

60004210210

Amartya Mishra

COMPS – C31

ML Experiment 1

ML - Experiment - 1

~~2 Weeks
11/11/2023~~

Aim: To perform Data cleaning & preprocessing on the Dataset & implement Linear Regression

Theory:

Data cleaning is the process of fixing or removing incomplete, incorrect, corrupted, incorrectly formatted, duplicate or incomplete data within a dataset.

While integrating multiple data sources, discrepancy is possible which makes outcomes & algorithms unreliable.

Steps to clean Data:

- Remove irrelevant data
- Remove duplicate & incomplete cases
- Fix structural errors
- Deal with missing Data
- Identify & review outliers
- Filter out outliers
- Encode categorical data
- splitting the Dataset
- feature scaling.

b. Linear Regression:

It is a special type of machine learning algorithm more specifically a supervised algorithm. It learns from the labelled dataset & maps the data points to the most optimized linear functions.

It computes the linear relationship between a dependent variable & one or more independent features.

Simple Linear Regression:

$$y_i = \beta_0 + \beta_1 x$$

y : dependent variable

x : Independent variable

β_0 : Intercept

β_1 : Slope

The goal of Algorithm is to find the best fit line equation that can predict the value based on independent variables.

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

$n \rightarrow$ NO. of datapoints

$y_i \rightarrow$ actual value

$\hat{y}_i \rightarrow$ Predicted value

Conclusion: Thus we understood steps of Data pre-processing & cleaning & thereafter implemented linear regression.

Implementation: Done in 3 ways, firstly on given data using direct formulas.

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

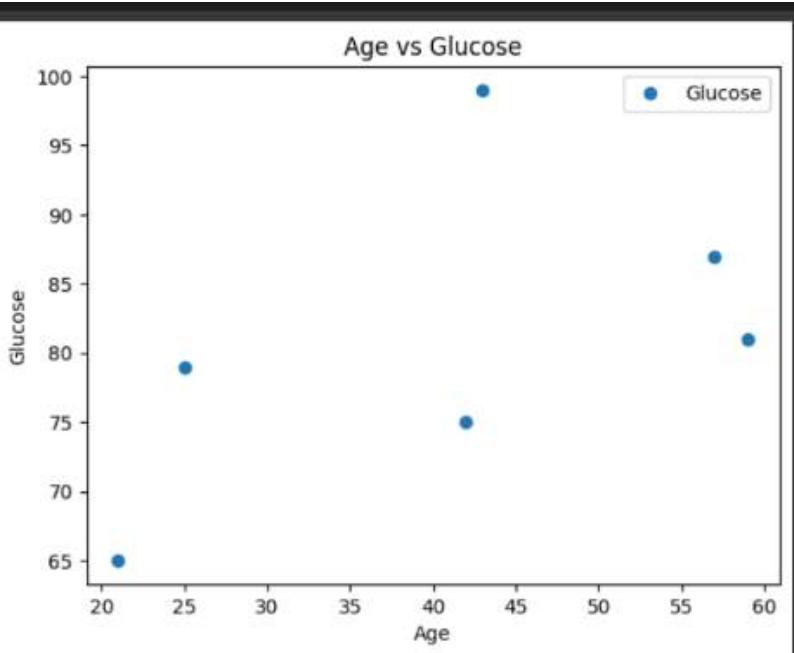
[ ] d = {'Age': [43, 21, 25, 42, 57, 59], 'Glucose': [99, 65, 79, 75, 87, 81]}

[ ] n = len(d['Age'])
print(n)

6

[ ] df = pd.DataFrame(data = d)

[ ] df.plot(x='Age', y='Glucose', style='o')
plt.xlabel('Age')
plt.ylabel('Glucose')
plt.title('Age vs Glucose')
plt.show()
```

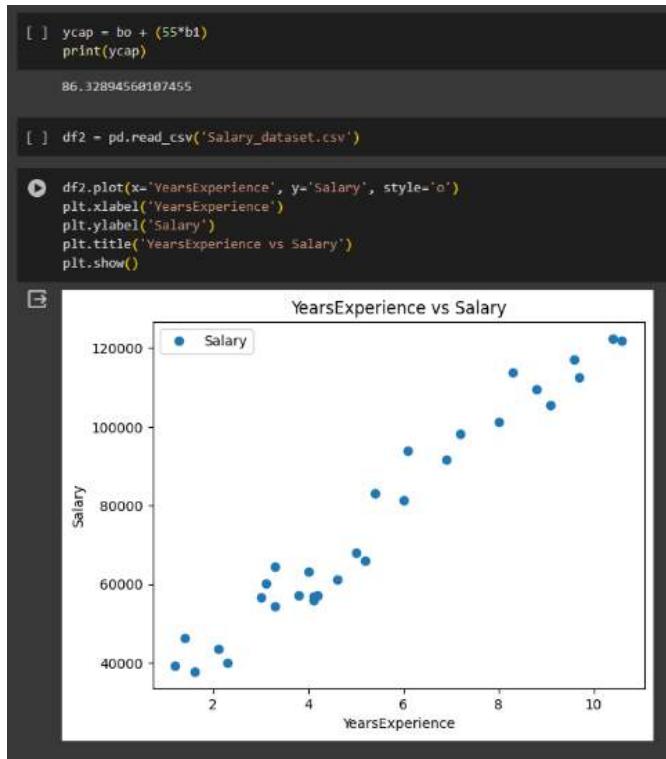


```
[ ] xs=ys=xy=x=y=0

[ ] for i in range(0,n):
    x += d['Age'][i]
    y += d['Glucose'][i]
    xs += d['Age'][i]**2
    ys += d['Glucose'][i]**2
    xy += d['Age'][i]*d['Glucose'][i]
```



Then on an imported dataset using inbuilt functions and libraries



```
[ ] from sklearn.model_selection import train_test_split
from sklearn.linear_model import LinearRegression
X = np.array(df2['YearsExperience']).reshape(-1, 1)
y = np.array(df2['Salary']).reshape(-1, 1)

[ ] X_train, X_test, y_train, y_test = train_test_split(X, y, test_size = 0.2)

