

22/02/2022

Machine Learning and Deep Learning

Unit - I

Introduction to Machine Learning

History of ML:

- 1) 1950 - Introduction of Turing Test by Alan Turing
Q → "Can Machines Think?"
- 2) 1952 - Samuel introduced first ML Program ie. The checkers game
- 3) 1957 - Frank Rosenblatt invented first neural network
- 4) 1967 - Pattern Recognition (KNN algorithm introduced)

25/02/2022

- 5) 1979 - Self Driving Cars → by Stanford university students.
- 6) 1982 - RNN - Recurrent Neural Network.
- 7) 1989 - Reinforcement Learning
- 8) 1995 - Random Forest and SVM (Support Vector Machine)
(accuracy of model > 95%)
- 9) 1997 - IBM deep blue (chess game)
- 10) 2006 - Netflix (Recommendation System)
- 11) 2010 - Kaggle Website
- 12) 2011 - IBM Watson (Jeopardy)
- 13) 2016 - Google Alpha Go

26/02/2022

Human Learning:

In a cognitive science, learning means gaining the information from observation.

There are 3 types of Learning. They are.

1) Learning under expert guidance

2) Learning by self knowledge gained from experts

3) Learning guided by knowledge

Machine learning:

It is a computer program which is said to learn from the experience E with respect to some class of task T and performance measure P, if its performance at task T as measured by performance P and improves with experience E.

Well Posed Problem:

Play checkers game

E → represents the experience to play the game.

T → represents the task of playing checkers.

P → Performance measure indicated by percentage of game won by the player.

1) Play checkers game
E → represents the experience to play the game.
T → represents the task of playing checkers.
P → Performance measure indicated by percentage of game won by the player.

2) Image Classification

E → taking the labelled data as input

T → to classify the image

P → to classify image based on percentage

(a)

E → Past data with images having labels

T → Task assigning class to new unlabelled data

P → The performance measure indicated by the percentage of image is correctly classified

How do Machines Learn?

To learn a machine, we need to take 3 steps with

1) Input Data:

Past data / information is utilized as a basis for future decision making.

a) Abstraction

The input data is represented in a broader way through the underlined environment.

b) Generalization

The abstracted representation is generalized to form a framework for making decision.

Process of ML:



Machine learning types:

i) Supervised Learning:

Supervised learning is an equation i.e. mathematical equation to solve a problem.

It consists of input and output data and a supervisor.

Input (Independent Variable)

Output (Dependent Variable)

Example: $y = f(x)$

Input: $\square, \triangle, \text{pentagon}$

Output: $\square \rightarrow \text{rectangle}$

$\triangle \rightarrow \text{triangle}$

$\text{pentagon} \rightarrow \text{pentagon}$

...

2) Unsupervised Learning:

Unsupervised learning is an equation i.e. mathematical equation to solve a problem.

It consists of input and output data and no supervisor.

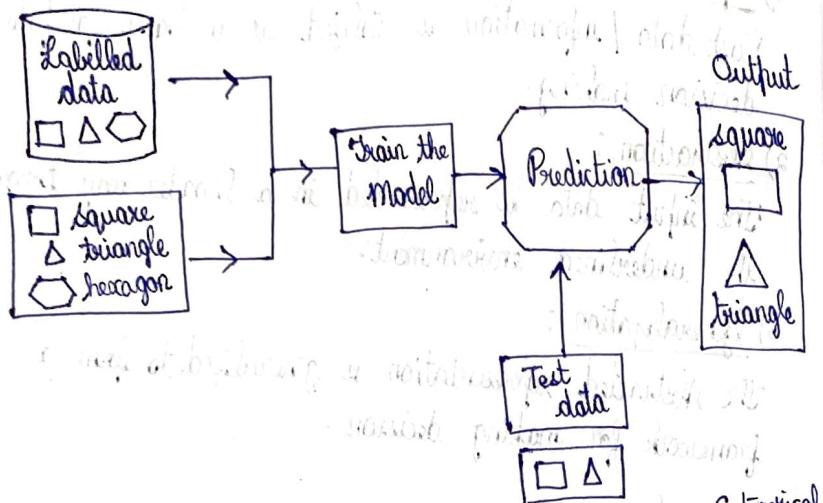
Input: $\square, \triangle, \text{pentagon}$

Output: $\square \rightarrow \text{rectangle}$

$\triangle \rightarrow \text{triangle}$

$\text{pentagon} \rightarrow \text{pentagon}$

...



Supervised learning → tumor detection (Yes/No) Categorical Value Data
 Sales of a product

Categorical Data → We use classification algorithms
 Continuous Data → We use Regression

- 1) Classification and 2) Regression
- * Supervised learning is a process of providing input data as well as correct output data to the machine learning model.
- * The aim of supervised learning is to find a mapping function to map the input variable x with the output variable y .
- * The input variable x is called as Independent Variable and output variable y is called as Dependent Variable.

Examples:

- 1) Predicting result of a game.
- 2) Sales Prediction.

3) Classifying set of emails as spam or not spam.

There are 2 different types of supervised learning:

1) Classification (and all make mistakes sometimes)

2) Regression.

Classification:

If we are trying to predict categorical and nominal variables.

Regression:

When we are trying to predict continuous or real value variables.

Classification Problems / Examples:

- 1) Image Classification
- 2) Prediction of a disease
- 3) WIN prediction of a game
- 4) Handwriting recognition

Regression is a statistical method to model the relationship between dependent variable (target variable) and independent variables (predictor).

It may be multi variable

Regression Problems / Examples:

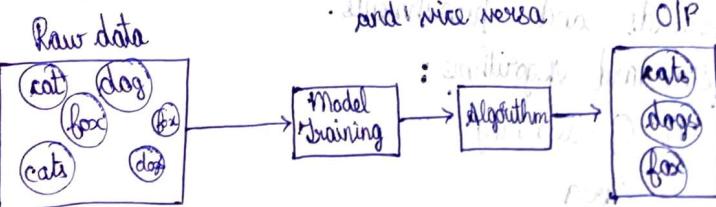
- 1) Sales / Price Prediction
- 2) Weather Forecasting

2) Unsupervised Learning:

Similarity Measures → Distance b/w clusters is more

→ less similar

and vice versa



- 1) Clustering → clustering similar data into clusters.
- 2) Association → Market Basket Analysis
- Common relation between the different input data.

Applications involved in Unsupervised Learning:

- * Unsupervised Learning contains 2 types
 - 1) Clustering
 - 2) Association

Clustering is a method of grouping the objects into clusters such that objects with most similarities remains into the group and has lesser distance (or) no similarities with objects and categorize them as per the presence and absence of those commonalities.

Association is used to find the relationship between variables in large database. It determines the set of items that occurs together in the dataset. Association rule makes marketing strategy more effective.

After finding Association Rule, we can use regression analysis to predict the sales.

Reinforcement Learning is based on rewards and punishments. Based on environment, agent must act. The model learns and updates itself based on the rewards and punishments.

Standard Algorithms:

- 1) Q-learning
- 2) Sarsa

Applications of Unsupervised learning:

- 1) Recommendation Systems
- 2) Anomaly Detection
- 3) Association Rule Mining
- 4) Similarity Detection
- 5) Product and Customer Segmentation.

Problems not to be solved by using ML:

- 1) Reasoning Power
- 2) Contextual Limitations
- 3) Scalability
- 4) Internal Working of Deep learning.

Applications of ML:

- 1) Banking Sector
- 2) Health Care
- 3) Fraud Detection
- 4) Self driving cars.
- 5) Image Recognition
- 6) Traffic Prediction
- 7) Sales Trading
- 8) Product Recommendations
- 9) Speech Recognition.

