

MACHINE LEARNING



SARASWATI Education Society's

SARASWATI College of Engineering

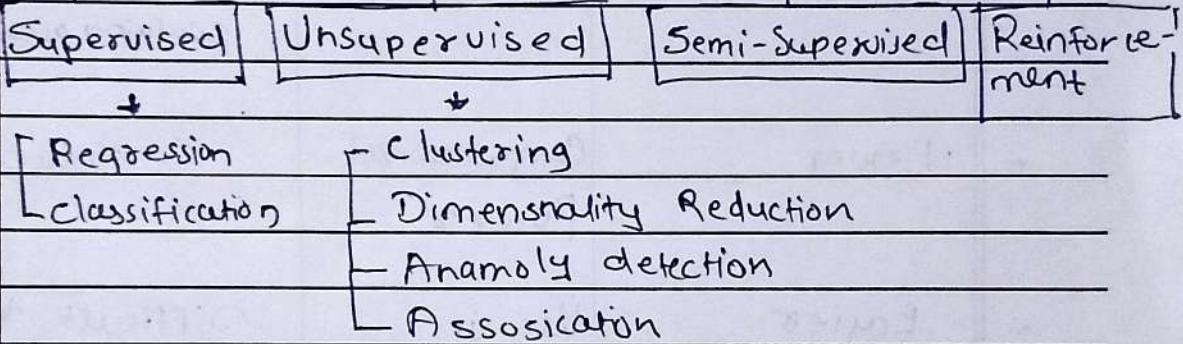
PAGE NO.: 1.

DATE: _____

Definition

Machine Learning is a field of computer science that uses statistical techniques to give computer systems the ability to "learn" with data, without being explicitly programmed.

Types of ML



Batch ML Vs Online ML

- Batch ML:
Also known as offline learning, involves training a machine learning model using entire dataset at once.

In Batch model is trained on a fixed dataset & Parameters are updated based on gradients computed from entire dataset.

- Online ML:
Also known as incremental learning or real-time learning, involves training a ML model continuously as new data become available over time.

In online learning, model is updated incrementally with each new data point or small batch of data.



Online learning libraries

River

vopal Rabbit

Offline Learning	Features	Online Learning
- less complex as model is constant	Complexity	Dynamic Comp. as the model keep evolving over time
- fewer	Computational Power	Continuous
- easier to implement	Use in production	Difficult to implement.
image classification, & anything in ML.	Applications	used in finance, health new data patterns.
Industry proven Tools, e.g. Scikit, TensorFlow, Pytorch, keras, Spark Mlib.	Tools	Active research, New pr tools. e.g. MOA, SAMOA, scikit-multiflow

SARASWATI College of Engineering

ML TYPES on the basis of Learning.

2 TYPES

Instance Based Model Based

1. Instance Based Learning

- It is also known as memory-based learning or lazy-learning
- here, the systems that learn the training examples by heart and then generalizes to new instances based on some similarity measures. It builds the hypotheses from the training instances.

Instance-based learning algorithms

1. K Nearest Neighbor (KNN)
2. Self-Organizing Map (SOM)
3. Learning Vector Quantization (LVQ)
4. Locally Weighted Learning (LWL)
5. Case-Based Reasoning ...etc.
etc.

2. Model Based Learning

- It involves creating a model that generalizes from the training data to make predictions on new, unseen data. This model can capture patterns, relationships, and features within the data to provide a way to understand and predict outcomes.

Model-based learning algorithms are:

decision trees, linear regression, neural networks, support vector machines etc.



SARASWATI College of Engineering

* Challenges in Machine Learning *

1. Data Collection
2. Insufficient Data / Labelled Data
3. Non Representative Data
4. Poor Quality Data
5. Irrelevant Features
6. Overfitting
7. Underfitting
8. Software Integration
9. Offline Learning / Deployment
10. Cost Involved

* Application of Machine Learning (B2B)

Sectors of ML.

1. Retail - Amazon / Big Bazaar
 - a. buying Pattern of user
 - b. deciding Positioning of Product
 - c. targeted marketing
 - d. seasonal Product etc.
2. Banking and Finance
 - a. Profile matching with past defaulters of loan.
 - b. Plans according to customers. etc.
3. Transport - OLA
 - a. Search Pricing according to time
 - b. Demand Forecasting
 - c.



4. Manufacturing - Tesla

- Fault occurrence sensors.
- Predictive Maintenance.

5. Consumer Internet - Twitter

- ~~tweets~~ Sentiment analysis
tweets

★ Machine Learning Development Life Cycle (ML DLC)

- Frame the Problem.
- Gathering Data.
- Data Preprocessing
- Exploratory Data Analysis
- Feature Engineering and Selection
- Model Training , Evaluation and Selection
- Model Deployment
- Testing
- Optimize



* Various Data Based Job Roles

1. Data Engineer
2. Data Analyst
3. Data Scientist
4. ML Engineer

1. Data Engineer Job Roles

- Scrape Data from the given sources.
- move / store the data in optimal servers / warehouses.
- Build data pipelines / APIs for easy access to the data.
- Handle databases / data warehouses.

2. Data Analyst Job Roles

- Cleaning & organizing Raw data.
- Analyzing data to derive insights.
- Creating data visualizations.
- Producing & maintaining reports
- Collaborating with teams / colleagues based on the insight gained.
- Optimizing data collection procedures

3. Data Scientist Job Roles

- Description: " A data scientist is someone who is better at statistics than any software engineer and better at software engineering than any statistician".
- Josh Wills on Quora
- It is a fullstack developer who can do all things like data engineering, data analytics & ML engineering



4. ML Engineer Job Role

- Deploying machine learning models to production ready environment
- Scaling and optimizing the model for production
- Monitoring and maintenance of deployed models.

* Tensor in Machine Learning

Tensors are mathematical objects that can be used to describe physical properties, like scalars & vectors, scalar is a zero rank tensor vector is a first rank tensor

Types of Tensors

- 0D Tensor / vector

Dimensions are zero

- 1D Tensor / vector

$[1, 2, 3, 4] \rightarrow$ 1D Tensor \therefore 1D Tensor
vector

2D Tensor / Matrices

Collection of vectors:

$$\begin{bmatrix} [1, 2, 3] \\ [4, 5, 6] \\ [7, 8, 9] \end{bmatrix} \Rightarrow 2D$$

ND Tensors

3D Tensor = $4 \times 3 \times 3$

4D Tensor = Vector of 3D tensor

5D Tensor = Metrics of 4D tensor

⋮



* Main Python Tools *

- Anaconda
- Jupyter Notebook
- Virtual Environment
- Kaggle Notebook
- Google Colab

End to End Project

Steps

0. Preprocess + EDA + feature selection
1. Extract input and output cols
3. Train test split
4. Train the model
5. Evaluate the model / model selection
6. Deploy the model

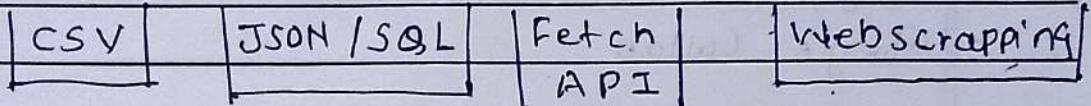
* How to Frame a Machine Learning Problem

1. Business Problem to ML Problem Conversion
2. Type of Problem
3. Current Solutions Availability
4. Getting Data
5. Metrics to measure
6. Online Vs Batch ?
7. Check Assumptions



* Gathering Data

Data Gathering



1) CSV - Comma Separated Values

Some important Commands and parameters

1. Importing pandas

```
import pandas as pd
```

2. Opening a local csv file

```
df = pd.read_csv('file-name.csv')
```

3. Opening a ext TSV file.

```
pd.read_csv('file-name.tsv', sep = '\t')
```

4. giving column names with

```
pd.read_csv('.csv', name = ['col1', 'col2'])
```

5. Choosing particular columns index column

```
pd.read_csv('.csv', index_col = 'col-name')
```

6. Header parameter

for starting first row in n no. of column.

7. Use cols parameter

import only required parameters.



8. squeeze parameters.

- Convert pandas df into series object.

9. Skiprows / nrows Parameter

- Skiprows for skipping unwanted rows with give index of number.

- nrows for importing n no. of rows

10. Encoding parameter

- change the encoding of character

11. Skip bad lines : Parser Error

12. dtypes parameter: Change datatype of column with this adding as dictionary.

13. Handling Dates

convert as object as "datetime 64 "

14. Convertors

Applying transformation to any features.

15. na-values parameter

assigning / specifying any values to na values

e.g. na_values = [-1, '0']

16.

16. Loading a huge dataset in chunks

chunksize = number

pd.read_csv('file.csv', chunksize=5000)



> Working with JSON / SQL

- JSON - JavaScript Object Notation

For loading JSON data

pd.read_json('data.json')

- SQL - Structured Query Language

For loading SQL data

- run SQL file in any database (e.g. mysql)

- install connector library

: pip install mysql-connector

- import mysql.connector

run query:

Conn = mysql.connector.connect(host='localhost',

user='root',

password=' ',

database='world')

it will return connection object
connection

* Now pass connection object on Function,

- pd.read_sql_query("SELECT * FROM city", conn)
↓ (object)
(table-name)



* Fetching Data from an API

API - Application Programming interface

These are the data pipelines for sharing data from one application to other

For loading data use requests module.

```
- import pandas as pd  
import requests  
response = requests.get('url')  
response.json()
```

Fetching data - using Web Scraping

for web scraping use beautiful soup library.

```
import pandas as pd  
import requests  
import bs4  
from bs4 import BeautifulSoup  
web = request.get('url').text  
now create object of library  
soup = BeautifulSoup(web, 'xml')
```



★ Understanding Your Data ★

D-19

1. Data size >> df.shape
2. Data head >> df.head()
3. Data sample >> df.sample(5) <Random>
4. Data type >> df.info()
5. missing values >> df.isnull().sum()
6. Maths summary of data >> df.describe()
7. Duplicate values >> df.duplicated().sum()
df.drop_duplicates()
8. corr betn. cols >> df.corr()
df.corr()['colname']

★ EDA - Univariate Analysis ★

D-20

Individual analysis of every columns.

1. Categorical Data.

a. Count plot.

>> sns.countplot(df['col-name'])

b. Pie Chart

>> df['col-name'].value_counts().plot(kind='pie')

2. Numerical Data.

a. Histogram >> plt.hist(df['col-name'])

b. Distplot >> sns.distplot(df['col-name'])

c. Boxplot >> sns.boxplot(df['col'])

des. df['col'].min(), .max(), .mean(),
df['col'].skew()



EDA - Bivariate and Multivariate

D-21

Exploratory Data Analysis

Type ② Bivariate

Type ③ Multivariate

The Bivariate & Multivariate analysis is easy with seaborn library.

① scatter plot (numerical - numerical)

sns.scatterplot (x, y, hue, style, size)

x = numerical - datavalue, y = numerical - datavalue

② bar plot (numerical - categorical)

sns.barplot (x = categorical, y = numerical, hue)

③ box plot (numerical - categorical)

sns.boxplot (x = categorical, y = numerical, hue)

④ Distplot (numerical - categorical)

sns.distplot (categorical - data, numerical, kde=)

⑤ HeatMap (categorical - categorical)

sns.heatmap (pd.crosstab ('categorical', 'numerical'))
(categorical)

pd.crosstab ('categorical-data', 'categorical-data')

⑥ ClusterMap (categorical - categorical)

sns.clustermap (pd.crosstab ('categorical'), 'categorical')



(7) Pair Plot

`sns.pairplot(dataframe, hue=)`

(8) Line Plot (numerical - numerical)

 $x = \text{time based number}$ $y = \text{numerical}$ `sns.lineplot(x=, y=)`

* Pandas Profiling *

Pandas Profiling is a library which provides a solution to understanding data with the help of generating reports for datasets that have numerous features.

- Installing Pandas Profiling

```
pip install pandas_profiling
```

- Importing:

```
from pandas_profiling import ProfileReport
```

- Make an object & add `dataframe`.`init`

```
report = ProfileReport(dataframe)
```

- Make html file:

```
report.to_file('output.html')
```

The One html file will be download in your folder.

Monday

FEATURE ENGINEERING

SARASWATI Education Society's

SARASWATI College of Engineering

Day 23

PAGE NO.: 16

DATE: 17/6/24

Feature Engineering

Feature Transformation	Feature Construction	Feature Selection	Feature Extraction
------------------------	----------------------	-------------------	--------------------

- Missing Value Imputation
- Handling Categorical Features
- Outlier Detection
- Feature Scaling

Feature Scaling (Standardization)

* When to use Standardization

Algorithms

1. k-Means
2. K-NN
3. PCA
4. ANN
5. Gradient Descent

Feature scaling is a technique to standardize the independent features present in the data in a fixed range

Standardization is also called as Z-score Normalization

$$x'_i = \frac{x_i - \bar{x}}{s}$$



Feature scaling - (Normalization)

Normalization

Normalization is a technique often applied as a part of data preparation for ml. The goal of normalization is to change the values of numeric columns in the dataset to use a common scale, without distorting difference in the range of values or losing information.

MinMax Scaling

Mean Normalization

MaxAbs Scaling

Robust Scaling

$$\text{MinMax Scaling} \quad x'_i = \frac{x_i - x_{\min}}{x_{\max} - x_{\min}}$$
$$\text{range} = [0, 1] \quad \min = 0, \max = 1.$$

Mean Normalization

$$\text{range} = [-1 \text{ to } 1] \quad x'_i = \frac{x_i - \bar{x}_{\text{mean}}}{x_{\max} - x_{\min}}$$

MaxAbs Scaling

$$\text{range} = [-1, 1] \quad x''_i = \frac{x_i}{|x_{\max}|}$$

used when sparse data

Robust Scaling

$$\text{range} = \text{IQR} \quad x'_i = \frac{x_i - \text{median}}{\text{IQR}}$$

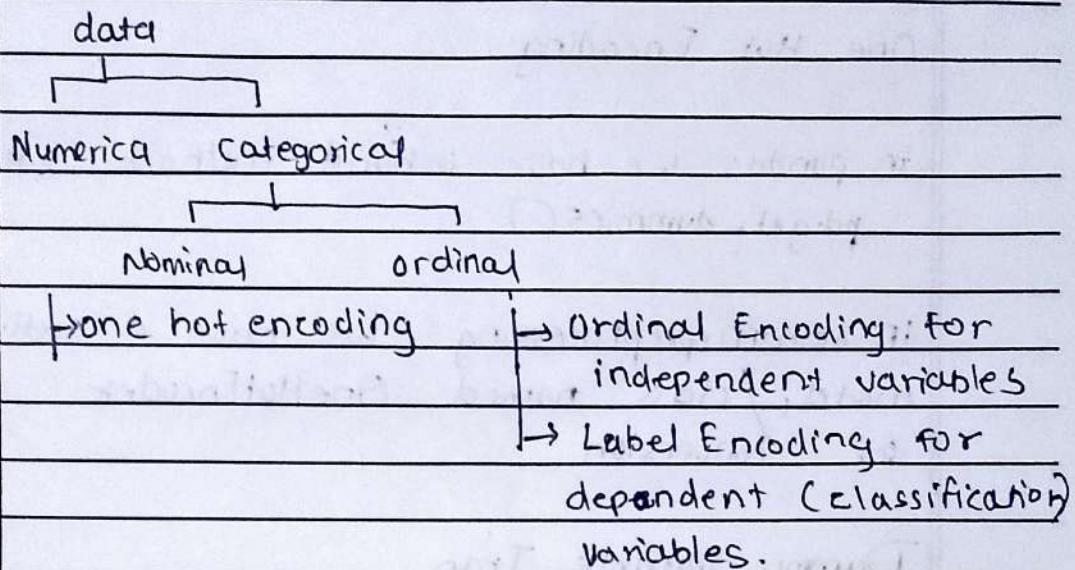
Used when Outliers on data



SARASWATI College of Engineering

Feature Transformation

2) Handling Categorical Features



Ordinal Data

Order matters, categories have weighted magnitude.

Ordinal Encoder \Rightarrow for categorical ordinal independent data (X)

Label Encoder \Rightarrow For categorical ordinal dependent variables (Y) output column is classification.



Nominal Categorical Data

To convert Nominal Data we use
One Hot Encoding

in pandas we have inbuilt method called
`pd.get_dummies()`

In `sklearn.preprocessing` we have a ~~cat~~
model / class named `OneHotEncoder`
for conversion

Dummy Variable Trap

The Dummy variable trap also called as
multicollinearity because it generates/causes
multicollinearity among the dummy or onehotEncoded
variables.

To reduce multicollinearity always drop
(n-1) columns from n columns which are
formed from encoding

Column Transformer

This is shortcut method for Encoding
multiple features.

Import.

```
from sklearn.compose import ColumnTransformer
transformer = ColumnTransformer ( transformers
= [ () , () , () ... ] , remainder=''' )
```



SK Learn · Pipelines

Pipelines chains together multiple steps so that the output of each step is used as input to the next step.

Pipelines makes it easy to apply the same preprocessing to train and test!

To use pipeline in ml we have library.

from: sklearn.pipeline import Pipeline, make_pipeline

Pipeline and make_pipeline these 2 classes do same work we can use anyone among them

Syntax:

Pipeline(steps: list of tuples, memory: str or object with joblib, verbose: bool)

e.g.

```
pipe = Pipeline([
    ('trf1', trf1),
    ('trf2', trf2)
], memory=None,
verbose=False)
```

Attributes:

1) named_steps: Bunch

Access the steps by name.



2) classes : ndarray of shape (n-classes)
The classes labels

3) n-features-in_ : int

No. of features seen during first step fit method

4) feature_names-in_ : ndarray of shape
(n-features-in_,)

Names of features seen during first step fit method

make_pipeline syntax:

sklearn.pipeline.make_pipeline(*steps : list of
estimator objects,
memory: default = None,
verbose: default = False)

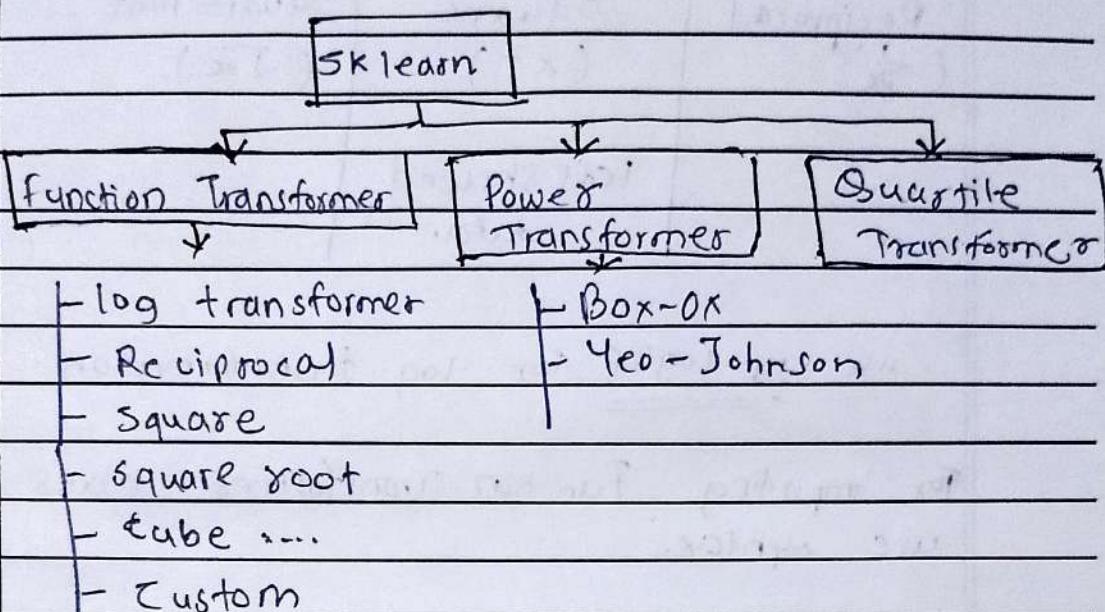
Example.

```
from sklearn.naive-bayes import GaussianNB  
from sklearn.preprocessing import StandardScaler  
from sklearn.pipeline import make_pipeline, Pipeline  
pipe = make_pipeline(StandardScaler(),  
                     GaussianNB())
```

```
pipe = Pipeline([('standardscalar', StandardScaler()),  
                ('gaussiannb', GaussianNB())])
```



* FUNCTION TRANSFORMER *



How To find if data is normal?

⇒ To find use

sns.distplot

pd.skew()

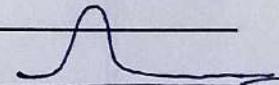
• Plot One [GOF plot]

* Log Transform

Take the log of the data values

when to apply?

⇒ on right skewed data.



when not to apply?

⇒ on -ve values.

why to use log transform.

⇒ convert the right skewed data into normal distribution. which is helpful for linear & logistic model.



Other transforms

Reciprocal

$$\left(\frac{1}{x}\right)$$

Square

$$(x^2)$$

Square root

$$(\sqrt{x})$$

left skewed
data.use np.log1p for log transformationfor importing Function Transformers class
use syntax.

from sklearn.preprocessing import FunctionTransformer



* Power Transformer *

Box - Cox
Transformer

Yeo - Johnson
Transform

Power Transformer is another class of sklearn which is used for transformed the data into normal distribution.

Box Cox Transform

The exponent here is

$$x_i^{(\lambda)} = \begin{cases} \frac{x_i^{\lambda} - 1}{\lambda} & \text{if } \lambda \neq 0, \\ \ln(x_i) & \text{if } \lambda = 0, \end{cases}$$

The exponent here is a variable called lambda(λ) that varies over the range of -5 to 5, and in the process of searching, we examine all values of λ .

Finally, we choose the optimal value (resulting in the best approximation to a normal distribution) for your variable.

Syntax for importing.

from sklearn.preprocessing import PowerTransformer
Now make an object of it.

Pt = PowerTransformer(method='box-cox')



* Yeo - Johnson Transform

if ..

$$x_i^{(\lambda)} = \begin{cases} [(\alpha_i + 1)^{\lambda} - 1] / \lambda & , \lambda \neq 0, \alpha_i \geq 0 \\ \ln(\alpha_i) + 1 & , \lambda = 0, \alpha_i \geq 0 \\ -[(-\alpha_i + 1)^{2-\lambda} - 1] / (2-\lambda) & , \lambda \neq 2, \alpha_i < 0 \\ -\ln(-\alpha_i + 1) & , \lambda = 2, \alpha_i < 0 \end{cases}$$

This transformation is somewhat of an adjustment to the Box-Cox transformation, by which we can apply it to negative numbers.

Syntax for import:

```
from sklearn.preprocessing import PowerTransformer
```

make an object of Yeo-Johnson Transform

```
pt = PowerTransformer()
```

(we don't need to specify "method" in Yeo-Johnson it is by default set on sqrt)



Encoding Numerical Features

Discretization
(Binning)

Binarization

* Discretization

Discretization is the process of transforming continuous variables into discrete variables by creating a set of contiguous intervals that span the range of the variable's values.

Discretization is also called binning, where bin is an alternative name for interval.

Why to use Discretization:

1. To handle outliers
2. To improve the value spread

* Types of Discretization *

Binning

Unsupervised

Binning

Equal width
(uniform)

Supervised

Binning

Decision Tree

Custom

Binning

Equal frequency
(quantile)

kmeans



* Unsupervised Binning *

① Equal width / Uniform Binning

$$\text{formula} = \frac{\text{Max} - \text{Min}}{\text{bins}}$$

Characteristics:

1. Handle Outliers for some spread
2. No change in spread.
3. Bins are equally distributed.

② Equal Frequency / Quantile Binning

$$\text{Intervals} = 10$$

Each interval contains 10% of total observations

Characteristics:

1. Handle Outliers
2. Change in value spread
3. Value is equally spread.

③ kMeans Binning

No of bins is decided by user k.

Use kmeans algorithm for Binning

6

Characteristics:

1. Adopts the distribution of the data
2. Cluster - Based Binning
3. .



* Custom Binning *

Encoding discretized variable.

SKLearn



kBinsDiscretizer()

Bins

Strategy

Encoding

→ Uniform

→ Ordinal

→ quantile

→ One hot encoding

→ Kmean

* Custom / Domain Based Binning

With the help of Pandas we can apply
Custom Based Binning.

These Custom Binning can be made as per
our requirement.

* Binarization

Here we convert the numerical values into
Binary value. (0, 1)s



Handling date time data

For doing operations on date & time we required to import necessary modules.

import datetime as dt

import time

in data if date and time is in object format then convert it into datetime64[ns] format for performing operations, using pandas function "to_datetime()"

e.g.

date['column'] = pd.to_datetime(date['column'])

(here, date is a dataframe)

There are different - different datetime functions in a datetime module.

import datetime as dt.

dt.time

dt.hour

dt.year

dt.minute

dt.month

dt.second

dt.days

dt.date

dt.month_name()

dt.day_name()

dt.dayofweek

dt.isocalendar().week

dt.quarter

dt.today()

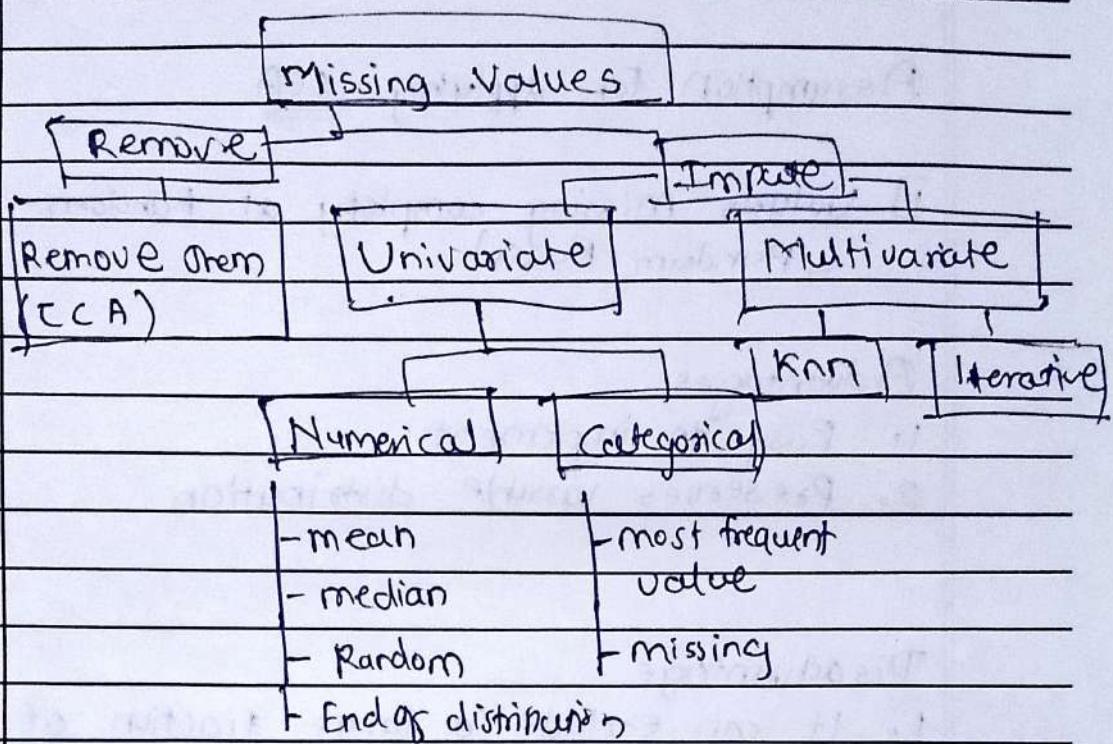
HANDLING Missing DATA

SARASWATI Education Society's

SARASWATI College of Engineering

PAGE NO.: 30

DATE: 27/6/2024



Complete Case Analysis

Complete-case analysis (CCA), also called "list-wise deletion" of cases, consists in discarding observations where values in any of the variables are missing.

Complete Case Analysis means literally analyzing only those observations for which there is information in all of the variables in the dataset.

Observation = Row

Variables = Column



Assumption for applying CCA

- 1] Values missing completely at Random locations.
(Random Rows)

Advantages.

1. Easy to implement
2. Preserves variable distribution

Disadvantage :

1. It can exclude a large fraction of the original dataset.
2. Excluded observations could be informative for the analysis.
3. When using our models in production, the model will not know how to handle missing data.

When to Apply CCA.

1. MCAR
2. missing values = [5% of data]



Handling Missing Numerical Data

Univariate

Imputation

Multivariate

Imputation

(2) Mean / median Imputation

Benefit

- 1] simple
- 2] ($5\% >]$ ~~seen~~ data values which are missing

Disadvantage

- 1] Distribution changes
- 2] Outliers
- 3] Covariance / Correlation changes

When to use

- 1] MCAR
- 2] $5\% <$

(3) Arbitrary Value Imputation

Replace missing Numeric value with any unusual value { -1, 99, 999 ...etc }

Benefit : Easy to apply

Disadvantage

- 1] PDF distribution
- 2] Variance / Covariance / Correlation changes

When to apply : When Data is not missing at Random position //



④ End of Distribution Imputation.

It is an extension of arbitrary value imputation

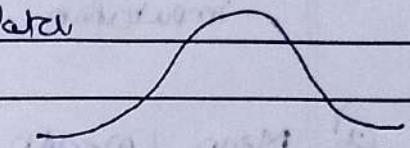
for Normally Distributed Data

it fill the missing values

with $(\text{mean} + 3\sigma)$

$(\text{mean} - 3\sigma)$

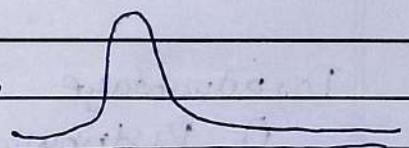
where σ is a sigma. S.D.



for Skewed Distributed Data

it fill the missing values

with [IQR proximity]



$$Q_1 - 1.5 * \text{IQR}$$

$$Q_1 = 25^{\text{th}}$$

$$Q_3 + 1.5 * \text{IQR}$$

$$Q_3 = 75^{\text{th}}$$

Benefit

1) easy to use.

Disadvantage:

1) PDF

2) Covariance

3) Variance

When to use

1) Data is Not Missing at Random

NMAR



* Handling Categorical Missing Data *

Most frequent Value
Imputation

Missing Category
Imputation

* Most Frequent Value Imputation

MFVI involves replacing missing values with most frequently occurring value (mode) in column.

To import syntax

```
from sklearn.impute import SimpleImputer  
make an object
```

```
imputer = SimpleImputer (strategy = 'most_frequent')  
imputer.fit_transform (x-train)
```

* Missing Category Imputation

MCII is a technique used primarily for categorical variables. It involves creating a new category (label) specifically for missing values.

To implement with sklearn. (syntax)

```
from sklearn.impute import SimpleImputer  
imputer = SimpleImputer (strategy = 'constant',  
fill_value = 'Missing')
```

```
imputer.fit_transform (x-train)
```

```
imputer.transform (x-test)
```



Random Sample Imputation

Missing Indicator

Automatic selection of parameter using Grid search CV.

1. Random Imputation

- ⇒ Preserves the variance of the variable
- b. Memory heavy for deployment, as we need to store the original training set to extract values from and replace the NA in coming observations.
- c. Well suited for linear models as it does not distort the distribution, regardless of the % of NA.

2. Missing Indicator

Replaces the Na with True & present value with False in new column. So model can improve accuracy.

import

```
from sklearn.impute import MissingIndicator  
mi = MissingIndicator()  
mi.fit(X-train)
```

Same method is inbuilt in SimpleImputer

```
si = SimpleImputer(add_indicator=True).
```

True One add-indicator parameter in SimpleImputer



3. Automatically - select - imputer - parameter using GridSearchCV.

GridSearchCV is a sklearn class in model-selection which is used for different value-parameter computation.

import

```
from sklearn.model_selection import train_test_split,  
GridSearchCV
```

* Multivariate Imputation technique:

- The KNN imputer is a technique used in multivariate imputation to fill missing values by considering their K-nearest neighbors.

import

```
class sklearn.impute.KNNImputer(*, missing_=  
    np.nan, n_neighbors=5,  
    weights='uniform', metric='minkowski',  
    add_indicator=False,  
    keep_empty_features=False)
```

- Iterative Imputer / MICE

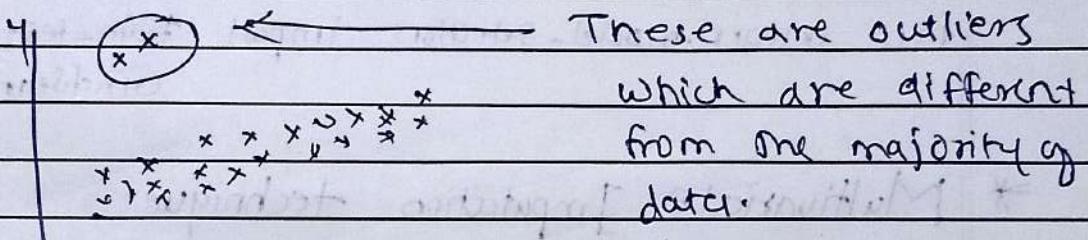
MICE stands for Multivariate Imputation by Chained Equations.

It is method for finding missing values by modeling each variable with missing data as a function of other variables. It is iterative process continues till convergence.



Outliers

Outliers are data points significantly different from the majority in a dataset. In ML outliers can impact model performance by skewing results or introducing noise. Identifying & handling outliers is crucial for building accurate & robust ML models.



Bad Impact of Outliers on ML algorithms.

- Linear
- Logistic
- AdaBoost
- Deep Learning

Non-effective of ML algorithms.

- Tree Based

How to Treat Outliers.

- ① Trimming (Removing directly)
- ② Capping (Setting boundaries)
- ③ Treat accordingly like if it's NaN value.
- ④ Do Discretization



How To Detect Outliers ?

1. Normal Distribution

if observation is $(\mu + 3\sigma) \leq (\mu - 3\sigma)$
out from these 2 Range's. then treat
that values as outliers.

2. Skewed Distribution

Detect using plotting box plot.

if No. is + less than $Q_1 - 1.5 IQR$
No. is greater than $Q_3 + 1.5 IQR$
then consider that value as outliers.

3. Percentile Based Approach

if observation is less than 2.5%

if observation is greater than 97.5%

then consider that value as outliers.

* Techniques for Outlier Detection & Removal

1. Z-score

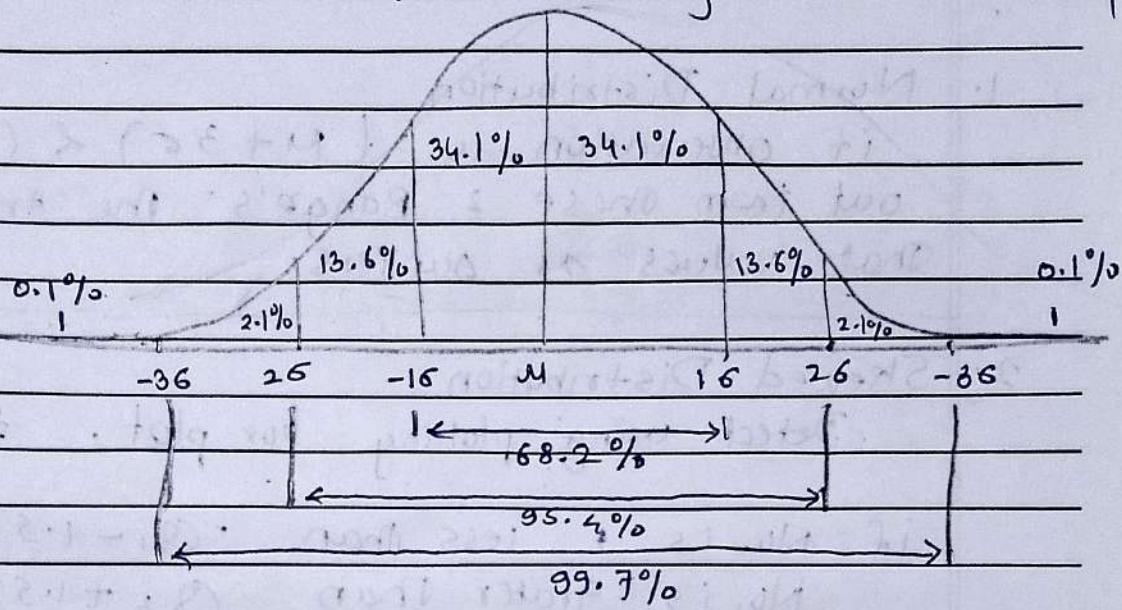
2. IQR based filtering

3. Percentile

4. Winsorization



Outlier Detection & Removal using Z-score method



Normal Curve (Bell curve)

Z score follows Normal curve distribution
Range of value distribution

$$\mu + \sigma \quad \} \quad 68\%$$
$$\mu - \sigma \quad \}$$

$$\mu + 2\sigma \quad \} \quad 95\%$$
$$\mu - 2\sigma \quad \}$$

$$\mu + 3\sigma \quad \} \quad 99\%$$
$$\mu - 3\sigma \quad \}$$

Outlier Treatment

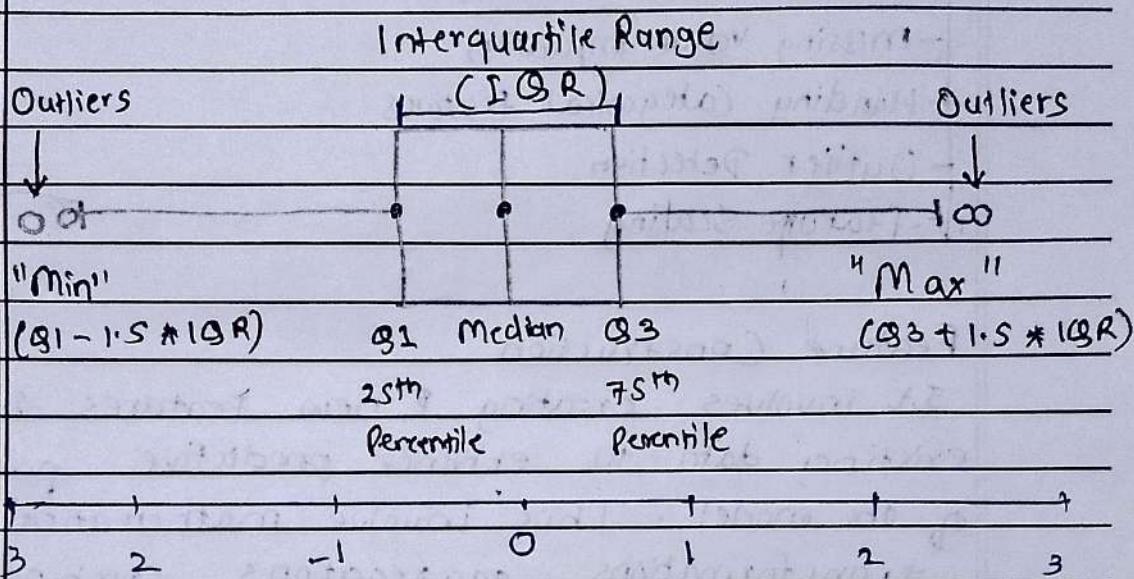
Trimming
directly remove the outliers

Capping
set an boundary to lower & upper limit for outliers



* Outlier Detection & Removal using IQR Method

The IQR method is used at the time when our data is skewed on right or left side



IQR Proximity Rule //

$$IQR = Q_3 - Q_1$$

$$75\text{th percentile} - 25\text{th percentile}$$

* Outlier Detection Using Percentile Method

for applying Percentile method there is no condition is mention we can apply on any type of data distribution.

In this method we set one equal amount of boundary on upper & lower limit in percentile.

most common is 1 percentile & 5 percentile.



Feature Engineering:

Feature Transformation Feature Construction Feature Selection Feature Extraction

- missing value imputation
- Handling Categorical features
- Outlier Detection
- feature Scaling

Feature Construction

It involves creating new features from existing data to enhance predictive power of the model. This involves mathematical transformations, aggregations, combinations of existing features, & domain-specific knowledge to generate features.

Feature Splitting

Feature splitting involves breaking down a single feature into multiple components, which can provide more granular information and potentially improve model performance.

Example:

① Splitting a timestamp into year, month, day, hour etc.

② Splitting text features

Splitting Categorical features into binary (dummy) variables.



Feature selection & Extraction in CDP

Curse of Dimensionality

Curse of features

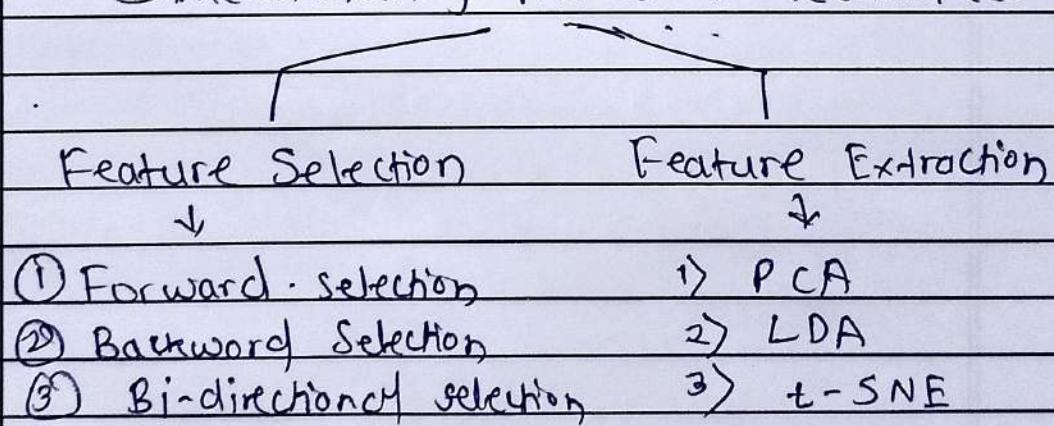
The curse of dimensionality is a phenomenon that occurs when analyzing & organizing data in high-dimensional spaces.

key aspects

1. Volume Increase : A dimensions increases the volume of the spaces grows exponentially.
2. Sparsity : Data points become sparse. The distance between data points increases.
3. Increased Computational Complexity : requires more computational resources to process.
4. Overfitting in ML models

To Reduce the Curse of Dimensionality we use

Dimensionality Reduction Techniques



Principle Component Analysis . (PCA)

PCA converts the higher dimensional data into best possible lower dimensional data.

Feature Extraction is an upgradation over a feature selection's limitations .

Benefits of PCA

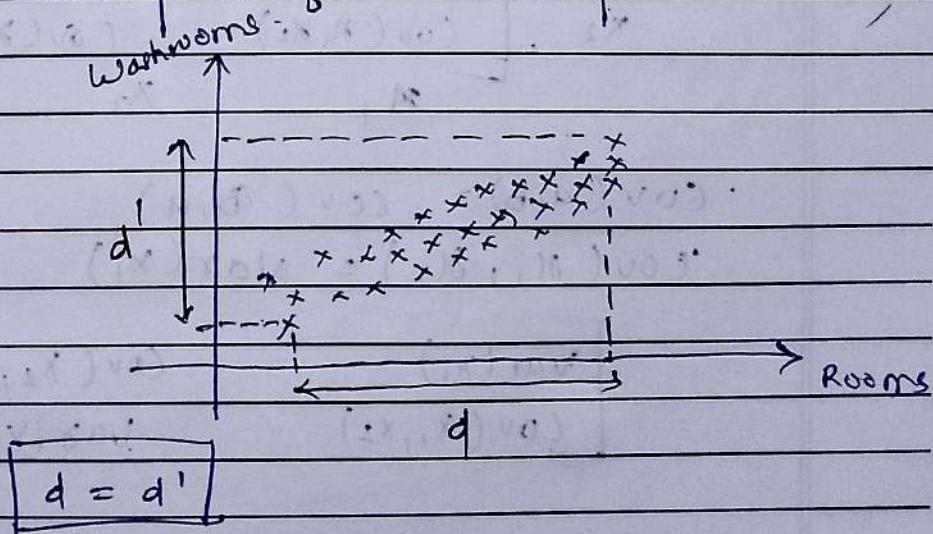
- 1) Faster execution of algorithm.
- 2) Visualization becomes very easy

PCA converts one higher dimensional data into 3D or 2D . if

e.g. 10 D dim. data \rightarrow PCA \rightarrow 3 D data

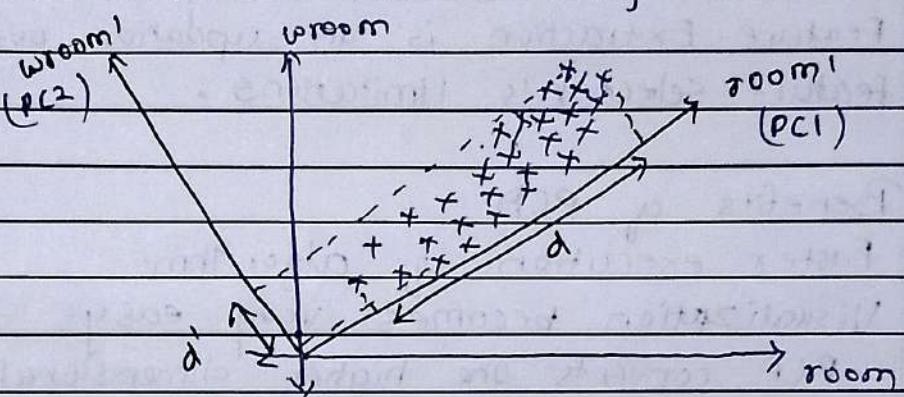
* If the data contains the columns in which variance is equally spreaded : then the feature selection technique - does not work at this time we use PCA . for better understanding see the below fig of equally variance spread & holding in 2 columns .

No. of rooms | No. of washrooms | Price.





According to the data spread the PC algo creates new axis at new plane which reduces the any one's dimension among 2 columns. The new axis columns, called as PC1 & PC2. see below fig.



The No. of PC $\leq n$

where n is no of original features in data.

* key points:

- Covariance and Covariance matrix
Co-variance is a relation between 2 columns

Covariance matrix

$$\begin{bmatrix} x_1 & \begin{bmatrix} \text{Cov}(x_1, x_1) & \text{Cov}(x_1, x_2) \\ \text{Cov}(x_2, x_1) & \text{Cov}(x_2, x_2) \end{bmatrix} \\ x_2 & \end{bmatrix}_{2 \times 2}$$

$$\text{Cov}(a, b) = \text{Cov}(b, a)$$

$$\text{Cov}(x_1, x_1) = \text{Var}(x_1)$$

$$\begin{bmatrix} \text{Var}(x_1) & \text{Cov}(x_2, x_1) \\ \text{Cov}(x_1, x_2) & \text{Var}(x_2) \end{bmatrix}_{2 \times 2}$$



Similarly for 3×3 , $\begin{matrix} & x & y & z \end{matrix}$

$$\begin{array}{c|ccc|c} & x & v_x & c_{x,y} & c_{x,z} \\ y & c_{y,x} & v_y & c_{y,y} & c_{y,z} \\ z & c_{z,x} & c_{z,y} & v_z & \\ \hline & x & y & z \end{array}$$

- Eigen Vectors.

These are special vector. The direction of eigen vectors is never changed while applying transformation. Only magnitude changes. Direction remains same in eigen vector. $A\vec{v} = \lambda\vec{v}$, always true.

- Eigen values.

These are the measurement of how much the eigen vector is stretched or shrunked after transformation.

Measures the amount of stretching or shrinking of eigen vector.

Problem:

Find a unit vector which maximizes the variance.

The largest eigenvector of covariance matrix always points into the direction of largest variance of data.



* STEPS TO SOLVE PCA

Step 1 : mean centring of the whole data

Step 2 : find covariance matrix

Step 3 : Find the eigen value & eigen vectors
for the covariance matrix. ($\lambda_1, \lambda_2, \lambda_3$
are eigen values)

Step 4 : After getting Principle Components
transform all the points from
higher dim space to lower w.r.t. Principle
Component.

for transforming points

firstly do one (dot) product of $U^T \cdot X$
with all points one by one on principle
Component and take transpose of the data
shape.

(1000, 3)	(2, 3)
↓ 4	↓ 3
Row	Col
vectors	
3D space	

1000, 3 (3, 2) (Transpose of
vector)

(1000, 2)



Simple Linear Regression

The Simple LR is a supervised learning method. The LR algorithm draws a line of best fit on a linear sort of data.

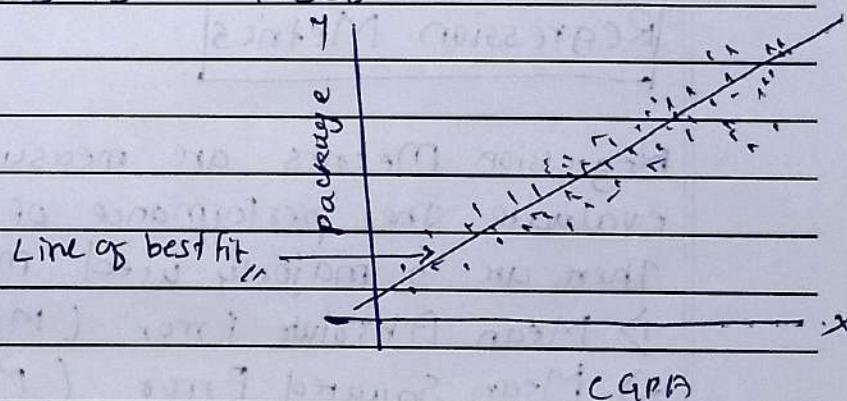
formula : $y = m \cdot x + b$

where y is dependent variable,

x is independent var.

m is weightage

b is offset



CGPA

import

from sklearn.model_selection

from sklearn.linear_model import LinearRegression
make object

lr = LinearRegression()

$$b = \bar{y} - m \bar{x}$$

$$m = \frac{n}{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}$$

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



Equation of Line of best fit.

$$\boxed{\bar{y} = \bar{m}\bar{x} + b}$$

Here, \bar{y} = mean of y_i

\bar{x} = mean of x_i

y_i = current value of y

x_i = current value of x

$y \rightarrow$

Regression Metrics

Regression Metrics are measures used to evaluate the performance of a regression model. There are 5 majorly used metrics.

1) Mean Absolute Error (MAE)

2) Mean Squared Error (MSE)

3) RMSE (Root Mean Squared Error)

4) R-squared (R^2)

5) Adjusted. (R^2)

* Loss Functions : It quantifies the difference between one predicted values and the actual values for a single data point. It is also called as Error function.

MAE, MSE & RMSE are comes under loss functions.



1. Mean Absolute Error (MAE)

measures the average magnitude of errors treating all errors equally regardless of their direction.

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

or

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

Advantages,

- 1) It has always a unit of dependent variable.
- 2) Robust to outliers.

Disadvantage

- Modulus graph is not differentiable

2. Mean Squared Error (MSE)

Measures average of squares of the errors

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

or

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

Advantages

- Differentiable
- use as loss function

Disadvantage

- Not Robust to Outliers.



3) Root Mean Square Error:

$RMSE = \sqrt{MSE}$ give higher weight to large errors

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

Advantage

Output is in same unit. \Rightarrow dependent var.

disadvantage

Robust to outliers

* Absolute

4) R² Score

It is also known as coefficient of determination or goodness of fit.

$$R^2 = 1 - \frac{SSR}{SSM}$$

where, SSR = Sum of Squared error Regression

SSM = Sum of squared error Mean

$$R^2 = 1 - \frac{\left[\sum_{i=1}^n (y_i - \hat{y}_i)^2 \right]_R}{\left[\sum_{i=1}^n (y_i - \bar{y})^2 \right]_M}$$

if $R^2 \rightarrow 0$ then bad

$R^2 \rightarrow 1$ then good

$R^2 \rightarrow$ true when $SSR > SSM$ very bad.



LPA ↑

Z(?)

Y_{mean}

Z(m)

← →

cgp4

Problem with R² score.

R² score increases with irrelevant features
also R² score fluctuates to high or in
increment when also when any unimportant
feature also added.

4) Adjusted R² score

Adjust the R² value based on the no of
predictor in the model.

$$\text{Adjusted } R^2 = 1 - \left[\frac{(1-R^2)(n-1)}{(n-1-k)} \right]$$

where, n → no. of rows

- k → total no of independent columns.



Multiple Linear Regression

Multiple LR Equation:

$$Y = \beta_0 + \beta_1 x_1 + \beta_2 x_2$$

The Multiple LR is an extension of simple LR that models the relationship between two or more independent variables & a single dependent variable by fitting linear eqⁿ to observed data.

Equation in general form:

$$Y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \dots + \beta_n x_n + \epsilon$$

where Y is independent var.

β_0 is intercept / offset

$\beta_1, \beta_2, \beta_3, \dots, \beta_n$ are coefficient for independent variables

x_1, x_2, \dots, x_p are independent variables

ϵ is a error term.

For understanding this algorithms in detail lets take an example:

We have data of 100 students containing 4 columns: CGPA | IQ | Graduate | LPA

(100, 4)	x_1	x_2	x_3	y
----------	-------	-------	-------	-----

is shape of data.

eqⁿ of hyperplane is -

$$\hat{y} = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \beta_3 x_3$$

wher. \hat{y} is y predicted



$$\hat{Y} = \begin{bmatrix} \hat{y}_1 \\ \hat{y}_2 \\ \vdots \\ \hat{y}_{100} \end{bmatrix} = \begin{bmatrix} \beta_0 & \beta_1 x_{11} & \beta_2 x_{12} & \beta_3 x_{13} \\ \beta_0 & \beta_1 x_{21} & \beta_2 x_{22} & \beta_3 x_{23} \\ \vdots & \vdots & \vdots & \vdots \\ \beta_0 & \beta_1 x_{1001} & \beta_2 x_{1002} & \beta_3 x_{1003} \end{bmatrix}$$

For n rows and m columns.

$$\hat{Y} = \begin{bmatrix} \hat{y}_1 \\ \hat{y}_2 \\ \vdots \\ \hat{y}_n \end{bmatrix} = \begin{bmatrix} \beta_0 & \beta_1 x_{11} & \beta_2 x_{12} & \beta_3 x_{13} & \dots & \beta_m x_{1m} \\ \beta_0 & \beta_1 x_{21} & \beta_2 x_{22} & \beta_3 x_{23} & \dots & \beta_m x_{2m} \\ \vdots & \vdots & \vdots & \vdots & \ddots & \vdots \\ \beta_0 & \beta_1 x_{n1} & \beta_2 x_{n2} & \beta_3 x_{n3} & \dots & \beta_m x_{nm} \\ & & & & & \beta_m x_{nm} \end{bmatrix}$$

$$= \begin{bmatrix} 1 & x_{11} & x_{12} & x_{13} & \dots & x_{1m} \\ 1 & x_{21} & x_{22} & x_{23} & \dots & x_{2m} \\ \vdots & \vdots & \vdots & \vdots & \ddots & \vdots \\ 1 & x_{n1} & x_{n2} & x_{n3} & \dots & x_{nm} \end{bmatrix} \begin{bmatrix} \beta_0 \\ \beta_1 \\ \vdots \\ \beta_m \end{bmatrix}$$

A or X

B

$$\hat{Y} = X \cdot B \quad \text{--- (1)}$$

where X is matrix & B is coefficient matrix

$$Y = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix} \quad e = Y - \hat{Y} = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix} - \begin{bmatrix} \hat{y}_1 \\ \hat{y}_2 \\ \vdots \\ \hat{y}_n \end{bmatrix}$$



$$e = \begin{bmatrix} y_1 - \hat{y}_1 \\ y_2 - \hat{y}_2 \\ \vdots \\ y_n - \hat{y}_n \end{bmatrix}$$

$$E = e^T e,$$

$$\left[(y_1 - \hat{y}_1) \quad (y_2 - \hat{y}_2) \dots (y_n - \hat{y}_n) \right] \begin{matrix} (y_1 - \hat{y}_1) \\ (y_2 - \hat{y}_2) \\ \vdots \\ (y_n - \hat{y}_n) \end{matrix} \begin{matrix} (1 \times n) \\ (n \times 1) \end{matrix}$$

$$(y_1 - \hat{y}_1)^2 + (y_2 - \hat{y}_2)^2 + (y_3 - \hat{y}_3)^2 + \dots + (y_n - \hat{y}_n)^2$$

$$\left[\sum_{i=1}^n (y_i - \hat{y}_i)^2 \right]$$

$$\because (A+B)^T = A^T + B^T$$

$$\therefore (A-B)^T = A^T - B^T$$

$$E = e^T e = (y - \hat{y})^T (y - \hat{y})$$

$$= (y^T - \hat{y}^T) (y - \hat{y})$$

$$= [y^T - (x\beta)^T] (y - x\beta)$$

$$\dots \text{from } (1) \boxed{\hat{y} = x\beta}$$

$$E = y^T y - y^T x\beta - (x\beta)^T y + (x\beta)^T x\beta$$

$$E = y^T y - 2y^T x\beta + \beta^T x^T x\beta$$

$$\text{From } \dots y^T x\beta = x\beta^T y$$

$$A^T B = B^T A$$



differentiate the eqn

$$\begin{aligned}\frac{dE}{d\beta} &= \frac{d}{d\beta} [Y^T Y - 2Y^T X \beta + \beta^T X^T X \beta] = 0 \\ &= 0 - 2Y^T X + \frac{d}{d\beta} [\beta^T X^T X \beta] = 0 \\ &= -2Y^T X + 2X^T X \beta^T = 0 \\ &= 2X^T X \beta^T = 2Y^T X\end{aligned}$$

$$\beta^T = \frac{Y^T X}{[X^T X]}$$

$$\beta^T = Y^T X (X^T X)^{-1}$$

$$(\beta^T)^T = [Y^T X (X^T X)^{-1}]^T$$

$$\beta = [(X^T X)^{-1}]^T (Y^T X)^T$$

$$\beta = [(X^T X)^{-1}]^T X^T Y$$

$$\beta = [(X^T X)^{-1} X^T Y]$$

$$\boxed{\beta = (X^T X)^{-1} X^T Y}$$

$$\begin{array}{c|cc} \beta_0 & [(m+1)*(m+1)] & [(m+1)*n] \\ \beta_1 & & [n*1] \\ \beta_2 & [(m+1)*n] & [n*1] \\ \vdots & & \\ \beta_m & & (m+1)*1 \end{array}$$

$$(m+1)*1 = (m+1)*1$$

$$\boxed{\beta = (X^T X)^{-1} X^T Y}$$

GRADIENT DESCENT

SARASWATI Education Society's



SARASWATI College of Engineering

PAGE NO.: 56

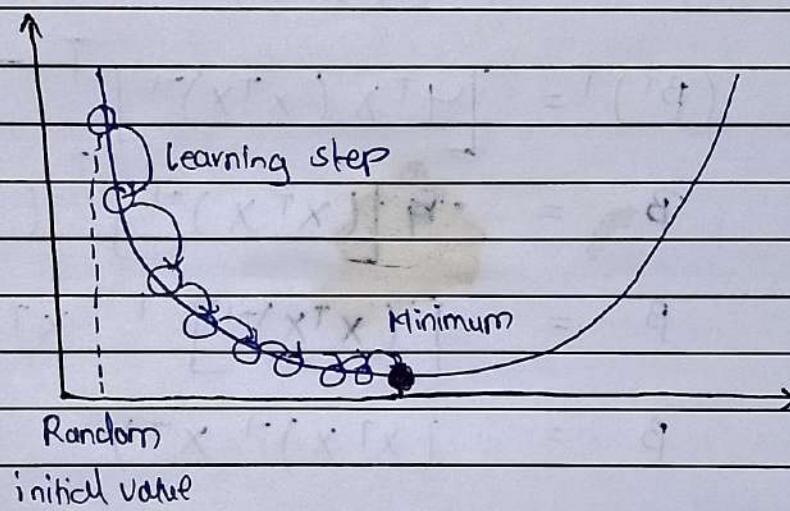
DATE:

Gradient : It is a derivative that defines the effects on outputs of the function with a little bit of variation in inputs

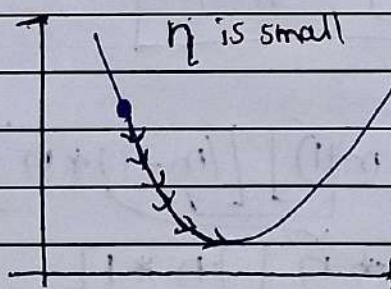
Gradient descent : is a method of unconstrained mathematical optimization.

It is a first order iterative algorithm for finding a local minimum of a differentiable multivariate function.

Cost



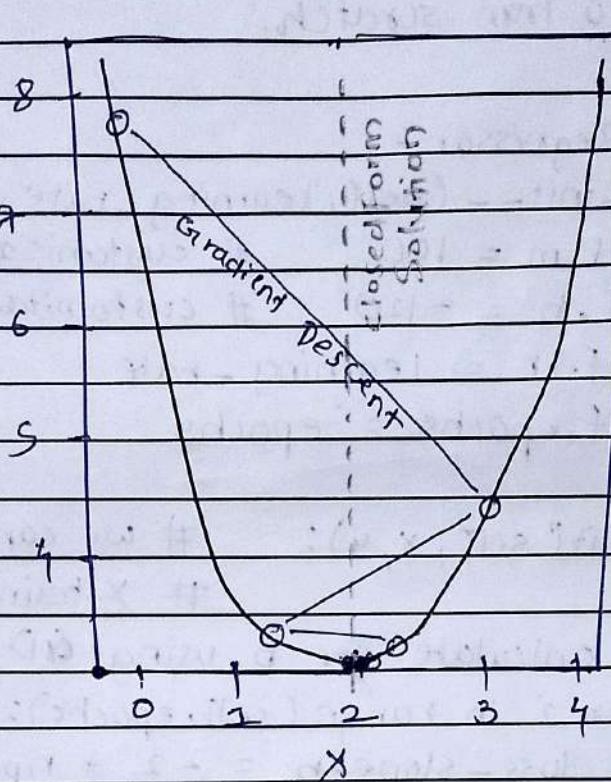
η is large



(η) is learning rate.

$$b_{\text{new}} = b_{\text{old}} - \text{slope}$$

$$b_{\text{new}} = b_{\text{old}} - \eta \text{slope}$$



$$\text{Gradient descent} = -2 \sum_{i=1}^n (y_i - mx_i - b)$$

Formula in code

for i in range (self.epochs):

$$\text{loss_slope_b} = -2 * \text{np.sum}(y - \text{self.m} * x.\text{ravel}() - \text{self.b})$$

$$\text{loss_slope_m} = -2 * \text{np.sum}(y - \text{self.m} * x.\text{ravel}() - \text{self.b}) * x.\text{ravel}()$$

$$\text{self.b} = \text{self.b} - (\text{loss_slope_b} * \text{self.lr})$$

$$\text{self.m} = \text{self.m} - (\text{loss_slope_m} * \text{self.lr})$$



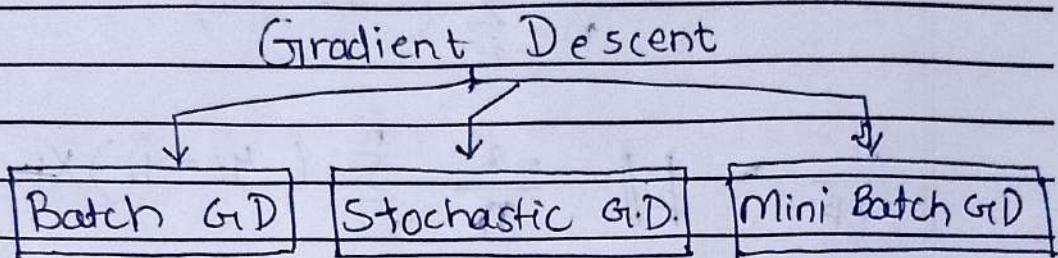
own Algo from scratch.

line:

```
1 class GDRegressor :  
2     def __init__(self, learning_rate, epochs):  
3         self.m = 100          # customizable  
4         self.b = -120         # customizable  
5         self lr = learning_rate  
6         self.epochs = epochs  
7  
8     def fit(self, X, y):      # we can replace with.  
9         # X-train, Y-train.  
10        # calculate m & b using. GD  
11        for i in range(self.epochs):  
12            loss_slope_b = -2 * np.sum((y -  
13                self.m * X.ravel() - self.b))  
14            loss_slope_m = -2 * np.sum((y -  
15                self.m * X.ravel() - self.b) * X.ravel())  
16            self.b = self.b - (loss_slope_b * self.lr)  
17            self.m = self.m - (loss_slope_m * self.lr)  
18            print(self.m, self.b)  
19  
20        def predict(self, X):  
21            return self.m * X + self.b  
22  
23    # make object.  
24    gd = GDRegressor(lr=0.001, epochs=100)  
25    gd.fit(X, y)      # replace with X-train, Y-train  
26    gd.  
27    y_pred = gd.predict(X_test)
```



* Types of Gradient Descent *



• Batch Gradient Descent

Batch gradient descent sums the error for each point in a training set, updating the model only after all training examples have been evaluated.

This process referred as a training epoch.

Batch GD usually produces a stable error gradient & convergence.

$$Y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \dots + \beta_n x_n$$

$$\beta_0 = \beta_0 - \eta \frac{\text{slope}}{\partial \beta_0} = \frac{\partial L}{\partial \beta_0}$$

$$\beta_1 = \beta_1 - \eta \frac{\text{slope}}{\partial \beta_1} = \frac{\partial L}{\partial \beta_1}$$

$$\beta_2 = \beta_2 - \eta \frac{\text{slope}}{\partial \beta_2} = \frac{\partial L}{\partial \beta_2}$$

$$\beta_n = \beta_n - \eta \frac{\text{slope}}{\partial \beta_n} = \frac{\partial L}{\partial \beta_n}$$

Loss function & its derivatives of each epochs -



Loss function for MSE

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

$$\frac{\partial L}{\partial \beta_1} = \frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i) x_{i1}$$

$$\frac{\partial L}{\partial \beta_2} = \frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i) x_{i2}$$

$$\frac{\partial L}{\partial \beta_3} = \frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i) x_{i3}$$

General formula of Loss function.

$$\frac{\partial L}{\partial \beta_m} = \frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i) x_{im}$$

Formula for intercept (β_0) derivative.

$$\frac{\partial L}{\partial \beta_0} = \frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)$$



- Stochastic Gradient Descent (SGD)

Stochastic gradient descent (SGD) runs a training epoch for each example within the dataset and it updates each training example's parameters one at a time. Since you only need to hold one training example, they are easier to store in memory. These make SGD faster - Its frequent updates can result in noisy gradients, but this can also be helpful in escaping the local minimum & finding global one.

Intercept (β_0)

$$\frac{\partial L}{\partial \beta_0} = -2 (y_i - \hat{y}_i)$$

where, $i = \text{idx}$

Coefficien. (β_1)

$$\frac{\partial L}{\partial \beta_1} = -2 (y_i - \hat{y}_i), x_{i1}$$

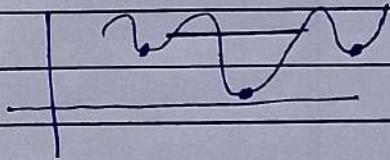
$\frac{\partial L}{\partial \beta_1} = -2 ((y_i - \hat{y}_i), x_i)$

When no of epochs are fixed in same size for batch & stochastic
Comparing with the speed
batch computes faster than stochastic

When to use Stochastic GD

1) Big data

2) Non convex function





- Mini-Batch Gradient Descent.

Mini-batch gradient descent combines concepts from both batch gradient descent and stochastic gradient descent. It splits the training dataset into small batches and performs updates on each of those batches. This approach strikes a balance between the computational efficiency of batch gradient descent & the speed of stochastic gradient descent.

$$\text{Min batch} = \text{Batch GD} + \text{Stochastic GD}$$



* Polynomial Regression

Polynomial regression is a form of regression analysis in which the relationship between the independent variable X and the dependent variable y is modeled as an n th degree polynomial in X .

Polynomial Regression is a special case of Linear Regression where we fit the polynomial equation on the data with a curvilinear relationship between the dependent and independent variables.

$$\text{Simple LR} \Rightarrow y = \beta_0 + \beta_1 x_1$$

$$\text{Multiple LR} \Rightarrow y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \dots + \beta_n x_n$$

$$\text{Polynomial LR} \Rightarrow y = \beta_0 + \beta_1 x_1 + \beta_2 x_1^2 + \dots + \beta_n x_1^n$$

Cost Function : CF of Polynomial is Mean Square Error,

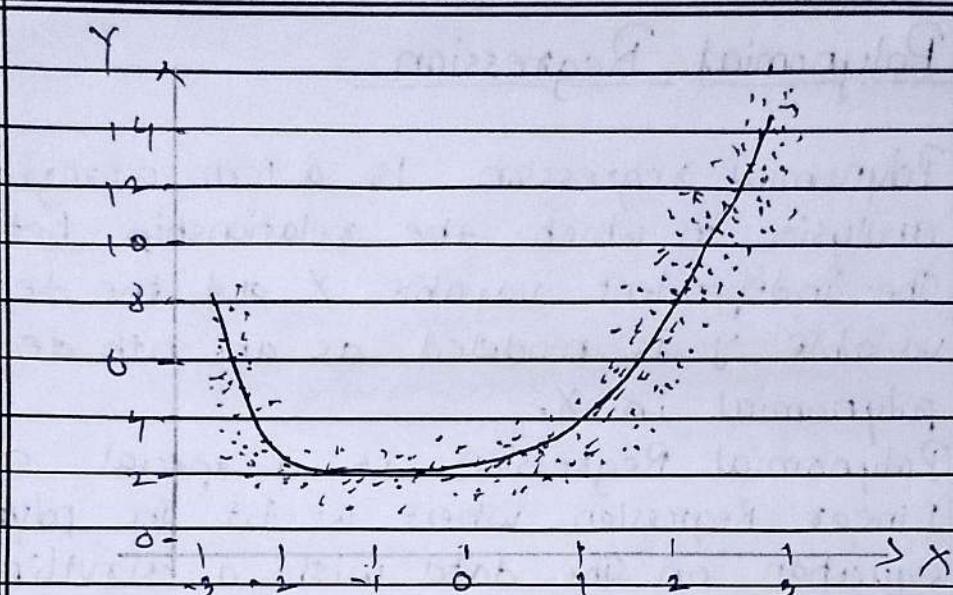
$$J = \frac{1}{n} \sum_{i=1}^n (\text{pred}_i - y_i)^2$$

Cost Function of Polynomial Reg.

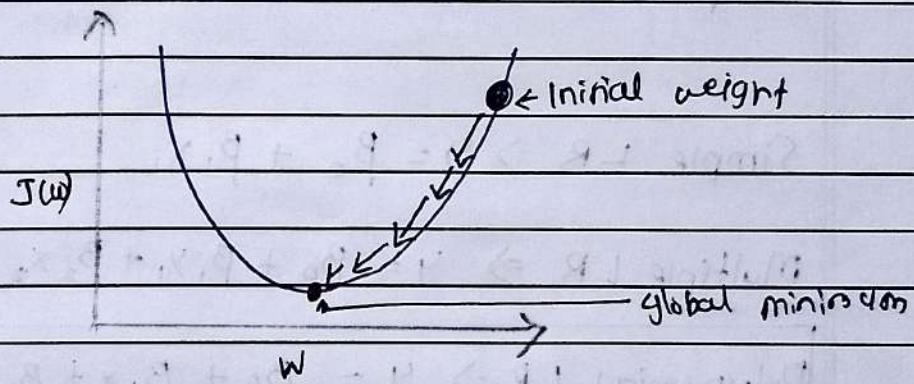
$$J = 1/n * \text{sum}(\text{square}(\text{pred} - y))$$

and one form,

$$J = 1/n * \text{sum}(\text{square}(\text{pred} - (\beta_0 + \beta_1 x + \beta_2 x^2 + \beta_3 x^3 + \dots)))$$



Polynomial regression.

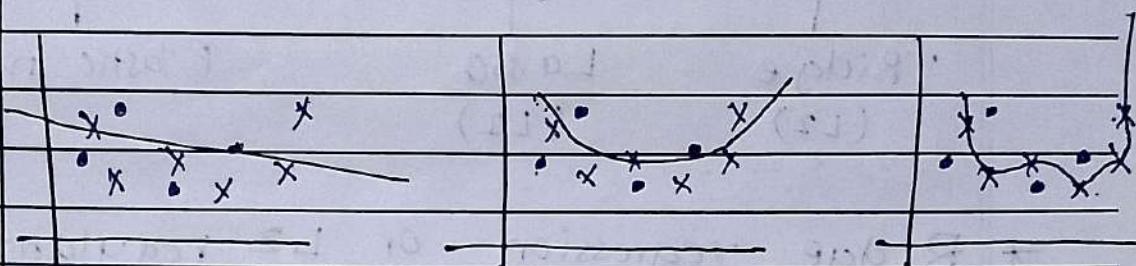


Gradient Descent for Polynomial Reg.



BIAS VARIANCE TRADE-OFF

Bias-variance tradeoff describes the relationship between a model's complexity. It tries to minimize the 2 sources of error which are bias & variance error.



Model

Model 2

Model 3

Bias

High

Mid
Low

Low

Variance

Low

Low

High

Fitting

Underfitting

best fit

Overfitting

Model

Simple

Sweet
Spot

Complex

- High variance means model makes strong assumptions about the data leads to underfitting.
- High variance means model is very sensitive to small fluctuations in the training data, leads to overfitting.

For Handling Bias Variance Trade-Off: 3

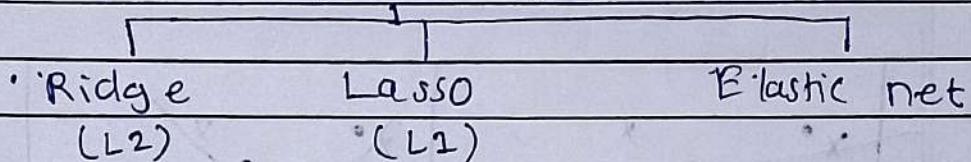
- 3 important techniques:
- 1) Bagging
 - 2) Boosting
 - 3) Regularization



Regularization

• Regularization is a set of methods for reducing overfitting in ML models.

Regularization



* Ridge regression - or L2 regularization

It penalizes high-value coefficients by introducing a penalty term in SSE loss fn.

- Using Ridge in code using sklearn class.

import.

from sklearn.linear_model import Ridge
make obj

R = Ridge(alpha = 0.01) # alpha is hyper-parameter.

train the model.

R.fit(x-train, y-train)

Predict.

R.predict(x-test).



Ridge Regression Mathematical Formulation.

$$\text{Loss Function.} = \sum_{i=1}^n (y_i - \hat{y}_i)^2 + \lambda m^2 \quad (L)$$

$$L = \sum_{i=1}^n (y_i - mx_i - b)^2 + \lambda m^2$$

$$\frac{\partial L}{\partial b} \quad \& \quad \frac{\partial L}{\partial m}$$

$$\frac{\partial L}{\partial b} = 0, \text{ from this.}$$

$$b = \bar{y} - m\bar{x} \quad \text{where } \bar{y} \rightarrow y\text{-mean} \\ \bar{x} \rightarrow x\text{-mean.}$$

Now.

$$\frac{\partial L}{\partial m} = \sum_{i=1}^n (y_i - mx_i - \bar{y} + m\bar{x})^2 + \lambda m^2 \\ (\text{replace } b \text{ by } \bar{y} - m\bar{x})$$

$$\frac{\partial L}{\partial m} = 2 \sum_{i=1}^n (y_i - mx_i - \bar{y} + m\bar{x})(-x_i + \bar{x}) + 2\lambda m = 0 \\ = -2 \sum_{i=1}^n (y_i - \bar{y} - mx_i + m\bar{x})(x_i - \bar{x}) + 2\lambda m = 0$$

$$= \lambda m - \sum_{i=1}^n [(y_i - \bar{y}) - m(x_i - \bar{x})](x_i - \bar{x}) = 0$$

$$= \lambda m - \sum_{i=1}^n (y_i - \bar{y})(x_i - \bar{x}) - m(x_i - \bar{x})^2 = 0$$

$$= \lambda m - \sum_{i=1}^n (y_i - \bar{y})(x_i - \bar{x}) + m \sum_{i=1}^n (x_i - \bar{x})^2 = 0$$

$$= \lambda m + m \sum_{i=1}^n (x_i - \bar{x})^2 = \sum_{i=1}^n (y_i - \bar{y})(x_i - \bar{x})$$

$$\therefore m = \frac{\sum_{i=1}^n (y_i - \bar{y})(x_i - \bar{x})}{\sum_{i=1}^n (x_i - \bar{x})^2 + \lambda}$$



$$m = \frac{\sum_{i=1}^n (y_i - \bar{y})(x_i - \bar{x})}{\sum_{i=1}^n (x_i - \bar{x})^2 + \lambda}$$

$$b = \bar{y} - m\bar{x}$$

loss function of N-Dim. Data.

$$L = (xw - y)^T (xw - y) + \lambda \|w\|^2$$

$$W = (x^T x)^{-1} x^T y$$

$$W = (x^T x)^{-1}$$

$$W = (x^T x + \lambda I)^{-1} x^T y$$

code for Ridge Reg. for N-dim. dataset.

class MyRidge:

def __init__(self, alpha=0.1):

self.alpha = alpha

self.coef_ = None

self.intercept_ = None

def fit(self, X_train, y_train):

X_train = np.insert(X_train, 0, 1, axis=1)

I = np.identity(X_train.shape[1])

I[0][0] = 0

result = np.linalg.inv(np.dot(X_train.T, X_train) + np.dot(self.alpha * I))

dot(X_train.T).dot(Y_train)



self.coef_ = result[1:]

self.intercept_ = result[0].

def predict(self, X-test):

return np.dot(self.coef_, X-test)

+ self.intercept_

Our class is created Now let, test it.

reg = MyRidge()

reg.fit(X-train, Y-train)

Y-pred = reg.predict(X-test)

r2-score(Y-test, Y-pred)

reg.coef_

reg.intercept_



Ridge Regression using Gradient Descent

Gradient Descent..

Loss function of linear reg.

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

Loss function in vector form. for linear reg.

$$L = (xw - y)^T (xw - y)$$

Ridge Reg. Loss fun -

$$L = \sum_{i=1}^n (y_i - \hat{y}_i)^2 + \lambda m^2$$

LF in Vector form:

$$L = (xw - y)^T (xw - y) + \lambda \|w\|^2$$

$$L = (xw - y)^T (xw - y) + \lambda \cdot w^T w$$

where,

$$x = \begin{bmatrix} 1 & x_{11} & x_{12} & x_{13} & \dots & x_{1n} \\ 1 & x_{21} & x_{22} & x_{23} & \dots & x_{2n} \\ 1 & x_{m1} & x_{m2} & x_{m3} & \dots & x_{mn} \end{bmatrix}, \quad y = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_m \end{bmatrix}$$

$$w = \begin{bmatrix} w_0 \\ w_1 \\ w_2 \\ \vdots \\ w_n \end{bmatrix} \rightarrow (n+1)$$



$$L = (xw - y)^T (xw - y) + \lambda w^T w$$

mul by $\frac{1}{2}$

$$L = \frac{1}{2} (xw - y)^T (xw - y) + \frac{1}{2} \lambda w^T w$$

$$L = \frac{1}{2} [w^T x^T x w - w^T x^T y - y^T w x + y^T y] + \frac{1}{2} \lambda w^T w$$

$$L = \frac{1}{2} [w^T x^T x w - 2w^T x^T y + y^T y] + \frac{1}{2} \lambda w^T w$$

$$\frac{\partial L}{\partial w} = \frac{1}{2} [2x^T x w - 2x^T y] + \frac{1}{2} 2x^T w$$

$$\frac{\partial L}{\partial w} = [x^T x w - x^T y + \lambda w] = \left(\frac{\partial L}{\partial w} \right)_{ii}$$



Key Points about Ridge Regression

1) How the coefficients get affected (β)

If λ value keep 0 then act as linear. neg

If λ value is increase from 0 to greater number all the coefficients will get shrinked get decreased towards zero but never turn zero.

2) Higher values are impacted more.

While increasing λ towards greater value (∞) the higher values will decreased faster than lower lesser values. but not closed to zero higher value = higher impact.

3) Bias Variance Trade off.

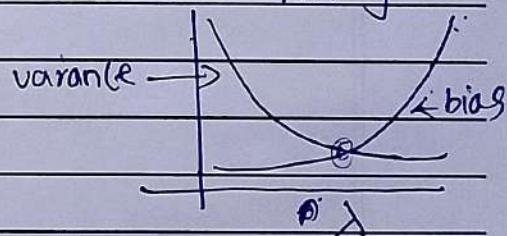
Bias & variance both affect by λ value.

if $\lambda \Rightarrow 0$ then bias ↓ overfitting.

$\lambda \Rightarrow$ greater then bias ↑ underfitting.

λ value should be

= low bias & variance



4) Impact on One Loss Function

as λ value increase the Loss function moves towards zero. but not become zero.

5) Why called Ridge

- Hard constraint Ridge regression is

- λ soln come on Ridge of circle / perimeter of circle for best results. best soln for MSE & λ value.



• LASSO REGRESSION (L1 Regularization)

Lasso regression is a regularization technique that penalizes high value, correlated coefficients. It introduces a regularization term (also called penalty term) into one (SS E) loss function. This penalty term is absolute value of the sum of coefficients.

Controlled in turn by hyperparameter lambda(λ), it reduces select feature weights to zero.

Lasso regression removes multicollinear feature from the model altogether.
(auto feature selection is applied)

Lasso Reg. Loss fun.

$$L = \sum_{i=1}^n (y_i - \hat{y}_i)^2 + [\lambda \cdot \|w\|_1] \rightarrow \text{Penalty term}$$

$$L = \text{MSE} + \lambda \|w\|_1$$

$$\text{where, } \|w\|_1 = [|w_1| + |w_2| + |w_3| + \dots + |w_n|]$$

→ import lasso for code from sklearn.

```
from sklearn.linear_model import Lasso
```

```
Lass_model = Lasso(alpha= 0.5)
```

```
Lass_model.fit(x_train, y_train)
```

```
Lass_model.predict(x_test)
```



Key point int Lasso Regression

- 1) How are coefficient affected?
 λ value \uparrow Then coefficient get shrinked.
unimportant values become zero
when λ value \uparrow
higher impact for decrementation of higher important values (coefficient).
- 2) Higher coefficient are affected more:
as λ value \uparrow Higher coefficient decreased fast & at last become zero.
- 3) Impact on Bias & Variance:
 λ value \uparrow :
overfitting \downarrow & Bias \uparrow & variance \downarrow
variance \uparrow & variance \downarrow .
- 4) Effect of Regularization on loss function
 λ value \uparrow increases
loss function above (curve) moves towards zero.
at the end it will reach zero.
- 5) Sparsity is created by lasso reg.



Lasso reg. & loss function.

$$L = \sum_{i=1}^n (y_i - \hat{y}_i)^2 + \lambda |m|$$

$$(b = \bar{y} - m\bar{x})$$

$$L = \sum_{i=1}^n (y_i - mx_i - \bar{y} + m\bar{x})^2 + 2\lambda|m|$$

$$\frac{\partial}{\partial m} = \sum_{i=1}^n (y_i - mx_i - \bar{y} + m\bar{x})^2 + 2\lambda m$$

$$= 2 \sum (y_i - mx_i - \bar{y} + m\bar{x})(-x_i + \bar{x}) + 2\lambda m$$

$$= -2 \sum [(y_i - \bar{y}) - m(x_i - \bar{x})] (x_i - \bar{x}) + 2\lambda m = 0$$

$$= - \sum [(y_i - \bar{y})(x_i - \bar{x}) - m(x_i - \bar{x})^2] + 2\lambda m = 0$$

$$= - \sum [(y_i - \bar{y})(x_i - \bar{x}) - m(x_i - \bar{x})^2] + \lambda m = 0$$

$$= - \sum (y_i - \bar{y})(x_i - \bar{x}) + m \sum (x_i - \bar{x})^2 + \lambda m = 0$$

$$\bullet m \sum (x_i - \bar{x})^2 = \sum (y_i - \bar{y})(x_i - \bar{x}) - \lambda$$

$$m = \frac{\sum (y_i - \bar{y})(x_i - \bar{x}) - \lambda}{\sum (x_i - \bar{x})^2}$$

for $m > 0$

$$m > 0$$

$$m = \frac{\sum (y_i - \bar{y})(x_i - \bar{x}) - \lambda}{\sum (x_i - \bar{x})^2}$$

for $m = 0$

$$m = \frac{\sum (y_i - \bar{y})(x_i - \bar{x})}{\sum (x_i - \bar{x})^2}$$

for $m < 0$

$$m = \frac{\sum (y_i - \bar{y})(x_i - \bar{x}) + \lambda}{\sum (x_i - \bar{x})^2}$$



Now for sparsity let's take an example.
assume - $(y_i - \bar{y}) = y$
 $(x_i - \bar{x}) = x$

for $m > 0$:

$$m = \frac{4x - \lambda}{x^2}$$

$$\text{supposed } 4x = 100 \\ \therefore x^2 = 50$$

$$m = \frac{100 - \lambda}{50}$$

$$\text{for } \lambda = 0 \quad m = \frac{100 - 0}{50} = \frac{100}{50} = 2. \\ \boxed{m=2}$$

$$\text{for } \lambda = 10 \quad m = \frac{100 - 10}{50} = \frac{90}{50} = \boxed{\frac{9}{5}}$$

$$\text{for } \lambda = 50 \quad m = \frac{100 - 50}{50} = \frac{50}{50} = \boxed{1},$$

$$\text{for } \lambda = 100 \quad m = \frac{100 - 100}{50} = \frac{0}{50} = \boxed{0},$$

when $\boxed{4x = \lambda}$ then $\boxed{m=0}$

when λ values goes in -ve then we
fall $m < 0$ eqn. which is $4x$.

$$m = \frac{4x + \lambda}{x^2} \text{ which again}$$

increase the m value to +ve greater so
algo stops at ~~that~~ p. 0



Now, for
 $m < 0$.

$$m = \frac{\sum (y_i - \bar{y})(x_i - \bar{x}) + \lambda}{\sum (x_i - \bar{x})^2}$$

$$m = \frac{x_4 + \lambda}{x^2} \quad \left. \begin{array}{l} x_4 = -100 \\ x^2 = 50 \end{array} \right\} \quad m = -2$$

$$\text{when } \lambda = 0, \frac{-100 + 0}{50} = -2 \quad \boxed{m = -2}$$

$$\lambda = 50, \frac{-100 + 50}{50} = -1 \quad \boxed{m = -1}$$

$$\lambda = 100, \frac{-100 + 100}{50} = 0 \quad \boxed{m = 0}$$

$$\lambda = 150, \frac{-100 + 150}{50} = \frac{50}{50} = 1 \quad \boxed{m = 1}$$

for $\lambda = 150$, m become 1 so, $m > 0$.
or $m \neq 0$.

so, we will use $m > 0$ formula.

$$m = \frac{x_4 - \lambda}{x^2} = \frac{-100 - 150}{50} = \frac{-250}{50} = \boxed{-5},$$

$\boxed{m = -5}$, so we can't go again for large value than -2 so algo stops at zero 0

2) why Ridge Reg not create sparsity.
Ridge Reg. formula.

$m = \frac{\sum (y_i - \bar{y})(x_i - \bar{x})}{(x_i - \bar{x})^2 + \lambda}$ in ridge reg.
 λ is at denominator
 so if it can't be zero but in lasso
 λ is in numerator so it can be
 zero cause denominator can't be
 zero //.



* ElasticNet Regression

- Elastic net regression essentially combines both ridge and lasso regression by inserting both the L1 & L2 penalty terms into the SSE loss function. Elastic net inset both penalty values into the cost function equation. Elastic net handles ~~not~~ addresses multicollinearity while also enabling feature selection.

Loss function

$$L = \sum (y_i - \hat{y}_i)^2 + a \|w\|_2^2 + b \|w\|_1$$

where a & b are independent.

$a = 0.5$ & $b = 0.5$. both have equal weightage.

$$\lambda = a+b, \text{ L1-ratio} = \frac{b}{a+b}$$

When

λ & L1-ratio are hyperparameters.

~~L1-ratio = 0.9. \Rightarrow 90% ridge & 10% lasso~~

~~90% ridge & 10% lasso~~

~~L1-ratio = 0.4. \Rightarrow 40% ridge & 60% lasso~~

$L1\text{-ratio} = 0.9 \Rightarrow 90\% \text{ lasso} \& 10\% \text{ ridge}$

$L1\text{-ratio} = 0.4 \Rightarrow 40\% \text{ lasso} \& 60\% \text{ ridge}$

implement elasticnet using sklearn.

From sklearn.linear_model import ElasticNet

reg = ElasticNet(alpha=0.005, L1-ratio=0.9)

reg.fit(X-train, Y-train)

reg.predict(X-test),

monday

LOGISTIC REGRESSION

SARASWATI Education Society's

SARASWATI College of Engineering

PAGE NO.: 79

DATE: 08/07/2024

Prerequisite for applying logistic regression.

- Data must be linearly separable

Perceptron Trick

$$\text{LOG Reg eqn } Ax + By + C = 0$$

$$Ax_1 + Bx_2 + C = 0$$

$$Ax_1 + Bx_2 + Cx_3 + d = 0$$

A, B, C

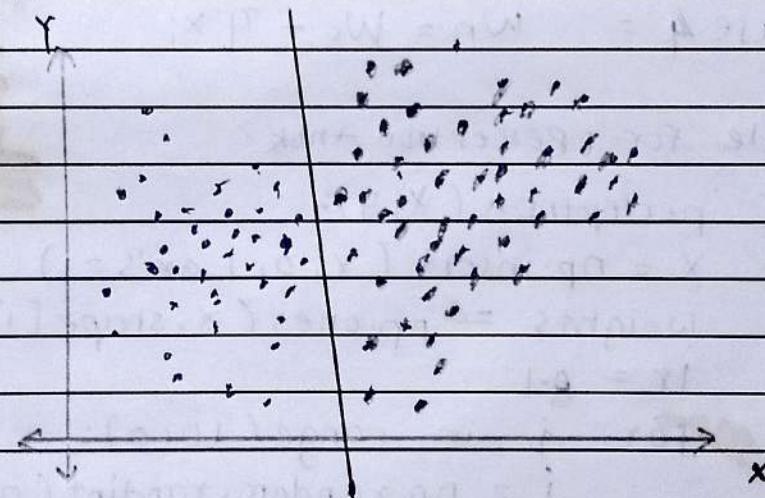
How to label regions.

If. eqn of line is to point is on

$Ax_1 + Bx_2 + C > 0$ = +v region

$Ax_1 + Bx_2 + C < 0$ = -v region

$Ax_1 + Bx_2 + C = 0$ = on the line



Alg:

$$Ax + By + C = 0$$

$$w_0 + w_1 x_1 + w_2 x_2 = 0, \quad w_0 = C, w_1 = A, w_2 = B$$

$$w_0 x_0 + w_1 x_1 + w_2 x_2 = 0$$

$$\sum_{i=0}^2 w_i x_i = 0 \quad \{w_0, w_1, w_2\} \begin{bmatrix} x_0 \\ x_1 \\ x_2 \end{bmatrix}$$



Algorithm. ①

if $x_i \in N$ and $\sum w_i x_i > 0$

$$w_n = w_0 - \eta x_i$$

if $x_i \in P$ and $\sum w_i x_i < 0$

$$w_n = w_0 + \eta x_i$$

~~where~~

where, $N = \text{Negative Point}$

$P = \text{Positive Point}$

Algorithm ②

$$w_n = w_0 + \eta (y_i - \hat{y}_i) x_i$$

$$\text{case 1} = w_n = w_0 \quad | \quad 0 \quad 1 \quad 1$$

$$\text{case 2} = w_n = w_0 \quad | \quad 0 \quad 0 \quad 0$$

$$\text{case 3} = w_n = w_0 + \eta x_i \quad | \quad 1 \quad 1 \quad 0$$

$$\text{case 4} = w_n = w_0 - \eta x_i \quad | \quad -1 \quad 0 \quad 1$$

code for Perceptron trick.

def perceptron(x, y):

$x = np.insert(x, 0, 1, axis=1)$

$weights = np.ones(x.shape[1])$

$lr = 0.1$

for i in range(1000):

$j = np.random.randint(0, 100)$

$y_hat = step(np.dot(x[j], weights))$

$weights = weights + lr * (y_{rj}) - y_hat) * x[j]$

return weights[0], weights[1:]



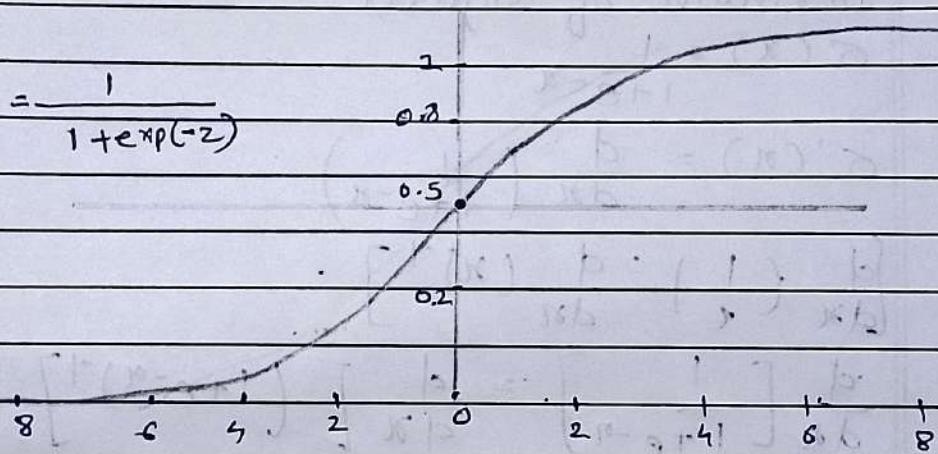
Logistic Regression.

For logistic reg. we have to make changes in our Algo. we can't make $(y_i - \hat{y}) = 0$ if $(y_i - \hat{y}) = 0$ then $w_n = w_0$ so eqn line will never change.

$$w_n = w_0 + \eta(y_i - \hat{y}_i)x_i$$

$$(y_i - \hat{y}_i) \neq 0$$

The Sigmoid Function



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

if $\sigma(z)$ is +ve then $0.5 > \sigma(z) \Rightarrow 1$

∴ if $\sigma(z)$ is -ve then $0.5 < \sigma(z) \Rightarrow 0$

$$m = -\frac{A}{B} - \text{coef}[0] - \text{coef}[1]$$

$$c = -\frac{C}{B} - \text{intercept} - \text{coef}[1]$$



$$L = - \sum_{i=1}^n -y_i \log(y_i) - (1-y_i) \log(1-y_i)$$

log loss error function / binary cross entropy

$$L = -\frac{1}{n} \sum_{i=1}^n y_i \log(y_i) + (1-y_i) \log(1-y_i)$$

Derivative of sigmoid:

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

$$\sigma'(x) = \frac{d}{dx} \left(\frac{1}{1+e^{-x}} \right)$$

$$\left[\frac{d}{dx} \left(\frac{1}{x} \right) = \frac{d}{dx} (x)^{-1} \right] \dots$$

$$\frac{d}{dx} \left[\frac{1}{1+e^{-x}} \right] = \frac{d}{dx} [(1+e^{-x})^{-1}]$$

$$= -\frac{1}{(1+e^{-x})^2} \cdot \frac{d}{dx} (1+e^{-x}) \quad \dots \text{(chain rule)}$$

$$= -\frac{1}{(1+e^{-x})^2} \cdot \frac{d}{dx} (e^{-x})$$

$$= -\frac{e^{-x}}{(1+e^{-x})^2} \cdot \frac{d}{dx} (-x) \Rightarrow \frac{e^{-x}}{(1+e^{-x})^2}$$

$$\therefore \frac{1 \cdot e^{-x}}{(1+e^{-x})(1+e^{-x})} = \frac{1}{1+e^{-x}} \times \frac{e^{-x}}{1+e^{-x}}$$

$$= \sigma(x) \left[\frac{e^{-x}}{1+e^{-x}} \right] \dots \left[\frac{1}{1+e^{-x}} = \sigma(x) \right]$$

$$= \sigma(x) \left[\frac{1+e^{-x}-1}{1+e^{-x}} \right] = \sigma(x) \left[\frac{1+e^{-x}}{1+e^{-x}} - \frac{1}{1+e^{-x}} \right]$$

$$\therefore \sigma(x) \left[1 - \sigma(x) \right] \Rightarrow$$

$$= [\sigma'(x) = \sigma(x) [1 - \sigma(x)]]$$



Gradient Descent for Logistic Reg.

Rows = m , Columns = n

1	2	3	\dots	n	y
x_{11}	x_{12}	x_{13}	\dots	x_{1n}	y_1

x_{21}	x_{22}	x_{23}	\dots	x_{2n}	y_2
----------	----------	----------	---------	----------	-------

x_{31}	x_{32}	x_{33}	\dots	x_{3n}	y_3
----------	----------	----------	---------	----------	-------

\vdots	\vdots	\vdots	\vdots	\vdots	\vdots
----------	----------	----------	----------	----------	----------

x_{m1}	x_{m2}	x_{m3}	\dots	x_{mn}	y_m
----------	----------	----------	---------	----------	-------

Coeff = $\{w_1, w_2, w_3, \dots, w_n\} \rightarrow (n+1)$

$$\sigma(w_0 + w_1 x_{11} + w_2 x_{12} + w_3 x_{13} + \dots + w_n x_{1n} + w_0) = y_1$$

$$\sigma(w_0 + w_1 x_{21} + w_2 x_{22} + w_3 x_{23} + \dots + w_n x_{2n} + w_0) = y_2$$

$$\hat{Y} = \begin{bmatrix} \hat{y}_1 \\ \hat{y}_2 \\ \vdots \\ \hat{y}_m \end{bmatrix} = \begin{bmatrix} \sigma(w_0 + w_1 x_{11} + w_2 x_{12} + \dots + w_n x_{1n}) \\ \sigma(w_0 + w_1 x_{21} + w_2 x_{22} + \dots + w_n x_{2n}) \\ \vdots \\ \sigma(w_0 + w_1 x_{m1} + w_2 x_{m2} + \dots + w_n x_{mn}) \end{bmatrix}$$

$$\hat{Y} = \begin{bmatrix} w_0 + w_1 x_{11} + w_2 x_{12} + \dots + w_n x_{1n} \\ w_0 + w_1 x_{21} + w_2 x_{22} + \dots + w_n x_{2n} \\ \vdots \\ w_0 + w_1 x_{m1} + w_2 x_{m2} + \dots + w_n x_{mn} \end{bmatrix}$$

$$\hat{Y} = \sigma \left(\begin{bmatrix} 1 & x_{11} & x_{12} & \dots & x_{1n} \\ 1 & x_{21} & x_{22} & \dots & x_{2n} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 1 & x_{m1} & x_{m2} & \dots & x_{mn} \end{bmatrix} \begin{bmatrix} w_0 \\ w_1 \\ w_2 \\ \vdots \\ w_n \end{bmatrix} \right)$$

$\hat{Y} = \sigma(XW)$



Loss fun.

$$L = -\frac{1}{m} \sum_{i=1}^m y_i \log(\hat{y}_i) + (1-y_i) \log(1-\hat{y}_i)$$

$$L = -\frac{1}{m} \left[\underbrace{\sum_{i=1}^m y_i \log(\hat{y}_i)}_A + \underbrace{\sum_{i=1}^m (1-y_i) \log(1-\hat{y}_i)}_B \right]$$

take, A:

$$\sum_{i=1}^m y_i \log(\hat{y}_i) < y_1 \log \hat{y}_1 + y_2 \log \hat{y}_2 + y_3 \log \hat{y}_3 + \dots + y_m \log \hat{y}_m$$

we also can write it as.

$$\begin{bmatrix} y_1 & y_2 & y_3 & \dots & y_m \end{bmatrix} \begin{bmatrix} \log \hat{y}_1 \\ \log \hat{y}_2 \\ \vdots \\ \log \hat{y}_m \end{bmatrix}$$

$$\begin{bmatrix} y_1 & y_2 & y_3 & \dots & y_m \end{bmatrix} \log \begin{bmatrix} \hat{y}_1 \\ \hat{y}_2 \\ \vdots \\ \hat{y}_m \end{bmatrix} = \hat{Y}$$

$$\boxed{y \log \hat{y} = y \log(\sigma(xw))} \quad \dots \hat{y} = \underline{\sigma(xw)}$$

$$\boxed{L = -\frac{1}{m} \left[y \log \hat{y} + (1-y) \log(1-\hat{y}) \right]}$$

$$\text{where } \hat{y} = \sigma(xw)$$



Loss function in matrix form.

$$L = -\frac{1}{m} \left[y \log(\sigma(wx)) + (1-y) \log(1-\sigma(wx)) \right]$$

• Applying Gradient descent

$$\frac{\partial L}{\partial w} = -\frac{1}{m} \left[y \log \hat{y} + (1-y) \log(1-\hat{y}) \right]$$

$$\frac{d}{dw} y \log \hat{y} \Rightarrow y \frac{d}{dw} \log \hat{y} \Rightarrow y \frac{d}{dL} (\hat{y})$$

$$\Rightarrow y \frac{d}{dL} \sigma(wx) \Rightarrow y \sigma(wx) [1 - \sigma(wx)] \frac{d}{dw} (wx)$$

$$\therefore \frac{d}{dw} y \log(1-\hat{y}) \Rightarrow \boxed{y(1-\hat{y})x} \quad \hat{y} = \sigma(wx)$$

$$\frac{d}{dw} (1-y) \log(1-\hat{y}) \Rightarrow (1-y) \frac{d}{dw} \log(1-\hat{y})$$

$$\Rightarrow \frac{(1-y)}{(1-\hat{y})} \cdot \frac{d}{dw} [1 - \hat{y}]$$

$$\therefore -\frac{(1-y)}{(1-\hat{y})} \cdot \frac{d}{dw} \sigma(wx) \therefore$$

$$\therefore -\frac{(1-y)}{(1-\hat{y})} [\sigma(wx) [1 - \sigma(wx)] \frac{d}{dw} (wx)]$$

$$\therefore -\frac{(1-y)}{(1-\hat{y})} \cdot \hat{y} (1-\hat{y}) x$$

$$\therefore \boxed{-y(1-y)x}$$

$$\frac{\partial L}{\partial w} = -\frac{1}{m} \left[y(1-\hat{y})x - \hat{y}(1-y)x \right]$$



$$= -\frac{1}{m} [4(1-4) - \hat{y}(1-4)]x$$

$$= -\frac{1}{m} [4 - 4\hat{y} - \hat{y} + \hat{y}\hat{y}]x$$

$$\Delta L = -\frac{1}{m} (4 - \hat{y})x$$

according to gradient descent-

$$w = w + \eta \left[\frac{1}{m} (4 - \hat{y})x \right] / /$$

where,

$$w = \begin{bmatrix} w_0 \\ w_1 \\ \vdots \\ w_n \end{bmatrix}, \quad x = \begin{bmatrix} 1 & x_{11} & x_{12} & \dots & x_{1n} \\ 1 & x_{21} & x_{22} & \dots & x_{2n} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 1 & x_{m1} & x_{m2} & \dots & x_{mn} \end{bmatrix}, \quad y = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_m \end{bmatrix}$$

(m+1), 1 m, n+1 m, 1

$$w = w + \frac{\eta}{m} (4 - \hat{y})x$$

↓ ↓ ↓ ↓
(n+1, 1) (m, 1) m, (n+1) (m, 1)

(1, m+1)

Accuracy and Confusion Matrix

23



SARASWATI Education Society's

SARASWATI College of Engineering

PAGE NO.: 87

DATE: 10/7/24

CLASSIFICATION METRICS

Accuracy

$$\text{Accuracy} = \frac{\text{No. of Correct Prediction}}{\text{Total No. of Predictions}}$$

implement

```
from sklearn.metrics import accuracy_score  
accuracy_score(y_true, y_pred)
```

Problem with Accuracy.

Accuracy never tells the Precision & Recall Trade-off.

- Imbalanced datasets - accuracy is misleading with imbalance dataset.
- Equal weightage for Errors. → accuracy treats all errors equally like False negative & False positive are not differentiable in accuracy score.

Confusion Matrix

		Prediction		P	
		True P	False	T P	F N.
Actual	True	Positive	Negative	T P	F N.
	False	Positive	Negative		
Predicted	True	True	False	F P	T N
	False	False	True		

I = Yes, TP = True Positive

O = No, TN = True Negative

FP = False Positive

FN = False Negative



accuracy from confusion matrix

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

Type I error

FP is also known as Type I error.

Type II error

FN is known as Type II error.

implement.

```
from sklearn.metrics import confusion_matrix  
confusion_matrix(y_true, y_pred)
```

Precision

Measures the proportion of true positive prediction among all positive

- Proportion of predicted Positive
- Not More Harmful than Recall

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

Recall

Measures the proportion of correctly predicted positive observation to all observations in actual class.

Cost of FN is high

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



F1 score

- It is used to penalize the lower value. It calculates the Harmonic mean for predictions.

$$\text{F1 score} = \frac{2 \times \text{Precision} \times \text{Recall}}{\text{Precision} + \text{Recall}}$$

* Multiclass Classification *

Here we calculate all classes precision.

• Precision

(1) Macro precision : for balanced classes take avg of all classes precision

(2) Weighted Precision : for unbalanced / weighted biased classes.

all precision values are multiplied by their class weight

$$\text{Class w.t.} = \frac{\text{each class actual value}}{\text{total items}}$$

for 3 portions.

$$\text{macro} = \frac{P_1 + P_2 + P_3}{3}$$

$$\text{Weighted} = \text{Class w.t.}_1 \times P_1 + \\ \text{Class w.t.}_2 \times P_2 + \text{Class w.t.}_3 \times P_3$$

• Recall

$$\text{Total Predicted class value} \\ \text{Total actual class value}$$

+ macro

+ weighted



F1-score

- F1-score is similar for all like
- we have to calculate F1-score individually for each class.

$$\text{Rq. Precision} = P_1, P_2, P_3$$

$$\text{Recall} = R_1, R_2, R_3.$$

$$F1\text{-score}_1 = \frac{2 P_1 \times R_1}{P_1 + R_1}$$

$$F1\text{-score}_2 = \frac{2 P_2 \times R_2}{P_2 + R_2}$$

$$F1\text{-score}_3 = \frac{2 P_3 \times R_3}{P_3 + R_3}$$

implementation using sklearn -

Macro.

```
from sklearn.metrics import precision_score,  
                           recall_score,  
                           f1_score.
```

```
precision_score(y-test, y-pred, average='macro')
```

same for Recall & f1-score.

Weighted

```
recall_score(y-test, y-pred, average='weighted')
```

same for Precision-score & f1-score.

All in One

```
from sklearn.metrics import classification_report  
print(classification_report(y-test, y-pred))
```



SARASWATI College of Engineering

Softmax Regression

Softmax regression (or multinomial logistic regression) is a generalization of logistic regression to the case where we want to handle multiple classes.

Softmax regression allows us to handle $y(1) \in \{1, 2, \dots, K\}$ where K is the number of classes.

$$\sigma(z)_i = \frac{e^{z_i}}{\sum_{j=1}^k e^{z_j}}$$
 where, k is class $\{1, 2, \dots, k\}$

e.g. Softmax can apply for 1st class of 3 classes

$$\sigma(z)_1 = \frac{e^{z_1}}{e^{z_1} + e^{z_2} + e^{z_3}}$$

$$\sigma(z)_i = \frac{e^{z_i}}{\sum_{j=1}^k e^{z_j}}$$
 $0 < \sigma < 1$

(classes = 0, 1, 2 etc.)
Yes, No, Opt etc.

Implementation

```
from sklearn.linear_model import LogisticRegression  
clf = LogisticRegression(multi_class='multinomial')  
clf.fit(X_train, Y_train)  
clf.predict(X_test)
```

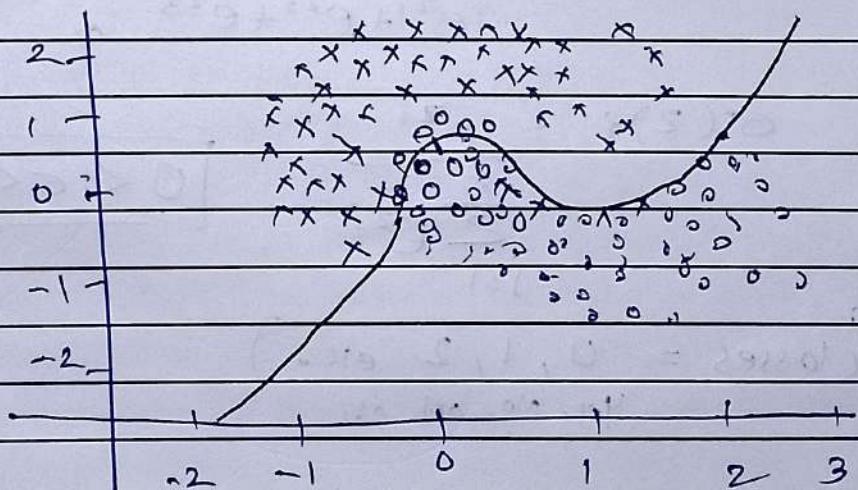


Polynomial Logistic Regression

Polynomial features are created by transforming the original input features into a new set of feature that include original features & their polynomial combination up to specific degree.

The transformation allows logistic regression to capture non-linear relationships between one input variables & the target variable.

Logistic regression with polynomial features is technique used for handling non-linear relationships in data. The idea is to transform input features into high-degree polynomial, to capture more complex relationship between variables.



|| DECISION TREE ||

SARASWATI Education Society's

SARASWATI College of Engineering

PAGE NO.: 93

DATE: 11/7/2024

A decision tree is a decision support hierarchical model that uses a tree-like model of decisions and their possible consequences, chance events' outcomes, resource costs & utility. It is one way to display an algorithm that only contains conditional control statements.

Example

Gender Occupation Suggestion

F Student PUBG

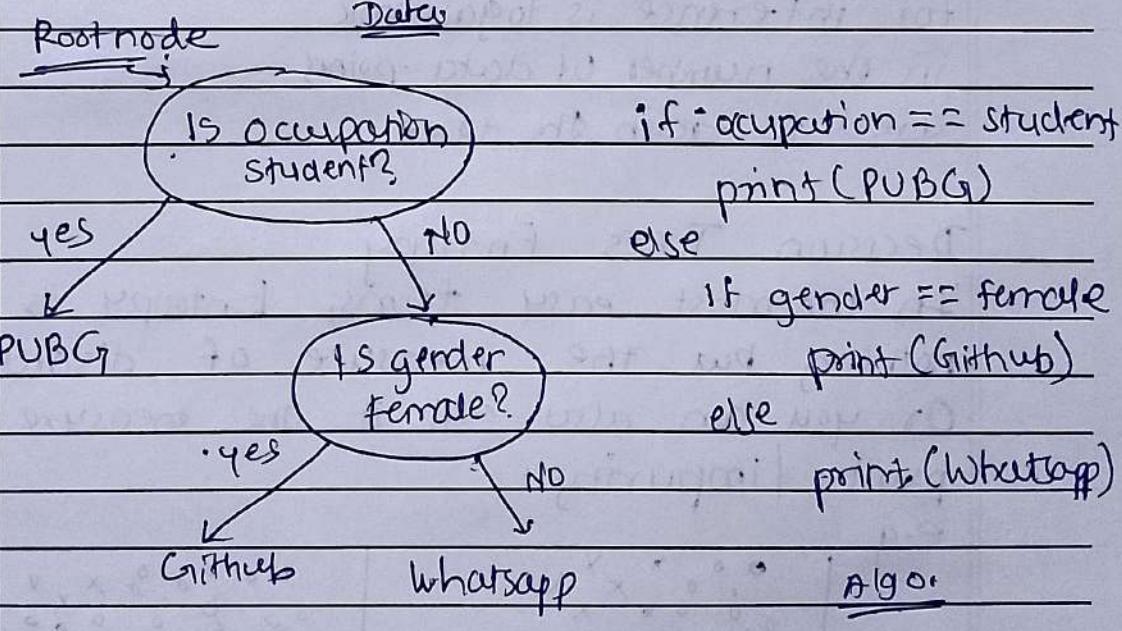
F Programmer Github

M Programmer Whatsapp

F Programmer Github

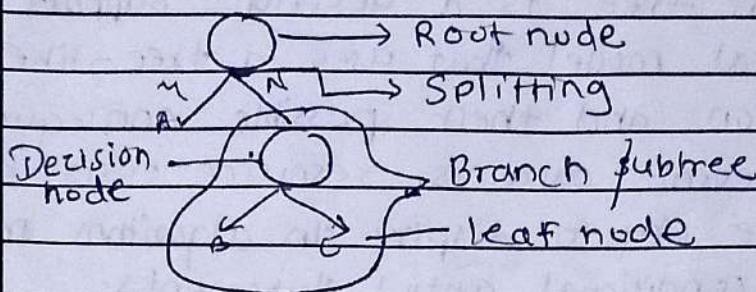
M Student PUBG

M Student PUBG





Terminology



Advantages

- Intuitive & easy to understand
- minimal data preparation is required

Disadvantages

overfitting

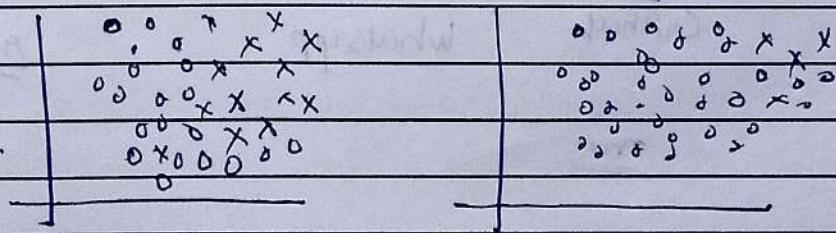
prior to errors for imbalanced datasets

The cost of using the tree for inference is logarithmic in the number of data points used to train the tree

Decision Trees Entropy

- In the most easy terms, Entropy is nothing but the measure of disorder. Or you can also call it one measure of purity / impurity.

e.g



High entropy

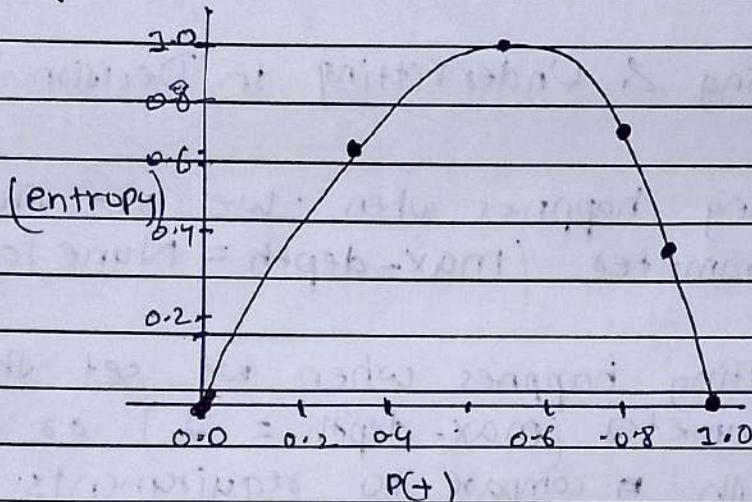
Low entropy



$$E(S) = \sum_{i=1}^c -P_i \log_2 P_i$$

where, P_i frequentist probability of an element (class 'i' in our data)

$$E(D) = -P_{yes} \log_2 (P_{yes}) - P_{no} \log_2 (P_{no})$$



* Information Gain *

Information Gain, is a metric used to train Decision Trees. Specifically ; this metric measures the quality of a split.