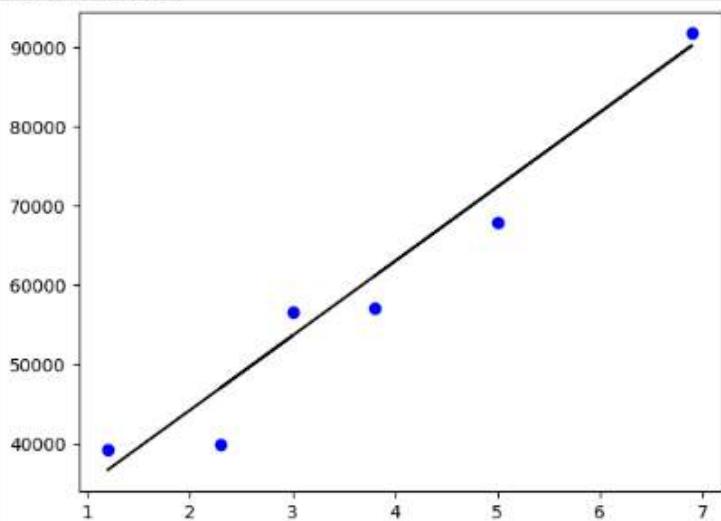
[ ] regr = LinearRegression()
regr.fit(X_train, y_train)
print(regr.score(X_test, y_test))

0.9450507958680934

[ ] y_pred = regr.predict(X_test)
plt.scatter(X_test, y_test, color ='b')
plt.plot(X_test, y_pred, color ='k')
print(X_test)
print(y_pred)
plt.show()
```

```
print(y_pred)
plt.show()
```

```
[[3.8]
 [6.9]
 [5. ]
 [2.3]
 [3. ]
 [1.2]]
[[61140.05516825]
 [90242.83467203]
 [72405.64722933]
 [47058.0650739]
 [53629.66044753]
 [36731.27234392]]
```



And then the same procedure is implemented from scratch:

```
#= ( ) b0 = 1
( ) b1 = 1
alpha = 0.01
{ x}
[ ] def find_ycap(x, b0, b1):
y = []
for i in range(len(x)):
y.append(b0 + b1*x[i])
return y

[ ] def J(x,y,ycap):
cost = 0
for i in range(len(y)):
cost += (ycap[i]-y[i])**2
cost = cost/len(x)
return cost

[ ] minitively
for i in range(10):
ycap = find_ycap(df2[["YearsExperience"]], (b0, b1))
cost = J(df2[["YearsExperience"]], df2[["Salary"]], ycap)
b0 = b0 - alpha * cost
b1 = b1 - alpha * cost
print(ycap)

[7500.656430486987, 17383.98883287328, 20866.32123474685, 35172.15234132791, 37654.404644186754, 45642.640853455465, 46785.814254791854, 46866.14661246462, 48056.34665746462, 54771.97769414655, 57954.71866681933, 58195.427674615521, ...]
[ ]
[ ] print(find_ycap(df2[["YearsExperience"]], b0, b1))

[1400.124283343719, 1600.286128340361, 11408.278377137003, 11408.479995230048, 11408.54881082525, 13408.875688113407, 11408.8562100811819, 11408.93125108846, 11408.93255408846, 11409.125803308065, 11409.38291326707, 11409.243435792
[ ]
[ ] Start coding or generate with AI.
```

60004210210

Amartya Mishra

COMPS – C31

ML Experiment 2

ML Experiment 2

~~Dr. Mitali A.~~

Aim: To Implement & understand the basics of logistic Regression.

Theory:

Logistic Regression is a supervised ML Algorithm used for Classification task where the goal is to predict the probability that an instance belongs to a given class or not.

It is a Statistical Algo which analyzes the relation between two factors

Types of logistic Regression:

- Binomial
- Multinomial
- Ordinal

$$Z = wX + b$$

X: Independent inputs

b: bias term (intercept)

w: Weights

Sigmoid function:

The output (Z) is mapped to get the probability

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

logistic regression Equation

$$P(x, b, w) = \frac{1}{1 + e^{-wx+b}}$$

Evaluating Model Methods:

- Accuracy
- Precision
- Recall
- F₁ Score

~~Conclusion~~ Thus, we implemented logistic Regression

Implementation:

```

[1] # Check for DataFrame and correct if necessary
if not isinstance(X_test, pd.DataFrame):
    X_test = X_test.values.reshape(-1, 1) # Assuming you want to use only the first column

# Convert labels to numerical representations if needed
le = LabelEncoder()
y_test = le.fit_transform(y_test)

# Normalize data (optional)
scaler = MinMaxScaler()
X_test_scaled = scaler.fit_transform(X_test.reshape(-1, 1))

# Fit logistic regression model
logr = LogisticRegression()
logr.fit(X_test_scaled, y_test)

# Create DataFrame
df = pd.DataFrame({'Study Hours': X_test.flatten(),
                    'Pass/Fail': y_test,
                    'Predicted': logr.predict(X_test_scaled)})

# Create scatterplot
df.plot.scatter(x='Study Hours', y='Pass/Fail', s=50, label='Actual')

# Generate x-values for prediction
x_values = np.linspace(X_test.min(), X_test.max(), 100).reshape(-1, 1)

# Normalize x-values for prediction (if normalized earlier)
if 'scaler' in locals():
    x_values_scaled = scaler.transform(x_values)
    y_values = logr.predict_proba(x_values_scaled)[:, 1]
else:
    y_values = logr.predict_proba(x_values)[:, 1]

# Plot logistic regression curve
plt.plot(x_values.flatten(), y_values, color='red', label='Logistic Regression Curve')

# Set labels and title
plt.xlabel('study Hours')
plt.ylabel('Pass/Fail')

```

```

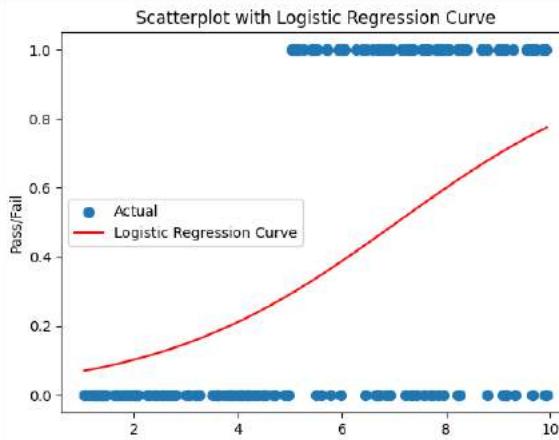
[1] # Plot logistic regression curve
plt.plot(x_values.flatten(), y_values, color='red', label='Logistic Regression Curve')

# Set labels and title
plt.xlabel('Study Hours')
plt.ylabel('Pass/Fail')
plt.title('Scatterplot with Logistic Regression Curve')

# Display legend
plt.legend()

# Show the plot
plt.show()

```



```

[1] w0 = 1
w1 = 3
alpha = 0.001

[1] def find_w(x, w0, w1):
    w = []
    for i in range(len(x)):
        w.append(w0 + w1*x[i])
    return w