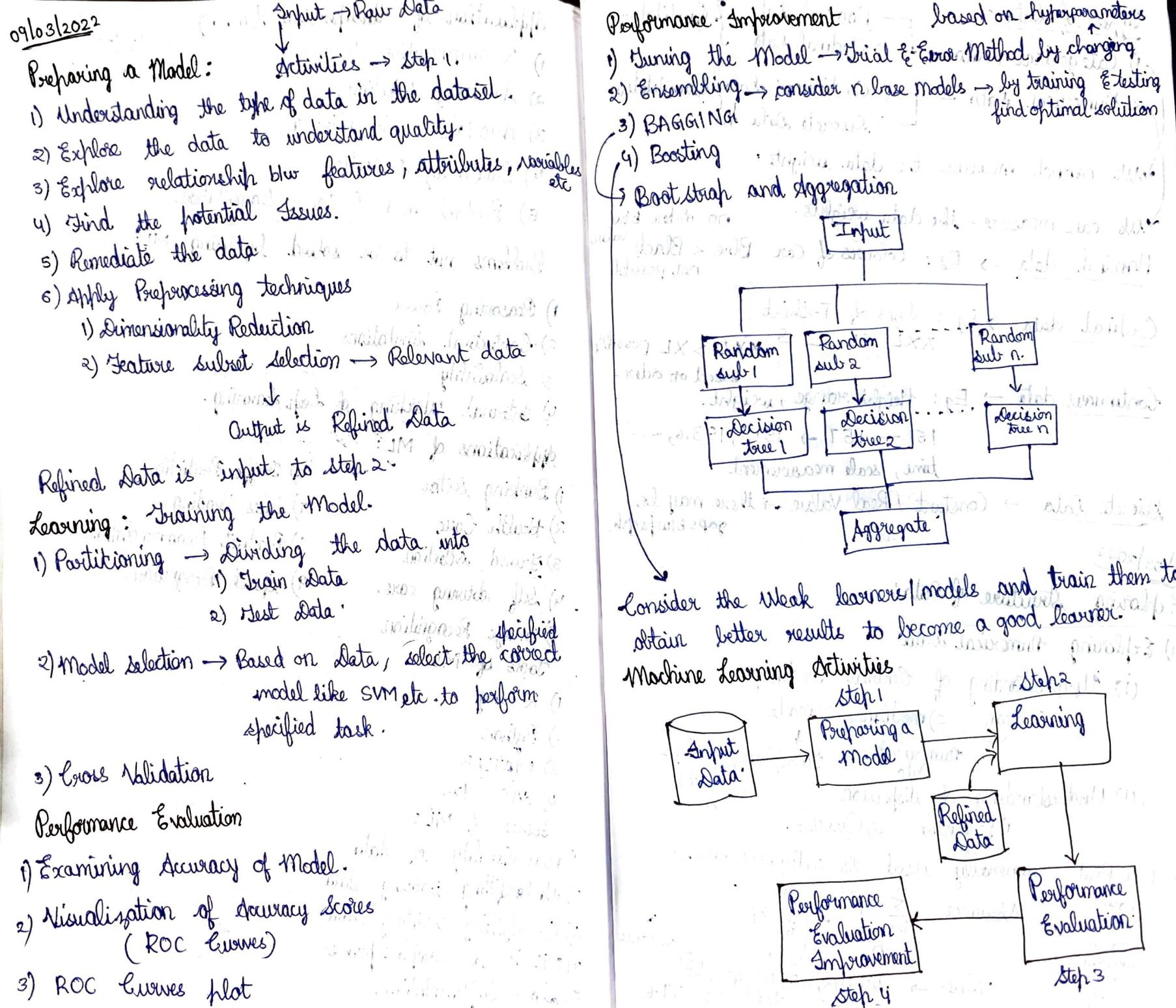
Tools of ML:

- 1) R lang
- 2) Python
- 3) MATLAB
- 4) SAS

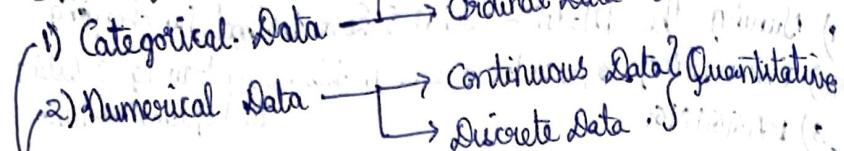
Issues of ML:

- 1) Poor Quality of data
- 2) Underfitting training data
- 3) Overfitting training data
- 4) ML is a complex process
- 5) Slow Implementation.

05/03/2022



Types of Data:



- We cannot measure the data weights.
 - We can measure the data weights.
- Nominal data → Eg: Colours of car Blue < Black values not possible

Ordinal data → Eg: Sizes of T-shirt
 $XXL < XL < L < S$ $XXL > XL$ possible based on order.

Continuous data → Eg: Height range, weight.
 $151 - 157 \rightarrow 151.1, 153.6, \dots$
 time, scale measurement

Discrete data → Constant / Real Value → there may be gaps b/w people.

11/03/2022

Exploring Structure of Data:

1) Exploring Numerical Data

(i) Understanding of Central Tendency
 1) mean 2) median 3) mode

Numerical Data To replace outliers

(ii) Understanding of dispersion
 1) Percentile 2) Quartile.

→ Box Plot generally used to categorize outliers

count
mean
median
variance
min
as%
50%
75%
max

$$\text{Variance} = \frac{\sum_{i=1}^n x_i^2}{n} - \left(\frac{\sum_{i=1}^n x_i}{n} \right)^2$$

Values → 44, 46, 48, 45, 47. ②
 34, 46, 59, 53, 52 79.6

1) Central Tendency

- mean
- median
- mode

25%, 50%, 75%

Q₁ Q₂ Q₃

25%, 25%, 25%, 25%

2) Data Dispersion

- Variance
- Standard Deviation
- Percentile
- Quartile

IQR - Inter Quartile Relationship

(ii) Plotting Numerical Data

use of Boxplot, histogram

Exploring Relation between variables

- 1) Scatter plot
- 2) Two-Way-Cross Tabulation x, y.

Data Quality and Remediation:

Best Quality → Standard Datasets.

