



Name Nidhi. A Std M.L Sec GD

Roll No. \_\_\_\_\_ Subject \_\_\_\_\_ School/College \_\_\_\_\_

School/College Tel. No. \_\_\_\_\_ Parents Tel. No. \_\_\_\_\_

Lab - 0

## I Method - 1: Initialising values directly into D.F

```
import pandas as pd
data = {'Name': ['Alice', 'Bob', 'Charlie', 'David'],
        'Marks': [25, 30, 45, 35],
        'USN': ['IBM177', 'IBM255', 'IBM240',
                'IBM203']}
```

```
df = pd.DataFrame(data)
print("Sample data:")
print(df.head())
```

OUTPUT

	Name	Marks	USN
1	Alice	25	IBM177
2	Bob	30	IBM255
3	Charlie	45	IBM240
4	David	35	IBM203

## 2. Method - 2: Importing datasets from sklearn dataset

```
from sklearn import datasets
import load_diabetes
diab = load_diabetes()
df = pd.DataFrame(diab.data, columns=diab.feature_names)
df['target'] = diab.target
print("Sample data:")
print(df.head())
```

Sample data:

age	sex	bmi	bp	s1	s2	s3	s4	s5	s6	target
0.038	0.050	0.061	0.02	-0.04	-0.03	-0.04	0.03	0.01	151	
-0.001	-0.044	-0.051	0.02	-0.08	-0.19	0.48	0.80	0.90	75	
0.085	0.050	0.044	-0.00	-0.04	0.18	0.34	0.40	-0.40	41	
0.005	-0.044	-0.011	-0.03	0.01	0.34	-0.18	0.44	0.44	206	

## 3) Method 3 Importing datasets from a specific .csv file

```
file-path = 'data.csv'
df = pd.read_csv(file-path)
print("Sample data")
print(df.head())
print("\n")
```

Sample data

ID	Name	Age	City
1	Alice	25	New York
2	Bob	30	Los Angeles
3	Charlie	45	Chicago
4	David	20	Houston
5	Eve	10	Washington DC

## ii) Method 4. Downloading datasets from Kaggle.

```
df = pd.read_csv("/content/dataset_of_diabetes.csv")
print(df.head())
```

ID	No-Patno	Gender	AGE	Weight	Gr	HbA1c
0	502	F	50	4.7	46	4.9
1	735	M	26	4.5	62	4.9
2	420	F	50	4.7	46	4.9
3	680	F	50	4.7	46	4.9
4	504	M	33	7.1	46	4.9

## TO DO

- ① Using the code, do exercise of "Stock Market Data Analysis":

```
import yfinance as yf  
import pandas as pd  
import matplotlib.pyplot as plt
```

```
tickers = ["HDFCBANK.NS", "ICICIBANK.NS",  
           "KOTAKBANK.NS"]
```

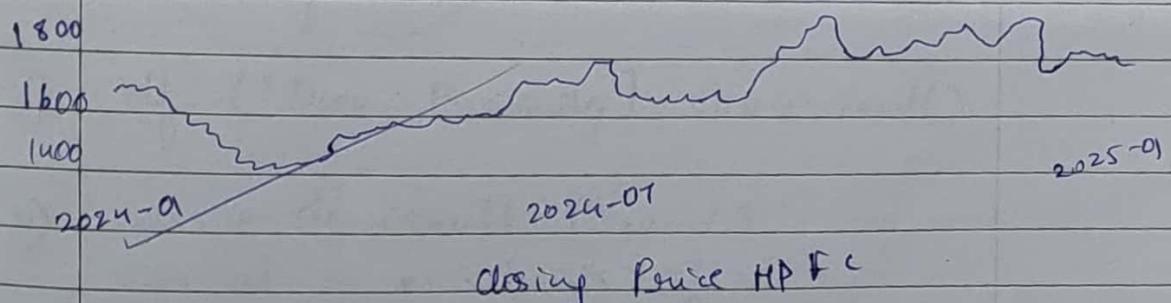
```
data = yf.download(tickers, start="2024-01-01",  
                   end="2024-12-31", group_by='ticker')
```

