) Support vectors machines

4) Kernel SVM

5) Naive Bayes

6) Decision Tree Classification

7) Random Forest Classification

Regression:

-> Relation between dependent & independent variables.

-> It is all about predicting a quantity.

-> used to predict the continuous values such as

(price, salary, age etc).

Types of Regression Algorithm:

- Simple Linear Regression
- Multiple Linear Regression
- Polynomial regression
- Support vector Regression
- Decision Tree regression
- Random Forest Regression

Difference (Regression & Classification)

Regression	Classification
① Must be continuous	① Must be discrete value
② Map the input value (x) with the continuous output variable (y).	② Map the input (x) with discrete output (y)
③ used with cont. data	③ used with discrete data
④ Find the best fit line	④ Divide dataset into all different classes
⑤ Used to solve regression problem.	⑤ Classification problem solve.

⑥ Can be divided into Linear & Non-Linear regression.

⑥ can be divided into binary classifier & multiclass classifier

⑦ Weather prediction, House price etc.

⑦ Speech recognition, Cancer cells etc.

⑧ What is a support vector?

⇒ Support vectors are simply the co-ordinates of individual observation.

→ It is

⑨ How to find the SVM for case in hand?

⇒

⑩ SVM Kernel Technique:-

It is a function that takes low dimensional input space & transforms it to a higher dimensional ~~page~~ space.

mostly useful in non-linear separation problem.

Types of SVM:

2 types:

① Simple SVM: used for linear regression & classification problems.

② Kernel SVM: Can add more features to fit a hyperplane instead of a two-dimensional space.

Kernel Function:

- ↳ Linear
- ↳ Polynomial
- ↳ Gaussian Radial Basis Function
- ↳ Sigmoid.

SVM Kernel technique:
It is a function that takes low dimensional input space & transform it to a higher dimensional space.

Linear

! biorep (1)

→ Linear kernel works really well when there are a lot of features.

→ Are faster than others.

→ Have fewer parameters to optimize.

$$f(x) = w^T x + b$$

w = weight
 x = data

b = linear coefficient

Polynomial

→ Not used in practice.

→ Not computationally efficient.

→ predictions are not accurate.

$$f(x_1, x_2) = (a + x_1^T x_2)^b$$

Gaussian Radial Basis Function (RBF)

→ usually choice for non-linear data.

$$f(x_1, x_2) = \exp(-\gamma \|x_1 - x_2\|^2)$$

γ = How much a single training point has on the other data point.

$\|x_1 - x_2\|$ = dot product between feature

Sigmoid!

→ Used in Neural Network.

$$\rightarrow F(x,y) = \text{tanh}(\alpha x + y + c)$$

α = weight vector

c = An offset value

v.v.Im (for both ML & AI)

Math

Step 1: (Determine the decision column)

play golf (14)	
yes	9
no	5

Step 2: Calculating entropy for the classes (play golf)

we know that, Entropy for two variables we know

$$\text{Entropy} = - \sum P(n) \cdot \log_2(P(n))$$

so, Entropy (5, 9) = $-\left\{ \frac{5}{14} \cdot \log_2\left(\frac{5}{14}\right) + \frac{9}{14} \cdot \log_2\left(\frac{9}{14}\right) \right\}$

$$= - \left(\frac{5}{14} \cdot (-1.485) + \frac{9}{14} \cdot (-0.848) \right)$$

$$= 0.911$$

$$\text{Entropy (play golf, temperature)} = P(\text{Hot}) \cdot E(\text{Hot}) + P(\text{Mild}) \cdot E(\text{Mild}) + P(\text{Cool}) \cdot E(\text{Cool})$$

$$\text{Entropy (play golf, humidity)} = P(\text{High}) \cdot E(\text{High}) + P(\text{Normal}) \cdot E(\text{Normal})$$

E (Normal)

$$\text{Entropy (play golf, wind)} = P(\text{Strong}) \cdot E(\text{Strong}) + P(\text{Weak}) \cdot E(\text{Weak})$$

Step-3:- (Calculating entropy for other attribute after split)