Issues

- 1) Incorrect Sample Set Collection
- 2) Errors in Data Collection

Data Remedies:

- 1) Handling Missing Values
- 2) Handling Outliers

Remove Outliers

- 1) Removing Outliers which are unnecessary.
- 2) Imputation → Assigning similar frequent values
- 3) Clipping ~5% ↔ 95% ~

2) Handling Missing Values

- 1) Eliminating records which are having missing values
- 2) Imputation of missing values
- 3) Estimating Missing Values.

Dimensionality Reduction:

It is a technique in which the conversion of higher dimension dataset to lower dimension dataset.

Prediction:

It is widely used to obtain a better fit predictive model while solving classification & regression values.

Feature Subset Selection:

It is a process to select the subset of relevant features and leaving out irrelevant features.

19/03/22

Unit - II

Selecting a Model:

- 1) Input Data → Raw data
- 2) Abstraction → Training data
- 3) Generalization → Testing data

$$y = f(x) + \epsilon$$

↓ ↓

Output Target independent
dependent function variable.

Target functions used:

- 1) Cost function → error function - tells us how bad our model (Only applicable to training dataset) is going to function.
- 2) Loss function → applicable on a particular datapoint
- 3) Objective function → We consider data and model

To select a model, we have 2 parts

- 1) Predictive 2) Descriptive → unsupervised learning

K Means Clustering
Market Based Analysis
Recommendation System
Pattern Recognition

Cost function:

It is also called as an error function which helps to measure the extent to which the model is going wrong in estimating the relationship between x and y . In that sense, cost function can tell how bad the model is performing. R^2 regression function can be used.

Loss function: It is similar to the cost function. The only difference is the loss function is defined on a data point while cost function is for entire training dataset.

Objective function:
 It is used to find most suitable solution to a point.
 We need to evaluate the quality and optimality of system.
 We can also call it goal function.

ML algorithms broadly divided into 2 types

- Supervised learning - primarily focussing on solving predictive problems.

- Unsupervised learning - primarily focussed on descriptive problem.

Predictive Models - They try to predict certain value using the values in an input dataset. It has a clear focus on what they want to learn and how they want to learn.

Examples: Predicting whether a transaction is fraud.
 a customer may move to another product
 win loss % of a cricket match.

- * The models which are used for prediction of target features of categorical values are called classification models.
- * The target feature is known as a class and the categories to which classes are divided into are called as levels.

Some of the popular classification models.

KNN, Naive Bayes, decision trees

- * The models which are used for prediction of numerical values of the target features of a data instance are called as regression models.

Examples:

- Prediction of revenue growth
- Prediction of rainfall amount
- Predicting sales of a product in coming years.

Descriptive Models:
 (Predictive) Descriptive models which group together similar data instances having a similar value of the different features are called clustering models.

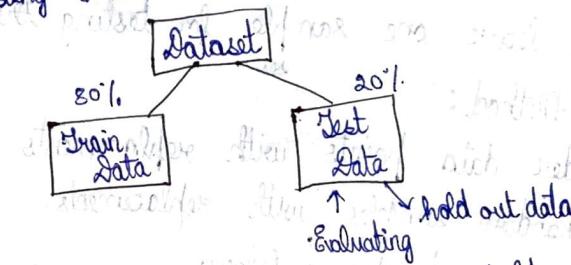
Examples: Customer grouping / segmentation based on social, demographic, ethnic factors.

The most popular model for clustering is K-Means.

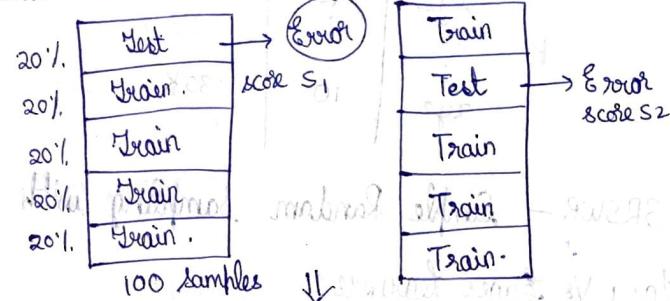
23/03/2022

Training of Model (for Supervised Learning):

- Hold Out Method:
 Used to divide training set into training data and testing data. So we hold out some part of training set i.e. testing data.



