

Experiment No.1

Name: Atharv Vijay Deshpande

PRN: 21410044

Batch : EN-3

Title: Implementation of Linear Regression Model

Part A: Implementation of Linear Regression Model on Synthetic data

Importing Libraries

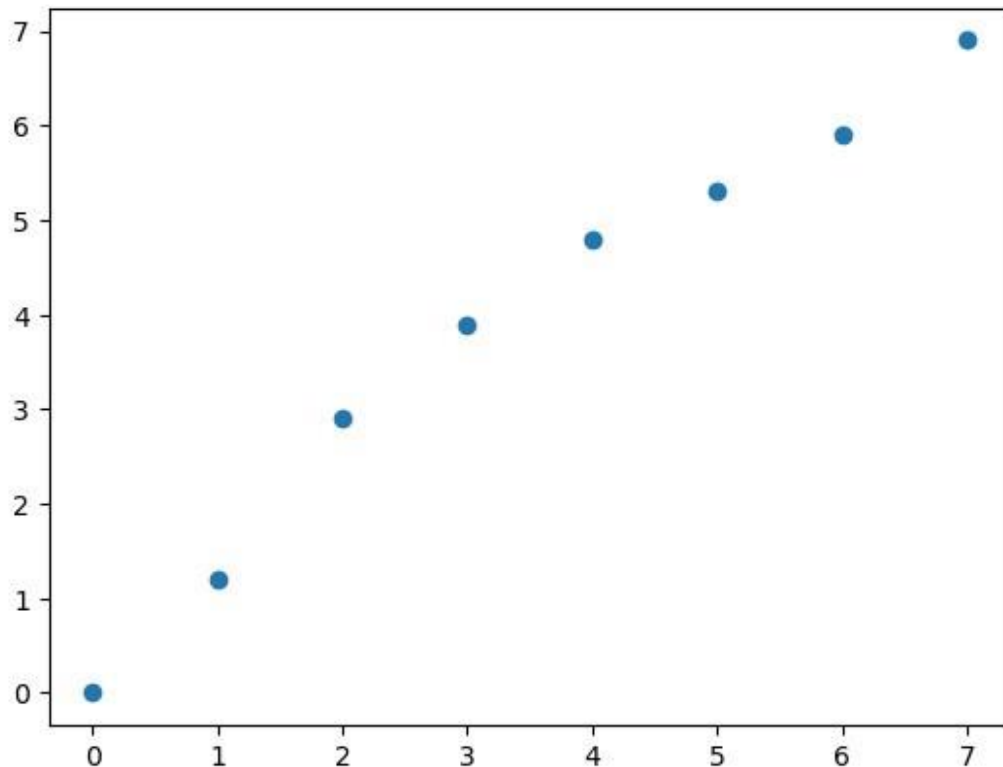
```
import numpy as np
import matplotlib.pyplot as plt
import pandas as pd
```

Creating Sample Data

```
#Independent Variable
x=np.array(np.arange(0,8))
array([0, 1, 2, 3, 4, 5, 6, 7])
#Dependent Variable
y=np.array([0,1.2,2.9,3.9,4.8,5.3,5.9,6.9])
array([0. , 1.2, 2.9, 3.9, 4.8, 5.3, 5.9, 6.9])
```

Visualizing the data

```
plt.scatter(x,y)
<matplotlib.collections.PathCollection at 0x7355ae7ef230>
```



Linear Regression Model

```
#Importing the LinearRegression model
from sklearn.linear_model import LinearRegression
#making an instance of LinearRegression
model=LinearRegression() #fitting the data to the model
model.fit(x.reshape(-1,1),y)#x should be 2d array while fitting the data

LinearRegression()
```

Model Parameters (Slope and Intercept)

```
print(f'Slope : {model.coef_[0]}') print(f'Intercept: {model.intercept_}')
Slope :0.9511904761904761
Intercept: 0.5333333333333332
```