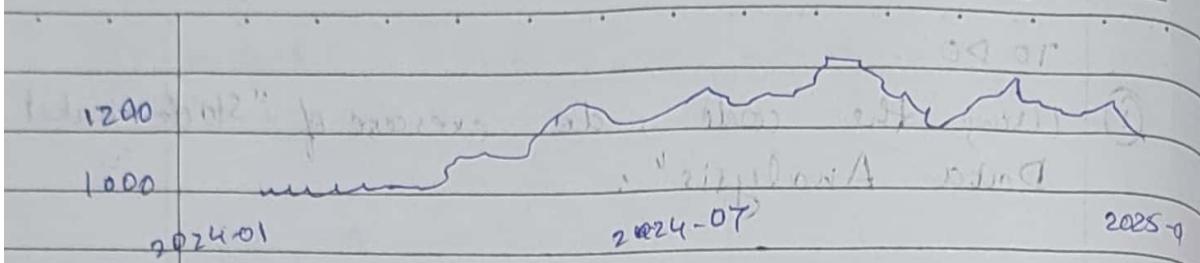
```
for ticker in tickers:  
    data.loc[:, (ticker, 'Daily Return')] =  
        data[ticker]['Close'].pct_change()
```

```
plt.figure(figsize=(12, 16))  
for i, ticker in enumerate(tickers):  
    plt.subplot(3, 1, i+1)
```

```
plt.tight_layout()  
plt.show()
```

```
plt.tight()  
plt.show()
```

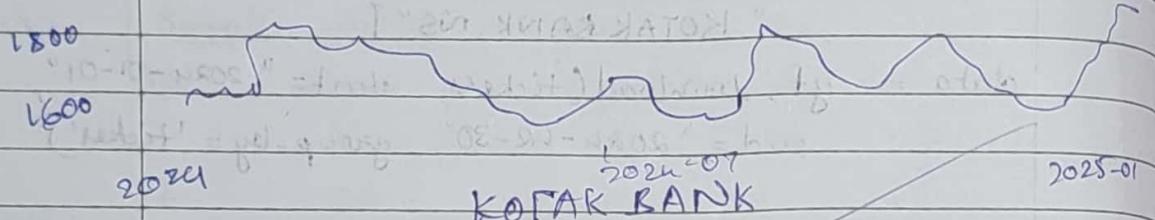




ICICI closing price

by 2024-07

KOTAK



KOTAK BANK

CLOSING PRICE

TC market price (short) short term

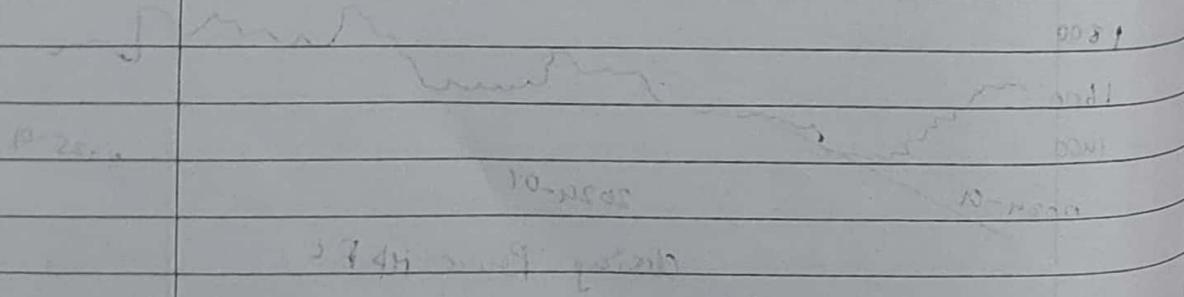
(medium term) long term  
03.03

Goal bank

((1.03) - 1.01) (1.03) - 1.01  
(1.03) - 1.01  
(1.03) - 1.01

((1.03) - 1.01)  
((1.03) - 1.01)

((1.03) - 1.01)  
((1.03) - 1.01)



KAIB-94

Diabetes

Point(df.isnull().sum())

0/7P

adult income

age

Class

income

dtype: int64

dtype: int64

There are no missing values.