- K-fold Cross Validation: Extension for hold out method.



avg. $K=5$
 The aggregate gives us the best accuracy / less error rate

20%	Train	Train	Train	Train	Test	→ Error S ₃
20%	Train	Train	Train	Train	Test	→ Error S ₂
20%	Train	Train	Train	Train	Test	→ Error S ₁
20%	Train	Train	Train	Train	Test	→ Error S ₅
20%	Train	Train	Train	Train	Test	→ Error S ₄

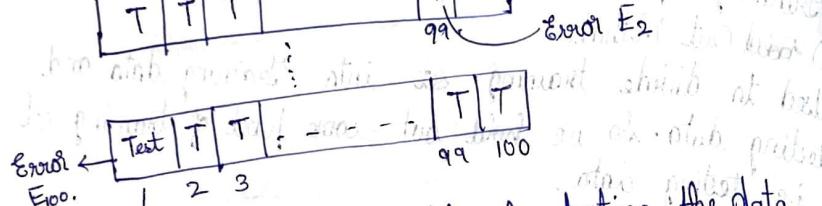
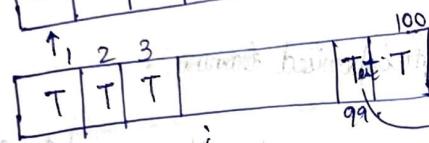
10 fold cross validation:

Here K value is 10.

leave one out cross validation:

Here $K = N = \text{Number of inputs}$

$$K = \frac{N}{100}$$



always we leave one sample for testing the data!

Bootstrap Method:

We consider data points with replacements

Random samples with replacements.

Name	Age	Salary
xyz	10	85K
xyz	10	65K
xyz	10	35K

SRSWR - Simple Random Sampling with Replacement.

Lazy Vs Eager Learners:

K-Nearest Neighbors Algorithm

Hold Out Method:

- The data is partitioned into training and testing data which is by holding back a part of the input data from validating the train model is called Hold Out Method.
- There are different hold out methods are used to improve the ML model by avoiding overfitting and underfitting of model.

K-fold Cross Validation:

- Data is divided into testing and training set.
- ML model is developed using the partition of data and then tested on the rest data.
- This process is repeated K-times with different random partitions to generate an average performance.

Bootstrap Method:

- It is a random sample conducted with replacements.
- While selecting the sample, the same can appear more than once.
- We can create more than one training data set from original dataset.
- It follows the technique SRSWR.
- s → simple, R → Random, S → Sampling with Replacement.

Lazy Vs Eager Learners:

- Lazy Learner:**
 - It simply stores training data that means memorizing data and wait to get test data.
 - The training time is less and predicted time is less.
- Example: KNN

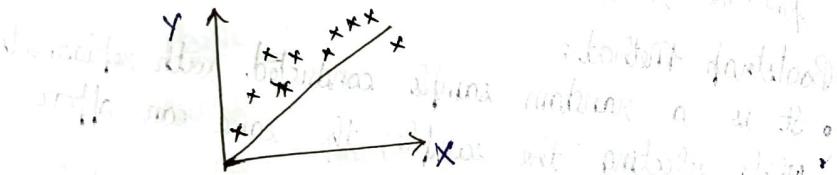
Eager Learners:

- It takes more time for training data
 - While new input is received it generalizes and gives the quick results.
- Example: SVM, Neural Networks, Decision Tree.

24/03/22

Model Representation:

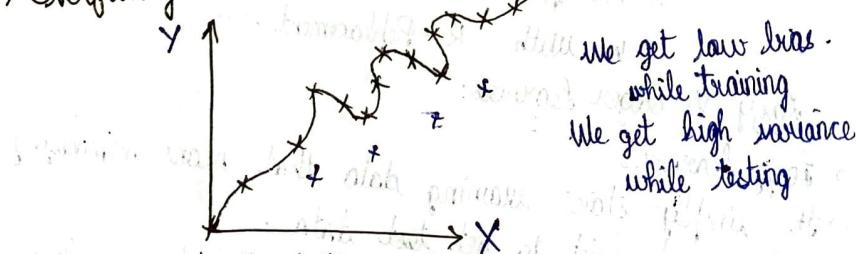
- 1) Underfitting: for low amount of data.
Model should be provided with refined input data.
- Bias - Average value b/w actual value and predicted value.
- Bias is high i.e. more erroneous, less variance.
- Low Variance and High Bias → indicate underfitting



To avoid underfitting

- 1) Consider more amount of data (i.e. train data).
- 2) Removing effective variables.

2) Overfitting: low Bias, high variance.



To avoid Overfitting

- 1) Apply K-fold cross Validation
- 2) Remove unnecessary data from data set

→ Best fit Model - Low bias, low variance.

Underfitting: It occurs when the model is not able to capture the underlying trend of data i.e. the model may not learn from the training dataset.

Overfitting: It occurs when the model is fully able to capture the underlying trend of data but not with test data.

Bias: It is an error used to find out the average difference between predicted and actual value.