Predicting the Values

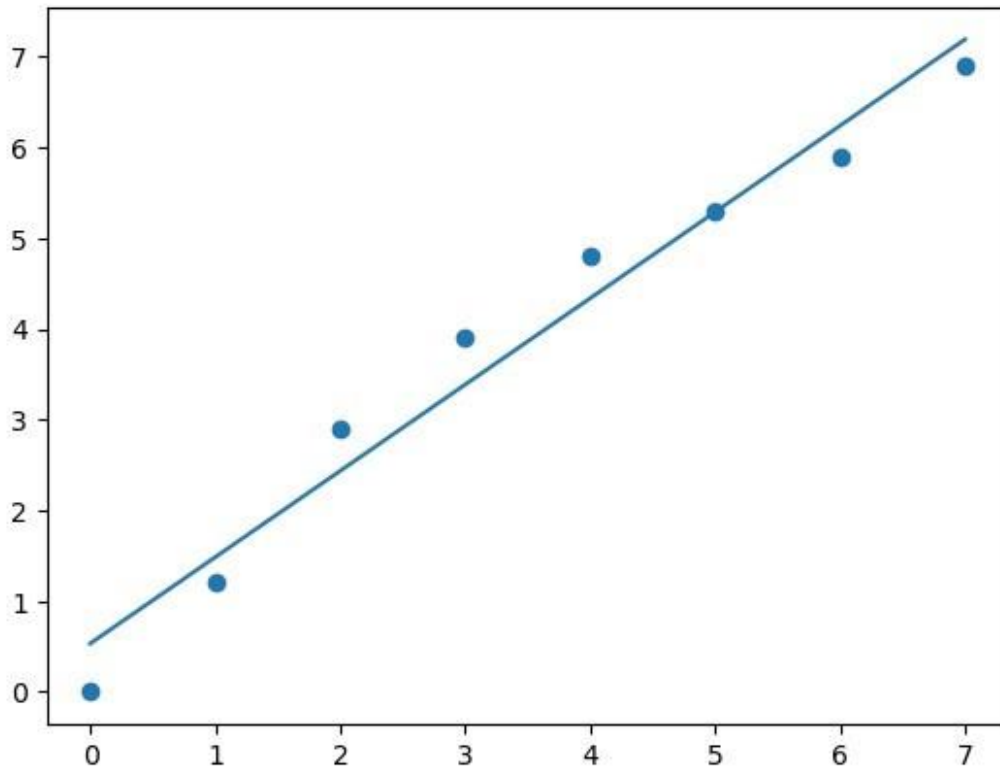
```
y_pred=model.predict(x.reshape(-1,1))
# y_pred=model.coef_[0]*x+model.intercept_ #y=mx+c equation of
regression line y_pred
```

```
array([0.53333333, 1.48452381, 2.43571429, 3.38690476, 4.33809524,  
       5.28928571, 6.24047619, 7.19166667])
```

Visualizing the model output

```
plt.scatter(x,y) plt.plot(x,y_pred)
```

```
[<matplotlib.lines.Line2D at 0x7355a601c770>]
```



Regression Performance metrics

1. Mean Absolute Error (MAE):

- **Definition:** MAE measures the average magnitude of errors in a set of predictions, without considering their direction. It is the average of the absolute differences between predicted and actual values.
- **Formula:**

$$MAE = \frac{1}{n} \sum_{i=1}^n |y_i - \hat{y}_i|$$

Where:

- y_i is the actual value, – \hat{y}_i is the predicted value, – n is the number of observations.
-

2. Mean Squared Error (MSE):

- **Definition:** MSE measures the average of the squares of the errors, giving more weight to larger differences between predicted and actual values.
- **Formula:**

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

Where:

- y_i is the actual value, – \hat{y}_i is the predicted value, – n is the number of observations.
-

3. R-squared (R^2) Score:

- **Definition:** R^2 , also called the coefficient of determination, represents the proportion of variance in the dependent variable that is predictable from the independent variables. It indicates how well the model explains the variability in the target variable.
- **Formula:**

$$R^2 = 1 - \frac{\sum_{i=1}^n (y_i - \hat{y}_i)^2}{\sum_{i=1}^n (y_i - \bar{y})^2}$$

Where:

- y_i is the actual value,
 - \hat{y}_i is the predicted value, – \bar{y} is the mean of the actual values, – n is the number of observations.
-

```
from sklearn.metrics import
mean_absolute_error,mean_squared_error,r2_score
print('Mean Absolute Error :',mean_absolute_error(y,y_pred)) print('Mean Squarred Error
:',mean_squared_error(y,y_pred)) print('R2 Value: ',r2_score(y,y_pred))
```

Mean Absolute Error : 0.36249999999999998
Mean Squarred Error : 0.14391369047619035
R2 Value: 0.9725878684807256

Part B:Implementation of Linear Regression Model on a dataset

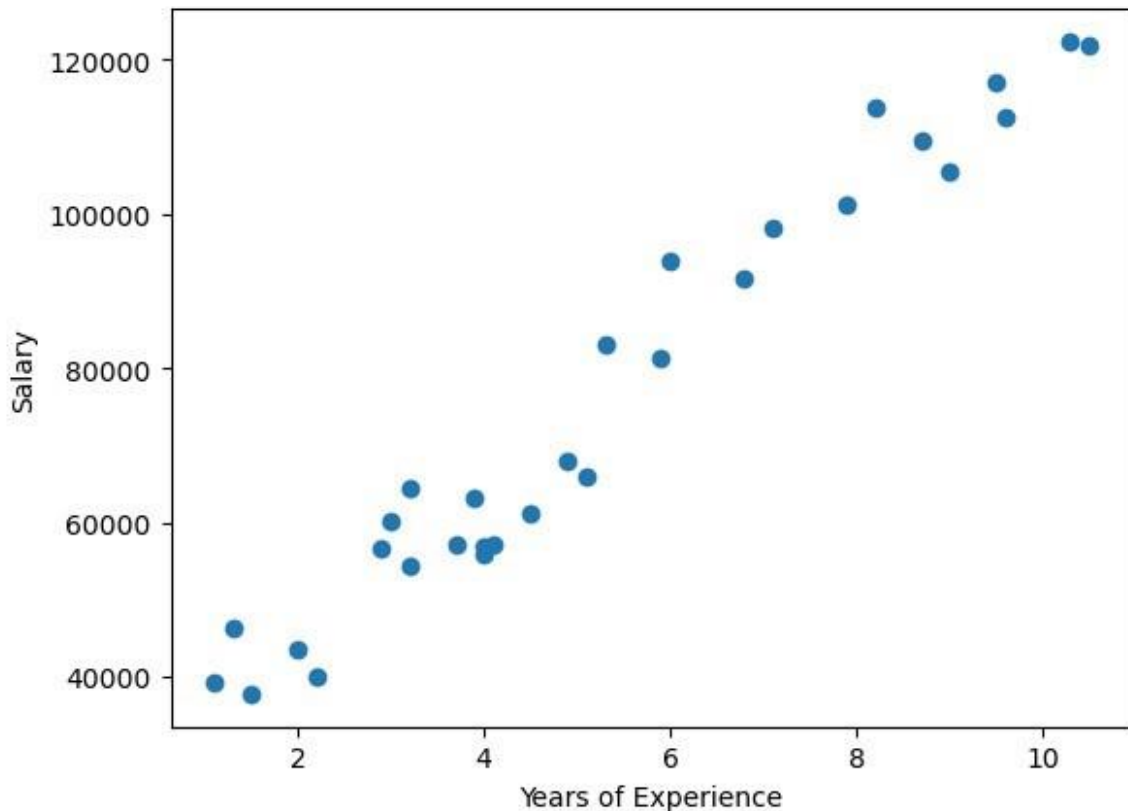
```
#pandas library for dataset handling and manipulation import pandas as pd
```