2. Gender &amp; Class (Diabetes)

↓

ordinal encoder

onehot encoder

ordinal encoder = Ordinal Encoder (categories: "F", "M")  
 "M" handle unknown = "use encoded value"  
 unknown\_value = -1

(1) Write python code consider filename as  
 to read csv

df = pd.read\_csv("housing.csv")

ii) display all columns  
 df.info()

iii) statistical information  
 df.describe()

iv) df[["Own Proximity"]].value.count()

v) miss = df.isnull().sum()  
 miss\_col = miss[miss > 0]  
 print(miss\_col)

Q i) which had missing value? How did you handle them

Diabetes:

columns: BMI BP, Glucose

they were replaced with median

ii) What categorical cols did you identify?  
How was it encoded?

Diabetes:

Gender Ethnicity,

Encoding Method.

Gender Binary (- Male 0 - Female)

Adult Income:

sex Work

iii) Min max scaling: Rescales data to fixed range.

standardization: Rescales data to have a mean of 0 & std

→ min max is used when a fixed range is given standardization does not work.

$\text{Scale} = \frac{x - \text{min}}{\text{max} - \text{min}}$

Iterative Dichotomiser (ID3) Algorithm.Algorithm

```
def ID3 (D, A):
```

if  $D$  is pure or  $A$  is empty:

return a leaf node with the majority class in  $D$ .

else:

$A_{best} = \arg\max(\text{InformationGain}(D, A))$

root = Node( $A_{best}$ )

for  $v$  in values( $A_{best}$ ):

$D_v = \text{subset}(D, A_{best}, v)$

child = ID3( $D_v, A - \{A_{best}\}$ )

root.add child(v, child)

return root.

ID3 algorithm is a decision tree algorithm used in machine learning and mining for classification tasks.

1. Prepare the Data:

- Each instance has a set of features (attribute)
- Each instance has a target class that you want to predict

2. Select the Best Attribute (Feature)

Check if the dataset is pure:

- If all instances in the dataset have the same class label, return that class as the result (leaf node)

3. Check if the dataset is empty:

- If empty, return most frequent class

4. Check if no more features to split  
 - If no, return the most freq.

5. Calculate  $I(G)$ :

$$I(G(DA)) = \text{Entropy}(D) - \sum_{V \in D} \frac{|D_V|}{|D|} \times \text{Entropy}(P_V)$$

$$\text{Entropy} = p(c_i) \times \log(p(c_i))$$

6. Choose feature with highest  $I(G)$   
 → is the splitting criterion

7. Split the dataset based on the selected feature.

8. Recursively apply ID3 to each subset

9. Prune the tree by removing branches that hold no value.

Example)	Weather	Play outside
	Sunny	Yes
	Rainy	No
	Overcast	Yes
	Sunny	Yes
	Rainy	No

$$G(\text{Play outside}) = - \left( \frac{3}{5} \log \frac{3}{5} + \frac{2}{5} \log \left( \frac{2}{5} \right) \right) \\ = 0.971$$

Weather	Ye	No	Entropy
Sunny	2	0	0.
Rainy	0	2	0.
Overcast	1	0	0

Weighted entropy

$$G(\text{Sweater}) = \left(\frac{2}{3} \times 0\right) + \left(\frac{2}{3} \times 0\right) + \left(\frac{1}{3} \times 0\right) = 0$$

$$\text{Information Gain} = 0.970 - 0 = 0.970.$$

Weather is the best attribute to split on

Weather		
Sunny	Rainy	Overcast
Yes	No	Yes

~~(to get weight of attribute)~~  
~~for each attribute calculate information gain~~  
~~10/3/25~~

~~(typical accuracy 0.7 to 0.8 in job)~~

~~1 = (0.970, 0.970, 0.970)~~

~~(classical aspect) the number~~

~~represent true 1~~

~~10/3/25 (above) represent the number~~

~~represent 0.970 with maximum weight & 0.970~~

~~represent 0.970 with minimum weight & 0.970~~

~~(the more weight a thing have more probability)~~

~~(if weight of a thing is zero)~~

17/03/25

Date \_\_\_\_\_  
Page \_\_\_\_\_

## Implementation of ID-3

```

import pandas as pd
import numpy as np

data = pd.read_csv("weather.csv")
df = pd.DataFrame(data)

def entropy(data):
    class_counts = data.value_counts()
    prob = class_counts / len(data)
    return -np.sum(prob * np.log2(prob))

def information_gain(df, feature, target):
    total_entropy = entropy(df[target])
    feature_value = df[feature].unique()
    weighted_entropy = 0
    for value in feature_values:
        subset = df[df[feature] == value]
        weight = len(subset) / len(df)
        weighted_entropy += weight * entropy(subset)
    return total_entropy - weighted_entropy

def id3(df, features, target):
    if len(df[target].unique()) == 1:
        return df[target].iloc[0]
    if not features:
        return df[target].mode()[0]
    gains = {feature: information_gain(df, feature, target) for feature in features}
    best_feature = max(gains, key=gains.get)
    tree = {best_feature: df}

```

```
for value in df[best_failure].unique():
    subset = df[df[best_failure] == value]
    tree[best_feature][value] = id3(subset, [f
        for f in features if f != best_feature],
        target)
return tree
```

features = ['Outlook', 'temperature', 'Humidity', 'Wind']
target = 'PlayTennis'
desc\_tree = id3(df, features, target)
print(desc\_tree).

### OUTPUT

```
{'Outlook': { 'sunny':
{ 'Humidity': { 'High': { 'No': 'No', 'Yes': 'Yes' } },
'Overcast': 'Yes', 'Rain': { 'Wind': { 'Weak': 'Yes',
'strong': 'No' } } } }}
```

I)

## Lab - 3 : End to End ML Project

### ① Get Data

```
from sklearn.datasets import fetch_california_housing
cal_house = fetch_california_housing(as_frame=True)
print(cal_house.DESCR)
print(cal_house.head())
```

### ② Discover & Visualise the data to gain insights

```
import matplotlib.pyplot as plt
import seaborn as sns

print(cal_house.frame.isnull().sum())

plt.figure(figsize=(8, 6))
sns.histplot(cal_house.frame['MedHouseVal'],
             kde=True)
plt.title('Dist')
plt.show()
```

### ③ Prepare the data for ML Algo

```
from sklearn.model_selection import train_test_split
from sklearn.preprocessing import StandardScaler
```

```
X = cal_house.frame.drop(['MedHouseVal'], axis=1)
y = cal_house.frame[['MedHouseVal']]
```

```
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=42)
```

scaler = StandardScaler()

X-train-scaled = scaler.fit\_transform(X-train)

X-test-scaled = scaler.transform(X-test)

④ Select a Model and Train it.

from sklearn.linear\_model import LinearRegression

from sklearn.metrics import mean\_squared\_error,  
r2\_score

linear-model = LinearRegression()

linear-model.fit(X-train-scaled, y-train)

y-pred-linear = linear-model.predict(X-test-scaled)

mse-linear = mean\_squared\_error(y-test, y-pred-linear)

r2-linear = r2\_score(y-test, y-pred-linear)

print(f'mse-linear: {mse-linear}, r2-linear: {r2-linear}')

