

## Datasets

### 1. Healthcare Dataset

→ This synthetic healthcare dataset has been created to serve as a valuable resource for data science, machine learning & data analysis enthusiasts. The inspiration behind this dataset is rooted in the need for practical and diverse healthcare data for research purposes.

Attributes: Name, Age, Gender, Blood Type, Medical Condition, Date of Admission, Doctor, Hospital, Insurance Provider, Billing amount, Room Number, Admission Type, Discharge Date, Medication, Test Results.

No. of Rows = 55501

No. of unique values = 49992

No. of Columns = 15

### 2. Workout & Fitness Tracker Dataset

→ It is designed to help analyse and predict workout efficiency based on user activity, health metrics and lifestyle factors.

Attributes: User ID, Age, Gender, Height (cm), Weight (kg), Workout Type, workout duration (mins), calories burned, heart rate (bpm), steps taken, distance (km), workout intensity, sleep hours, water intake (liters), daily calories intake, resting heart rate (bpm), VO2 Max, Body Fat (%), mood Before workout, mood After workout.

No. of Rows = 10000

No. of columns = 21

### 3. Fraud Detection Transactions Dataset

This dataset is designed to help data scientists and machine learning enthusiasts develop robust fraud detection models.

Attributes : transaction-id, user-id, transaction-amount, transaction-type, timestamp, account-balance, device-type, location, merchant-category, ip-address-flag, previous-fraudulent-activity, daily-transaction-count, avg-transaction-amount-7d, failed-transaction-count-7d, card-type, card-age, transaction-distance, authentication-method, risk-score, is-weekend, fraud-label.

No. of total rows = 50001

No. of columns = 21.

### 4. Diabetes Dataset

This dataset includes various health parameters, lifestyle habits and genetic predispositions that contribute to diabetes risk.

Attributes: Age, Pregnancies, BMI (Body Mass Index), glucose, Blood Pressure, HbA1c (Hemoglobin A1c level), LDL (low-density lipoprotein), HDL (High-Density Lipoprotein), Triglycerides, waist Circumference, Hip circumference, WHR (Waist-to-Hip Ratio), Family History, Diet Type, Hypertension, Medication Use, Outcome

No. of rows = 9539

No. of columns = 17

## 5. Student Performance & Learning Style Dataset

This dataset provides insights into how different style habits, learning styles, & external factors influence student performance.

Attributes: student-id, age, gender, study-hours-per-week, preferred-learning-style, online-courses-completed, participation-in-discussions, assignment-completion-rate(1:1), Exam-score(1:1), Attendance-Rate(1:1), use-of-educational-tech, self-reported-stress-level, time-spent-on-social-media, sleep-hours-per-night, final-grade

No. of rows = 10001

No. of columns = 15

88  
3/03/25

## Lab-2

## Linear Regression

- Linear Regression is a statistical method used to model the relationship between a dependent variable and one or more independent variables.
- Linear Regression is also a type of supervised machine-learning algorithm that learns from the labelled datasets and maps the data points with most optimized linear functions which can be used for prediction on new datasets.
- It computes linear relationship between the dependent variable & one or more independent features by fitting a linear equation with observed data.

## Advantages:

- Simplicity & ease of implementation
- Interpretability
- Computationally Efficient
- Baseline Model
- Feature importance

Types of Linear Regression

Simple LR

Multiple LR

⇒ Single LR: It is the simplest form of LR & it involves only one independent variable & one dependent variable. The equation for simple LR is:

$$y = \beta_0 + \beta_1 x$$

$y \rightarrow$  dependent variable       $\beta_0 \rightarrow$  Intercept

$x \rightarrow$  independent variable       $\beta_1 \rightarrow$  Slope

Code:

```
import numpy as np  
import pandas as pd  
import matplotlib.pyplot as plt
```

```
df = pd.read_csv("experience-salary.csv")
```

```
feature_col = df.columns[0]
```

```
target_col = df.columns[1]
```

```
X = df[feature_col].values
```

```
Y = df[target_col].values
```

```
mean_X = np.mean(X)
```

```
mean_Y = np.mean(Y)
```

```
mean_XY = np.mean(X*Y)
```

```
mean_X2 = np.mean(X**2)
```

```
b1 = (mean_XY - mean_X * mean_Y) / (mean_X2 - mean_X**2)
```

```
b0 = mean_Y - b1 * mean_X
```

```
print("Intercept(b0): {b0:.4f}")
```

```
print("Slope(b1): {b1:.4f}")
```

```
print("Linear Regression Equation: {}target_col{} =  
{b0:.2f} + {b1:.2f}*{}feature_col{}")
```

~~plt.scatter(x, y, color="blue", marker="o", s=30,  
label="Actual Data")~~

~~plt.plot(x, b0+b1\*x, color="red", label="Regression Line  
(Mean-Based)")~~

plt.xlabel('feature-col')

plt.ylabel('target-col')

plt.title("Linear Regression using mean-based formula")

plt.legend()

plt.grid(True)

plt.show()

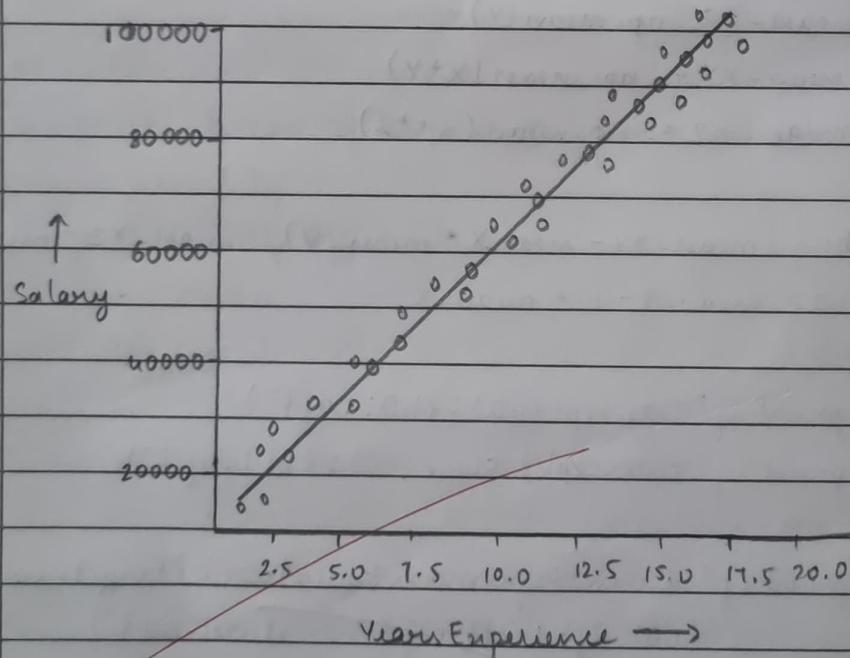
Output:

Intercept ( $b_0$ ): 8775.4214

Slope ( $b_1$ ): 4301.3422

Linear Regression Eq: Salary = 8775.42 + 4301.34 \* YearsExperience

LR using mean-based Formula



86  
10/21/20

Lab-3

## Housing [Chapter-2; End to End ML]

The script focuses on housing price prediction using ml techniques in Python. It follows the end-to-end ML pipeline from data loading to model evaluation. Below is a structured breakdown:

### 1. Data Loading & Exploration

- Loads the California housing dataset from a csv file (`housing.csv`) which contains 20,640 records (rows) and 10 attributes (columns).
- Displays basic statistics using `.head()`, `.info()` & `.describe()`.
- Analyzes categorical feature distribution (`ocean_proximity.value_counts()`)
- Plots histograms to visualize numerical feature distributions.

### 2. Data Splitting

- Uses random shuffling (split-train-test) and hash-based splitting (split-train-test-by-id) for reproducibility.
- Implements stratified sampling based on income categories to maintain proportional representation.

### 3. Data Visualization

- Scatter plots illustrate housing locations based on longitude & latitude.
- Color-coded visualizations highlight population density.

and median house value.

- Scatter matrix explores relationships between key features.

#### 4. Feature Engineering

- Creates new derived features:

- \* Rooms per household =  $\text{total\_rooms} / \text{households}$
- \* Bedrooms per room =  $\text{total\_bedrooms} / \text{total\_rooms}$
- \* Population per household =  $\text{population} / \text{households}$

- Handles missing values using median imputation

#### 5. Correlation Analysis

- Computes Pearson's correlation coefficient between numerical features & median house value.

#### 6. Data Preprocessing

- Uses OneHotEncoder for categorical variables (ocean-proximity).

- Applies pipelines to automate preprocessing:
  - \* Numerical pipeline: Handles missing values, feature engineering, & standardization.
  - \* Categorical pipeline: Applies one-hot encoding

- Uses ColumnTransformer to combine numerical & categorical processing.

## 7. Model Training & Evaluation

- Trains linear regression to predict house prices
- Evaluates performance using root mean squared error (RMSE)
- Performs cross-validation for better performance estimation.

✓ 113/25

Lab 4

## 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 in 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 : 1.5174

Temperature : 1.5567

Humidity : 1.0000

Windy : 0.9852

Information Gain for each feature:

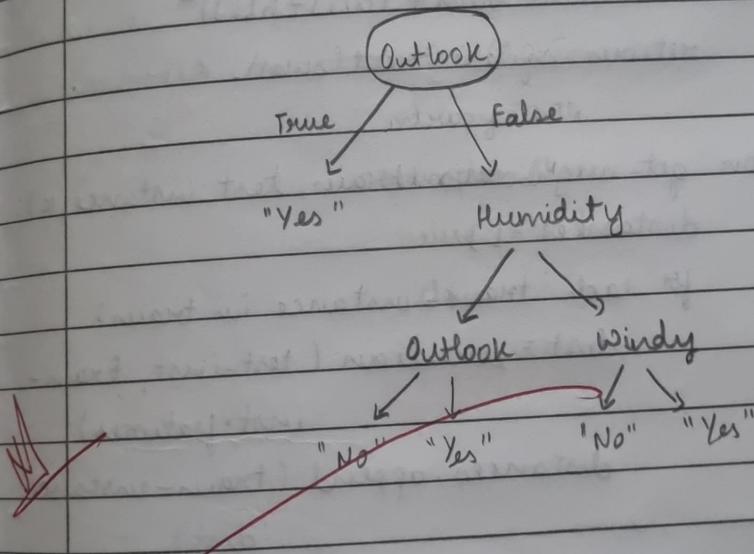
Outlook : 0.2467

Temperature : 0.0292

Humidity : 0.1518

Windy : 0.0481

Root Node : Outlook



## 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 euclidian(a, b):

sum = 0

for i=1 to length(a):

sum = sum + (a[i] - b[i])<sup>2</sup>

return square-root(sum).

function get-neighbours(train, test-instance, k):

distances = []

for each train-instance in train:

dist = euclidian(test-inst, train-inst.features)

distances.append((train-instance, 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 occurrences of each class in neighbours

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 use is iris.csv

feature 1 = sepal length

feature 2 = sepal width

Predicted class for [5.1, 3.5, 1.4]: Iris-versicolor.

## 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-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):
    returns sum (a[i] × b[i]) for all
```