The information gain is based on the decrease in entropy after a data-set is split on an attribute. Constructing a decision tree is all about finding attribute that returns the highest information gain.

$$\text{Info Gain} = E(\text{Parent}) - \left\{ \begin{array}{l} \text{Weighted Average} \\ * E[\text{Children}] \end{array} \right\}$$



Weighted Entropy = $P * E$

* Gini Impurity *

$$G.I = 1 - (P_y^2 + P_N^2)$$

* Overfitting & Underfitting in Decision Trees.

Overfitting happens when we set the hyperparameter (max_depth = None) or high

Underfitting happens when we set the hyperparameter (max_depth = 1) or very less with respect to requirements.

* Regression Trees *

Used when non-linear relationship in dataset is present

For Decision Tree visualization

Use library called dtreeviz.

install

pip install dtreeviz

import

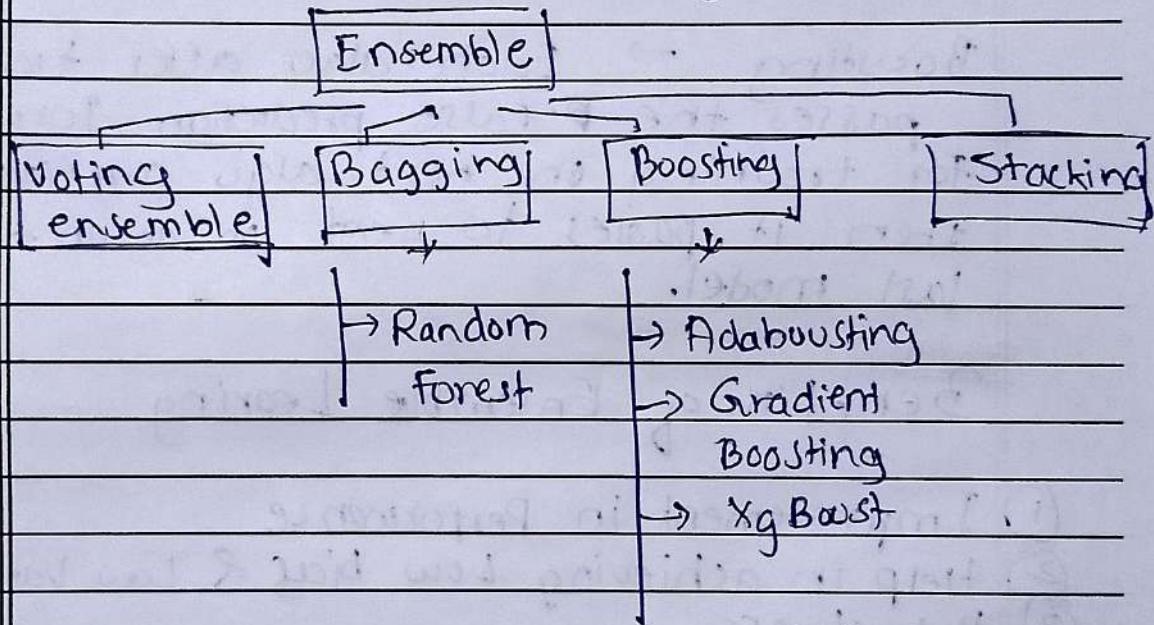
from dtreeviz.trees import *



Ensemble methods use multiple learning algorithms to obtain better predictive performance than could be obtained from any of the constituent learning algo.

Ensemble learning trains two or more ml algorithms to a specific classification or regression task.

Type of Ensemble Learning



1. Voting Ensemble

- Base models are different Algo's and same Data.

2. Stacking

- Same as Voting but additionally we give weights to every model's output and another last model is trained with results + additional weights of



remaining models.

Bagging \rightarrow Bootstrapped Aggregation.

- Based models are same but
the data is different for each model.

Random Forest \rightarrow it is special case
of bagging.

Boosting \rightarrow Each algo after training
passes the false prediction forward
for training the next algo correctly
then it passes to next on and till
last model.

Benefits of Ensemble Learning

- (1) Improvement in Performance
- (2) Help in achieving Low bias & Low Variance
- (3) Robustness



SARASWATI College of Engineering

Voting Ensemble

Ensemble voting is a popular technique in ensemble learning where the final prediction is made by combining the predictions of multiple individual models.

It can be used for classification & regression problems.

Voting Classifier

The concept is to build a single model that learns from various models and predicts output based on their aggregate majority and predicts output based on their of votes for each output class, rather than building separate specialized models & determining the accuracy for each of them.

2 different types of voting classifier

• Hard Voting: The predicted output class is a bi-class with the highest majority of votes,

• Soft Voting: The average probabilities of the classes determine which one will be the final prediction.



To import voting.

>> from sklearn.ensemble import VotingClassifier

Hard voting-

>> vc = VotingClassifier(estimators = [('lr', LogisticRegression()), ('rf', RandomForestClassifier()), ('knn', KNeighborsClassifier()), Voting = "hard")]

Soft Voting

vc = VotingClassifier(estimators = estimators, Voting = "soft")

Voting Regressor

It combines or 'ensembles' multiple regression models & overperforms the individual models present as its estimators. ML models as $m_1, m_2, m_3 \dots m_n$. each model will produce a prediction P_n for given data I . If we pass thru Voting Regressor then final prediction will be Average of all Predictions.

$$\text{Pruting} = \frac{1}{n} \sum_{n=1}^N P_n \Rightarrow \text{Simple Avg}$$

$$\text{Pruting} = \sum_{n=1}^N (w_{tn} \cdot P_n) \Rightarrow \text{Weighted Avg}$$

where w_{tn} is custom w.t.

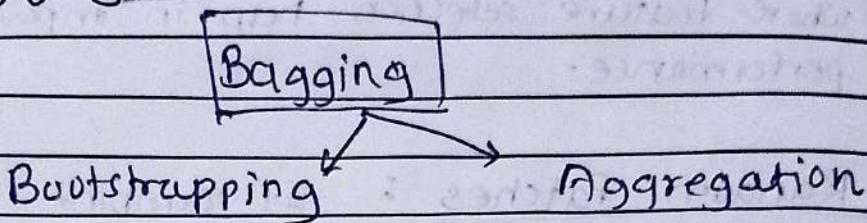
Implementation:

>> from sklearn.ensemble import VotingRegressor

>> vr = VotingRegressor(estimators)



Bagging.



Bagging also known as bootstrap aggregation, is commonly used to reduce variance within a noisy data set.

In bagging, a random sample of data in a training set is selected with replacement. After generating several data samples, these weak models are then trained independently. Depending on the type of task - regression or classification,

There are 3 different ways for performing Bagging

1. Pasting
2. Random Subspace
3. Random Patches

1. Pasting : It is sampling without replacement. each subset is created by randomly sampling the training data without duplicating instances.

use case: When dataset is large enough, & risk of overfitting due to duplicated instances in bootstrapped samples is a concern.

2. Random Subspace : It involves training each base model on a random subset of the features rather than the instances.



Use case - dealing with high dimensional data where feature selection help in improving model performance.

3. Random Patches : It combine both The random subspace & pasting methods. It involves training each base model on a random subset of both instance & features.

Use case - beneficial when dataset is large & high-dimensional.

Bagging Classifier

Is an ensemble meta-estimator that fits base classifier each on random subsets of original dataset & then aggregate their individual predictions to form final prediction.

```
class sklearn.ensemble.BaggingClassifier( estimator=  
    None, n_estimators=10, max_samples=1.0, max-  
    features=1.0, bootstrap=True, bootstrap_features=  
    False, oob_score=False, warm_start=False,  
    n_jobs=None, random_state=None, verbose=0)
```

Implementation of code

```
>> from sklearn.ensemble import BaggingClassifier  
>> bag = BaggingClassifier( estimator = DecisionTree  
    .Classifier(), n_estimators = 500, max_samples=  
    0.5, bootstrap = True, random_state = 42)  
>> bag.fit(X-train, Y-train)  
>> bag.predict(X-test)
```



Bagging Regressor

A bagging regressor is an ensemble estimator that fits base regressors each on random subset of the original dataset & then aggregate their individual predictions (either by voting or averaging) to form a final prediction.

```
class sklearn.ensemble.BaggingRegressor (estimator=None,  
n_estimators=10, max_samples=1.0,  
max_features=1.0, bootstrap=True,  
bootstrap_features=False, oob_score=False,  
warm_start=False, n_jobs=None,  
random_state=None, verbose=0)
```

implementation of code

```
>> from sklearn.ensemble import BaggingRegressor  
>> bag = BaggingRegressor (random_state=1)  
>> bag.fit(x_train, y_train)  
>> bag.predict(x_test)
```

RANDOM FOREST

SARASWATI Education Society's

SARASWATI College of Engineering

PAGE NO.: 104

DATE: 13/7/2024

Random forests or random decision forests is an ensemble method for classification, regression and other task that operates by constructing a multitude of decision trees at training time. For classification tasks, the mean output of the random forest is the class selected by most trees.

For regression tasks, the mean or average prediction of the individual trees is returned.

Random forests converts the Highbias-High Variance HBHV model.

LowBiasHigh Variance algorithm into

LowBias Low variance algorithm

LBHV converts LBLV

Difference Between Bagging & Random Forest
Bagging -

- ① In Bagging we can use multiple algorithms (DecisionTree, SVM, KNN, etc.) it is generalize method.
- ② in bagging feature / column sampling the No. of columns choosed in starting are always fixed for all trees creation ③ Column Sampling is at tree level
Columns not shuffle they are fixed.
e.g if 2 columns are selected Col 3, & Col 5
So in all model it will be fixed.



Random Forest

- (1) In Random Forest we only use Decision Tree model / Algorithm It is focused method.
- (2) Column Sampling at node level
- (3) Column / feature Sampling the no. of columns choose in starting is fixed and the internal columns are shuffle in each node.
eg. if 2 columns are selected col2, & col4 after 1 node again randomly columns will shuffle now it will choose col3, col1 for next not if will go on & on..

• Random Forest Classifier.

```
class sklearn.ensemble.RandomForestClassifier (n_estimators=100,
criterion='gini', max_depth=None, min_samples_split=2,
min_samples_leaf=1, max_features='sqrt', max_leaf_nodes=None,
min_impurity_decrease=0.0, bootstrap=True,
oob_score=False, n_jobs=None, random_state=None,
verbose=0, warm_start=False, ccp_alpha=0.0,
class_weight=None, max_samples=None,
monotonic_cst=None)
```

implementation in code

```
>> from sklearn.ensemble import RandomForestClassifier
rf = RandomForestClassifier (n_estimators=500,
    random_state=42, max_features=2)
rf.fit (x-train, y-train)
rf.predict (x-test)
```



RandomForestRegressor

is a meta estimator that fits a no. of decision tree regressors on various sub-samples of the dataset & uses averaging to improve the predictive accuracy.

Syntax.

```
class sklearn.ensemble.RandomForestRegressor(n_estimators=100, *, criterion='squared_error', max_depth=None, min_samples_split=2, min_samples_leaf=1, min_weight_fraction_leaf=0.0, max_features=1.0, max_leaf_nodes=None, min_impurity_decrease=0.0, bootstrap=True, oob_score=False, n_jobs=None, random_state=None, verbose=0, warm_start=False, ccp_alpha=0.0, max_samples=None, monotonic_cst=None)
```

Implementation using sklearn.

```
>> from sklearn.ensemble import RandomForestRegressor  
>> rfr = RandomForestRegressor(n_estimators=1000)  
>> rfr.fit(x_train, y_train)  
>> rfr.predict(x_test)
```



* Out of Bag Evaluation OOB - Score

Random forests do not require validation dataset. Most random forests use a technique called out-of-bag-evaluation (OOB evaluation) to evaluate the quality of the model. OOB evaluation treats the training set as if it were on the test set of a cross-validation.

Over $\sim 33\%$ of training examples are omitted by decision tree while choosing data with replacement.

Implementation:

```
from sklearn.ensemble import RandomForestClassifier  
rfc = RandomForestClassifier(oob_score=True)  
rfc.fit(x-train, y-train)  
rfc.predict  
rfc.oob_score_
```

* Feature Importance using Random Forest & Decision Trees

It also known as variables or attributes, are individual measurable properties or characteristics of the phenomena being observed.

Formula for Decision trees feature importance.

$$\begin{aligned} f_{ik} &= \sum_{j \in \text{node split on feature } k} n_j \\ &= \frac{\sum_{j \in \text{all nodes}} n_j}{\sum_{j \in \text{all nodes}}} \end{aligned}$$



$$n_i = \frac{N-t}{N} \left[\text{impurity} - \left(\frac{N-t-r}{N-t} * \text{right-impurity} \right) - \left(\frac{N-t-l}{N-t} * \text{left-impurity} \right) \right]$$

where, n_i = node importance

$N-t$ = No. of rows to Node

N = Total No. of rows present in data

impurity = 'gini'

$N-t-r$ = Node to right

$N-t-l$ = Node to left.

right-impurity = impurity on right node

left-impurity = impurity on left node

$$\leq n_i$$

$$f_{ik} = \frac{\sum_{j \in \text{node split on feature } k}}{\sum_{j \in \text{all nodes}} n_j}$$

Implementation of Random forest feature importance

⇒ from sklearn.ensemble import RandomForestClassifier

rf = RandomForestClassifier()

rf.fit(X, y)

rf.feature_importances_

⇒ from sklearn.ensemble import DecisionTreeClassifier

clf = DecisionTreeClassifier()

clf.fit(X, y)

clf.feature_importance_

↑

Implementation of Decision Tree importance



ADABoost

AdaBoost Adaptive Boosting is used in conjunction with many other algs to improve performance. The output of the other learning algorithms ('Weak learners') is combined into a weighted sum that represents the final output of the boosted classifier.

$$\alpha_m = \frac{1}{2} \ln \left(\frac{1 - \text{error}}{\text{error}} \right)$$

$$\alpha_m = \frac{1}{2} \ln \left(\frac{1 - E_m}{E_m} \right)$$

where E_m = weighted error rate of the weak classifier.

α_m = weight of model m.

$$E_m = \frac{\sum_{i \neq k \in m} (x_i) w_i^{(m)}}{\sum_{i=1}^N w_i^{(m)}}$$

For misclassified

$$\text{new_wt} = \text{curr_wt} * e^{2\alpha}$$

For correctly classified

$$\text{new_wt} = \text{curr_wt} * e^{-2\alpha}$$



SARASWATI College of Engineering

AdaBoost Classifier

An AdaBoost classifier is a meta-estimator that begins by fitting a classifier on the original data set, and then fits subsequent classifiers on modified versions of the original data set.

Syntax:

```
class sklearn.ensemble.AdaBoostClassifier( estimator=  
    None, *, n_estimators=50, learning_rate=1.0,  
    algorithm='SAMME.R', random_state=None )
```

Implementation in sklearn:

```
>> from sklearn.ensemble import AdaBoostClassifier  
abc = AdaBoostClassifier()  
abc.fit(X,y)  
abc.predict(X-test)
```

Bagging

Vs Boosting

①	Types of model used	HBLV model used
	L BH V	used
	model used	
②	full grown decisiontree	Shallow grown decisiontree
③	Parallel training	Sequential (one-by-one)
④	equally weightage of base learners	biased or unequal weightage of base learners.

K-Means Clustering

SARASWATI Education Society's

SARASWATI College of Engineering

PAGE NO.: 111

DATE: 16/7/2029

K-means clustering

Is a unsupervised learning method of vector quantization, that aims to partition n observations into k clusters in which each observation belongs to the cluster with the nearest mean (cluster centers or centroid), serving as a prototype of the cluster.

Steps for achieving k-means clusters.

1) Decide n clusters



2) Initialize n centroids



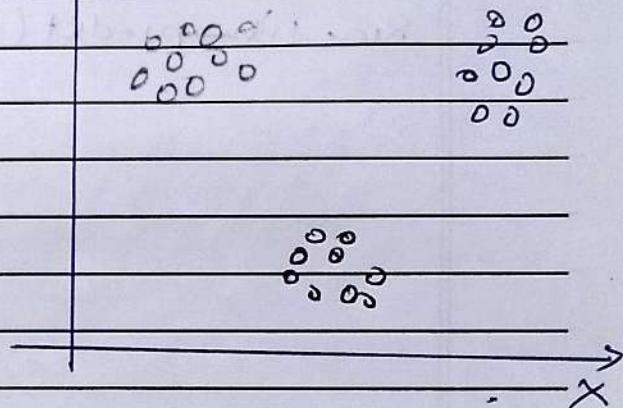
3) Assign cluster



4) Move centroid



5) Finish



To find k or decide number of clusters we follow (Elbow Method)

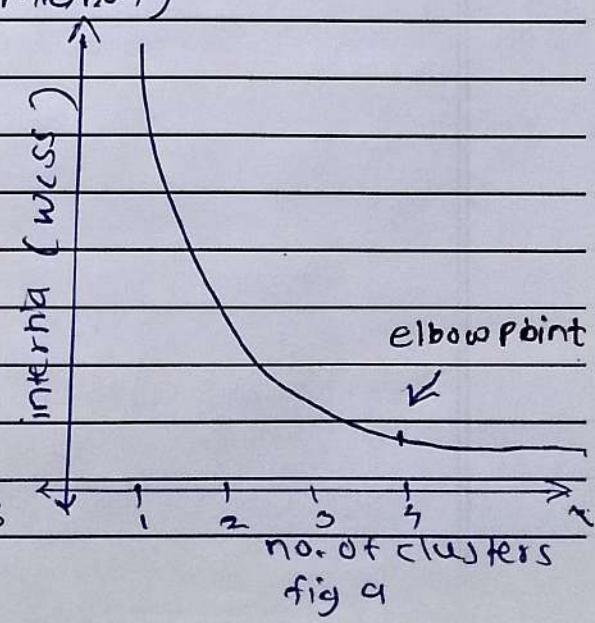
= WCSS \Rightarrow Within clusters

Sum of Squared distance

$$WCSS = \sum_{i=1}^n d_i^2$$

A the point whe WCSS got stable we choose no. of cluster

in fig a. 4 clusters is chosen





Syntax:

```
class sklearn.cluster.KMeans(n_clusters = 8, *,  
    init = 'k-means++', n_init = 'auto', max_iter = 300,  
    tol = 0.0001, verbose = 0, random_state = None,  
    copy_x = True, algorithm = 'lloyd')
```

Implementation in sklearn:

```
>> from sklearn.cluster import KMeans  
>> km = KMeans(n_clusters = 4)  
>> km.fit(x)  
>> km.predict(x-new) OR  
km.fit_predict(x)
```

Gradient boosting is a ML technique based on boosting in a functional space, where the target is pseudo-residuals rather than the typical residuals used in traditional boosting. It gives a prediction model in the form of an ensemble of weak prediction models.

Gradient boosting is used with decision trees of a fixed size as base learners.

* Algorithm *

Input: training set $\{(x_i, y_i)\}_{i=1}^n$, a diff. loss function $L(y, f(x))$, no of iterations M .

1. Initialize $f_0(x) = \arg \min_f \sum_{i=1}^n L(y_i, f)$.

2. For $m=1$ to M :

- (a) For $i=1, 2, \dots, N$ compute

$$\gamma_{im} = - \left[\frac{\partial L(y_i, f(x_i))}{\partial f(x_i)} \right]_{f=f_{m-1}}$$

- (b) Fit a regression tree to the targets γ_{im} giving terminal regions R_{jm} , $j=1, 2, \dots, J_m$.

- (c) For $j=1, 2, \dots, J_m$ compute

$$\gamma_{jm} = \arg \min_{x_i \in R_{jm}} \sum L(y_i, f_{m-1}(x_i) + r).$$

- (d) Update $f_m(x) = f_{m-1}(x) + \sum_{j=1}^{J_m} \gamma_{jm} I(x \in R_{jm})$.

3. Output $f(x) = f_M(x)$.



The Loss function that we have used in the above algo is Least square.

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

(MSE)

Syntax. for Regressor.

```
class sklearn.ensemble.GradientBoostingRegressor(  
    loss='squared_error', learning_rate=0.1,  
    n_estimators=100, subsample=1.0, criterion='friedman_mse',  
    min_samples_split=2, min_sample_leaf=1,  
    min_weight_fraction_leaf=0.0, max_depth=3,  
    min_impurity_decrease=0.0, init=None, random_state=None,  
    max_features=None, alpha=0.9, verbose=0,  
    max_leaf_nodes=None, warm_start=False,  
    validation_fraction=0.1, n_iter_no_change=None,  
    tol=0.0001, ccp_alpha=0.0)
```

implementation in sklearn.

```
from sklearn.ensemble import GradientBoostingRegressor  
reg = GradientBoostingRegressor(random_state=0)  
reg.fit(X_train, y_train)  
reg.predict(X_test)
```



Syntax for classifier.

```
class sklearn.ensemble.GradientBoostingClassifier(*,
loss='log-loss', learning_rate=0.1, n_estimators=100,
subsample=1.0, criterion='friedman-mse',
min_samples_split=2, min_sample_leaf=1,
min_weight_fraction_leaf=0.0, max_depth=3,
min_impurity_decrease=0.0, init=None, random_state=
None, max_features=None, verbose=0, max_leaf_
nodes=None, warm_start=False, validation_fraction=
0.1, n_iter_no_change=None, tol=0.0001, ccp_alpha=0)
```

implementation

```
From sklearn.ensemble import GradientBoostingClassifier
```

```
clf = GradientBoostingClassifier(n_estimators=100,
learning_rate=1.0, max_depth=
1, random_state=0)
```

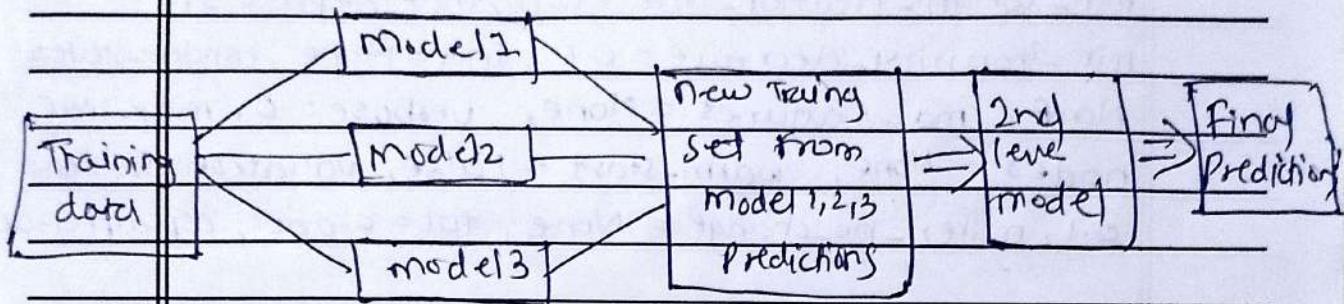
```
clf.fit(X-train, Y-train)
```

```
clf.predict(X-test)
```



STACKING

Stacking is a way to ensemble multiple classifications or regression model. These The point of stacking is to explore a space of different models for the same problem.