① $E(\text{playgolf, outlook})$

② $E(\text{playgolf, temperature})$

③ $E(\text{playgolf, Humidity})$

④ $E(\text{playgolf, windy})$

Entropy for two variables we know

$$\text{Entropy}(S, T) = \sum P(e) \cdot E(e)$$

$$\text{Entropy} = \left\{ \left(\frac{20}{40} \right) \cdot \left(\frac{20}{20} \right) + \left(\frac{20}{40} \right) \cdot \left(\frac{20}{20} \right) \right\} = 1$$

✓ ① $E(\text{playgolf, outlook}) = P(\text{sunny}) \cdot E(\text{sunny}) +$

$P(\text{Rainy}) \cdot E(\text{Rainy}) + P(\text{Overcast}) \cdot E(\text{Overcast})$

② $E(\text{playgolf, temperature}) = P(\text{Hot}) \cdot E(\text{Hot}) + P(\text{mild}) \cdot E(\text{mild}) + P(\text{Cool}) \cdot E(\text{Cool})$

③ $E(\text{playgolf, Humidity}) = P(\text{High}) \cdot E(\text{High}) + P(\text{Normal}) \cdot E(\text{Normal})$

④ $E(\text{playgolf, windy}) = P(\text{True}) \cdot E(\text{True}) + P(\text{False}) \cdot E(\text{False})$

E(Playgolf outlook) calculation:

		Playgolf		Total
		yes	no	
Outlook	Sunny	3	2	5
	Rainy	2	3	5
	Overcast	4	0	4

$$E(\text{sunny}) = E(3, 2)$$

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

$$= - \{ 0.442 + 0.53 \}$$

$$E(\text{Rainy}) = E(2, 3)$$

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

$$= 0.972$$

$$E(\text{overcast}) = E(4, 0)$$

$$= - \left\{ \frac{4}{4} \cdot \log_2 \left(\frac{4}{4} \right) + \frac{0}{4} \cdot \log_2 \left(\frac{0}{4} \right) \right\}$$

$$= (0 \cdot 0 + 0 \cdot 0) = 0$$

$$\therefore E(\text{Play golf, Outlook}) = \frac{5}{14} \cdot E(3, 2) + \frac{4}{14} \cdot E(1, 0) + \frac{5}{14} \cdot E(2, 3)$$

$$= \frac{5}{14} \times 0.972 + \frac{4}{14} \times 0 + \frac{5}{14} \times 0.972$$

$$= 0.35 + 0 + 0.35$$

$$E(\text{Play}) = 0.7$$

Now, $P(\text{Play golf, Temperature})$ Calculation :-

	Playgolf		Total
	Yes	No	
Mild	4	2	6
Cool	3	1	4
Hot	2	2	4

$$E(\text{Mild}) = E(4, 2)$$

$$= -\left(\frac{4}{6} \log_2 \left(\frac{4}{6}\right) + \frac{2}{6} \log_2 \left(\frac{2}{6}\right)\right)$$

$$= - (0.39 + 0.53)$$

$$= 0.92$$

$$E(\text{Cold}) = E(3,1)$$

$$= \frac{3}{4} \cdot \log_2\left(\frac{3}{4}\right) + \frac{1}{4} \cdot \log_2\left(\frac{1}{4}\right)$$

$$= -(0.31 + 0.5)$$

$$= 0.81$$

$$E(\text{Hot}) = E(2,2)$$

$$= 2 \cdot \frac{1}{4} \cdot \log_2\left(\frac{2}{4}\right) = 2 \cdot \frac{1}{4} \cdot \log_2\left(\frac{1}{2}\right)$$

$$= - (0.5 + 0.5) = -1$$

$$= 1$$

$$\therefore E(\text{play golf, Temperature}) = E(\text{Mild}) \cdot P(\text{Mild}) + E(\text{Cold}) \cdot$$