Import the dataset

```
df=pd.read_csv('archive/Salary_Data.csv') df.head()
YearsExperience Salary
0 1.1 39343
1 1.3 46205
2 1.5 37731
3 2.0 43525
4 2.2 39891
df.shape (30, 2)
```

Visualizing the data

```
plt.scatter(df['YearsExperience'],df['Salary'])
plt.xlabel('Years of Experience') plt.ylabel('Salary')
Text(0, 0.5, 'Salary')
```



Separating data into Independent/Features Variable and Dependent/Target Variable

```
X=np.array(df['YearsExperience']) Y=df['Salary']
```

Splitting the data into train and test set

```
from sklearn.model_selection import train_test_split
X_train,X_test,Y_train,Y_test=train_test_split(X,Y,test_size=0.2,random_state=0)
```

`train_test_split(X, Y, test_size=0.2, random_state=0)` splits the dataset `X` (features) and `Y` (labels) into training (80%) and testing (20%) subsets. The `test_size=0.2` determines the proportion of the data for testing, while `random_state=0` ensures reproducibility of the split. It returns four subsets: `X_train`, `X_test`, `Y_train`, and `Y_test`.

Linear Regression Model

```
#Importing the LinearRegression model
from sklearn.linear_model import LinearRegression
#making an instance of LinearRegression
model=LinearRegression() #fitting the data to the model
model.fit(X_train.reshape(-1,1),Y_train)#X should be 2d array while fitting the data
```

```
LinearRegression()
```

Model Parameters (Slope and Intercept)

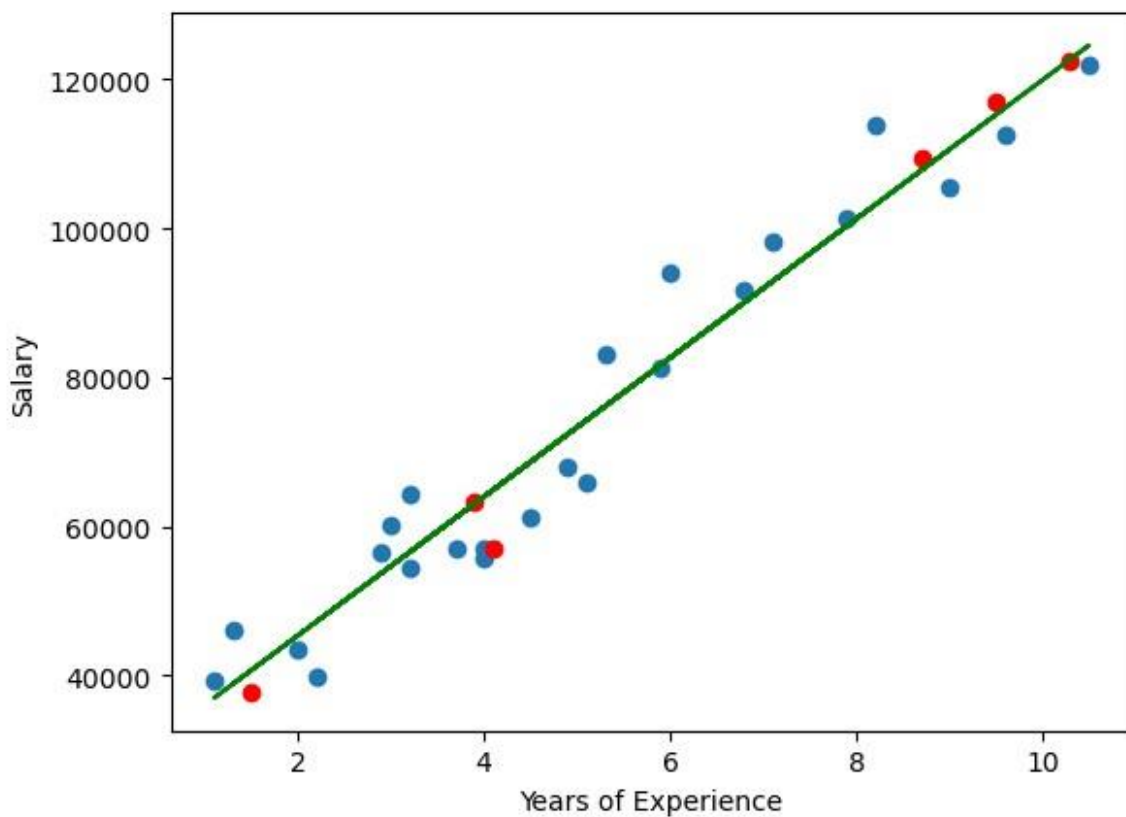
```
print(f'Slope : {model.coef_[0]}') print(f'Intercept: {model.intercept_}')  
Slope :9312.575126729187  
Intercept: 26780.099150628186
```

Predicting the Salary

```
Y_pred=model.predict(X_test.reshape(-1,1)) Y_pred  
array([ 40748.96184072, 122699.62295594, 64961.65717022,  
63099.14214487,  
115249.56285456, 107799.50275317])
```

Visualizing the model

```
model_eq=model.coef_[0]*X_train+model.intercept_  
plt.scatter(X_train,Y_train)#training data plt.scatter(X_test,Y_test,c='red')#test data  
plt.plot(X_train,model_eq,c='green')#linear regression model plt.xlabel('Years of Experience')  
plt.ylabel('Salary')  
Text(0, 0.5, 'Salary')
```



Regression Model Performance metrics

```
from sklearn.metrics import  
mean_absolute_error,mean_squared_error,r2_score  
print('Mean Absolute Error :',mean_absolute_error(Y_test,Y_pred)) print('Mean Squarred Error  
:',mean_squared_error(Y_test,Y_pred))  
print('R2 Value: ',r2_score(Y_test,Y_pred))  
Mean Absolute Error : 2446.1723690465064  
Mean Squarred Error : 12823412.298126562  
R2 Value: 0.988169515729126
```

Conclusion:

Merits and Demerits of Linear Regression

Merits:

1. **Simplicity:** Linear regression is easy to understand and interpret, making it a good starting point for regression analysis.
2. **Efficiency:** Computationally efficient, especially for smaller datasets, and works well when the relationship between dependent and independent variables is linear.
3. **Low Variance:** With fewer parameters, it is less likely to overfit, especially in cases with fewer features.
4. **Good interpretability:** The coefficients in linear regression offer clear interpretations of the relationship between variables (i.e., the change in the dependent variable for a unit change in the independent variable).
5. **Widely used:** Suitable for many real-world problems, and often provides good results when assumptions are met.

Demerits:

1. **Assumes Linearity:** Linear regression assumes a linear relationship between the features and target variable, which might not always hold in real-world data.
2. **Sensitive to Outliers:** Outliers can have a significant impact on the regression model, leading to poor predictions.
3. **Multicollinearity:** If the independent variables are highly correlated, the model's coefficients can become unstable and difficult to interpret.
4. **Limited Flexibility:** It is not capable of capturing complex relationships in the data that are non-linear.
5. **Requires Assumptions:** Linear regression assumes normality of errors, homoscedasticity (constant variance of errors), and no multicollinearity, which might not be satisfied in real datasets.

Experiment No.2

Name: Atharv Vijay Deshpande

PRN: 21410044

Batch: EN-3

Title: Implementation of Logistic Regression Model

Part A: Implementation of Logistic Regression Model on Synthetic data

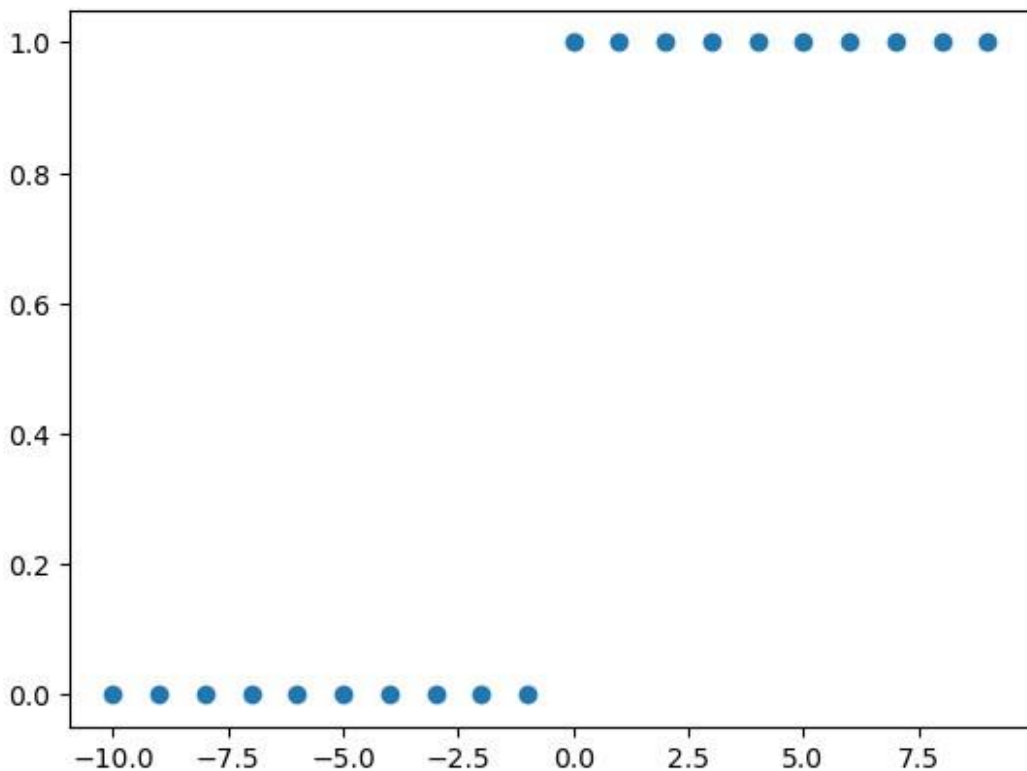
Importing Libraries

```
import numpy as np
import pandas as pd
import matplotlib.pyplot as plt
import seaborn as sns
import warnings
warnings.filterwarnings('ignore')
```