⑤ Fine Tune Your Model

from sklearn.metrics import root\_mean\_squared\_error

import numpy as np

y-pred = model.predict(X-test)

rmse = root\_mean\_squared\_error(y-test, y-pred)

print(f'rmse: {rmse}')

rmse  
y-test  
y-pred

## Lab - 4.

### Simple Linear Regression

function linear-reg(x, y, n) {

m = 0

b = 0

n = length(x)

for i ← 1 to n:

pred = [i, x[i]]

for j ← 0 to n-1:

pred[j] = m \* [j] + b

append pred to pred

error = {}

for j ← 0 to n-1:

error[j] = pred[j] - y[j]

error.append(error)

cost =  $(1/n) * \sum(\text{error}^2)$  for error in error

m-gradient =  $\frac{2}{n} \times \sum(\text{error} \cdot x[j])$  for j ← 0 to n-1

b-gradient =  $\frac{2}{n} \times \sum(\text{error})$

m = m - m-gradient \* g

b = b - b-gradient \* g

return m, b,

logistic linear Regressionfunction logistic\_regression ( $x, y, a, num$ ): $m = \text{length}(x[0])$  # features $n = \text{length}(x)$ weights = array of 0 to length  $m$ 

bias = 0

for  $i$  1 to num:

pred = {}

for  $j \leftarrow 0$  to  $n-1$ :lin = dot\_product(weights,  $x_j$ )  
+ bias

pred = sigmoid(lin)

predicted.append(pred)

errors = {}

for  $j \leftarrow 0$  to  $n-1$ :error = pred[i] -  $y[i]$ 

errors.append(error)

cost =  $-(\frac{1}{n}) * \sum(y[i] * \log(\text{prediction}[i]))$   
 $+ (1 - y[i]) * \log(1 - \text{pred}[i])$   
for  $j$  from 0 to  $n-1$ )weight\_grad =  $(\frac{1}{n}) * \text{dot\_product}(\text{errors}, x)$ bias\_grad =  $(\frac{1}{n}) * \sum(\text{errors})$ weights = weight -  $\alpha * \text{weight\_grad}$ bias = bias -  $\alpha * \text{bias\_grad}$ 

return weights, bias

## Algo linear Regression

1. Initialise  $m \& b = 0$
2. Make prediction  $m x + b$ .
3. Find error by  $y - \text{pred}$ .
4. Update  $m \& b$  using gradient descent
5. Repeat 4 steps for num iterations

## Algo logistic Regression

1. ~~Let~~  $m$  be the length of features initialized to 0
2.  $n = \text{no. of datapoints}$
3. Initialize weight & bias to 0.
4. Make prediction  $\text{weight} \cdot x + \text{bias}$
5. ~~Please~~ Use sigmoid function for pred.
6. Find errors using  $\text{pred} - y$   $\forall i$
7. Update weight & bias using gradient descent
8. Repeat steps 4 & return weight & bias

## Algo Multilinear Regression

~~Same as linear but with multiple features.~~

## Multilinear Regression

function multilinear-reg ( $x, y, \alpha, \text{num}$ ):

$$m = \text{length}(x[0])$$

weights = array of 0 of length m

bias = 0

$$n = \text{length}(x)$$

for  $i \leftarrow 1$  to num:

$$\text{pred}[] = \{ \}$$

for  $j \leftarrow 0$  to  $n-1$ :

$$pr = \text{dot-product}(\text{weight}, x[i:j]) +$$
  
$$\text{bias}$$

pred.append(pr)

$$\text{errors}[] = \{ \}$$

for  $j \leftarrow 0$  to  $n-1$ :

$$\text{error} = \text{pred}[j] - y[j]$$

errors.append(error)

$$\text{cost} = \frac{1}{n} \times \text{sum}(y[j] \times \log(\text{pred}[j]) +$$
  
$$(1 - y[j]) \times \text{error}^2 \text{ in}$$
  
$$\text{error in errors})$$

$$\text{weight\_grad} = \frac{1}{n} \times \text{dot prod}(\text{error} \times)$$

$$\text{bias\_grad} = \frac{1}{n} \times \text{dot prod}(\text{sum(error)})$$

$$\text{weights} = \text{weights} - \eta \times \text{weight\_grad}$$

$$\text{bias} = \text{bias} - \eta \times \text{bias\_grad}$$

between weight & bias.