Variance: It is a measure of how far the set of data points are spread out from their mean value. It means to find the difference of deviation from the actual value.

$$\text{Variance} = \frac{\sum_{i=1}^n (x_i - \bar{x})^2}{N-1}$$

\bar{x} - mean of all DPs
N - No. of DPs.

- Underfitting can be avoided by giving more training data and reducing the effective features.
- Overfitting is avoided by K-fold cross validation, Hold back of validation dataset, Removing the nodes with no predictive use.

Evaluating the Performance of a Model:

Classification Metrics

Confusion Matrix:

$n = \text{Total no. of Data points}$	Actual; No	Actual; Yes	Precision
Predicted - No	TN	FP	Recall
Predicted - Yes	FN	TP	

Predicted

	<u>Actual</u>	
Match would be win (T)	Match Won (P) \rightarrow (TP)	
Match win (F)	Match lost (P) \rightarrow (FP)	
Match lost (F)	Match win (N) \rightarrow (FN)	
Match lost (T)	Match lost (N) \rightarrow (TN)	

Accuracy: $\frac{TP + TN}{TP + FP + TN + FN}$

$$TN = 9, \quad FP = 4, \quad FN = 2, \quad TP = 85$$

$$A = \frac{94}{100} = 0.94 \Rightarrow 94\% \text{ accuracy}$$

Misclassification Rate / Error Rate:

$$\text{Error Rate} = \frac{FP + FN}{TP + TN + FP + FN}$$

- Confusion Matrix is a performance measurement.
- It gives the summary on prediction of classification.
- It is an $n \times n$ matrix used to compare actual and predicted values.

Type I Error \rightarrow FP

Type II Error \rightarrow FN

→ Accuracy defines how often the model predicts correct output.

→ Misclassification error is an error rate defined how often the model gives wrong prediction.

Precision:

It defines number of correct outputs provided by the model.

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

Recall:

It defines out of total +ve classes how our model predicted correctly. It must be as high as possible for best performance of model.

$$\text{Recall} = \frac{TP}{TP + FN} \quad (\text{Sensitivity})$$

Specificity:

It measures the proportion of negative examples which have been correctly classified.

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

F-Measure:

$$F\text{-Measure} = \frac{2 * \text{Recall} * \text{Precision}}{\text{Recall} + \text{Precision}}$$

ROC Curve:

ROC stands for Receiver Operating Characteristic curve.
In this we have True Positive Rate (TPR) = $\frac{TP}{TP + FN}$

$$\text{False Positive Rate (FPR)} = \frac{FP}{FP + TN}$$

30/03/2022

Performance Measure for Regression:

- MAE - Mean Absolute Error

Average differences b/w Predicted and (Average) values.

$$\text{MAE} = \frac{1}{n} \sum_{i=1}^n |y_i - \hat{y}_i|$$

↓
Predicted value.
↓
Actual Data Point.

It is one of the best method/metric for regression.
It is the average of difference between original value and predicted value.

② Root mean square error (RMSE):

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

Eg: $y_i = (10, 5)$ $\hat{y}_i = (8, 6)$

 $MAE = \frac{1}{2} [|10-8| + |5-6|] = \frac{3}{2} = 1.5$
 $RMSE = \sqrt{\frac{1}{2} [(10-8)^2 + (5-6)^2]} = \sqrt{\frac{1}{2} (4+1)} = \sqrt{\frac{5}{2}} = 1.582$

③ R^2 Error:

Residual \rightarrow Error

$$R^2 \text{ Error} = \frac{SST - SSE}{SST}$$

$$= 1 - \frac{SSE}{SST} \rightarrow \text{Sum of square error}$$

$$\frac{SSE}{SST} \rightarrow \text{Sum of square Total}$$

$$SSE = \sum (Y_{\text{actual}} - Y_{\text{pred}})^2$$

$$SST = \sum (Y_{\text{actual}} - Y_{\text{mean}})^2$$

$$R^2 \text{ Error} = 1 - \frac{\sum (Y_{\text{actual}} - Y_{\text{pred}})^2}{\sum (Y_{\text{actual}} - Y_{\text{mean}})^2}$$

Range of R^2 error is (0, to 1)

④ Adjusted R^2

$$\text{Adjusted } R^2 = 1 - \frac{(1-R^2) \cdot (N-1)}{(N-P-1)}$$

↓
f² value
↓
Total no. of value
↓
Predicted value.