Creating Sample Data

```
X=np.array(np.arange(-10,10))
Y=np.array([0,0,0,0,0,0,0,0,0,0,1,1,1,1,1,1,1,1]) plt.scatter(X,Y)
```

<matplotlib.collections.PathCollection at 0x704701423e30>



Sigmoid function

```
def sigmoid(x):    return 1/(1+np.exp(-x))
```

Sigmoid Function in Logistic Regression

The **sigmoid function** is used in logistic regression to map predicted values to probabilities between 0 and 1. It is crucial for binary classification tasks, as it converts the linear output into a probability.

Sigmoid Function Formula:

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

Where:

- (z) is the input (usually the linear combination of features and weights),
- (e) is Euler's number (approximately 2.718).

Explanation:

- The sigmoid function takes any real number (z) and outputs a value between 0 and 1.
- When (z) is very large and positive, sigma(z) approaches 1, indicating a high probability of the positive class.
- When (z) is very large and negative, sigma(z) approaches 0, indicating a low probability of the positive class.
- When (z = 0), (sigma(z) = 0.5), representing a 50% probability.

This function allows logistic regression to model probabilities, which is essential for making decisions in binary classification.

Logistic Regression Model

```
from sklearn.linear_model import LogisticRegression
log=LogisticRegression()
log.fit(X.reshape(-1,1),Y)

LogisticRegression()
```

Predicting the Values

```
Y_predicted=log.predict(X.reshape(-1,1))
print(Y_predicted)

[0 0 0 0 0 0 0 0 0 1 1 1 1 1 1 1 1 1]
```

Linear Regression Model on Same data

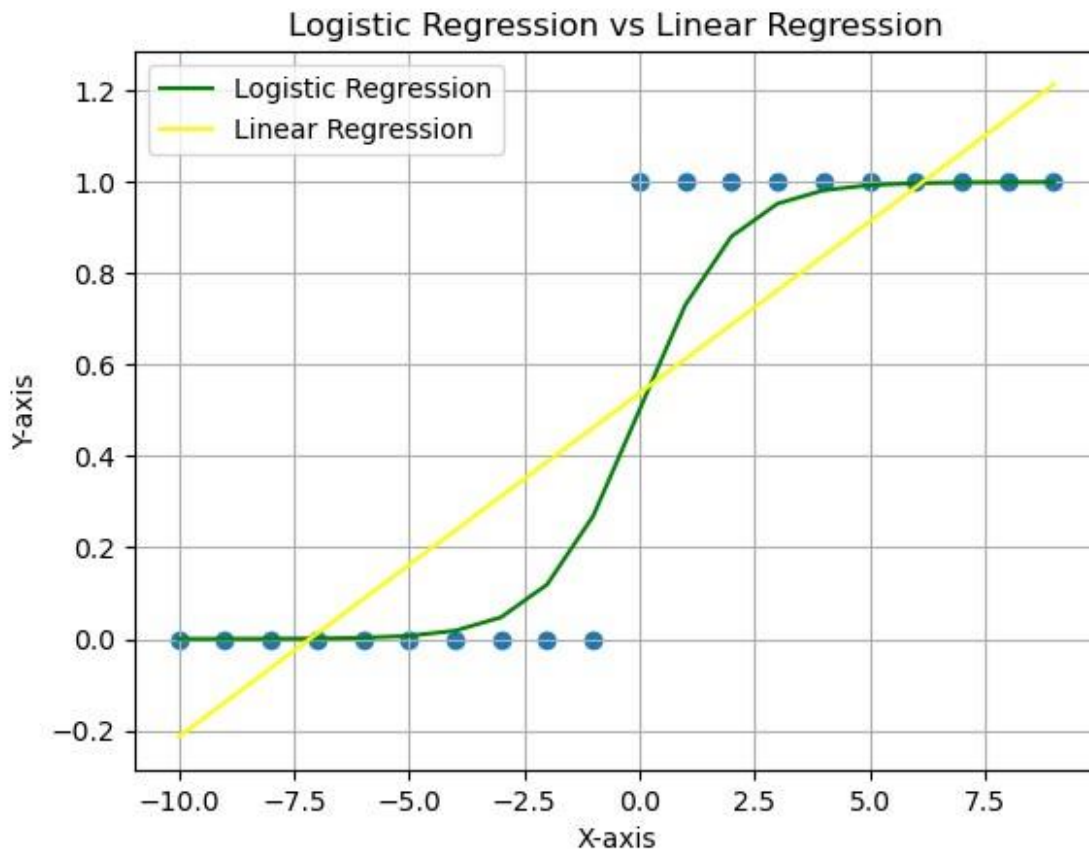
```
from sklearn.linear_model import LinearRegression
lr=LinearRegression()
lr.fit(X.reshape(-1,1),Y)

LinearRegression()

lrp=lr.predict(X.reshape(-1,1))
```