$$P(\text{Cold}) + E(\text{Hot}) \cdot P(\text{Hot})$$

$$= \frac{6}{14} \cdot 0.92 + \frac{4}{14} \cdot 0.81 + \frac{4}{14} \cdot 1$$

$$= 0.231 + 0.286$$

$$= 0.517$$

$$= 0.907$$

$E(\text{Playgolf, Humidity}) :-$

	Playgolf		Total
	Yes	No	
High	3	4	7
Normal	6	1	7

$$\therefore E(\text{High}) = \frac{3}{7} \cdot \log_2\left(\frac{3}{7}\right) + \frac{4}{7} \cdot \log_2\left(\frac{4}{7}\right)$$

$$= - (0.524 + 0.461) = 0.985$$

$$\therefore E(\text{Normal}) = E(6, 1)$$

$$= \left(\frac{6}{7}\right) \cdot \log_2\left(\frac{6}{7}\right) + \left(\frac{1}{7}\right) \cdot \log_2\left(\frac{1}{7}\right)$$

$$= - (0.191 + 0.401)$$

$$= 0.592$$

$$E(\text{Playgolf, Humidity}) = \frac{7}{14} \cdot 0.985 + \frac{7}{14} \cdot 0.592$$

$$= 0.4925 + 0.296$$

$$= 0.7885$$

R·E (playgolf, windy) :-

	playgolf		Total
	Yes	No	
True	3 6	3 2	6
False	6	2	8

① $E(\text{True}) = E(3, 3)$

$$= \frac{3}{6} \cdot \log_2\left(\frac{3}{6}\right) + \frac{3}{6} \cdot \log_2\left(\frac{3}{6}\right)$$

$$= 0.5 + 0.5 = 1$$

② $E(\text{False}) = E(6, 2)$

$$= \frac{6}{8} \cdot \log_2\left(\frac{6}{8}\right) + \frac{2}{8} \cdot \log_2\left(\frac{2}{8}\right)$$

$$= 0.811 + 0.5$$

$$= 0.811$$

∴ $E(\text{playgolf, windy}) = \frac{6}{14} \times 1 + \frac{8}{14} \times 0.811$

$$= 0.43 + 0.463$$

$$= 0.893$$

Step-4 (Calculating Info gain)

① $\text{Gain}(\text{playgolf}, \text{outlook}) = \text{Entropy}(\text{playgolf}) -$

$$= 0.94 - 0.7$$

$$= 0.24$$

② $\text{Gain}(\text{playgolf}, \text{Temp}) = \text{Entropy}(\text{playgolf}) -$

$$= 0.94 - 0.907$$

$$= 0.033$$

③ $\text{Gain}(\text{playgolf}, \text{Humidity}) = \text{Entropy}(\text{playgolf}) -$

$$= 0.94 - 0.79$$

$$= 0.15$$

④ $\text{Gain}(\text{playgolf}, \text{windy}) =$

$$= \text{Entropy}(\text{playgolf}) - \text{Entropy}(\text{playgolf}, \text{windy})$$