[1] def calculate_jw(x, w, y):
    temp = []
    for i in range(len(x)):
        temp.append(y[i]-1/(1+math.e**(-w[i])))
    return sum(temp)

[1] x_data['Study Hours']
y_data['Pass/Fail']

for i in range(125):
    w = find_w(x, w0, w1)
    jw = calculate_jw(x, w, y)
    prevW0=w0
    prevW1=w1
    w0 = w0 + alpha*(jw)
    w1 = w1 + alpha*(jw)
    if w0==prevW0 or w1==prevW1:
        break

[1] pred=[]
for i in w:
    print(1/(1+math.e**-i),end=" ")
    if(1/(1+math.e**-i)>=0.5):
        pred.append(1)
    else:
        pred.append(0)
print()

```

```

[1]: x_data['Study_Hours']
y_data['Pass/Fail']

for i in range(125):
    x = (float(xc), float(yc))
    y = calculate_3d(x, xc, yc)
    predprob.append(y)
    predval.append(int(y))
    xc = xc + 0.05*(x[0])
    yc = yc + 0.05*(x[1])
    if abs(predval[i] - predprob[i]) <= 0.01:
        break

[2]: pred[]
for i in w:
    print(i[0].lstrip('0'), end=" ")
    if i[1].lstrip('0') == '1':
        pred.append(i[0])
    else:
        pred.append(i[1])
print()
print(pred)

0.3094223325423 -0.5740026214029609 0.404612953318331764 0.40300906207533 0.2213499495792160 0.2213499495792160 0.19803227969439224 0.533068892786214 0.402670159511033 0.4046223117048653 0.12093613610544 0.532099151500
[0, 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, 31, 32, 33, 34, 35, 36, 37, 38, 39, 40, 41, 42, 43, 44, 45, 46, 47, 48, 49, 50, 51, 52, 53, 54, 55, 56, 57, 58, 59, 60, 61, 62, 63, 64, 65, 66, 67, 68, 69, 70, 71, 72, 73, 74, 75, 76, 77, 78, 79, 80, 81, 82, 83, 84, 85, 86, 87, 88, 89, 90, 91, 92, 93, 94, 95, 96, 97, 98, 99, 100]

[3]: conf_matrix = metrics.confusion_matrix(pred, y)
print(conf_matrix)
metrics.accuracy_score(pred,y)*100

[[779 110]
 [ 43 68]]
68.7
```

60004210210

Amartya Mishra

COMPS – C31

ML Experiment 3

60004210210
Amaranya Mishra
COMPS - C31

ML Experiment - 3

Aim: To Implement CART Algorithm.

11/3/24

Theory:

CART (Classification & Regression tree) is a variation of the decision tree algorithm. It can handle both classification & regression tasks. It is a predictive algorithm used in ML & it explains how the target variables' values can be predicted based on other parameters.

CART Algorithm

- 1) Tree Structure
- 2) Splitting criteria

$$gini = 1 - \sum_{i=1}^k (p_i)^2$$

- 3) Pruning

Advantages

- 1) Results are Simplistic
- 2) Trees are non-parametric & non-linear
- 3) Trees implicitly perform feature selection

Disadvantages

- 1) Overfitting
- 2) High variance
- 3) Low bias.

Conclusion: Thus we Implement CART Algorithm.

Implementation:

ML Exp3 CART:pylab - Colab

File Edit View Insert Runtime Tools Help

+ Code + Test

Comment Share

Connect Colab AI

```
import numpy as np
import pandas as pd

df = pd.read_csv('content/social_network_ads.csv')
df.head()

Gender Age EstimatedSalary Purchased
0 1 19 19000 0
1 1 33 20000 0
2 0 26 33000 0
3 0 27 67000 0
4 1 19 75000 0

[ ] X = df.loc[:, [0, 1, 2]].values
y = df.loc[:, -1].values
from sklearn.model_selection import train_test_split
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.3, random_state=42)

[ ] from sklearn.tree import DecisionTreeClassifier
from sklearn.metrics import accuracy_score
scores = []
clf = DecisionTreeClassifier(max_depth=3, criterion='gini')
clf.fit(X_train, y_train)
y_pred = clf.predict(X_test)
accuracy = accuracy_score(y_true=y_test, y_pred=y_pred)
scores.append(accuracy)
print(scores)
for y in y_pred:
    print(y, end=" ")
print()
for y in y_test:
    print(y, end=" ")
```

ML Exp3 CART.ipynb

```
[1]: outlook temp humidity wind decision
0: Sunny Hot High Weak 0
1: Sunny Hot High Strong 0
2: Overcast Hot High Weak 1
3: Rain Mid High Weak 1
4: Rain Cool Normal Weak 1
5: Rain Cool Normal Strong 0
6: Overcast Cool Normal Strong 1
7: Sunny Mid High Weak 0
8: Sunny Cool Normal Weak 1
9: Rain Mid Normal Weak 1
10: Rainy Mid Normal Strong 0
11: Overcast Mid High Strong 1
12: Overcast Hot Normal Weak 1
13: Rain Mid High Strong 0

[2]: def calc_gini_for_attribute(class_name, col_idx, target_col='decision'):
    total_count = len(df[(df[col_idx] == class_name) & (df[target_col] == 1)])
    count_of_0 = len(df[(df[col_idx] == class_name) & (df[target_col] == 0)])
    prob_of_1 = count_of_1 / total_count
    prob_of_0 = count_of_0 / total_count
    gini_val = (prob_of_1 * (1 - (prob_of_0)**2))
    return gini_val, total_count
calc_gini_for_attribute('Sunny', 'outlook', target_col='decision')
(0.48, 5)

[3]: col1 = 'outlook'
```

ML Exp3 CART.ipynb

```
[1]: [Sunny, 'Overcast', 'Rain']

[2]: cols = ['outlook', 'temp', 'humidity', 'wind']
gini_dict = {}
for col in cols:
    print(col)
    gini_val = 0
    for value in list(df[col].unique()):
        gini_val, var_count = calc_gini_for_attribute(value, col, df)
        print(f'For attr: {value}, Value = {gini_val}')
        gini_val_col += var_count/len(df) * gini_val
    print(f'Grand Gini for attr: {col}')
    print(f'{gini_dict[col]} = {round(gini_val_col, 3)}')

[3]: outlook
For attr: Sunny, Value = 0.48
For attr: Overcast, Value = 0.0
For attr: Rain, Value = 0.48
0.365

[4]: humidity
For attr: Hot, Value = 0.5
For attr: Mid, Value = 0.40444444444444445
For attr: Cool, Value = 0.375
0.44