Visualizing the model output of Logistic and Linear Regression

```
plt.scatter(X,Y)
plt.plot(X,sigmoid(X),c='green',label='Logistic Regression')
plt.plot(X,lrp,c='yellow',label='Linear Regression') plt.xlabel('X-axis') plt.ylabel('Y-axis')
plt.title('Logistic Regression vs Linear Regression') plt.grid() plt.legend()
plt.show()
```



Logistic Regression Performance metrics

```

from sklearn.metrics import confusion_matrix, accuracy_score, precision_score, recall_score, f1_score

conf_matrix = confusion_matrix(Y, Y_predicted) print("Confusion Matrix:\n", conf_matrix)

accuracy = accuracy_score(Y, Y_predicted) print("Accuracy:", accuracy)
Confusion Matrix:
[[10  0]
 [ 0 10]]
Accuracy: 1.0

precision = precision_score(Y, Y_predicted) recall =
recall_score(Y, Y_predicted) f1 = f1_score(Y, Y_predicted)

print("Precision:", precision) print("Recall:",
recall) print("F1 Score:", f1)

Precision: 1.0
Recall: 1.0
F1 Score: 1.0

```

Part B: Implementation of Logistic Regression Model on a dataset

#pandas library for dataset handling and manipulation import pandas as pd

Import dataset

```

df=pd.read_csv('archive/Social_Network_Ads.csv') df.head()
  Age  EstimatedSalary  Purchased
0   19         19000         0
1   35         20000         0
2   26         43000         0
3   27         57000         0
4   19         76000         0

df.shape
(400, 3)
df.isna().sum()#checking for null values
Age      0
EstimatedSalary  0

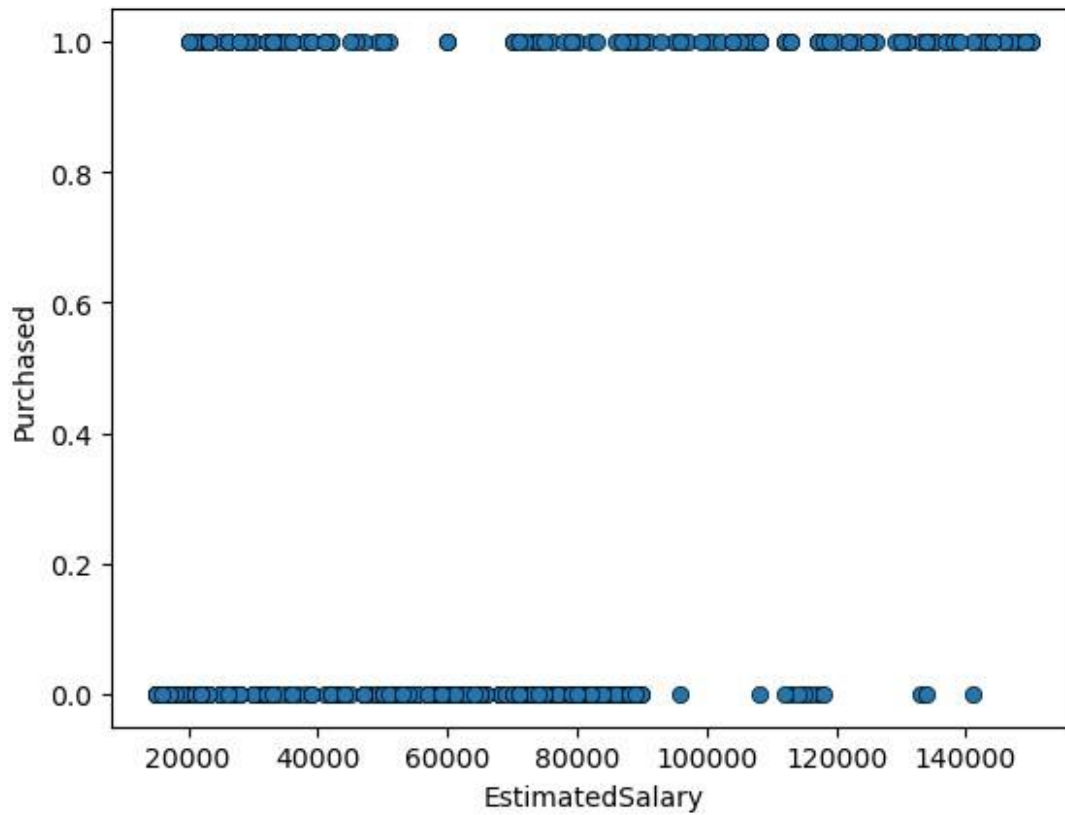
```

```
Purchased      0  
dtype: int64
```