## K-Neighbour

## K-Nearest Neighbours

Algorithm.

1. Choose the number of neighbours (k):

2. Calculate the distance:

For each instance in training dataset,  
calculate dist b/w training dataset  
with test data.  $\text{dist} = \sqrt{(x-x')^2 + (y-y')^2}$

3. Sort distances in ascending order

4. Select k nearest neighbours.

5. For classification

Take a majority vote among the 'k'  
neighbours to determine predicted label.

6. For regression.

Take an average of k neighbour  
to return predicted value

7. Return the prediction.

Pseudocode

function kNN(train, test, k):

    distances = []

    for each (features, label) in train:

        distance = cal-dis(features, test)

        distances.append((dist, label))

    distances.sort(key=lambda x: [10])

    k-n-n = distances[:k]

```
if classifier has:  
    vote = {}  
    for neighbor in k-n-n:  
        label = neighbor[1]  
        if label in votes:  
            votes [label] += 1  
        else:  
            votes [label] = 1  
predicted-label = max (votes, key=votes.get)  
return predicted-label
```

```
else:  
    total = 0  
    for neighbor in k-n-n:  
        total += neighbor[1]  
predicted-val = total / k  
return predicted-val.
```

function cal-dist (features1, features2):

```
return sqrt (sum ((f1-f2)**2 for f1, f2 in  
zip (features1, features2)))
```

### SVM

#### Algorithm

- ① Input Data
- ② Choose Kernel Function
- ③ Formulate the Optimization Problem
- ④ Solve Optimization Problem
- ⑤ Make Predictions
- ⑥ Output the model

function svm (train, labels, kernel, c):

n = length (train)

alpha = array of zeros of size n

b = 0

w = array of zeros of no. of features

for i from 0 to n-1:

if alpha[i] > 0:

b = labels[i] - dot\_prod(w, feature[i])

break

return (w, b)

function predict (model, test):

(w, b), model

des\_val = dot\_prod(w, test) + b

if des\_val >= 0:

return +1;

else

return -1

function dot\_prod (vect1, vect2):

result = 0

for i from 0 to len(vect1)-1:

res += vect1[i]\*vect2[i]

return res

## i) Random Forest Ensemble Method

1. Load dataset  $D$  with features  $X$  and target  $Y$
2. Preprocess the data:
  - Handle missing value
  - Encode categorical variables
  - Split dataset into training set  $D_{train}$  and test set  $D_{test}$ .
3. Initialise RandomForest with parameters:
  - n\_estimators (no. of trees)
  - max\_depth (max depth of trees)
  - random\_state (seed for replicability)
4. For each tree  $t=1$  to  $n_{estimators}$ :
  - Randomly select a bootstrap sample of data points from  $D_{train}$
  - For each node in the tree:
    - i) Randomly select a subset of feature
    - ii) Calculate Gini Impurity or Entropy for each feature and split based on the best value
      - Gini Impurity:  $Gini = 1 - \sum (p_i)^2$
      - Entropy:  $Entropy(S) = - \sum (p_i * \log_2(p_i))$
    - Split data at the best feature & continue until max depth is reached or stopping condition are met.
5. After all trees are built, make predictions:
  - a) For each test instance  $x_{test}$ , aggregate

predictions from all trees using majority voting

i) Prediction = mode (prediction from all tree)

6. Evaluate the model:

Calculate accuracy, precision, recall, F1 score.

7. Return Model performance

Boosting Ensemble Method.