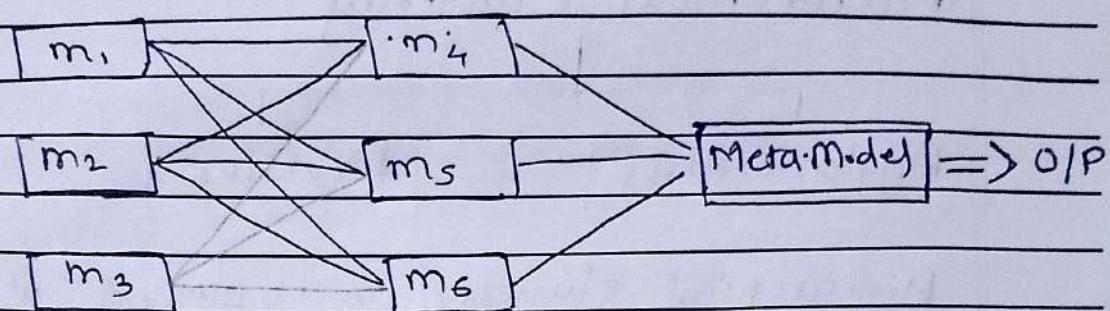


There are two or three methods for stacking:

① **Stacked Model**: An ensemble of models in which the meta-model is trained on out-of-fold predictions made by one base models during k-fold cross validation.

② **Blended Model**: An ensemble of models in which the meta-model is trained on predictions made by the base models on a holdout dataset (e.g. the validation dataset).

③ **Multi Layer Stacking**: This create a hierarchical structure of models where each layer's models are trained on predictions of the previous layer.



Syntax

```
class sklearn.ensemble.StackingClassifier (estimators,
final_estimator = None, *, cv = None, stack_method =
'auto', n_jobs = None, passthrough = False,
verbose = 0)
```

implementation using sk learn

```
from sklearn.ensemble import StackingClassifier
estimators = [ ('rf', RandomForestClassifier(n_estimators =
10, random_state = 42)),
('knn', KNeighborsClassifier(n_neighbors = 10)),
('gbdt', GradientBoostingClassifier())]
```

```
from sklearn.linear_model import LogisticRegression
```

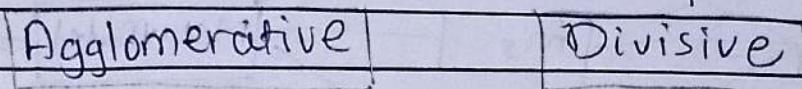
```
cif = StackingClassifier (estimators = estimators,
final_estimator = LogisticRegression(),
cv = 10)
```

```
cif.fit (X_train, y_train)
```

```
cif.predict (X_test)
```



Hierarchical Clustering



Hierarchical clustering is a method of cluster analysis that seeks to build a hierarchy of clusters. Strategies are of 2 ways.

- Agglomerative Clustering: This is a "bottom-up" approach: Each obs. starts in its own cluster, & pairs of clusters are merged as one moves up the hierarchy.
- Divisive Clustering: This is "top-down" approach: All observations start in one cluster, & splits are performed recursively as one moves down the hierarchy.

The merges & splits are determined in a greedy manner. The results of hierarchical clustering are usually presented in a dendrogram.

Algorithm:

1. Initialize the Proximity Matrix
2. Make each point a cluster
3. Inside a loop
 - a. Merge the 2 closest clusters
 - b. Update the Proximity Matrix
4. Until only one cluster is left.



Types of Agglomerative Clustering

1. Min (single-link) \Rightarrow 'single'

'single' uses the minimum of distances betⁿ all observations of the two sets.

2. Max (complete link) \Rightarrow 'complete' or 'maximum'
uses the maximum distance betⁿ all observations of the two sets.

3. Average \Rightarrow 'average'

uses the average of the distance of each obs of the two sets.

4. 'ward'

minimizes the variance of clusters being merged

syntax

```
class sklearn.cluster.AgglomerativeClustering(n_clusters=  
= 2, *, metric='euclidean', memory=None, connectivity=  
None, compute_full_tree='auto', linkage='ward',  
distance_threshold=None, compute_distances=False)
```

implementation

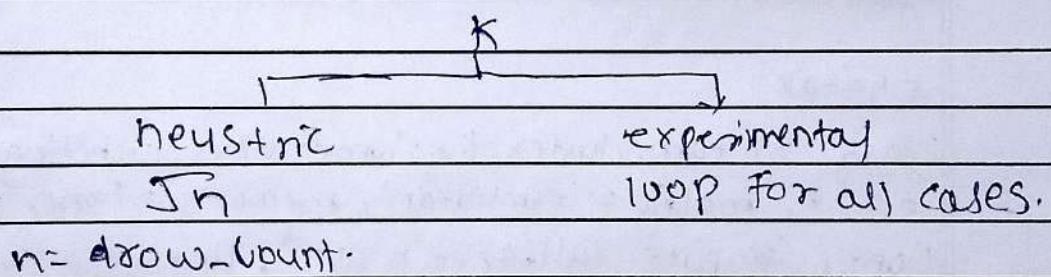
```
from sklearn.cluster import AgglomerativeClustering  
agc = AgglomerativeClustering()  
agc.fit(X)  
agc.labels_
```

K-nearest neighbors algorithm (K-NN) is a non-parametric supervised learning method. It is used for classification & regression.

K-nn algorithm assumes the similarity between the new case / data & available cases & put the new case into the category that is most-similar to the available categories.

K-nn alg. stores all the available data & classifies a new data point based on the similarity means when new data appears then it can be easily classified into a well suited category by using K-NN algorithm.

* How to select data K.



n = draw_count.

* Overfitting & Underfitting.

if $K=1$ or very low then \Rightarrow Overfitting

if $K=n$ or near to n then \Rightarrow Underfitting

n = no of rows.

* Limitations

- 1) Large dataset \Rightarrow Lazy learning technique
- 2) High dim data \Rightarrow Curse of dimensionality.
- 3) Outliers

- 4) Non-homogeneous Scales (feature scales)
- 5) Imbalanced dataset
- 6) Inference and not for prediction uses.

Syntax

```
class sklearn.neighbors.KNeighborsClassifier(  
    n_neighbors = 5, *, weights = 'uniform', algorithm =  
    'auto', leaf_size = 30, p = 2, metric = 'minkowski',  
    metric_params = None, n_jobs = None)
```

Implementation

```
from sklearn.neighbors import KNeighborsClassifier  
neigh = KNeighborsClassifier(n_neighbors = 3)  
neigh.fit(X, y)  
neigh.predict(x-test)
```



Assumptions of Linear Regression

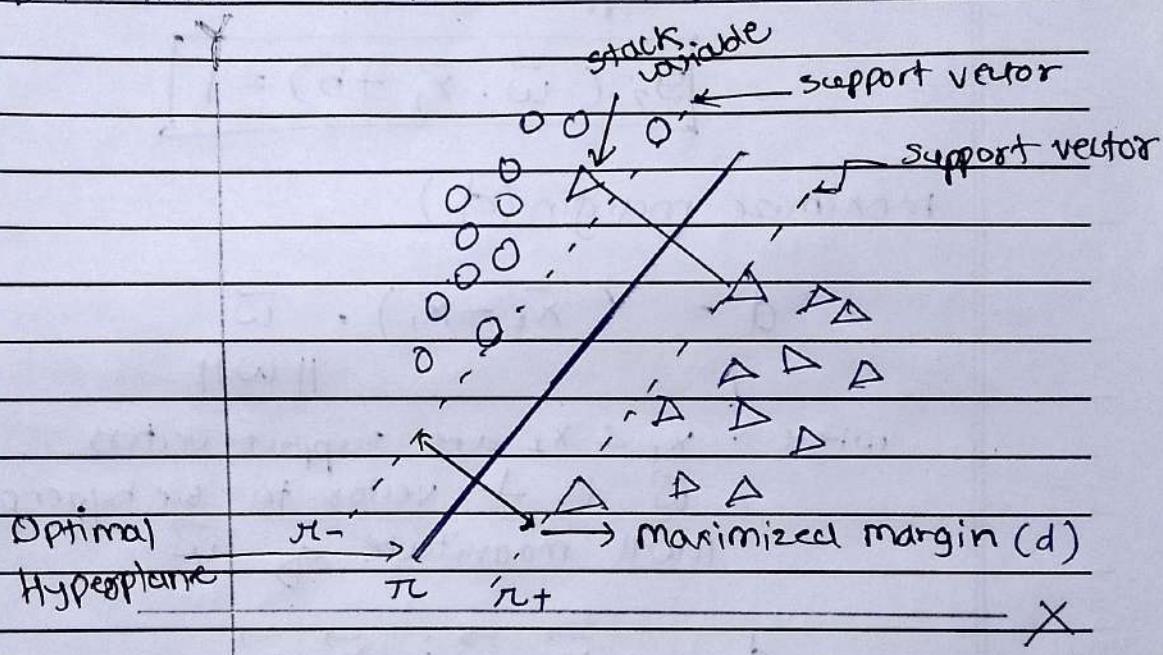
1. Linear Relationship between input & output
2. No Multicollinearity
3. Normality of Residual
4. Homoscedasticity
5. No Autocorrelation of Errors

- 1 → Find linearity b/w each feature with target column.
- 2 → All input features should independent in nature.
 - calculate variance inflation factor for checking multicollinearity
 - calculate one co-relation b/w columns for checking multicollinearity
- 3 → Error should be normally distributed
- 4 → Having the same scatter the ~~feature~~ spread should uniform on residual values.
- 5 → Should NO positive auto-correlation of Residuals.



Support Vector machines (SVM) are supervised max-margin models with associated learning algorithms that analyze data for classification & regression analysis.

SVM distinguish between two classes by finding the optimal hyperplane that maximizes the margin between the closest data points of opposite classes. SVM algorithm can handle both linear & nonlinear classification tasks.



Let's derive decision rule of SVM.
we kno.

$$\vec{w} \cdot \vec{u} > c$$

$$\text{or } \vec{w} \cdot \vec{u} - c > 0$$

$$\vec{w} \cdot \vec{u} + b > 0$$

for any x_i

$$y = \begin{cases} +1 & \text{if } \vec{w} \cdot \vec{x}_i + b \geq 0 \\ -1 & \text{if } \vec{w} \cdot \vec{x}_i + b < 0 \end{cases}$$

y = Predicted value.



assume $\vec{w}^T \vec{x} + b = 1$ for $\pi +ve$
 $\vec{w}^T \vec{x} + b = -1$ for $\pi -ve$

Optimization function

$$\begin{cases} \vec{w} \cdot \vec{x}_i + b \geq 1 \\ \vec{w} \cdot \vec{x}_i + b \leq -1 \end{cases} \quad \begin{cases} y_i (\vec{w} \cdot \vec{x}_i + b) \geq 1 \\ y_i (\vec{w} \cdot \vec{x}_i + b) \geq -1 \end{cases}$$

$$y_i (\vec{w} \cdot \vec{x}_i + b) \geq 1$$

for support vectors.

$$y_i (\vec{w} \cdot \vec{x}_i + b) = 1$$

maximize margin (d)

$$d = (\vec{x}_2 - \vec{x}_1) \cdot \frac{\vec{w}}{\|\vec{w}\|}$$

where \vec{x}_1 & \vec{x}_2 are support vectors

\vec{w} is \perp vector to our hyperplane.
 $\|\vec{w}\|$ magnitude of \vec{w} .

$$d = \frac{\vec{x}_2 \cdot \vec{w} - \vec{x}_1 \cdot \vec{w}}{\|\vec{w}\|} \dots \textcircled{①}$$

$$y_2 (\vec{w} \cdot \vec{x}_2 + b) = 1$$

$$1 (\vec{w} \cdot \vec{x}_2 + b) = 1$$

$$\vec{w} \cdot \vec{x}_2 = (1-b) \dots \textcircled{②}$$

similarly for x_1

$$-1 (\vec{w} \cdot \vec{x}_1 + b) = 1$$

$$\vec{w} \cdot \vec{x}_1 = -b-1 \dots \textcircled{③}$$

replace one value in eqn $\dots \textcircled{④}$



$$\frac{\vec{x}_1 \vec{w} - \vec{x}_2 \vec{w}}{\|\vec{w}\|} = \frac{(-b) - (-b-1)}{\|\vec{w}\|}$$

$$= \frac{1}{\|\vec{w}\|}$$

$d = 2$	$\dots \quad \textcircled{3}$
$\ \vec{w}\ $	

To maximize one eqⁿ $\textcircled{3}$.

$$\underset{(\vec{w}^*, b^*)}{\operatorname{argmax}} \frac{1}{2} \|\vec{w}\|^2 \text{ such that } y_i(\vec{w} \cdot \vec{x}_i + b) \geq 1 \quad y_i(\vec{w} \cdot \vec{x}_i + b) \geq 1$$

Hard Margin optimization function.

Hard Margin SVM: Is a type of SVM that tries to find a linear separator (hyperplane) that perfectly separates the data points of different classes without any errors

Constraint: requires data to be perfectly separable.

Soft Margin SVM: It allows some data points to be misclassified or to lie within the margin. It introduces a penalty for misclassified points to handle cases where data is not perfectly separable.

Constraint: It uses a regularization parameter (C) to control the trade-off between



achieving a larger margin and having a lower error rate. A smaller C allows more misclassifications, while a larger C tries to classify all training examples correctly.

The optimize function for soft margin is

$$\max f(x) \Leftrightarrow \min \frac{1}{2} f(x)$$

$$\underset{(w^*, b^*)}{\operatorname{argmin}} \quad \frac{\|w\|}{2} + C \sum_{i=1}^n \xi_i$$

subject to:

$$y_i(w \cdot x_i + b) \geq 1 - \xi_i \quad \text{for all } i = 1, 2, \dots, N$$

$$\xi_i \geq 0 \quad \text{for all } i = 1, 2, \dots, N$$

where,

w is weight vector

b is bias term

(zeta) $\xi_i(x_i)$ are slack variables that allow some misclassification

C is regularization parameter controls trade off betw max one margin & min one misclassification error.

y_i is class label of i^{th} e.g.

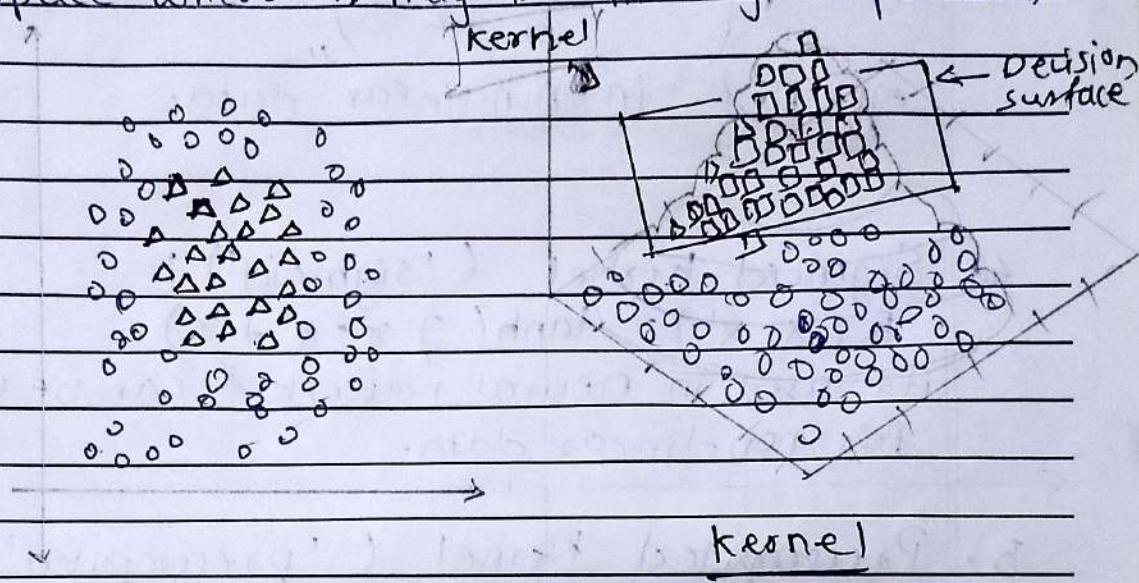
x_i is feature vector of i^{th} training e.g.

N is no of training. e.g.



Kernel Trick in SVM

The "Kernel Trick" is a method used in SVMs to convert data (that is not linearly separable) into a higher-dimensional feature space where it may be linearly separated.



different types of kernels used in SVM.

1. Linear kernel ('linear')
 - It is dot product of two vectors used for linearly separable data. $k(x, x') = x \cdot x'$
2. Polynomial kernel ('poly'):
 - is a non-linear kernel that represents the similarity of vectors in feature space over polynomial of original variables $k(x, x') = (\gamma x \cdot x')^d$
 - γ is scaling parameter, d is constant, d is degree of polynomial.
 - γ is scaling parameter
 - it is useful for scenarios where the relationship betⁿ class labels & attributes is non-linear.



3. Radial Basis Function (RBF) kernel ('rbf'):

- formula $K(x, x') = \exp(-\gamma \|x - x'\|^2)$

$$\gamma = \frac{1}{2\sigma^2}$$

$$K(x, x') = \exp\left(-\frac{\|x - x'\|^2}{2\sigma^2}\right)$$

effective in non-linear data.

4. Sigmoid kernel ('sigmoid'):

$$K(x, x') = \tanh(\gamma x \cdot x' + \gamma)$$

its used in neural network & can be used
for non-linear data.

5. Precomputed kernel ('precomputed'):

- requires manually compute the kernel
matrix before training SVM.

- kernel matrix should be an $n \times n$ matrix
where n is no. of samples.

6. Callable kernel

custom kernel function as a callable
fun should take arguments x & y , & return
kernel matrix.

e.g. def custom_kernel(x, y):
 return np.dot(x, y.T)



SVM syntax & parameters.

```
class sklearn.svm.SVC(*, C=1.0, kernel='rbf',
                      degree=3, gamma='scale', coef0=0.0,
                      shrinking=True, probability=False, tol=0.001,
                      cache_size=200, class_weight=None, verbose=False,
                      max_iter=-1, decision_function_shape='ovr',
                      break_ties=False, random_state=None)
```

implementation.

```
>> from sklearn.svm import SVC
>> model = SVC(kernel='linear', C=1)
>> model.fit(X,y) # X-train, y-train
>> model.predict(X-test)
```



|| NAIVE BAYES CLASSIFIER ||

SARASWATI Education Society's

SARASWATI College of Engineering

PAGE NO.: 130

DATE: 21/7/2029

Condition Probability

A, B:

$$P(A|B) = \frac{P(A \cap B)}{P(B)}, \quad \{P(B) \neq 0\}$$

Independent Events

A, B

$$P(A \cap B) = P(A) * P(B)$$

Mutually Exclusive Event

A, B

$$P(A \cap B) = 0$$

$$P(A|B) = \frac{n(A \cap B)}{n(B)}$$

$$P(A|B) = \frac{P(A \cap B)}{P(B)}$$

$$P(A) = \frac{n(A \cap B)}{n(B)}$$

$$P(A|B) = 0$$

$$P(A|B) = P(A)$$

Bayes Theorem

We know, Conditional probability

$$P(A|B) = \frac{P(A \cap B)}{P(B)} \quad \text{--- (1)}$$

$$A \cap B = B \cap A$$

$$\text{If } P(B|A) = P(B \cap A) = \frac{P(A \cap B)}{P(A)} \quad \text{--- (2)}$$

$$P(A \cap B) = P(B|A) * P(A) \quad \text{--- (3)}$$

$$P(A|B) = \frac{P(A \cap B) * P(A)}{P(B)} \Rightarrow \begin{array}{l} \text{prior} \\ \Downarrow \\ \text{Evidence} \end{array}$$

Posterior

given $\{P(B) \neq 0\}$



Naive Bayes Classifier

naive Bayes classifiers are a family of linear "probabilistic classifiers" which assumes that the features are conditionally independent, given the target class.

Naive Bayes is a conditional probability model. It assigns probabilities $p(C_k | x_1, \dots, x_n)$ for each of the K possible outcomes or classes C_k given a problem instance to be classified, represented by a vector $x = (x_1, \dots, x_n)$ encoding n features.

$$P(C_k | x) = \frac{P(C_k) P(x | C_k)}{P(x)}$$

$$\text{Posterior} = \frac{\text{Prior} \times \text{likelihood}}{\text{evidence}}$$

Types of Naive Bayes Classifiers

- Gaussian Naive Bayes (GaussianNB): This is used with Gaussian distributions - i.e. normal distributions - & continuous variables. This model is fitted by finding the mean and standard deviation of each class.

$$f(x) = \frac{1}{\sqrt{2\pi}} e^{-\frac{1}{2} \left(\frac{x-\mu}{\sigma}\right)^2}$$

- Multinomial Naive Bayes (MultinomialNB): This type assumes that features are from multinomial distributions. This variant is useful when using discrete data, such as frequency.



Counts, & it is typically applied within NLP.

- Bernoulli Naive Bayes (BernoulliNB) : This classifier is used with Boolean variables. That is, variables with two values, i.e. True & False or 1 & 0 .

Implementation .

```
from sklearn.naive_bayes import GaussianNB  
gnb = GaussianNB()  
gnb.fit ( X-train, y-train)  
gnb.predict ( X-test)
```

XGBoost (eXtreme Gradient Boosting) is an open-source software library which provides a regularizing gradient boosting framework for C++, Java, Python, R, Julia, Perl, Scala. It works on Linux, Microsoft Windows & macOS. It runs on a single machine, as well as the distributed processing frameworks Apache Hadoop, Apache Spark, Apache Flink & Dask.

XGBoost initially started by Tianqi Chen as a part of Distributed ML community group (DMLC).

Developed by Tianqi Chen from University of Washington.

Features.

Performance

Speed

Flexibility

1. Flexibility

- Cross Platform support
- Multiple Language Support
- Integration with other libraries & tools
- Support all kinds of ML problems.



2. Speed

- Uses Parallel Processing in tree building stage
- Uses Optimized Data Structures
- Cache Awareness
- Out of Core Computing - sequential
- Distributed Computing - parallel
- GPU Support

3. Performance

- Regularized Learning Objective
- Handling Missing Values
- Sparsity Aware Split Finding
- Efficient Split Finding
(Weighted Quantile Sketch)
+ Approximate Tree Learning
(Convert continuous var into discrete)
uses histogram based training using weighted quantile sketch.
- Tree Pruning

Comparing with Other Boosting Algorithms

XGBoost vs. AdaBoost

AdaBoost Emphasizes incorrect predictions using weights. Correct predictions reduce weight; incorrect prediction increase weight.
less effective than XGBoost in accuracy & speed.



XGBoost vs. CatBoost

CatBoost is another gradient boosting framework. It specializes in handling categorical features without preprocessing. perform well out-of-the-box with less hyperparameter tuning.

Like XGBoost, CatBoost has built in support for handling missing data.

CatBoost especially useful for datasets with many categorical features.

XGBoost vs. LightGBM

LightGBM (Light Gradient Boosting Machine)

LightGBM was developed by microsoft in 2016

most decision tree learning algorithm grow trees depth-wise. it uses leaf-wise tree growth strategy. like XGBoost, LightGBM exhibits fast model training speed & accuracy & performance well with large datasets.