Visualizing the data

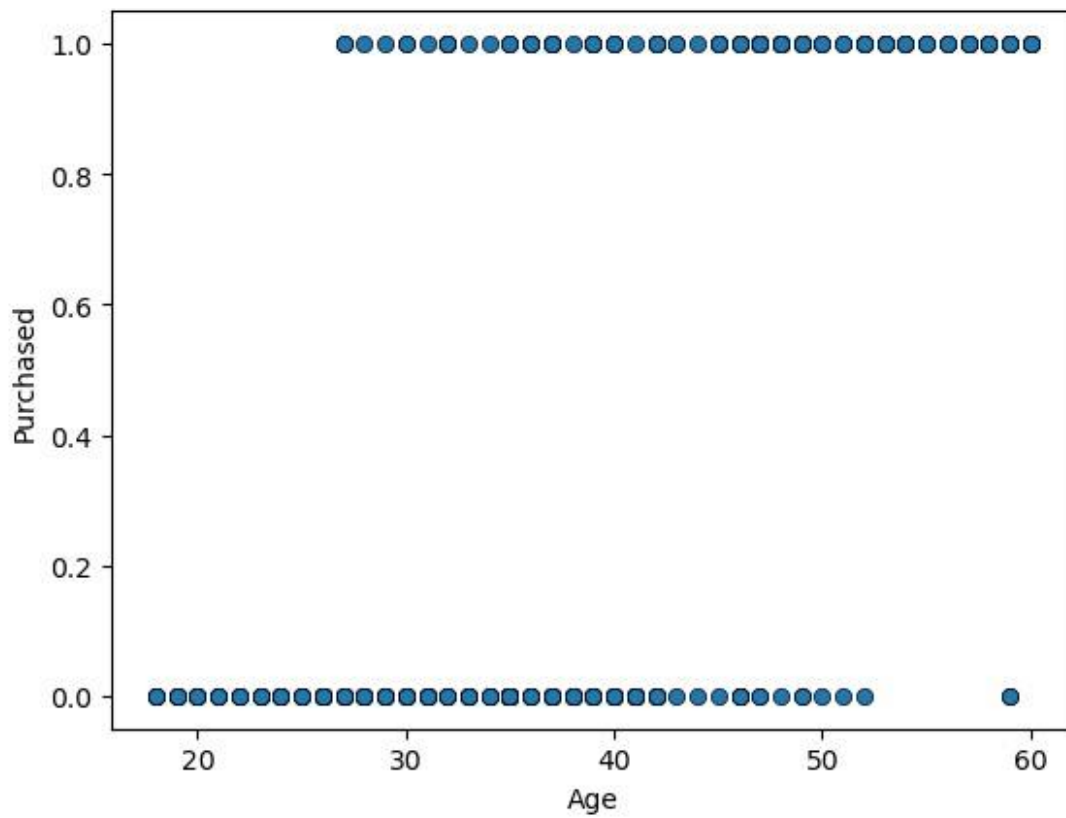
```
sns.scatterplot(x=df['EstimatedSalary'], y=df['Purchased'], edgecolor='black')
```

```
<Axes: xlabel='EstimatedSalary', ylabel='Purchased'>
```

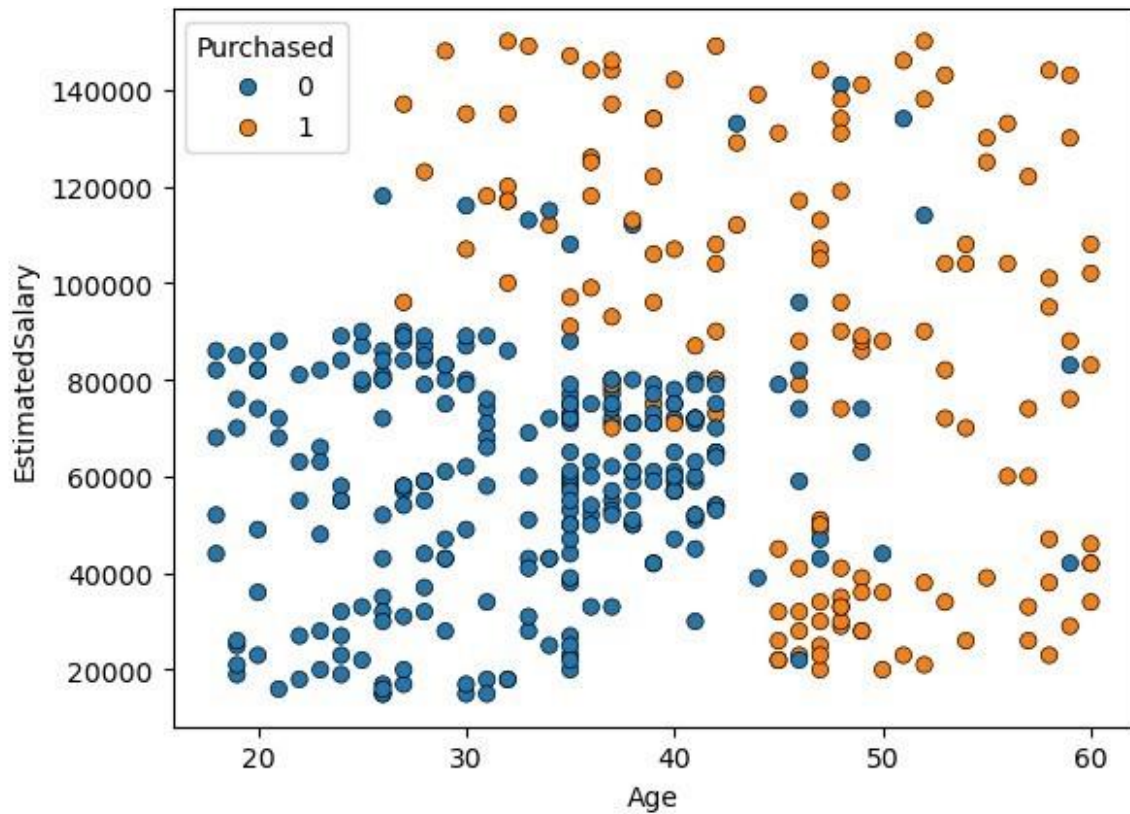


```
sns.scatterplot(x=df['Age'], y=df['Purchased'], edgecolor='black')
```

```
<Axes: xlabel='Age', ylabel='Purchased'>
```



```
sns.scatterplot(x=df['Age'], y=df['EstimatedSalary'], hue=df['Purchased'], edgecolor='black')  
<Axes: xlabel='Age', ylabel='EstimatedSalary'>
```