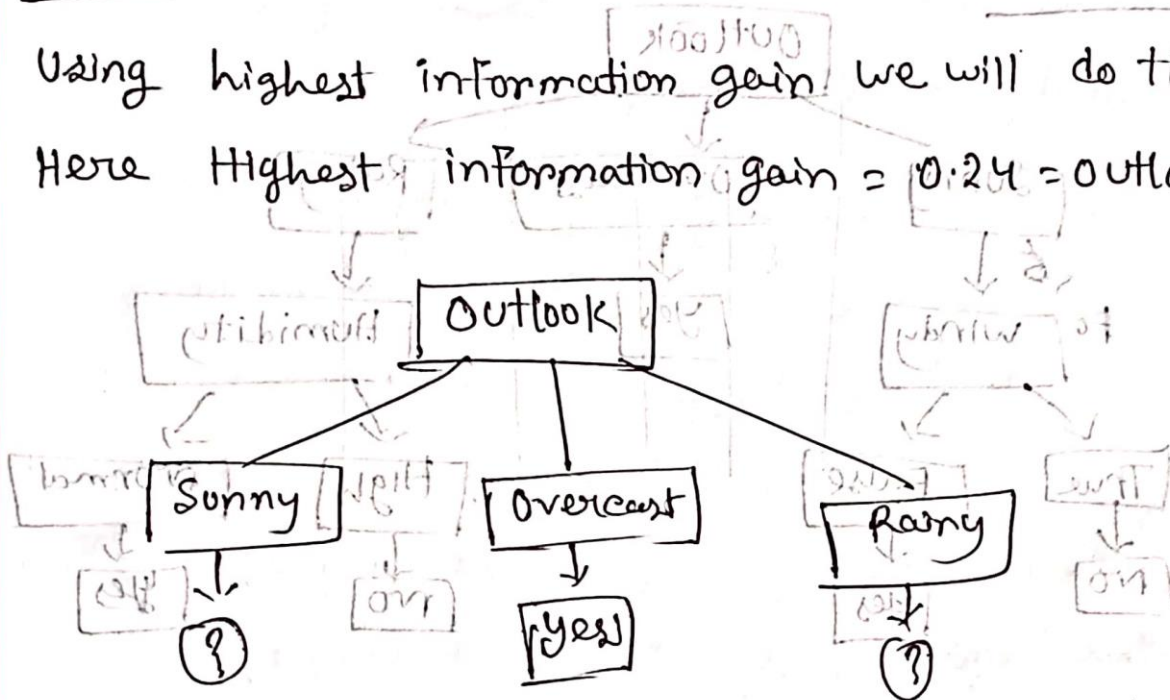
$$= 0.94 - 0.893$$

$$= 0.047$$

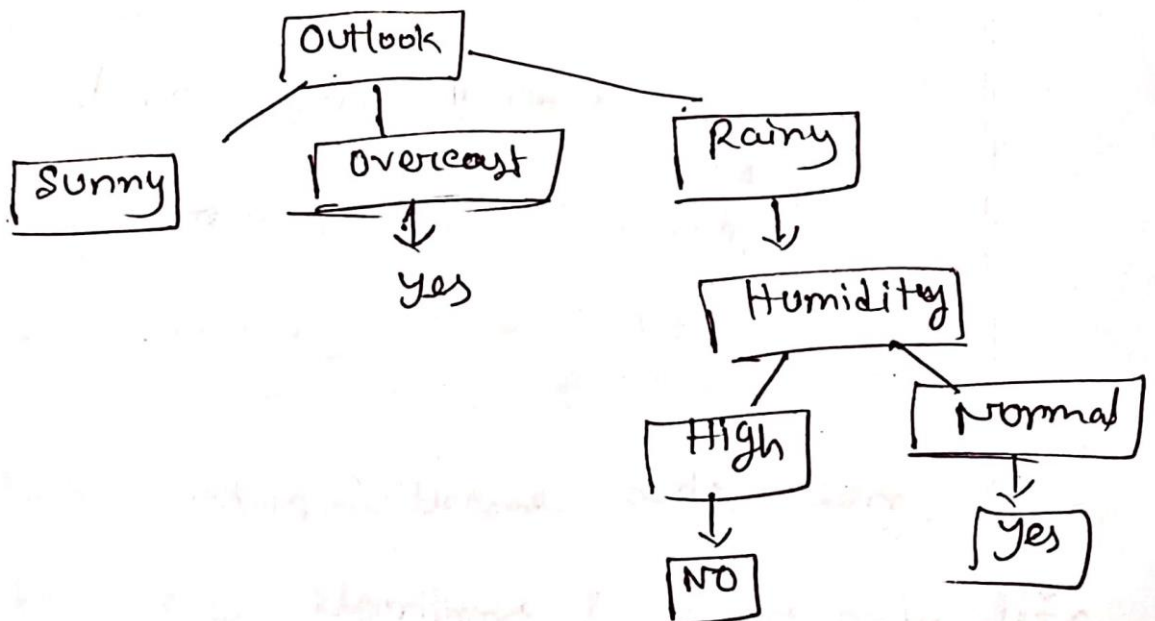
Step-5:- (perform the first split)

Using highest information gain we will do that.

Here Highest information gain = 0.24 = outlook



Step-6:-



Step-7:-