1. Load dataset  $D$  with features  $X$  and target  $Y$ .
2. Preprocess the data:
  - a. Handle missing values
  - b. Encode categorical variables
  - c. Split dataset into training set  $D_{train}$  and test set  $D_{test}$ .
3. Initialise the weight of all training samples equally.  
 $w_1 = w_2 = \dots = w_N = 1/N$   
 where  $N$  is the no of training sample
4. For each iteration  $t=1$  to  $n$  estimator.
  - a. Train a weak classifier  $h_t$  using the weighted dataset  $D_{train}$  with sample weights  $w_i$ .
  - b. Calculate the error of the classifier  

$$\epsilon_t = \left( \sum_{i=1}^N w_i \mathbb{I}(h_t(x_i) \neq y_i) \right) / \sum w_i$$

where  $\mathbb{I}$  is the indicator function
  - c) Compute the classifier weight ( $\alpha_t$ ):  

$$\alpha_t = 0.5 * \log((1 - \epsilon_t) / \epsilon_t)$$
  - d) Update the weights of misclassified points:  

$$w_i' = w_i * \exp(\alpha_t * \mathbb{I}(h_t(x_i) \neq y_i))$$
  - e) Normalize the weight so that they sum to 1:  

$$w_i = w_i' / \sum w_i'$$

5. After all iteration, combine weak classifiers:

a) Final prediction for test point

$$x = \text{sign}(\sum_{t=1}^T h_t(x))$$

6. Evaluate the model

a) Calculate accuracy, precision, recall, F1-score

7. Retrain model performance

Lab - 8

## 1. K-means Clustering Algorithm.

1. Load dataset  $D$  with data-points  $x_i$ , where  $x_i \in \mathbb{R}^d$  ( $d$ -dimensional space)

## 2. Preprocess the data:

- Handling the missing value
- Standardize the data

3. Initialize  $k$  centroids randomly from the dataset:

$c_1, c_2, \dots, c_k$

## 4. Repeat until convergence (centroids no longer change):

a. For each data point  $x_i$ :

- Calculate the distance to each centroid  $c_j$ :

$$\text{distance}(x_i, c_j) = \|x_i - c_j\|_2$$

- Assign point  $x_i$  to the closest centroid:

$$\text{Cluster}(x_i) = \arg \min_j \text{distance}(x_i, c_j)$$

b. Update centroids by calculating the mean of points in each cluster:

$$c_j = (1/|S_j|) * \sum_{x_i \in S_j} x_i$$

where  $S_j$  is the set of points assigned to centroid  $c_j$ .

5. Convergence: If centroids do not change significantly b/w iterations, stop.

6. Output.

a. Final centroids  $C_1, C_2, \dots, C_k$

b. Cluster assignments for each data point

7. Optionally, calculate clustering metrics  
(e.g., silhouette score, inertia).

Half 9

## Principle Component Analysis

1. Load data  $D$  with data points  $x_i$ , where  $x_i \in \mathbb{R}^d$  ( $d$ -dimensional space)
2. Preprocess the data:
  - a. Handle missing values
  - b. Standardize the data by subtracting mean of each feature  

$$x_j = (x_j - \mu_j) / \sigma_j$$

where  $\mu_j$  is the mean of feature  $j$  &  $\sigma_j$  is the standard deviation
3. Compute the covariance matrix  $\Sigma$  of standardized data:  

$$\Sigma = (1/N) * X^T * X$$

where  $X$  is the data matrix (each row is a data point)
4. Perform eigenvalue decomposition on the covariance matrix  $\Sigma$ .  

$$\Sigma * v = \lambda * v$$

where  $\lambda$  are the eigenvalues, and  $v$  are the eigenvectors
5. Sort the eigenvalues in descending order and select the top  $k$  eigenvectors
6. Project the data onto the  $k$  eigenvectors (PCA).  

$$X_{\text{new}} = X * V$$

where  $V$  is the matrix of the top  $k$  eigenvectors

Output:

Transformed dataset  $X_{\text{new}}$  (reduced dimensionality)

Optionally, visualize the reduced data ( $k=2$  or  $k=3$ ).

