

Theory for Viva

Practical 1:-

1. Data wrangling?

→ Process of cleaning, transforming & preparing raw data in structure suitable for analysis.

Big data → large collection of dataset & cannot process using traditional techniques.

Numpy, Pandas

Matplotlib, Seaborn,

Scikit learn

Label encoding:-

It converts the label in numeric form as most of library are compatible with the numerical data.

1) Numpy:- It is the fundamental library used in python. you can perform basic mathematical operation and also operation of array [add, multiply and basic operation like sorting & indexing]
 slice, index

2) Pandas:-

Pandas is an open source library providing high performance easy-to-use data structure and analysis tools.

It can used for manipulating the data frame.

- Update, delete, add column to data frame
- Handling & Identifying missing values
- Plotting.

3) Scikit learn = open source library featuring algorithms like random forest, k-means clustering, SVM etc.

It is use for classification: spam detect

Regression (Linear & Logistic)
model selection (for giving accuracy & score).

Q) matplotlib lib :- matplotlib lib is an Data visualization library use to plot different visualization technique such as charts, graphs, maps & plots.

Q) Seaborn :- Seaborn is advance version of matplotlib. It has better design and specialized visualization technique. It is high level library provide the simplified function.

Technique to handle missing data
fill in using fillna & interpolation
dropna
replace

Practical 2 :-

Data transformation :- refers to ETL i.e. Extract transform load.

↓
Smoothing -> noises are removed

Aggregation -> Presenting data in suitable format.

Generalization -> low level -> high level

Normalization -> Converting data variable in given range to similar data

2- Scale

min-max

Decimal scaling.

Method to detect outliers

Box plot + visual

Z-score :-

Scatter plot

Inter quartile method

$x = [1, 2, 2, 3, 4, 4, 4, 5]$

mode $\Rightarrow 4$

unimodal \rightarrow one mode

bimodal - 2 modes. (two values appear most often same high freq.)

multimodal, no mode

describe() \rightarrow mean, median, max, min, count, 25%, 75%.

Practical 4:- Concept of linear regression

\rightarrow It is algorithm used for the prediction of the target variable based on independent variable or feature variable. used in supervised learning task.

It has linear relationship betwⁿ feature variable & target variable

It gives the mean square error betwⁿ actual value and predicted value which tells accuracy of model

In eqⁿ $y = mx + b + e$

- a. Gradient descent \rightarrow method of linear regression use to find optimal coefficient for regression model by minimizing difference betwⁿ predicted & actual value.
- b. Least square method

Supervised and Unsupervised learning:-

- ① In supervised learning \rightarrow algorithm learns from labeled data or known data where each training set consist features & corresponding variable. based on this, it make on new pattern patterns insight.
Even unseen data base.

- ② Unsupervised learning :- algorithm learn from a unlabeled data which has no target variable to predict. It clusters similar data pts. target goal is to find pattern and insights in data without any guidance.

Interpolate:- Interpolate is method or function used to fill the missing value in dataset by the estimations.
method

1) Linear:- estimate missing value bas on adjacent non-miss value

2) Polynomial:- use polynomial fn on non-missing value to estimate value

Practical 5

What is logistic Regression:-

Sigmoid
function

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

logistic regression is an algorithm used for binary classification, which has only two possible outcomes true/false Yes/no 1/0. It predicts the probability of occurrence of data point of logistic case on particular class.

method of logistic regression

- ① binary logistic
- ② multinomial logistic regression:-

- linear regression give continuous output \swarrow, \searrow
- logistic regression give constant output $\swarrow \rightarrow$ sigmoid
- stock prices, house price
- Predicting where patient has cancer or not

Confusion matrix:- It's a matrix of table used to give performance of model
True-ve, True+ve, False+ve, False-ve.

Standard score \rightarrow Remove mean & unit the so
* as mean can be zero center

Precision \rightarrow measure of accuracy of test instance
Recall \rightarrow measure of ability to correctly identify the -ve instance.

Practical 6:- Naive Bayes base of Bayes Probability
which ass all feature are independent of
each other and calculate probability of
feature occur in classes

Practical 8 \rightarrow matplotlib!

Practical 9 \rightarrow 8 10

Pugplot - 1 dimensional representation.

(11) \rightarrow

what is Impala:-

It's an open source parallelly processing SQL
query engine data store in apache-hadoop
clusters.

all to run SQL queries on big data
providing fast process.

(12) Scala \rightarrow Scalable programming language
which has syntax similar to the
java.

It's Object oriented & "fine" lang

Everything in scala is object & operation performed
are method call.

Hadoop :- Open source framework for
distributed storage which store
documents in chunk on different
computer system & process large
data set.

It is ~~fine~~ Scalable

It is used for big data analysis. & Data source
variety.

MapReduce \rightarrow It is a model used in Hadoop
which processes and analyzes large data set
in parallel clusters. - break task in

Small - task & distribute across node in cluster
: