

Dataset-1Lab-0Global Blood Group Distribution
worldwide Dataset.attributes:

- 1) Country : Where the person having blood group is from.
 - 2) Population : Number of people
 - 3) Blood Groups :
[O+, A+, B+, AB+, O-, A-, B-, AB-]
- Missing values in blood groups:
O-, A-, B-, AB-
- Total unique values - 126

Dataset-2Metro Transit Data

This dataset was created for analyzing metro data across Indian cities.

No. of rows: 226

No. of columns = 10

City → Name of the city

Country → Name of the country

Name → Name of the metro system

Service opened → Opening date of metro

Last expanded - Year when the metro system was last expanded.

Lines → Total number of metro lines

System length → Length of the metro length in Km.

Dataset - 3

Startup Growth & Funding Trends

No. of columns = 12

No. of rows = 500

Attributes: Many of groups need

already for business modelling

Startup Name : The name of the startup

Industry - The sector in which the start-up operates.

Funding Rounds - The total number of funding rounds

Valuation - The startup's post-money

Revenue - The estimated revenue

Employee - The number of employees working with start-ups.

Market Share - The percentage of market share that the

start-up has captured

PBA Profitable - A binary indicator

(1 = Profitable, 0 = Not Profitable)

Dataset - 4

Titanic Dataset - EDA & Logistic Regression

No. of rows = 831

No. of columns = 10 to 12 columns

1. Survival prediction: To build a logical regression
 2. Data cleaning and preprocessing:
To perform data cleaning
 3. EDA: Visualizing patterns.
- Name - person's name
- Sex → either
- | | |
|--------|---|
| Female | ↑ |
| Male | ↓ |
- Age : age of a person.

Ticket → which type of ticket.

Fare → Price of the fair

Data-5

Mars rover datasets

id → unique identifier for each image

sol → Martian sol (day) when the image was captured

camera-name → rover's camera

img-src → URL link to the image

earth-date : The earth date corresponding to the martian sol.

rover-status : current operational status

→ active	↓
complete	

landing-date → Date when it was landed.

launch-date → Date when launched from earth.

No. of rows: 856 values

No. of columns: 10 columns

Linear Regression in Machine Learning

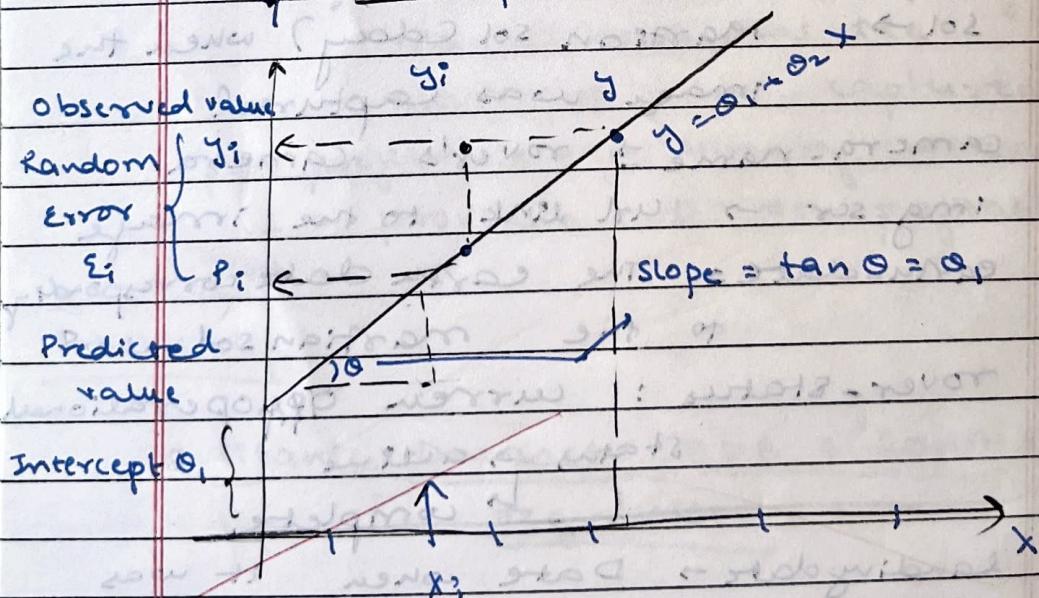
Linear regression is a statistical method used to model the relationship between a dependent variable and one independent variable.

- It is a model-based learning and also a type of supervised machine learning

Regression: try to find the relationship between variables.

Equation for linear expression:

$$\hat{y} = \beta_0 + \beta_1 x + \epsilon$$



$$\text{Slope} = \frac{\text{Change in } y}{\text{Change in } x} = \frac{\Delta y}{\Delta x}$$

- Linear regression learns from labelled datasets

code for linear regression:

```

import pandas as pd
import numpy as np
// load dataset
df = pd.read_csv("data_for_lr.csv")
df = df.iloc [:50] // first 50 points
// first row & second column
x = df.iloc[:, 0].values
y = df.iloc[:, 1].values
// least square formula
x_mean, y_mean = np.mean(x),
np.mean(y)
// slope - m & intercept - b
m = np.sum((x - x_mean) * (y - y_mean)) /
np.sum((x - x_mean) ** 2)
b = y_mean - m * x_mean
// Predict values
y_pred = m * x + b
mse = np.mean((y - y_pred) ** 2)
print(f"Slope (m): {m},\nIntercept (b): {b},\nMSE: {mse}")
// Plotting
plt.scatter(x, y, color='blue',
label='Data points')

```

plt. plot (x, y-pred, color = 'red',
label = "Regression Line")

plt. xlabel ("X values")

plt. ylabel ("Y label")

plt. legend ()

plt. title ("Linear Regression")

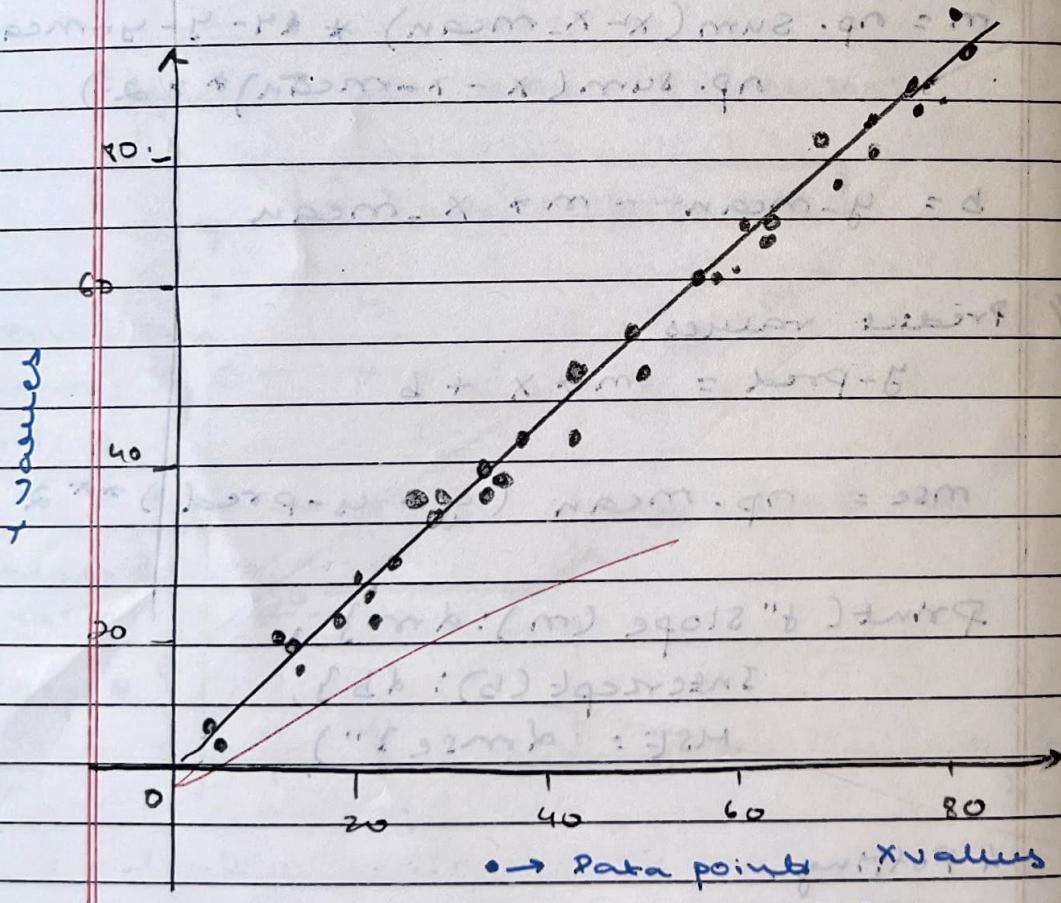
plt. show ()

Output:

Slope (m): 0.9926528446814438

Intercept (b): -0.0738029436

MSE: 6.42113898574697



formula: $B_1 = \frac{\bar{xy} - \bar{x} \cdot \bar{y}}{\bar{x^2} - (\bar{x})^2}$

$B_0 = \bar{y} - B_1 \bar{x}$

8/13/25

Lab - 2

End-to-End Machine Learning Project.

Step-1

Download the dataset in CSV format
i.e. housing.tgz → housing.csv

attributes: 1) longitude

2) latitude

indexes: 20640 3) housing-median-age numerical

4) total_rooms

data

5) total_bedrooms

↓

6) population

can use regression

7) households

8) median_income

9) median_house_value

10) ocean_proximity. } string/text

can only be categorized.

value

(Object-type)

Step-2

Load the dataset, for quick glance;

housing.head() → displays first 5 rows
with attributes.

housing.info() → Gives info about whole
dataset

Step-3

Categorize based on ocean_proximity

housing.describe() → gives mean, count,
std, max, min, 25%, 75%, 50% of
all attributes. [null values are ignored]

wood planks & (1.0 angles)

red pitholw w202
was given

For

use histogram for each attributes.
bins \Rightarrow 50 } divides the data into 50-equal width.

Step-4

create a Test set

Data Snooping: It occurs when model is trained or tested on data that was used in preprocessing (decision making process)

Test-data: test size

Training data = Total dataset - Test Data [160/16]

iloc[] \rightarrow specific rows.

np.random.permutation()

using unique ID's & immutable identifier, or use

hash function.

random-state & using

split-train-test()

Visualization Methods:

1) Histogram:

- hist() \rightarrow visualize the data based on histograms

2) using scatter plots

using scatter plots between different types of attributes

(alpha = 0.1) \rightarrow density based

Scatter plotting for visualization

3) Correlation:

using `.corr()` for attributes.

Step - 5

Data cleaning:

- 1) remove corresponding districts
- 2) get rid of whole attribute
- 3) set values to some value

use `dropna()`, `fillna()`, `drop()`

Step - 6

Select and Train a model.

use regression models.

for validation use `.train-test-split()`

English

26/03/21

Lab - 04

Decision Tree ID3 Classification

Pseudocode.

function ID3 (Data, features, target):

If all instances in Data have the same class:

Return the class label

If features is empty:

Return the most common class
on Data

BestFeature \leftarrow feature with highest
information gain in Features

Tree \leftarrow create a root node labeled
BestFeature

For each value v in BestFeature:

subset \leftarrow subset of Data where

BestFeature = v

If subset is empty:

Add a leaf node with the
most common class in Data

else:

childNode \leftarrow ID3 (subset,

features - {BestFeature}, target)

Add childNode to Tree

Return Tree

Entropy for each feature:

outlook : 1.5774

temperature : 1.5567

humidity : 1.000

windy : 0.9852

Information gain for each feature:

outlook : 0.2469

temperature : 0.0292

humidity : 0.1818

windy : 0.0481

Root Node (Best split): outlook.

outlook <= 0.5
entropy = 0.94
samples = 14
value = [5, 9] class = 1

entropy = 0.0
samples = 4
value = [0, 4]
class = 1

humidity <= 0.5
entropy = 1.0
samples = 10
value = [5, 10] class = 0

outlook <= 0.5
entropy = 0.722
samples = 5
value = [5, 10] class = 0

windy <= 0.5
entropy = 0.722
samples = 5
value = [1, 4] class = 0

entropy = 1.0
samples = 2
value = [1, 1] class = 0

entropy = 0.0
samples = 5
value = [3, 0] class = 0

entropy = 1.0
samples = 2
value = [1, 1] class = 0

entropy = 0.0
samples = 3
value = [0, 3] class = 1

Lab - 5

KNN Algorithm.

Function split (data, 0.2) :

split_index $\leftarrow \text{length}(\text{data}) \times (1 - 0.2)$

train = first split_index elements of data

test = remaining elements of data

return train, test.

function euclidean (a, b) :

sum = 0

for i = 1 to length(a) :

sum = sum + $(a[i] - b[i])^2$

return square-root (sum).

function get_neighbours (train, test_instances, k) :

distances = []

for each train_instances in train

dist = euclidean (test_instance,

0.5 * train_instances.features)

distances.append (train_instances, dist)

distances.sort()

neighbours \leftarrow first k elements of distances

return neighbours.

function predict (train, test-instance, k):
 neighbours \leftarrow get-neighbours
 (train, test-instance, 3).
 count \leftarrow occurrences of each class in
 neighbours.
 return class with max(count).

junction KNN()

train-test = split dataset (dataset=0.7)
 $k = 3$

predictions = []
for each test-instance in test-set:

predicted-class = predict (train,
 test-instance.features, k)

predictions.append (predicted-class)

Data set used is Iris.csv

feature 1 = species. Sepal.length()

feature 2 = Sepal width

Predicted class for [5.7, 3.5, 1.7]:

Iris - virginica

SVM Algorithm:

function load-dataset (filename):

dataset = []

open(file)

skip header

for each row in file:

features = first 2 feature values

label = 1 if class is 'Iris-setosa'
else 0

dataset.append([features, label])

return dataset

function train-test-split (dataset, test-ratio):

split-under = length(dataset) * test-ratio

train-set = dataset[:split-under]

test-set = dataset[split-under:]

return train-set, test-set.

function dot-product (a, b):

return sum (a[i] * b[i]) for all

function train-sum (train, learning-rate, lambda-param, epochs)

wt [0, 0]

b ← 0

for each epoch in epoch:

for each data point (x, y) in train:

if ($y \times (\text{dot-product})(w, x) + b \geq 1$):

 update $w = w - \text{learning-rate} \times x$

$\times 2 \times \text{lambda-param} \times w$

else:

 update $w = w - \text{learning-rate} \times 2 \times \text{lambda-param} \times w - y \times x$

 update $b = b + \text{learning-rate} \times y$

end if

return w, b .

) algorithm = k nearest neighbor

function predict(x, w, b):

 return 1 if dot-product(w, x) + $b > 0$ else -1

filename = 'content/iris.csv'

dataset = load-dataset(filename)

train, test = train-test-split(dataset)

$w, b = \text{train-sum}(train)$

plot-sum(train, w, b)

for each point in test:

 prediction = predict(point.features, w, b)

 print(Actual | Predicted).

Lab-6

Random forest classification.

Random forest classifier ($x\text{-train}$, $y\text{-train}$,
 T , max_depth , min_samples_split):

$n\text{-samples}$, $n\text{-features} = x\text{-train.shape}$

for $t=1$ to T :

bootstrap-idx = random.sample(
range($n\text{-samples}$), $n\text{-samples}$).

$x\text{-bootstrap} = x\text{-train}[bootstraps_idx]$

$y\text{-bootstrap} = y\text{-train}[bootstraps_idx]$

Train decision Tree:

tree = DecisionTree($\text{max_depth} = \text{max_depth}$,
 $\text{min_samples_split} = \text{min_samples_split}$)

trees.append(tree)

return trees

Predict($x\text{-test}$):

prediction = []

for tree in trees:

tree_predictions = tree.predict($x\text{-test}$)

predictions.append(tree_predictions)

final_predictions = majority_vote
(predictions)

return final_predictions

majority_vote(predictions):

return mode(predictions)

Lab-7

Boosting Ensemble

1. Initialize weights for each training sample:
2. For each iteration ($t = 1$ to T):
 - a. Train a weak classifier (e.g. a decision stump) using the weighted samples.
 - b. calculate the classifier

$$\alpha_t = 0.5 \times \log \left[\frac{(1 - \text{error})}{\text{error}} \right]$$

If error is high, alpha will be low, meaning less importance for this classifier.
 - c. calculate the classifier's error:

$$\text{error} = (\text{sum of weights of misclassified samples}) / (\text{sum of all sample weights})$$
 - d. update the sample weight
 - increase the weight of misclassified sample
 - ~~Decrease the weight of correctly classified samples.~~
 - ~~Normalize the weight to make sure they can sum to 1~~

3. After all iterations, combine the weak classifiers:

- each classifier contributes based on its importance (α)

- final prediction =

: (sign sum of all classifiers)

4. Use the combined classifier to make predictions on new data.

Lab - 8

k-means clustering.

1. Initialize K clusters' centroids [randomly selected data points].
2. Repeat until convergence:
 - a. Assign step:
 - for each data point, calculate the distance to each centroid.
 - Assign each data point to the closest centroid [cluster].
 - b. Update step:
 - For each cluster, compute the mean of all data points assigned to that cluster.
 - update the centroid to this mean.
3. Repeat the process until the centroids do not change significantly (convergence)
4. The final centroids represent the cluster centers, and each data point belongs to one of the clusters.

Lab-9

Input:

data matrix X with shape (n-samples, n-features)

Desired number of components k

Step 1 - center the data

For each feature column in X :

Subtract the means of that column

or at string for each value

Step 2: compute the covariance matrix

$$\text{cov} = (1 / (\text{n-samples})) \cdot X \cdot \text{transposed}$$

Step 3: compute eigenvalues and eigenvectors
of the covariance matrix [eigenvalues,
eigenvectors] = eig(cov)

Step 4: sort eigenvectors by decreasing

now as the eigenvalues largest

sort eigenvectors so that those with
highest eigenvalues come first

Step 5: select top k eigenvectors

Select the first k eigenvectors to form
the projection matrix W

Step 6: Transform the original data

$$X_{\text{reduced}} = X \cdot W$$

Output: X_{reduced} (data in k dimension)