[5]: temp
For attr: Weak, Value = 0.175
For attr: Strong, Value = 0.5
0.35
```

```
MLExp3 CART.ipynb - Colaboratory
```

```
# colab.research.google.com/drive/1oxcXzaprzbLobsE/NZdtpKaB4
```

```
File Edit View Insert Runtime Tools Help
```

```
+ Code + Text
```

```
Search
```

```
(x)
```

```
Comment Share
```

```
Connect Delete All
```

```
||| gini_dict
```

```
{'outlook': 0.543, 'temp': 0.44, 'humidity': 0.367, 'wind': 0.429}
```

```
[ ] def calc_gini(cols, df):
    gini_dict = {}
    for col in cols:
        gini_for_attr = 0
        for value in list(df[col].unique()):
            gini_val, var_count = calc_gini_for_attribute(value, col, df)
            gini_for_attr += var_count / len(df) * gini_val
        gini_dict[col] = round(gini_for_attr, 3)
    return gini_dict
```

```
calc_gini(cols, df)
```

```
['outlook': 0.543, 'temp': 0.44, 'humidity': 0.367, 'wind': 0.429}
```

```
[ ] def get_sel_attr():
    cols = list(df.columns)
    attr_gini = calc_gini(cols, df)
    min = 10
    for col in cols:
        if attr_gini[col] < min:
            min = attr_gini[col]
            sel_attr = col
    return sel_attr
```

```
[ ] def split_df(sel_attr, df, father):
    global id
    print(sel_attr)
    list_of_unique_values = list(df[sel_attr].unique())
    for value in list_of_unique_values:
        if check_termination(df[df[sel_attr] == value]):
            if not check_compatibility([sel_attr] * len(value)):
                id += 1
                final_tree.append({'id': id, 'data': df[df[sel_attr] == value], 'cond': sel_attr + " == " + value, 'children': [0, 0], 'isRoot': False, 'isLeaf': True, 'father': father})
            else:
                print("Cannot terminate when [sel_attr] is " + value)
        else:
            print("Cannot terminate when [sel_attr] is " + value)
```

```
ML Exp3 CART.ipynb - Colab: x +  
File Edit View Insert Runtime Tools Help Saving...  
Comment Share   
+ Code Text  
  
[ ] def split_dt(df=df_attr, df=df_attr):  
    print("df_attr:  
    \n", df_attr)  
    list_of_unique_values = list(df_attr.unique())  
    for value in list_of_unique_values:  
        if check_terminating(df=df_attr, sel_attr == value):  
            print("Cannot terminate when (sel_attr) is ", value)  
            df_attr = df_attr.drop([value])  
        else:  
            final_tree.append([sel_id, df_attr[(df_attr[sel_attr] == value)], "cond"], sel_attr + " is " + value, [children], isRoot=False, isLeaf=True, father=df_attr)  
    else:  
        print("Cannot terminate when (sel_attr) is (value)")  
        df_attr = df_attr.drop([value])  
        new_attr = df_attr[sel_attr == value].drop([sel_attr], axis=1)  
        print("attr removed from header")  
        final_tree.append([sel_id, df_attr[(df_attr[sel_attr] == value)], "cond"], sel_attr + " is " + value, [children], isRoot=False, isLeaf=True, father=df_attr)  
        new_best_attr = get_best_attr(new_df)  
        print("New attr: ", new_best_attr)  
        split_dt(new_best_attr, new_df, 10)  
    return [value]  
  
[ ] def check_terminating(df):  
    df_length = len(df)  
    zero_count = len(df[df['decision'] == 0])  
    one_count = len(df[df['decision'] == 1])  
    higher = zero_count/df_length if zero_count/df.length > one_count/df.length else one_count/df.length  
    print(higher)  
    if (higher > 0.9):  
        return True  
    return False  
  
[ ] df = pd.read_csv('iris.csv')  
df['Species'] = df['Species'].map({'Iris-setosa': 0, 'Iris-versicolor': 1, 'Iris-virginica': 2})  
df['SepalLength'] = df['SepalLength'].map(lambda x: round(x))  
df['SepalWidth'] = df['SepalWidth'].map(lambda x: round(x))  
df['PetalLength'] = df['PetalLength'].map(lambda x: round(x))  
df['PetalWidth'] = df['PetalWidth'].map(lambda x: round(x))  
df['Outlook'] = df['Outlook'].map({'Sunny': 0, 'Rainy': 1, 'Cloudy': 2})  
df['Humidity'] = df['Humidity'].map(lambda x: round(x))  
df['Wind'] = df['Wind'].map(lambda x: round(x))  
  
# Check if Outlook is Sunny  
# Can't terminate when Outlook is Sunny  
# Can Attr: humidity  
# humidity
```

```
+ Code + Test
File Edit View Insert Runtime Tools Help Edit changes saved
Comment Share
Connect Colab AI
[ ] id = 0
final_tree = [{"id": 0, "data": "df", "cond": "None", "children": [], "isRoot": true, "isLeaf": false, "father": null}
split_of("outlook", df, 0)
    .itself
0.0
    Cannot terminate when outlook is Sunny
    New Attr: humidity
    humidity
    1.0
        terminating when humidity is High
    1.0
        terminating when humidity is Normal
    1.0
        terminating when outlook is Overcast
    0.0
        Cannot terminate when outlook is Rain
        New Attr: wind
        wind
        1.0
            terminating when wind is Weak
        1.0
            terminating when wind is Strong
        {success: False}

[ ] for obj in final_tree:
    print(f"obj: {obj['id']} --- cond: {obj['cond']} --- Father: {obj['father']}")