Unsupervised learning metrics:

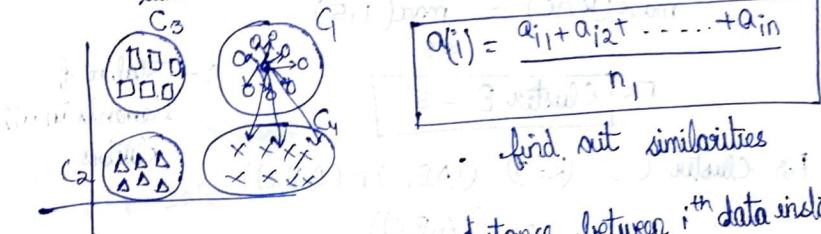
i) Clustering:

Internal Evaluation \rightarrow Euclidean, Manhattan Distances

$$\text{ii) Silhouette Coefficient} = \frac{b(i) - a(i)}{\max\{a(i), b(i)\}}$$

Internal evaluation methods generally measures the cluster quality. Based on homogeneity of data, belongs to same cluster and heterogeneity of data, belong to different clusters.

a(i) - average distance between ith datapoint and all other data instances belonging to same clusters.



$$a(i) = \frac{a_{i1} + a_{i2} + \dots + a_{in}}{n}$$

find out similarities

b(i) - the lowest average distance between ith data instance and data instances of all other clusters.

$$b_{14}(\text{average}) = \frac{b_{14}(1) + b_{14}(2) + \dots + b_{14}(n)}{n_4}$$

cluster 1 cluster 4

$$b_{12}(\text{average}) = \frac{b_{12}(1) + b_{12}(2) + \dots + b_{12}(n)}{n_2}$$

$$b_{13}(\text{average}) = \frac{b_{13}(1) + b_{13}(2) + \dots + b_{13}(n)}{n_3}$$

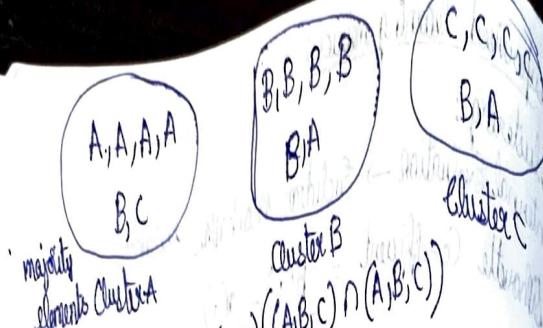
$$\star b(i) = \min \{ b_{12}(\text{average}), b_{13}(\text{average}), b_{14}(\text{average}) \}$$

find out dissimilarities.

External Evaluation:

$$(1) \text{ Purity} = \frac{1}{n} \sum_k \max_{l \in L} (K \cap S_l) \rightarrow \text{Target class label name.}$$

↓
Total No. of classes present
in the cluster



$$\text{For Cluster A} = (K=3) ((A|B|C) \cap (A|B|C))$$

$$= (A|B|C).$$

$$\max(K \cap C) = \max(4|1|1) = 4$$

Cluster A = 4

$$\text{For Cluster B} = (K=2) ((A|B) \cap (A|B))$$

$$= (A|B)$$

$$\max(K \cap C) = \max(1|5) = 5$$

Cluster B = 5

$$\text{For Cluster C} = (K=3) ((A|B|C) \cap (A|B|C))$$

$$= (A|B|C)$$

$$\max(K \cap C) = \max(1|4) = 4.$$

Cluster C = 4

$$\text{Purity} = \frac{1}{18} (4+5+4) = \frac{13}{18} = 0.72$$

We assign a label to each cluster based on most frequent class in it then the purity becomes the number of correctly matched class and cluster labels divided by number of total data points.

$$\text{Purity} = \frac{\text{Cluster}(A+B+C)}{\text{Total}}$$

The Purity value increases when the number of clusters increases.

$n \rightarrow$ Total number of data instances
 $K \rightarrow$ Number of distinct classes in a cluster
 $C \rightarrow$ Class / Cluster label.

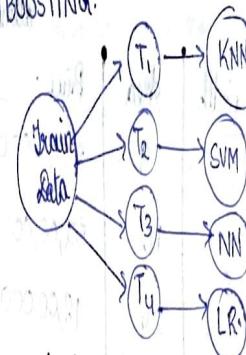
Improving the Performance of a Model:

i) Model Parameter Tuning

ii) Ensembling

(i) BAGGING

(ii) BOOSTING



06/04/2022, 07/04/2022

Basics of Feature Engineering:

→ To find relevant attributes in a dataset

→ It is a preprocessing step → Construction

i) Feature transformation → Extraction

ii) Feature subset selection.

Feature Engineering is a process of translating the data into features such that these features are able to represent the dataset more efficiently.

It is an important preprocessing step which improves the performance of a model.

It consists of 2 variants

i) Feature Transformation

ii) Feature subset selection

Feature Transformation:

Transforming the data (structured/unstructured) into new set of features which represents the underlying problem which ML is trying to solve.

(i) Feature Construction:

	apt-length	apt-width	Price
3	3	3	23,00,000
2	5	5	52,00,000
6	4	4	18,00,000
3	1	1	30,00,000

↓ Construction of length × width = area attribute

	apt-length	apt-width	area	Price
3	3	3	9	23,00,000
2	5	5	10	52,00,000
6	4	4	24	18,00,000
3	1	1	3	30,00,000

Encoding Categorical Variables (Ordinal)

Eg: Predicting the win loss for an athlete

Age	City of Origin	Parent Athlete	Chance of Win
18	City A	yes	yes
20	City B	no	yes
23	City C.	yes	no
19	City B	yes	yes

Age	City A		City B		City C.		
	yes	No	yes	No	yes	No	
18	1	0	0	1	0	1	0
20	0	1	0	0	1	1	0
23	0	0	1	1	0	0	1
19	0	1	0	0	1	0	0

Since dimensions increases, space usage increases so we drop irrelevant features.

Encoding Categorical (Ordinal) files :

Low value → high priority
 $A=1, B=2, C=3, D=4$

Maths	Science	Grade	Maths	Science	Grade
78	75	B	78	75	2
56	62	C	56	62	3
89	90	A	89	90	1
91	95	A	91	95	1
66	78	D	66	78	4

Converting Numerical (Continuous) Features to Categorical Features

For Eg: Prediction of House Pricing Category

apt-area	apt-price	apt-price	Ordinal
4720	23,60,000	Medium	2
2430	12,15,000	Low	1
4368	21,84,000	Medium	2
3969	19,84,500	Low	1
6142	30,71,000	Medium	2
7912	39,56,000	High	3

Range : $x < 20L \rightarrow \text{Low}$

$20L < x < 35L \rightarrow \text{Medium}$

$x > 35L \rightarrow \text{High}$

If we assume Low = 1, Medium = 2, High = 3. Then
 Text → Tokenization → Variable Count → Normalization
 file with
 text

12/04/2022

Feature Extraction:

Principal Component Analysis:

- 1) Standardization
- 2) Covariance Matrix
- 3) Eigen Vectors & eigen values
- 4) Principal Components
- 5) Reducing dimensions

i) Standardization:

$$z = \frac{\text{value} - \text{mean}}{\text{SD}} = \frac{x - \mu}{\sigma}$$

Covariance Matrix:

$$\text{Cov}(a, b) = \frac{1}{n-1} \sum_{i=1}^n (a_i - \bar{a})(b_i - \bar{b})$$

$$\text{Covariance Matrix} = \begin{bmatrix} \text{Cov}(a, a) & \text{Cov}(a, b) \\ \text{Cov}(b, a) & \text{Cov}(b, b) \end{bmatrix}$$

Eigen Values and Eigen Vectors:

$$\begin{bmatrix} 0.7351 \\ 0.6778 \end{bmatrix} \begin{bmatrix} 0.6778 \\ 0.7351 \end{bmatrix}$$

Principal Components:

Singular Value Decomposition (SVD):

$$A = \begin{bmatrix} \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot \end{bmatrix}_{m \times n} = U \Sigma V^T$$

U, V are orthogonal & unitary
 Σ is a singular matrix.

U columns \rightarrow left singular vectors of A
 V columns \rightarrow right singular vectors of A .
 Σ — values of A .

$$A = \begin{bmatrix} 3 & 1 & 1 \\ -1 & 3 & 1 \end{bmatrix}$$

$U \rightarrow$ 1) Eigen vectors of $A A^T$
 2) Orthogonalize (divide with distance)

$V \rightarrow$ 2) $A^T A$.

13/01/22

Feature Subset Selection:

Issues in huge data

- Performance of model decreases.

A	B	C	D

F	O

Advantages:

- faster
- efficient

Key Drivers of feature selection:

feature relevance

supervised learning

Predictor Variables

- irrelevant
- weakly redundant

In Unsupervised learning \rightarrow no labelled data

Test data \rightarrow grouping into clusters \Rightarrow based on irrelevant feature elimination ie relevant features.

\rightarrow potentially redundant

Measures to be taken for feature relevance

i) Supervised learning:

ii) Mutual information:

$$MI(C, F) = H(C) + H(F) - H(C, F)$$

$$= - \sum_{i=1}^t P(i) \log_2 (P(F-i))$$

i. Overall feature selection Process:

Steps Involved:

* Feature generation

Take an empty set, and insert the attributes one by one in forward direction

$$n \rightarrow 2^n$$

i) Forward Selection

2) Backward Selection

3) Bidirectional