XGBoost for Regression

For split finding in Regression 2 algorithms are used.

- (1) Exact Greedy Algorithm for Split finding
- (2) Split finding
- (3) Approximate Algorithm for Split finding

XGBoost Objective function contains loss function & regularization term. it tells diff. b/w actual value & predicted values.

XGBoost expects to have base learners which are uniformly bad at one remainder so that when all the predictions are combined, bad predictions cancels out & better ones sum up to form final good predictions.

For building XGBoost tree for regression.

Step 1: Calculate the similarity scores. it helps in growing the tree.

$$\text{Similarity Score} = \frac{(\text{sum of residuals})^2}{\text{No. of residuals} + \lambda}$$

Step 2: Calculate the gain to determine how to split the data.

$$\begin{aligned} \text{Gain} = & \text{Left tree (Similarity score)} + \\ & \text{Right. tree (Similarity score)} - \\ & \text{Root (Similarity score)} \end{aligned}$$

$$[L_t + R_t - Root]$$



Step 3: Prune the tree by calculating the difference between Gain & gamma

$$\boxed{\text{Gain} - \text{gamma}}$$

if result is positive no. then do not prune
if result is negative, then prune & again subtract gamma from the next Gain value way up the tree.

Step 4: Calculate output value for the remaining leaves

$$\text{Output value} = \frac{\text{Sum of residuals}}{\text{No. of residuals} + \lambda}$$

1st model is avg of all x_i (feature values)
(base model)

$$\text{Similarity Score} = \frac{(\sum \text{Residuals})^2}{N + \lambda} //$$

$$\text{Output Value} = \frac{\sum \text{Residuals}}{N + \lambda} //$$

where N is no of residuals.



Loss Function / Objective function.
 Loss fun regularization para

$$L = \sum_{i=1}^n L(y_i, \hat{y}_i) + \mathcal{R}(f_2(x_i))$$

$$\mathcal{R}(f) = \gamma T + \frac{1}{2} \lambda \|w\|^2$$

where, $T \Rightarrow$ no. of leaf nodes

$w \Rightarrow$ weight of leaf

$\gamma, \lambda \Rightarrow$ regularization Parameters

Derivation

Stage I

$$\boxed{m_1} \text{ mean} \rightarrow L = \sum_{i=1}^n L(y_i, \hat{y}_i)$$

$$f_1(x_i)$$

Stage II

$$\boxed{m_1} \quad \boxed{m_2}$$

mean dtree

$$L = \sum_{i=1}^n L(y_i, \hat{y}_i) + \mathcal{R}(f_2(x_i))$$

$$f_1(x_i) \quad f_2(x_i)$$

$$\hat{y}_i = f_1(x_i) + f_2(x_i)$$

$$L = \sum_{i=1}^n L(y_i, f_1(x_i) + f_2(x_i)) + \mathcal{R}(f_2(x_i))$$

Stage III

$$\boxed{m_1} \quad \boxed{m_2} \quad \boxed{m_3}$$

mean dt₁ dt₂

$$f_1(x_i) \quad f_2(x_i) \quad f_3(x_i)$$

$$L = \sum_{i=1}^n (y_i, f_1(x_i) + f_2(x_i) + f_3(x_i)) + \mathcal{R}(f_3(x_i))$$

Stage t

$$\boxed{m_1} \quad \boxed{m_2} \dots \dots \boxed{m_t}$$

mean dt₁ dt_t

$$f_1(x_i) \quad f_2(x_i) \quad f_t(x_i)$$

$$\underset{(t)}{L} = \sum_{i=1}^n L(y_i, f_1(x_i) + f_2(x_i) + \dots + f_t(x_i)) + \Omega(f_t(x_i))$$

$$\underset{(t)}{L} = \sum_{i=1}^n L(y_i, \hat{y}_i^{(t-1)} + f_t(x_i)) + \Omega(f_t)$$

At any stage t

$$\underset{(t)}{L} = \sum_{i=1}^n L(y_i, \hat{y}_i^{(t-1)} + f_t(x_i)) + \Omega(f_t)$$

we have to find ~~the~~ such output for $f_t(x_i)$ which will minimize the whole objective function

but as our $\Omega(f_t(x_i))$ faces one Taylor Series (zig-zag curve) & we need to approximate the $\sum_{i=1}^n L(y_i, \hat{y}_i^{(t-1)} + f_t(x_i))$

The solution is apply Taylor series is technique which do the approximation using polynomial.

Taylor series equation

$$f(n) \approx f(a) + f'(a)(n-a) + f''(a)(n-a)^2. \dots \quad ?!$$

$$x \rightarrow \hat{y}_i^{(t-1)} + f_t(x_i) \quad a \rightarrow \hat{y}_i^{(t-1)}$$

from deriving our simplified objective function is

$$\begin{aligned} \underset{(t)}{L} \approx \sum_{i=1}^n & \left[L(\hat{y}_i^{(t-1)}, \hat{y}_i^{(t-1)} + f_t(x_i)) + g_i f_t(x_i) + \frac{1}{2} h_i f_t''(x_i) \right] \\ & + \Omega(f_t(x_i)) \end{aligned}$$

where

$$g_i = \frac{\partial L(y_i, \hat{y}_i^{(t-1)})}{\partial \hat{y}_i^{(t-1)}}$$

(gradient)

$$h_i = \frac{\partial^2 L(y_i, \hat{y}_i^{(t-1)})}{\partial \hat{y}_i^{(t-1)2}}$$

(hessian)

We can remove one constant term.

$$L(y_i, \hat{y}_i^{(t-1)})$$

Now, objective function is

$$L^{(t)} = \sum_{i=1}^n \left[g_i f_t(x_i) + \frac{1}{2} h_i f_t^2(x_i) \right] + \lambda \sum_j w_j^2$$

Now, expand $\sum_j w_j^2$ to \Rightarrow

$$\sum_{i=1}^n \left[g_i f_t(x_i) + \frac{1}{2} h_i f_t^2(x_i) \right] + \gamma I + \frac{1}{2} \lambda \sum_{j \neq i} w_j^2$$

Now, replace.

~~w_j~~

$$f_t(x_i) \rightarrow w_j$$

w_j is weight.

$$\sum_{i=1}^n \rightarrow \sum_{j=1}^n$$

a_i and replace

$$\sum_{i=1}^n \rightarrow \sum_{j=1}^n \sum_{i \in S_j}$$

$$a_1 + a_2 + \dots + a_n = \\ a_3 + a_1 + a_2 + a_4 + \dots + a_n$$

SARASWATI College of Engineering

so our objective function becomes.

$$L^{(t)} = \sum_{j=1}^T \left[\sum_{i \in I_j} g_i w_j + \frac{1}{2} \sum_{i \in I_j} h_i w_j^2 \right] + \gamma T + \frac{1}{2} \lambda \sum_{j=1}^T w_j^2$$

$$L^{(t)} = \sum_{j=1}^T \left[\sum_{i \in I_j} g_i w_j + \frac{1}{2} \sum_{i \in I_j} h_i w_j + \frac{1}{2} \lambda w_j^2 \right] + \gamma T$$

$$L^{(t)} = \sum_{j=1}^T \left[\sum_{i \in I_j} g_i w_j + \frac{1}{2} \left(\sum_{i \in I_j} h_i + \lambda \right) w_j^2 \right] + \gamma T$$

differentiat. w.r.t. w_j

$$\frac{\partial L^{(t)}}{\partial w_j} = \sum_{i \in I_j} g_i + \left(\sum_{i \in I_j} h_i + \lambda \right) w_j = 0$$

$$w_j = \frac{-\sum_{i \in I_j} g_i}{\sum_{i \in I_j} h_i + \lambda}$$

gradient
↓
heisian

we know

$$g_i = \frac{\partial L}{\partial y_i^{(t-1)}} (y_i, \hat{y}_i^{(t-1)})$$

diff. w.r.t. $\hat{y}_i^{(t)}$

$$\frac{\partial L}{\partial \hat{y}_i^{(t)}} = \sum_{i=1}^n (y_i - \hat{y}_i)$$

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

↑
↓

11



and

$$h_i = \frac{\sigma^2 L(y_i, \hat{y}_i < t-1)}{\sum_{j=1}^n (\hat{y}_j - y_j)^2}$$

diff. w.r.t \hat{y}_i

$$\frac{\partial L}{\partial \hat{y}_i} = \sum_{j=1}^n (y_j - \hat{y}_j) \Rightarrow \sum_{j=1}^n (y_j - \hat{y}_j)$$

$$\Rightarrow \sum_{j=1}^n I_{ij}$$

Now put the differentiated values
in w_j . It becomes -

$$w_j = - \sum_{i \in I_j} (y_i - \hat{y}_i < t-1)$$

$$w_j = \frac{\sum_{i \in I_j} I_{ij}}{\sum_{i \in I_j} I_{ij} + \lambda}$$

$$w_j = \frac{\sum_{i \in I_j} (y_i - \hat{y}_i < t-1)}{\sum_{i \in I_j} I_{ij} + \lambda}$$

w_j = sum of residuals.

NO. of residuals + λ

* XGBoost Output value of for
Regression.

Output value for classification

$$w_j = \frac{\text{sum of residual}}{\sum_{i \in I_j} p_i (1-p_i) + \lambda}$$

where $p_i \rightarrow$ Predicted probability at previous timestamp.

$$w_j = \frac{\text{sum of residual}}{\sum_{i \in I_j} p_i (1-p_i) + \lambda}$$

where $p_i \rightarrow$ Predicted Probability at previous time-stamp.

* Similarity Score

Derivation

We have to find the minimum loss

$$L^{(t)} = \sum_{j=1}^T \left[\left(\sum_{i \in I_j} g_i \right) w_j + \frac{1}{2} \left(\sum_{i \in I_j} h_i + \lambda \right) w_j^2 \right] + \gamma T$$

$$w_j = - \frac{\sum_{i \in I_j} g_i}{\sum_{i \in I_j} h_i + \lambda}$$

$$L^{(t)}(q) = \sum_{j=1}^T \left[\frac{\left(\sum_{i \in I_j} g_i \right) \left(-\sum_{i \in I_j} g_i \right)}{\sum_{i \in I_j} h_i + \lambda} + \frac{1}{2} \left(\sum_{i \in I_j} h_i + \lambda \right) \frac{\left(\sum_{i \in I_j} g_i \right)^2}{\sum_{i \in I_j} h_i + \lambda} \right] + \gamma T$$

$$L^{(t)}(q) = \sum_{j=1}^T \left[- \frac{\left(\sum_{i \in I_j} g_i \right)^2}{\sum_{i \in I_j} h_i + \lambda} + \frac{1}{2} \cdot \frac{\left(\sum_{i \in I_j} g_i \right)^2}{\sum_{i \in I_j} h_i + \lambda} \right] + \gamma T$$



$$= T \sum_{i=1}^T \left[-\frac{1}{2} \frac{(\sum_{i \in I} g_i)^2}{\sum_{i \in I} h_i + \lambda} \right] + \gamma T$$

$$\underline{\underline{L^{(t)}(q) = -\frac{1}{2} \sum_{j=1}^T \frac{(\sum_{i \in I_j} g_i)^2}{\sum_{i \in I_j} h_i + \lambda} + \gamma T}}$$

Similarity Score

Tree building using Similarity Score

Similarity score of Parent.

$$L^{(t)}(q_p) = -\frac{1}{2} \frac{(\sum_{i \in I_p} g_i)^2}{\sum_{i \in I_p} h_i + \lambda} + \gamma$$

Similarity score of Child

$$L^{(t)}(q_c) = -\frac{1}{2} \left[\frac{(\sum_{i \in I_L} g_i)^2}{\sum_{i \in I_L} h_i + \lambda} + \frac{(\sum_{i \in I_R} g_i)^2}{\sum_{i \in I_R} h_i + \lambda} \right] + 2\gamma$$

where I_L for left node I_R for right node

$$L^{(t)}(q_p) - L^{(t)}(q_c) = -\frac{1}{2} \frac{(\sum_{i \in I_p} g_i)^2}{\sum_{i \in I_p} h_i + \lambda} + \gamma - \left(-\frac{1}{2} \left[\frac{(\sum_{i \in I_L} g_i)^2}{\sum_{i \in I_L} h_i + \lambda} + \frac{(\sum_{i \in I_R} g_i)^2}{\sum_{i \in I_R} h_i + \lambda} \right] + 2\gamma \right)$$

take $\frac{1}{2}$ common & $\gamma - 2\gamma = -\gamma$

$$L_{\text{split}} = \frac{1}{2} \left[\frac{(\sum_{i \in I_L} g_i)^2}{\sum_{i \in I_L} h_i + \lambda} + \frac{(\sum_{i \in I_R} g_i)^2}{\sum_{i \in I_R} h_i + \lambda} - \frac{(\sum_{i \in I_p} g_i)^2}{\sum_{i \in I_p} h_i + \lambda} \right] - \gamma$$

Loss reduction after the fit //

Density-based spatial clustering of applications with noise (DBSCAN)

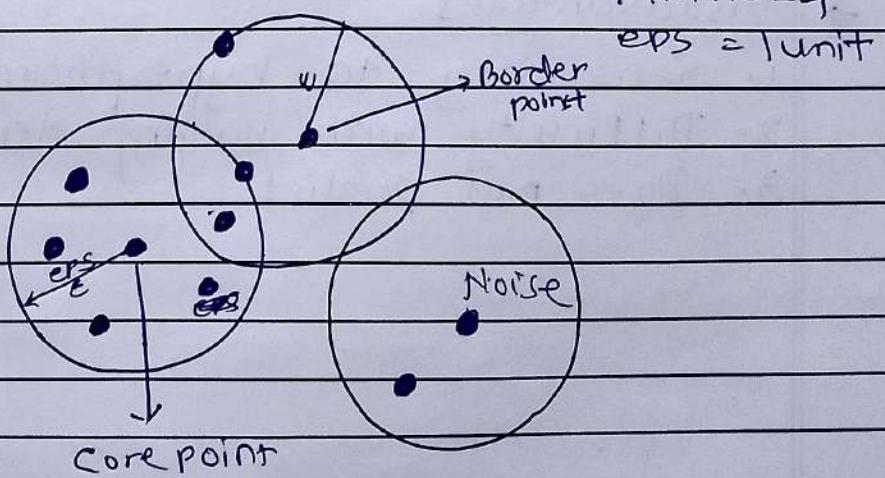
It is a density-based clustering non-parametric algorithm. Given a set of points in some space, it groups together points that are closely packed & marks as outliers points that lie alone in low-density regions.

In this algorithm we have 3 types of data points.

Core Point: A point is a core point if it has more than minPts points within eps .

Border Point: A point which has fewer than minPts within eps but it is in the neighborhood of a core point.

Noise or Outlier: A point which is not a core or border point.





Steps in Algorithm.

Step 1: Identify all points as either core point, border point or noise point.

Step 2: for all of the unclustered core points

a: Create new cluster

b: add all the points that are unclustered & density connected to the current point into this cluster.

Step 3: for each unclustered border point assign it to the cluster of nearest core point

Step 4: leave all the noise points as it is.



Advantages

1. Robust to outliers
2. No need to specify clusters
3. Can find arbitrary shaped clusters
4. Only 2 hyperparameter to tune



Disadvantages

1. Sensitivity to hyperparameters
2. Difficulty with varying density clusters
3. Does not predict



Syntax Parameters

```
class sklearn.cluster.DBSCAN(eps=0.5, *,  
    min_samples=5, metric='euclidean', metric-  
    params=None, algorithm='auto', leaf_size=30,  
    p=None, n_jobs=None)
```

metric : str, or callable, default='euclidean'

algorithm : { 'auto', 'ball-tree', 'kd-tree', 'brute' }

* Implementation.

```
>> from sklearn.cluster import DBSCAN  
>> cluster = DBSCAN(eps=3, min_samples=2)  
>> cluster.fit(X)  
>> cluster.labels_
```

Imbalanced data

SARASWATI Education Society's

SARASWATI College of Engineering

PAGE NO.: 148

DATE: 24/7/2024

Imbalanced data is a common problem in ML, which brings challenges to feature selection, class separation & evaluation & result in poor model performance.

It can be found in countless of applications in finance, healthcare, and public sectors.

Example

bank employee responsible for detecting the validity of credit card transaction.

a) Normal b) Fraudulent

most transaction are normal - but 0.1% are likely to fraudulent

Normal (0) = 99%

Fraudulent (1) = 0.1%

This is imbalanced data.

* Challenges with imbalanced data,

① The majority class - class with highest no of samples.

② The minority class - class with the lowest no of samples.

③ The class ratio - ratio betw minority class & majority class

* Sampling techniques to handle imbalanced data.

1. Random Under-Sampler (RU)

Down samples the majority class by randomly selecting instances to match one minority class.



Pros : → Reduction in bias
→ Faster training

Cons : = Info loss leading to underfitting
= Sampling Bias

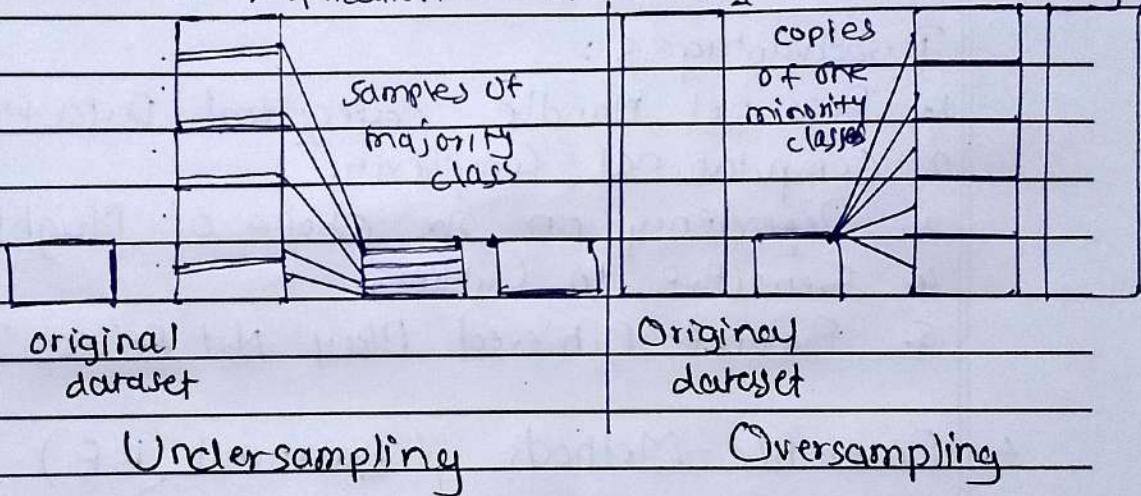
2. Random Over-Sampler (RO)

Upsamples the minority class by randomly duplicating instances until class size are balanced.

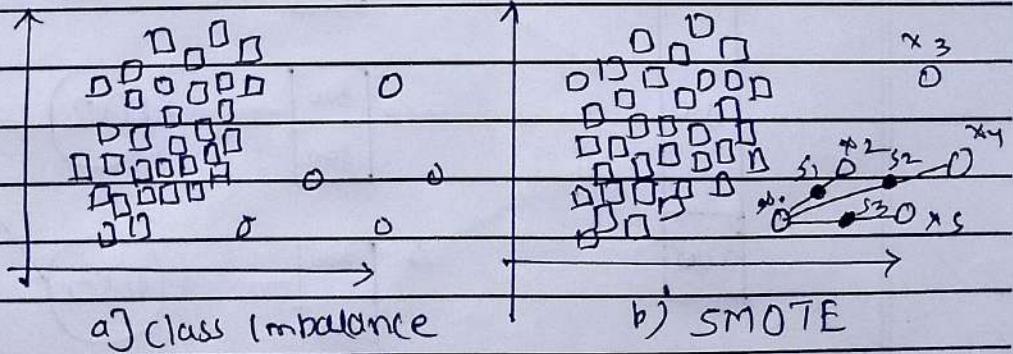
Pros : → avoids data deletion | Reduce bias

cons : → Increased size

→ Duplication of data may cause overfitting.



3. SMOTE



Steps to apply SMOTE

Approach: Generates new instances for the minority class to avoid exact duplication.

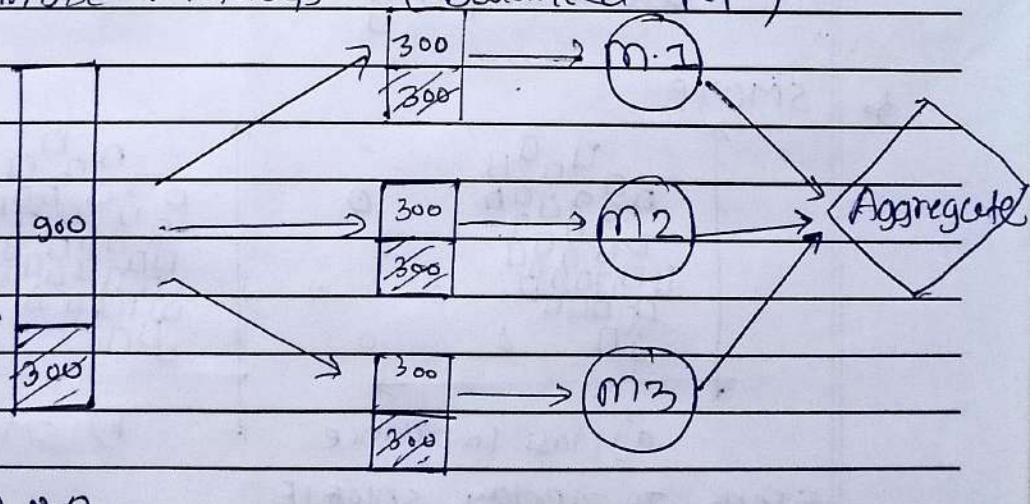
Algorithm:

- Step 1: Train KNN on minority class obs. And each obs 5 closest neighbours.
- Step 2: Select examples from minority class at random
- Step 3: Select a neighbour of each ex. at random
- Step 4: Extract a random no. between 0 & 1
- Step 5: Cal the new ex. as = original sample
 → factor * (original sample - neighbour)
- Step 6: The final dataset consist of the original dataset + newly created examples.

Disadvantages :

1. Does not handle Categorical Data well
2. Computational Complexity
3. Dependency on one choice of Neighbors
4. Sensitive to Outliers
5. Balance Achieved May Not Reflect True Nature

4. Ensemble Methods (Balanced RF)



Approach: Modifies the traditional RF algo to handle class imbalance by using balanced bootstrapping.



Algorithm.

Step 1: Each tree is trained on a bootstrap sample that contains an equal no of instances from each other.

Step 2: This done by randomly under-sampling the majority class &/or over-sampling the minority class.

5. Cost Sensitive Learning

Is an approach within ML that considers varying costs associated with diff types of errors. This method diverges from traditional approaches by introducing a cost matrix, explicitly specifying the penalties like class weight.

Class weights in Cost-Sensitive Learning
class weights adjust the importance of different classes in imbalanced datasets:

It increases model sensitivity to minority class by assigning higher weights to them.

Implementation

1. Manual: set weights manually based on class importance.
2. Automatic: use library features like in sklearn. There is a set parameter for every differentiable Loss function models. "class_weight" for automatic weight calculation.