```
function train-sum (train, learning-rate,
lambda-param, epochs)
    w ← [0, 0]
    b ← 0
    for epoch in epochs:
        for each data point (x, y) in train:
            if (y × (dot-product (w, x) + b)) ≥ 1:
                update w = w - learning-
                rate × 2 ×
                lambda-param ×
                w
```

else :

update  $w = w - \text{learning\_rate} \times (2x)$

$\lambda\text{ambda\_param} = w - y \times x$

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

-1

filename = "content / iris.csv"

dataset = load-dataset (filename)

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

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

plot-summ (train,  $w, b$ )

for each point in test:

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

print (Actual, Predicted)

W

Lab - 6

Random Forest

Random Forest Classifier ( $x$ -train,  $y$ -train,  $T$ , max-depth, min-samples-split):

trees = []

$n$ -samples,  $n$ -features =  $x$ -train.shape

for  $t=1$  to  $T$ :

bootstrap-idn = random.sample(range(n-samples), n-samples)

$x$ -bootstrap =  $x$ -train[bootstrap-idn]

$y$ -bootstrap =  $y$ -train[bootstrap-idn]

tree = Decision Tree (max-depth = max-depth,  
min-samples-split = min-samples-split)

tree.fit( $x$ -bootstrap,  $y$ -bootstrap)

trees.append(tree)

return trees.

Predict ( $x$ -test):

predictions = []

for tree in trees:

tree-predictions = tree.predict( $x$ -test)

predictions.append(tree-predictions)

final-predictions = majority-vote(predictions)

return final-predictions

majority-vote(predictions):

return mode(predictions)

## Lab-7

### Boosting Ensemble [Adaboost]

1. Initialize weights for each training sample:
  - set each sample's weight to  $1/N$  (where  $N$  is the no. of samples).
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's error:
    - $\text{Error} = (\text{sum of weights of misclassified samples}) / (\text{sum of all sample weights})$
  - c. calculate the classifier's importance (alpha):
    - $\text{Alpha} = 0.5 * \log((1 - \text{Error}) / \text{Error})$
    - If Error is high, alpha will be low, meaning less importance for this classifier.
  - d. Update the sample weights:
    - 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 (alpha)
  - Final predictions = signs (sum of all classifiers)

4. use the combined classifier to make prediction on new data.

## k-Means Clustering Algorithm

1. Initialize k cluster centroids (randomly selected data points).

2. Repeat until convergence :

a. Assign step :

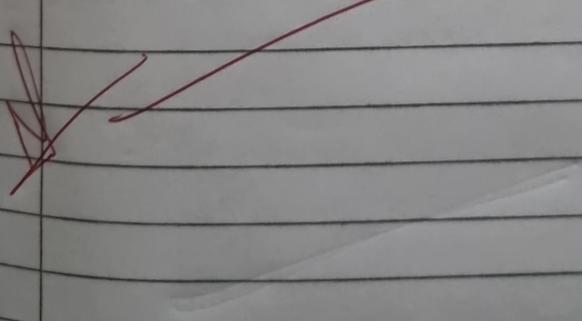
- For each data point, calculate the distance to each centroid.
- Assign each data point to the closest centroid (cluster).

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.



## Logistic Regression

1. Initialize weights  $w$  & bias  $b$  to small values  
(e.g., zeros)
2. For  $i=1$  to  $n$ -iter :
  - a. Compute linear combination :  $z = X \cdot w + b$
  - b. Apply sigmoid function :  $y_{\text{pred}} = 1 / (1 + \exp(-z))$
  - c. Compute binary cross-entropy loss :
$$\text{Loss} = -\left(\frac{1}{n}\right) * \sum [y * \log(y_{\text{pred}}) + (1-y) * \log(1 - y_{\text{pred}})]$$
  - d. Compute gradients :
$$dw = \left(\frac{1}{n}\right) * X \cdot T * (y_{\text{pred}} - y)$$

$$db = \left(\frac{1}{n}\right) * \sum (y_{\text{pred}} - y)$$
  - e. Update weights and bias :
$$w = w - \alpha * dw$$

$$b = b - \alpha * db$$
3. Output the final weights & bias.

## PCA (Principal Component Analysis)

1. standardize the dataset (mean = 0, std = 1)

for each feature in  $X$ :

subtract mean & 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 eigenvectors of the covariance matrix

$$\text{eig-vals}, \text{eig-vecs} = \text{eig}(\text{cov})$$

4. Sort eigenvectors by decreasing eigenvalues (variance)  
sort eig-vecs according to sorted eig-vals in decreasing order.

5. Select the top  $K$  eigenvectors

$$W = \text{eig-vecs}[:, :K]$$

6. Project the data onto the new subspace

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

## Multi-linear Regression

1. Initialize weight vector  $w$  & bias  $b$  to zeros
2. For  $i=1$  to  $n$ -iter:
  - a. Compute predictions:
$$y_{\text{pred}} = x \cdot w + b$$
  - b. Compute loss (Mean Squared Error):
$$\text{loss} = (1/n) * \sum (y_{\text{pred}} - y)^2$$
  - c. Compute gradients:
$$dw = (2/n) * x \cdot T * (y_{\text{pred}} - y)$$
$$db = (2/n) * \sum (y_{\text{pred}} - y)$$
  - d. Update weights & bias:
$$w = w - \alpha * dw$$
$$b = b - \alpha * db$$
3. Return the final weights  $w$  & bias  $b$ .