obj: 0 --- cond: None --- Father: null
obj: 1 --- cond: outlook is sunny --- Father: 0
obj: 2 --- cond: humidity is High --- Father: 0
obj: 3 --- cond: humidity is Normal --- Father: 0
obj: 4 --- cond: outlook is Overcast --- Father: 0
obj: 5 --- cond: outlook is Rain --- Father: 0
obj: 6 --- cond: wind is Weak --- Father: 0
obj: 7 --- cond: wind is Strong --- Father: 0
obj: 8 --- cond: wind is Strong --- Father: 0
```

60004210210

Amartya Mishra

COMPS – C31

ML Experiment 4

6004210210
Amartya Mishra
COMPS - C31

ML Experiment 4

Aim: To implement principle component Analysis (PCA)

~~TOPIC~~

Theory:

PCA works on the condition that while the Data in a higher dimensional space is mapped to data in lower dimension space the Variance of the data in the lower dimensional space should be maximum.

PCA is a statistical procedure that uses an orthogonal transformation that converts a set of co-related variables to a set of uncorrelated variables

It is an unsupervised learning algorithm used to examine the iterations among a set of variables

Steps

I) Standardization

$$z = \frac{x - \bar{x}_1}{\sigma}$$

II) Covariance Matrix Computation

$$\text{cov}(x_1, x_2) = \sum_{i=1}^n \frac{(x_i - \bar{x}_1)(x_{i2} - \bar{x}_2)}{n-1}$$

III) Compute Eigen values & Eigen Vectors of covariance matrix to identify principal components.

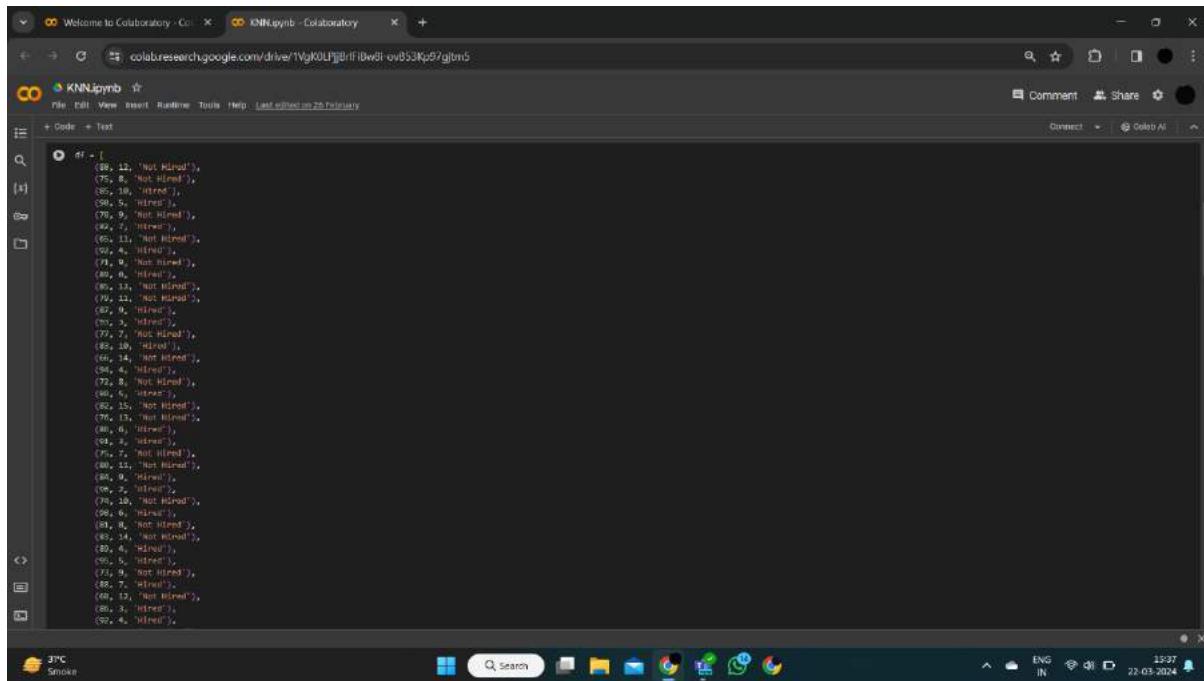
$$AX = \lambda X$$

$$AX - \lambda X = 0$$

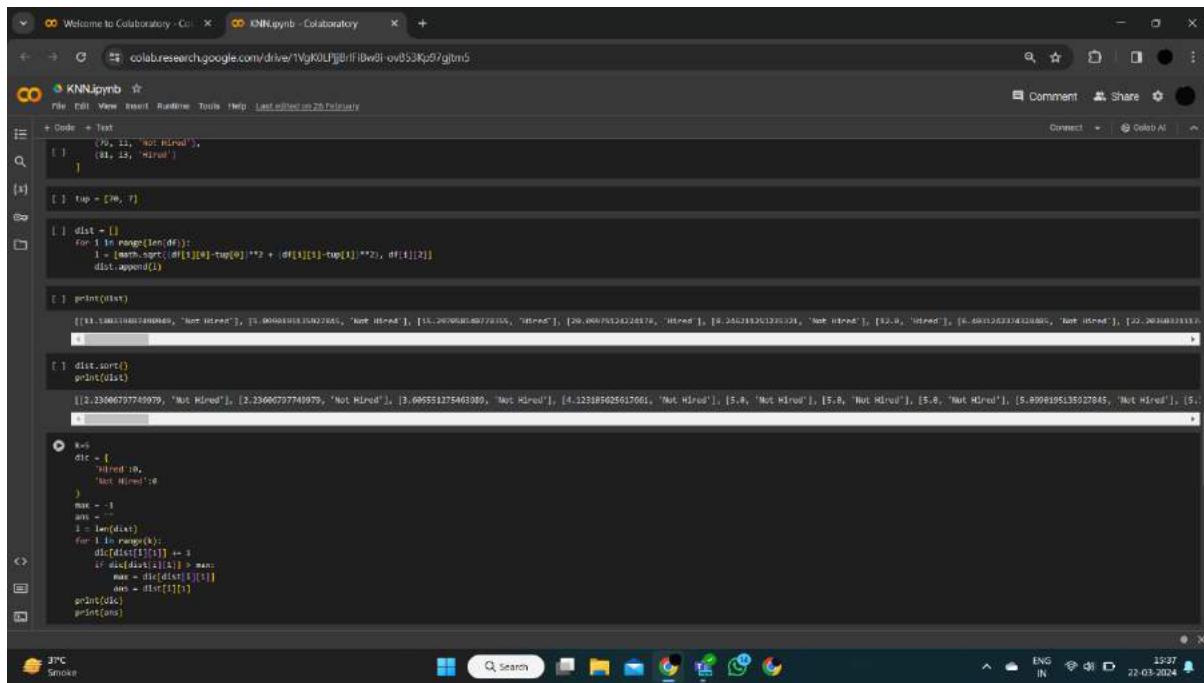
$$(A - \lambda I)X = 0$$

Conclusion : Thus we implement PCA

Implementation:



```
df = [
    (88, 12, 'Not Hired'),
    (89, 2, 'Not Hired'),
    (88, 18, 'Hired'),
    (90, 5, 'Hired'),
    (78, 9, 'Not Hired'),
    (84, 7, 'Hired'),
    (85, 13, 'Not Hired'),
    (80, 4, 'Hired'),
    (79, 3, 'Not Hired'),
    (80, 13, 'Hired'),
    (79, 11, 'Not Hired'),
    (82, 9, 'Hired'),
    (83, 5, 'Hired'),
    (93, 7, 'Not Hired'),
    (85, 18, 'Hired'),
    (66, 14, 'Not Hired'),
    (84, 4, 'Hired'),
    (72, 8, 'Not Hired'),
    (86, 1, 'Hired'),
    (86, 15, 'Not Hired'),
    (78, 10, 'Hired'),
    (80, 6, 'Hired'),
    (89, 2, 'Hired'),
    (76, 7, 'Not Hired'),
    (80, 13, 'Not Hired'),
    (84, 9, 'Hired'),
    (85, 2, 'Hired'),
    (89, 14, 'Not Hired'),
    (85, 4, 'Hired'),
    (85, 5, 'Hired'),
    (83, 10, 'Not Hired'),
    (88, 7, 'Hired'),
    (86, 13, 'Not Hired'),
    (88, 3, 'Hired'),
    (82, 4, 'Hired')
]
```



