

1) Obesity Prediction Dataset :

Understanding lifestyle factors that contribute to obesity can help with early intervention, health recommendations, and machine learning applications in healthcare.

Attributes :

- 1) Gender No. of columns : 14
- 2) Age No. of rows : 2111
- 3) Height
- 4) Weight
- 5) family history with overweight
- 6) FARE → frequent high-calorie food consumption
- 7) FCVC - Vegetable consumption
- 8) NCP - No. of meals/day
- 9) CAEC - Frequency of consuming food b/w meals
- 10) Smoke
- 11) H2O - daily water intake
- 12) SCC - If person monitors their calorie intake
- 13) FAF - Physical activity
- 14) CALC - Alcohol consumption

2) Road Accident Survival Dataset

Analyzing the impact of Age, Speed, and safety measures on accident survival.

Attributes :

- 1) Age No. of Rows : 76
- 2) Gender No. of columns : 6
- 3) Speed of impact
- 4) Helmet Used
- 5) Seatbelt Used
- 6) Survival

3) Indian Pima (Hope)3) Diabetes dataset

Includes various health parameters, lifestyle habits that contribute to diabetes risk.

Attributes :

- | | |
|-------------------------|-----------------------|
| 1) Age | 11) Hip Circumference |
| 2) Pregnancies | 12) MR |
| 3) BMI | 13) Family History |
| 4) Glucose | 14) Diet Type |
| 5) Blood Pressure | 15) Hypertension |
| 6) LDL | 16) Medication Use |
| 7) HDL | 17) Outcome |
| 8) HbA1c | |
| 9) Triglycerides | No. of Rows : 9538 |
| 10) Waist circumference | No. of columns : 17 |

4) Titanic Dataset

- | | |
|-------------------------------|--|
| 1) Passenger ID | No. of columns : 10
No. of rows : 891 |
| 2) Survived | |
| 3) Pclass | |
| 4) Name | |
| 5) Sex | |
| 6) Age | |
| 7) SibSp - Siblings / spouses | |
| 8) Parch - Parent child | |
| 9) Ticket | |
| 10) Fare | |

5) Customer Personality Analysis

Helps businesses to better understand its customers and make it easier for them to specify products according to the specific needs, behaviours.

Attributes

- | | |
|-------------------|---------------------|
| 1) ID | No. of rows : 2240 |
| 2) Year_Birth | No. of columns : 10 |
| 3) Education | |
| 4) Marital_Status | |
| 5) Income | |
| 6) Kidhome | |
| 7) TeenHome | |
| 8) D6_Customer | |
| 9) Recency | |
| 10) Complaint | |

~~Q3~~
3/3/25

LINEAR REGRESSION

Linear regression is a supervised machine learning algorithm that is used to model the relationship between two variables by fitting a straight line to the data. The goal is to find the line that best describes how one variable (the independent variable, or "x") influences another value (the dependent variable, or "y").

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

↓ ↓ ↓
 output Coefficients input
 error

- ~~Y~~ y - is the dependent variable
- x - is the independent variable
- β_0 - is the intercept
- β_1 - is the slope.

Linear Regression Code

```
import numpy as np
```

```
import pandas as pd
```

```
import matplotlib.pyplot as plt
```

```
np.random.seed(42)
```

```
x_values = np.random.uniform(0, 50, 50)
```

```
y_values = 3 * x_values + 5 + np.random.normal(0, 10, 50)
```

```
def LR_formular(x: list[int], y: list[int]):
```

```
x_values = x_values.reshape(-1)
```

```
N = len(x_values)
```

```
sum_x = np.sum(x_values)
```

$$\text{sum}_x = np \cdot \text{sum}(\text{x-value})$$

$$\text{sum}_y = np \cdot \text{sum}(\text{y-value})$$

$$\text{sum}_{xy} = np \cdot \text{sum}(\text{x-value} * \text{y-value})$$

$$\text{sum}_{x^2} = np \cdot \text{sum}(\text{x-value} ** 2)$$

$$b_1 = (\text{sum}_{xy})$$

$$x\text{-mean} = \text{sum}_{xy} / \text{len}(x)$$

$$x\text{-mean} = \text{sum}_x / \text{len}(x)$$

$$y\text{-mean} = \text{sum}_y / \text{len}(y)$$

$$x^2\text{-mean} = \text{sum}_{x^2} / \text{len}(x)$$

$$b_1 = (y\text{-mean} - x\text{-mean} * y\text{-mean}) / (x^2\text{-mean} - x\text{-mean} ** 2)$$

$$b_0 = y\text{-mean} - b_1 * x\text{-mean}$$

return b_0, b_1

def predict($b_0: \text{int}, b_1: \text{int}, x: \text{int}$):
 return $b_0 + b_1 * x$

~~orange~~ = orange = [x for x in range(1, 31)]

year = [random(2 * x + 10 + random.uniform(-5, 5))
 for x in range(30)]

points(x=orange)

points(y=year)

$b_0, b_1 = \text{LR-formula}(x=orange, y=year)$

y = predict($b_0, b_1, 32$)

points(x=orange, y=y)

plt.plot(x=orange, y=year, color='blue', label='Data')
plt.plot(x=orange, y=y, color='red', label='Model')

plt.plot (x, y, 'red') for x in
xarray, color = 'red', label
= 'Regression line')

plt.xlabel ('X values')

plt.ylabel ('Y values').

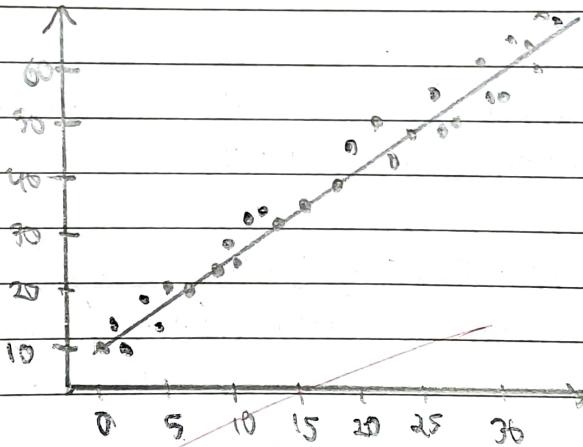
plt.title ('Linear Regression')

plt.legend()

plt.show()

INPUT:

[1, 2, 3, 4, 5, 6, 7, 8, 9, 10, 11, 12, 13, 14, 15, 16, 17, 18, 19,
20, 21, 22, 23, 24, 25, 26, 27, 28, 29, 30]
[12.36, 15.09, 13.28 ...]



880
101210

LAB - 2

Date _____
Page _____

END-TO-END MACHINE LEARNING PROJECT

The dataset contains 20640 entries and 10 columns.

The data is split 80-20 with 80% training and 20 test data.

1) Data Loading and Exploration :

Dataset : housing.csv was loaded using pandas

Exploration : .head(), .info(), .describe() used to inspect dataset structure.

Visualisation : Histograms and scatter plots were created to understand feature distributions.

2) Data Splitting : Implemented random sampling and stratified sampling based on income categories to ensure representative test and training sets.

3) Data cleaning and pre-processing

- * Handling missing values using Simple Imputer to fill missing values with median values

- * Feature Engineering : Created new features such as rooms_per_household.

- * Categorical Encoding : Applied Ordinal Encoder and One-Hot Encoder for categorical data transformation.

4) Data Transformation Pipeline

- * Used Scikit-learn Pipelines to automate pre-processing :

- Numerical Features : Imputation, feature scaling

- Categorical Features : One-Hot encoding

5) Model Training and Evaluation :

Linear Regression :-

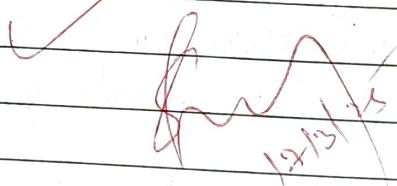
- i) Model Training using Scikit - Learn
 - housing - prepared : Preprocessed feature matrix after applying transformations
 - housing labels : Target variable
 - fit() : Trains the model by finding the best-fitting line that minimizes the error.

ii) Making Predictions

The models predictions were compared with actual values to check performance

iii) Evaluating the Model

Root Mean Squared Error (RMSE) a common metric for regression models is used



DECISION TREE - ID3 CLASSIFICATION

Dataset Used: PlayTennis.csv from Kaggle Datasets

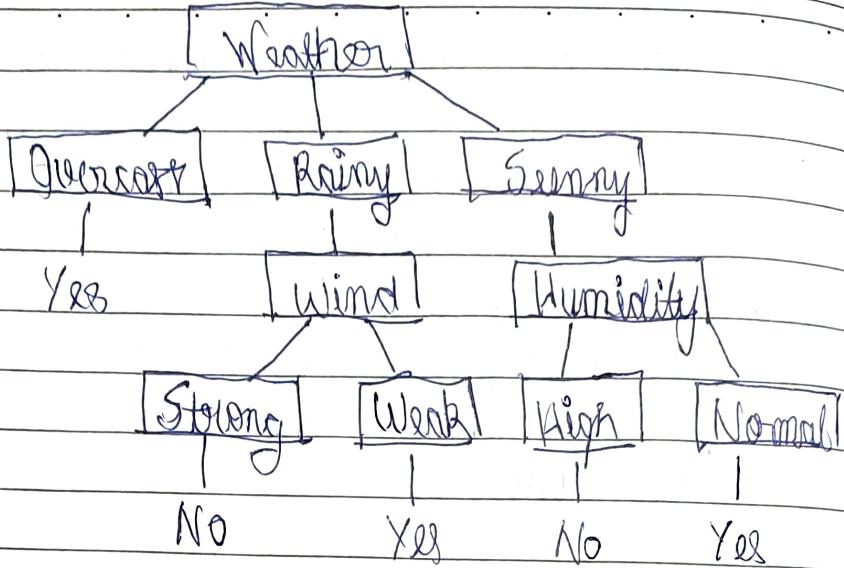
Number of columns: 5

Number of rows: 200

ID3 Algorithm:

function ID3 (Data, Target Attribute):

- 1) If all instances in Data have the same class
 - a) Return the class.
- 2) If Attributes is empty
 - a) Return the most common class in Data
- 3) Select Best Attribute with the highest information gain
- 4) Create a decision node with Best Attribute
- 5) for each value v in Best Attribute:
 - a) Split Data where Best Attribute = v
 - b) Add a branch to the decision node:
 - i) If subset is empty:
 - Add a leaf with the most common class
 - ii) Else:
 - Add the result of ID3 (Subset, Target, Attribute - Best Attribute)
- 6) Return the decision tree



One Prediction that was made :-

- 'Weather': 'Sunny'
- 'Temperature': 'Cool'
- 'Humidity': 'High'
- 'Wind': 'Strong'

Prediction : No

Decision Tree explanation in this case:-

- 1) If the weather is Overcast, the decision is Yes (Direct classification)
- 2) If the weather is rainy, the next check is wind.
 - * If the wind is Strong, decision is No
 - * If the wind is Weak, decision is Yes
- 3) If the weather is sunny, the next check is humidity
 - * If the humidity is high, the decision is No
 - * If the humidity is normal, the decision is Yes.

Information Gain Formulas

$$S(T) = \sum_{v \in \text{values}(S)} \frac{|S_v|}{|S|} S(T_v)$$

$S(T)$: Entropy of the whole dataset

v = Possible values of attribute X

S_v = Subset of S where attribute $X = v$

$|S_v|$ = Number of instances in subset S_v

$|S|$ = Total number of instances in S .

LAB-4

Date _____
Page _____

KNN Algorithm :

function split (data, 0.2) :

 split_index ← length (data) × (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_instance, k) :

 distances = []

 for each train_instance in train :

 dist = euclidean (test_instance, train_instance)
 [features]

 distances.append (train_instance, dist)

distances.sort ()

 neighbours ← first R elements of
 distances

gettrain_neighbours

function predict (train, test_instance, k) :

 neighbors ← get_neighbours (train, test_instance, 3)

 count occurrences of each class in
 neighbors

 return class with max(count)

function `knn()`:

`train, test = split_dataset (dataset, 0.2)`

`K=3`

`predictions = []`

for each test instance in test set:

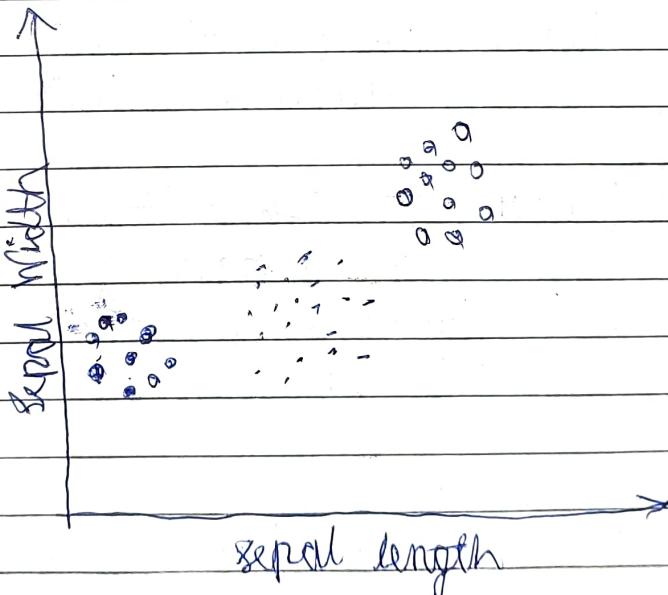
`predicted_class = predict (train, test instance features, k)`

`predictions.append (predicted_class)`

Dataset used is Iris.csv :-

Consists of 2 features = sepal length
sepal width

Multiclass target variable = Iris setosa
Iris versicolor
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 'versicolor'

dataset.append ([features, label])

return dataset

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

split_index = length (dataset) * (1 - test_ratio)

train_set = dataset [:split_index]

test_set = dataset [:split_index]

return train_set, test_set

function dot-product (a, b):

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

function train (train, learning_rate, lambda, param, epochs)

w \leftarrow [0, 0]

b \leftarrow 0

for epoch in epochs:

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

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

update $w = w - \text{learning-rate} \times 2x (y - 1) / \lambda$

param $\times w$

else :

update $w = w - \text{learning-rate} \times (x \text{ (symbol) - prediction}) \times w \cdot y \cdot x$

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

end if

return w, b

function predict(x, w, b) :

return 1 if dot product(w, x) + $b \geq 0$, else 0

filename = "content/iris.csv"

dataset = load_dataset(filename)

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

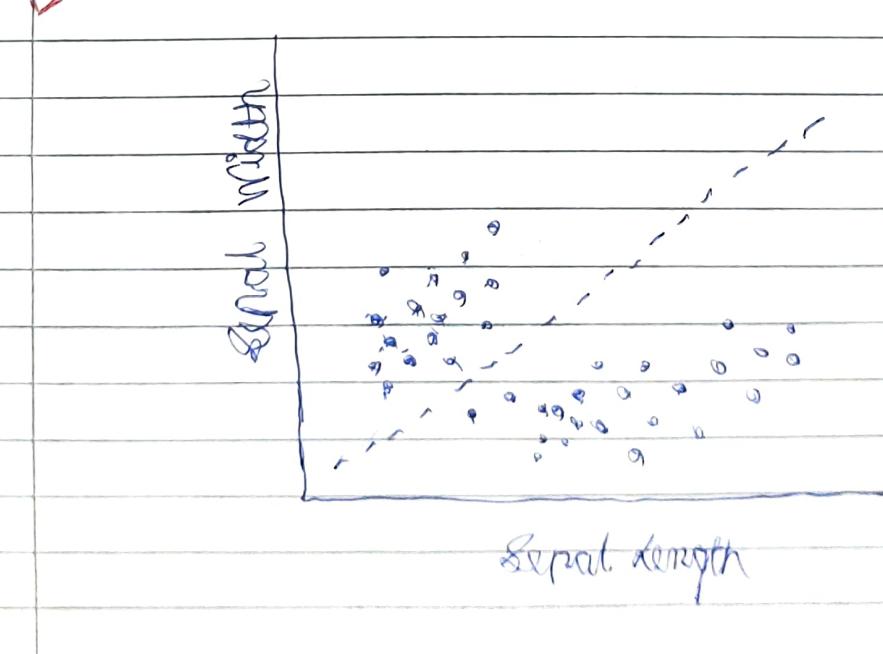
$w, b = \text{decision_sum}(\text{train})$

plot-sum(train, w, b)

for each point in test:

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

print("Actual, predicted")



LAB - 5

RANDOM FOREST :-

- 1) Load the dataset
- 2) Split the dataset into training and testing

function train test split (data, test size):
 test len = int (len(data) * test size)
 data_copy = shuffle(data)
 return training and test set.

- 3) GINI INDEX (groups, class labels)
 Initialize total variance as sum of the variances of all groups.
 Initialize gini = 0
 For each group in groups:
 if the group skip
 Calculate proportion of each class
 In the group
- 4) Test-split (index, value, dataset)
 - Initializing empty left and right groups
- 5) Get-split (dataset)
 Get the unique class labels from the dataset
- 6) To-terminal (group)
 - Return the most common class label from the group.

- 7) **SPLIT** (node, max_depth, min_size, current_depth)
 Extract left and right groups from node groups
- 8) **BUILD TREE** (train_data, max_depth, min_size)
 Initialize root node using GET_SPLIT on train_data
- 9) **PREDICT** (tree, row)
 If $\text{row}[\text{tree}, \text{index}] < \text{tree.value}$
 If tree.left is a dictionary (subtree).
 Else left
 Else return tree.left
 Else:
 if tree.right is a dictionary, decide right
 Else return tree.right.



LAB-6

Boosting Procedure

function AdaBoost:

 Input : X (Training features), y (Training labels), n_clf (Number of classifiers)

 Output : Trained AdaBoost model

Initialize an empty list 'clf' to store the decision stumps

Initialize sample weights:

$$\text{weights} = [1 / \text{number of samples}]$$

for each classifier in the range of n_clf :

 Initialize best classifier to store the best decision stump

 Initialize min_error to a large value

 For each feature in the dataset

 for each unique threshold in that feature

 for each polarity (+ or -)

 Make predictions based on the feature

 Calculate classifier weight (alpha):

$$\alpha_{clf} = \alpha * \log(1 - \min_error) / \max_error$$

 Make predictions using 'best_classifier' and update sample weights

Function predict:

Input: X , Trained AdaBoost model
 Initializing cif vector to store prediction of each classifier
 Get classifier prediction and multiply by classifier weight.

Main Program:

Call `AdaBoost.fit()` on the training set
 Call `AdaBoost.predict()` on the test set

Calculate accuracy: point accuracy

Plot accuracy versus no. of classifiers.

K MEANS PSEUDOCODE

function KMeans($x, k, \text{max_iters} = 100, \epsilon = 10^{-4}$)

Input: X (data points, matrix of size $m \times n$ samples, n features)

k : Number of clusters

~~Step 1: Randomly initialize the centroids~~

~~Step 2: Repeat for max iterations or until convergence.~~

~~Step 3a: Assign each point to nearest centroid~~

~~Calculate the distance from each data point to each centroid~~

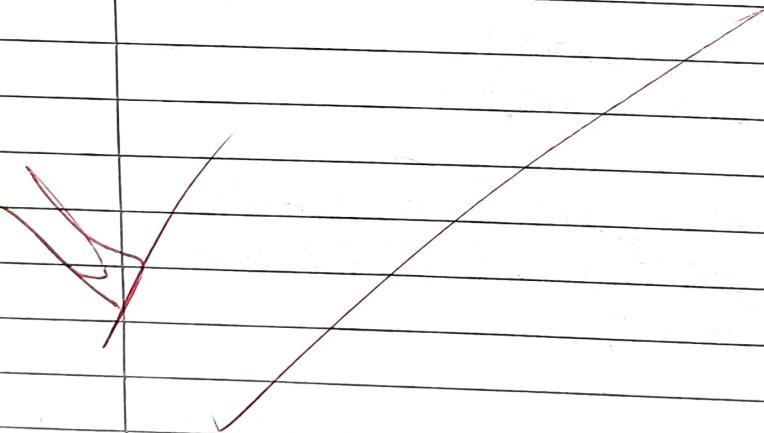
~~Step 3b: Update the centroids~~

~~Step 3c: Check for convergence~~

~~If the centroids don't change significantly break the loop.~~

Step 3 d : If convergence is not reached
update centroids and repeat.

Between the centroids and the label of the
data points



PCA (Principal Component Analysis)

- 1) Standardize the dataset (mean = 0, std = 1)
for each feature, in X :
Subtract mean and divide by standard deviation
- 2) Compute the covariance matrix of the standardized data.
 $\text{cov} = (X \cdot T + X) / (\text{n-samples} - 1)$
- 3) Compute eigenvalues and eigen vectors of the covariance matrix.
 $\text{eig_vals, eig_vecs} = \text{eig}(\text{cov})$
- 4) Sort eigen vectors by decreasing eigen values
(Workingly) Sort eig-vects according to sorted eig-val in decreasing order.
- 5) Select the top n eigen vectors.
 $w = \text{eig_vecs}[:, :k]$
- 6) Project the data onto the new subspace
 $X_{\text{reduced}} = X \cdot w$