```

df = [
    (79, 11, 'Not Hired'),
    (81, 13, 'Hired')
]

top = [80, 7]

dist = []
for i in range(len(df)):
    l = [math.sqrt((df[i][0]-top[0])**2 + (df[i][1]-top[1])**2), df[i][2]]
    dist.append(l)

print(dist)
[[13.188070000000002, 'Not Hired'], [1.8000000000000002, 'Not Hired'], [15.207000000000002, 'Not Hired'], [10.280000000000002, 'Not Hired'], [8.250200000000001, 'Not Hired'], [12.0, 'Hired'], [4.4810203308120802, 'Not Hired'], [12.20304022111111, 'Not Hired']]

dist.sort()
print(dist)
[[2.23606707740979, 'Not Hired'], [3.23606707740979, 'Not Hired'], [3.695551275461089, 'Not Hired'], [4.123185025917061, 'Not Hired'], [5.0, 'Not Hired'], [5.0, 'Not Hired'], [5.0, 'Not Hired'], [5.8998195138927845, 'Not Hired'], [5.1, 'Not Hired']]

K=5
dic = {}
'Hired':0,
'Not Hired':0
max = -1
ans = ''
l = len(dist)
for i in range(k):
    dic[dist[i][1]] += 1
    if dic[dist[i][1]] > max:
        max = dic[dist[i][1]]
        ans = dist[i][1]
print(dic)
print(ans)
```

```
# KNN.pyrb
# File Edit View Insert Runtime Tools Help [Last edited on 26 February]
# Comment Share
# Connect Colab AI
# + Code + Text
# 
# K=5
# dict = {
#     'Hired':0,
#     'Not hired':1
# }
# max = -1
# ans = ''
# 2 = len(df)
# for i in range(0,2):
#     dist=df[i][1][1] += 3
#     if dist>dict[1][1] > max:
#         max = dist[1][1]
#         ans = dist[1][1]
# print(ans)
# print(ans)

# ['Hired': 0, 'Not hired': 1];
# Not hired

# [ ] def acc(tp, df):
#     dist = []
#     for i in range(len(df)):
#         l = math.sqrt((df[0][0]-tp[0])**2 + (df[1][1]-tp[1])**2 + (df[2][2]-tp[2])**2)
#         dist.append(l)
#     dist.sort()
#     k=1
#     dict = {
#         'Hired':0,
#         'Not hired':1
#     }
#     max = -1
#     ans = ''
#     l = len(dict)
#     for i in range(k):
#         dist[i][1][1] += 1
#         if dict[dist[i][1][1]] > max:
#             max = dict[dist[i][1][1]]
#             ans = dist[i][1][1]
#     return ans==tp[2]

# JPC
# Smoke
# Welcome to Colaboratory - Colab
# KNNApyrb - Colaboratory
# colab.research.google.com/drive/1Vgk0UPjBfIbw8l-ov053Kp97gjmS
# File Edit View Insert Runtime Tools Help [Last edited on 26 February]
# Comment Share
# Connect Colab AI
# + Code + Text
# 
# K=1
# max = -1
# ans = ''
# l = len(dict)
# for i in range(k):
#     dist[i][1][1] += 1
#     if dict[dist[i][1][1]] > max:
#         max = dict[dist[i][1][1]]
#         ans = dist[i][1][1]
# return ans==tp[2]

# [ ] testdata = [
#     (96, 2, "Hired"),
#     (96, 16, "Not hired"),
#     (86, 8, "Hired"),
#     (95, 7, "Hired"),
#     (77, 13, "Not hired"),
#     (82, 11, "Hired"),
#     (85, 14, "Not hired"),
#     (66, 16, "Not hired"),
#     (90, 5, "Hired"),
#     (77, 12, "Not hired"),
#     (84, 10, "Not hired"),
#     (80, 3, "Hired"),
#     (73, 13, "Not hired"),
#     (89, 8, "Hired"),
#     (94, 7, "Hired")
# ]
# yes=0
# no=0
# for tp in testdata:
#     if acc(tp, df):
#         yes += 1
#     else:
#         no += 1
# print(yes, no)
# print("Accuracy:", yes/(yes+no))

# 14 1
# Accuracy: 0.9333333333333333
```

60004210210

Amartya Mishra

COMPS – C31

ML Experiment 5

60004210210
Amartha Mishra
COMPS - C31

ML Experiment 5

Ques 2)

Aim: To implement KNN Algorithm

Theory:

K-nearest Neighbor (KNN) Algo w/ a supervised ML Algo Employed to take tackle classification & Regression models problems.

Distance metrics used in KNN

1) Euclidean dist

$$d(x, y_i) = \sqrt{\sum_{j=1}^d (x_j - y_{ij})^2}$$

2) Manhattan Distance

$$d(x, y) = \sum_{i=1}^n |x_i - y_i|$$

3) Minkowski Distance

$$d(x, y) = \left(\sum_{i=1}^n (x_i - y_i)^p \right)^{1/p}$$

choosing the value of k

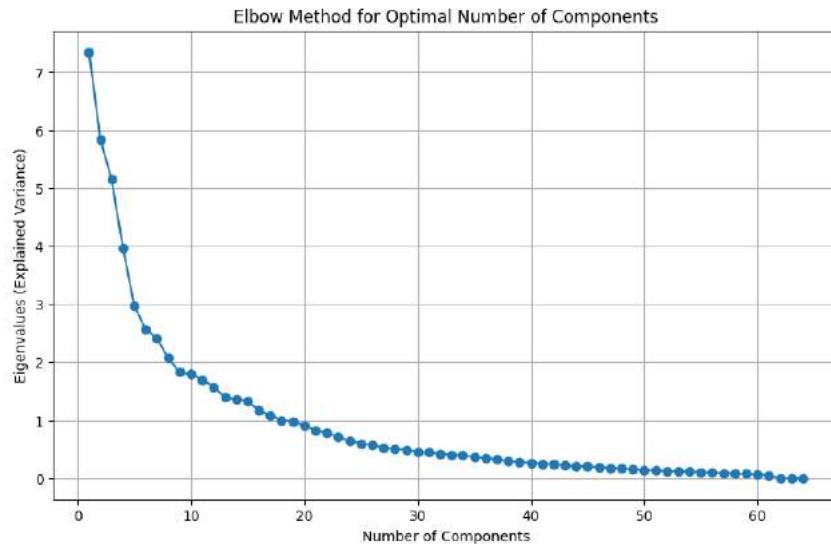
The value of k depends on input data

If input data has more outliers, a higher value of k is

Conclusion: Thus we implemented & understood KNN Algorithm.

Implementation:

```
import numpy as np
import matplotlib.pyplot as plt
import pandas as pd
from sklearn.datasets import load_digits
from sklearn.decomposition import PCA
from sklearn.preprocessing import StandardScaler
digits_data = load_digits()
X = digits_data.data
scaler = StandardScaler()
X_scaled = scaler.fit_transform(X)
pca = PCA()
X_pca = pca.fit_transform(X_scaled)
eigenvalues = pca.explained_variance_
plt.figure(figsize=(10, 6))
plt.plot(range(1, len(eigenvalues) + 1), eigenvalues, marker='o',
linestyle='--')
plt.title('Elbow Method for Optimal Number of Components')
plt.xlabel('Number of Components')
plt.ylabel('Eigenvalues (Explained Variance)')
plt.grid(True)
plt.show()
optimal_num_components = 10
X_reduced = X_pca[:, :optimal_num_components]
df_reduced = pd.DataFrame(X_reduced, columns=[f'PC{i}' for i in
range(1, optimal_num_components + 1)])
df_reduced['target'] = digits_data.target
df_reduced.to_csv('reduced_digits_dataset.csv', index=False)
print("Digits Wine dataset saved successfully.")
```



```
▶ print(X_pca)
[[ 1.91421366e+00 -9.54501571e-01 -3.94603482e+00 ... -0.00000000e+00
  0.00000000e+00  8.24385469e-15]
 [ 5.88980330e-01  9.24635800e-01  3.92475494e+00 ...  8.37680823e-16
 -1.88056170e-17  3.08670512e-17]
 [ 1.30203906e+00 -3.17188827e-01  3.02333293e+00 ...  4.30699370e-16
  1.79508547e-17  2.15555254e-17]]
```

60004210210

Amartya Mishra

COMPS – C31

ML Experiment 6

60004210210
Amartya Mishra
COMPS - C31

ML - Experiment 6

Aim: To implement Back propagation *(RPNK)*

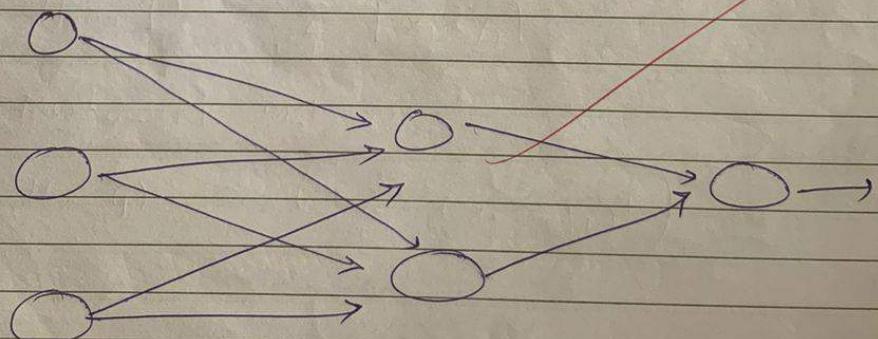
Theory:

Back propagation is an algorithm that back propagates the errors from the output nodes to input nodes.

It is widely used algorithm for training feed forward neural networks.

It computes the gradient of loss function with respect to the network weights.

It is very efficient, rather than naively directly computing the gradient concerning each weight.



Parameters:

x = inputs training vector (x_1, x_2, \dots, x_n)

t = target vector (t_1, t_2, \dots, t_n)

s_n = error at output unit

δ_j = error at hidden layer

α = learning rate

v_{0j} = bias of hidden unit j

conclusion: Thus, we implemented back propagation

Implementation:

```
import numpy as np
class NeuralNetwork:
    def __init__(self, input_size, hidden_size, output_size):
        self.input_size = input_size
        self.hidden_size = hidden_size
        self.output_size = output_size
        # Initialize weights and biases
        self.weights_input_hidden = np.random.randn(self.input_size,
                                                    self.hidden_size)
        self.bias_input_hidden = np.random.randn(1, self.hidden_size)
        self.weights_hidden_output = np.random.randn(self.hidden_size,
                                                    self.output_size)
        self.bias_hidden_output = np.random.randn(1, self.output_size)
    def sigmoid(self, x):
        return 1 / (1 + np.exp(-x))
    def sigmoid_derivative(self, x):
        return x * (1 - x)
    def forward(self, inputs):
        self.hidden_input = np.dot(inputs, self.weights_input_hidden) +
        self.bias_input_hidden
        self.hidden_output = self.sigmoid(self.hidden_input)
        self.final_input = np.dot(self.hidden_output,
        self.weights_hidden_output) + self.bias_hidden_output
        self.final_output = self.sigmoid(self.final_input)
        return self.hidden_output, self.final_output
    def backward(self, inputs, targets, learning_rate):
        error = targets - self.final_output
        delta_output = error * self.sigmoid_derivative(self.final_output)
        delta_hidden = np.dot(delta_output, self.weights_hidden_output.T) *
        self.sigmoid_derivative(self.hidden_output)
        self.weights_hidden_output += np.dot(self.hidden_output.T,
        delta_output) * learning_rate
        self.bias_hidden_output += np.sum(delta_output, axis=0,
        keepdims=True) * learning_rate
        self.weights_input_hidden += np.dot(inputs.T, delta_hidden) *
        learning_rate
        self.bias_input_hidden += np.sum(delta_hidden, axis=0,
        keepdims=True) * learning_rate
        return error
    def train(self, inputs, targets, learning_rate):
        hidden_output, final_output = self.forward(inputs)
        error = self.backward(inputs, targets, learning_rate)
        print("Output of hidden layer:")
        print(hidden_output)
        print("Output of output layer:")

        print(final_output)
        print("Error found:")
        print(error)
        print("Updated weights after 1 iteration:")
        print("Weights from input to hidden layer:")
        print(self.weights_input_hidden)
        print("Weights from hidden to output layer:")
        print(self.weights_hidden_output)
```

```
dataset = pd.read_csv('reduced_digits_dataset.csv')
inputs = dataset.drop(columns=['target']).values
targets = dataset['target'].values.reshape(-1, 1)
input_size = inputs.shape[1]
output_size = len(np.unique(targets))
hidden_size = 3
nn = NeuralNetwork(input_size, hidden_size, output_size)
nn.train(inputs, targets, learning_rate=0.1)
```

60004210210

Amartya Mishra

COMPS – C31

ML Experiment 7

MI Experiment 7

60004210210
Amartya Mishra
COMPS - C31

Aim: To Implement Support Vector Machine.

Theory:

Support Vector machine (SVM) is used for linear/non linear classification, regression & even outlier detection.

It is a supervised MI Algorithm.

The main objective of the SVM is to find the optimal hyperplane in a N-dimensional space that can separate the data points in different classes in the feature space.

The hyperplane tries that the margin between the closest points of different classes should be as maximum as possible.

The dimension of the hyperplane depends upon the number of features.

Mathematical Intuition

$$\mathbf{w}^T \mathbf{x} + b = 0$$

$$d_i = \frac{\mathbf{w}^T \mathbf{x}_i + b}{\|\mathbf{w}\|}$$

$$\hat{y} = \begin{cases} 1 & : w^T x + b \geq 0 \\ 0 & : w^T x + b < 0 \end{cases}$$

Types of SVM

- 1) Linear
- 2) Non linear

Advantages

- 1) Effective in high dimensional cases
- 2) Its memory is efficient as it uses a subset of training points in the decision function called support vectors
- 3) Different kernel functions can be specified for the decision function & its possible to specify custom kernels.

Conclusion : Thus we implemented SVM.

Implementation:

```
from sklearn.datasets import load_breast_cancer
from sklearn.model_selection import train_test_split
from sklearn.naive_bayes import GaussianNB
from sklearn.metrics import accuracy_score
data = load_breast_cancer()
X = data.data
y = data.target
X

X_train, X_test, y_train, y_test = train_test_split(X, y,
test_size=0.3, random_state=42)
classifier = GaussianNB()
classifier.fit(X_train, y_train)
y_pred = classifier.predict(X_test)
accuracy = accuracy_score(y_test, y_pred)
print(accuracy)

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

60004210210

Amartya Mishra

COMPS – C31

ML Experiment 8

ML - Experiment 8

60004210210

Amaranya Mishra
COMPS - C3I

Aim: To Implement Bayesian classification

Theory :

Bayesian Classification is a probabilistic approach to learning & inference based on different views of what it means to learn from Data in which probability is used to represent uncertainty about the relation ship being learnt.

It is used to determine the Probability of hypothesis with a prior knowledge

It depends on conditional probability

$$\text{The Bayes thm is } P(A|B) = \frac{P(B|A) \times P(A)}{P(B)}$$

$P(A|B)$ Probability of B given A happened (Posterior)

$P(B|A)$ Prob. of B given A occurs (likelihood)

$P(A)$ prob. of A Prior prob. (Prob. of hypothesis)

$P(B)$ Prob. of B Posterior^{marginal} Prob. (Prob. of evidence)

Conclusion: Thus we implement Bayesian classification.

Implementation:

```
from sklearn.datasets import load_breast_cancer
from sklearn.model_selection import train_test_split
from sklearn.svm import SVC
from sklearn.metrics import accuracy_score
data = load_breast_cancer()
X = data.data
y = data.target
X

X_train, X_test, Y_train, Y_test = train_test_split(X, y,
test_size=0.3, random_state=42)
classifier = SVC(kernel=' ')
classifier.fit(X_train, Y_train)
Y_pred = classifier.predict(X_test)
accuracy = accuracy_score(Y_test, Y_pred)
print("Accuracy:", accuracy)

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

MI - Assignment 2

60004210210

Amartya Mishra
COMPS - C31

~~Roll No.~~
~~Amartya~~
~~PF~~

- 1) No - kmeans & gaussian mixture model (GMM) will generally not produce same cluster for a given dataset / data.

K-Means

It is a centroid based clustering Algorithm. It works by initially assigning a fixed No. of centroids at random location.

Then it iteratively arranges the data point to the clustered centroid.

Gaussian Mixture Model

GMM is probabilistic clustering model that assumes the data is generated by a massive gaussian distribution.

It uses expectation minimization.

Hence K-means & GMM make different assumption & hence produce different cluster appropriation. The choice of Algorithm depends on the characteristics of the data.

2) Hidden Markov Model is a statistical model that is used to describe the probabilities relationship between a sequence of observation & a sequence of hidden states.

Applications:

- 1) Speech recognition:
can model statistical properties of words or sentences to recognize patterns.
- 2) Gesture Recognition:
used to model different gestures based on observed movements.
- 3) Bioinformatics
Protein structure prediction, sequence alignment & model various biological sequences.
- 4) Robotics & Autonomous System.
Localization & mapping
- 5) Independent Component Analysis
It is used to separate mixed signals into their original independent components.
It assumes that input is a combination of sources. By finding the linear transformation that maximizes the statistical independence of the components.

ICA can be used to extract components
It leverages non gaussian nature of signals.

- 4) Deep Neural Network: They are NN based
are composed of multiple layers of nodes which
attempt to model high level abstraction in
data.
It consists of IIP layer, multiple hidden layers
& output layer.
Each set of layer contains nodes for computation
~~They are capable of learning features that
can be used for various task.~~

(At)

Machine learning Assignment 2

Q1

a) Video Surveillance :

Application :

17/07/24

Machine learning can be used in video surveillance for various tasks such as object detection, activity recognition & anomaly detection.

It can help automatically detect & track objects of interest, identifying suspicious behaviour & alerting security personals in real time.

There many such more application in almost all fields of industry from face detection, motion detection to automatic tag detection etc.

Suitable ML Technique:

Convolution Neural Networks (CNN's) are commonly used for such task due to their ability to effectively analyze spatial features in image or frames of videos.

b) Sentiment Analysis:

→ Application :

Sentiment Analysis involves determining the

Sentiments or opinion Expressed in text.

It could be extracted from:

- Social Media Post
- Product reviews
- customer feed backs
- Such technology can be used by businesses to understand customer opinions or analyze user feed back for product improvement.

Suitable Methods:

- NLP - Natural language processing models.
- RNN - Recurrent Neural Networks.
- Transformer Models : EBPTA.

(c) Image recognition :

Application:

It involves identifying & classifying objects or patterns with image
It is used in various Applications

- Medical diagnosis
- Autonomous vehicles
- facial Recognition

Suitable MI Techniques:

Convolution Neural Network - CNN are most suitable techniques for image recognition tasks