Separating data into Feature Variables and Class

```
X=df[['Age']]  
Y=df[['Purchased']]
```

Splitting the data into train and test set

```
from sklearn.model_selection import train_test_split
```

```
X_train,X_test,Y_train,Y_test=train_test_split(X,Y,test_size=0.2,random_state=0)
```

Logistic Regression Model

```
from sklearn.linear_model import LogisticRegression
log=LogisticRegression() log.fit(X_train,Y_train) LogisticRegression()
```

Predicting the Class for Test data

```
0 0 1 0 0 0 0 1 0 0 1 0 1 1 0 0 0 1 1 0 0 1 0 0 1 0 0 0 1 0 0 0 0 0 0
0 1
```

```
0 0 0 0 1 1]
```

```
Result= pd.DataFrame({
    'Actual': Y_test,
    'Predicted': Y_pred
})
```

Result

	Actual	Predicted
132	0	0
309	0	0
341	0	0
196	0	0
246	0	0
14	0	0
363	0	0
304	0	0
361	1	1
329	1	1

[80 rows x 2 columns]

Performance Evaluation of the model

```
Y_pred=log.predict(X_test) print(Y_pred)
[0 0 0 0 0 0 0 0 0 1 0 0 0 0 0 0 0 0 1 0 0 1 0 1 0 1 0 0 0 0 0 0 1 0 0
0 0
```



```

from sklearn.metrics import confusion_matrix, accuracy_score, precision_score,
recall_score, f1_score
conf_matrix = confusion_matrix(Y_test, Y_pred)
print("Confusion Matrix:\n", conf_matrix)

accuracy = accuracy_score(Y_test, Y_pred)
print("Accuracy:", accuracy)
precision = precision_score(Y_test, Y_pred)
recall = recall_score(Y_test, Y_pred)
f1 = f1_score(Y_test, Y_pred)

print("Precision:", precision)
print("Recall:", recall)
print("F1 Score:", f1)

Confusion Matrix:
[[57  1]
 [ 4 18]]
Accuracy: 0.9375
Precision: 0.9473684210526315
Recall: 0.8181818181818182
F1 Score: 0.8780487804878049

```

Conclusion:

Merits and Demerits of Logistic Regression

Merits

1. **Simplicity and Interpretability:**
 - Logistic regression is easy to implement and interpret. The coefficients represent the log odds of the dependent variable, making it straightforward to understand the influence of predictors.
2. **Efficiency:**
 - It is computationally efficient and performs well on smaller datasets. The model can be trained quickly compared to more complex algorithms.
3. **Probabilistic Output:**
 - Logistic regression provides probabilities for class membership, allowing for nuanced decision-making. This is particularly useful in applications where uncertainty needs to be quantified.
4. **Works Well with Linearly Separable Data:**
 - The algorithm performs well when the classes are linearly separable, meaning that a linear decision boundary can effectively separate the classes.
5. **Feature Scaling Not Required:**
 - Unlike other algorithms, logistic regression does not require normalization or standardization of features, simplifying preprocessing.

Demerits

1. **Assumes Linear Relationship:**

- Logistic regression assumes a linear relationship between the independent variables and the log odds of the dependent variable, which may not hold in realworld data.
- 2. **Sensitivity to Outliers:**
 - The model can be sensitive to outliers, which may disproportionately influence the fitted model. This can lead to misleading interpretations.
- 3. **Limited to Binary Classification:**
 - While logistic regression can be extended to multiclass classification (using techniques like one-vs-all), it is inherently designed for binary outcomes.
- 4. **Overfitting:**
 - In cases with many features or multicollinearity among predictors, the model may overfit the training data, leading to poor generalization on unseen data.
- 5. **Assumes Independence of Features:**
 - Logistic regression assumes that the features are independent of each other. In cases of multicollinearity, the results can be unreliable.

Experiment No.3

Name: Atharv Vijay Deshpande

PRN: 21410044

Batch: EN-3

Title: Implementation of Polynomial Regression

Part A: Implementation of Polynomial Regression Model on Synthetic data

Importing Libraries

```
import numpy as np
import pandas as pd
import matplotlib.pyplot as plt
import seaborn as sns
import warnings
warnings.filterwarnings('ignore')
```

Creating and Visualizing Sample Data

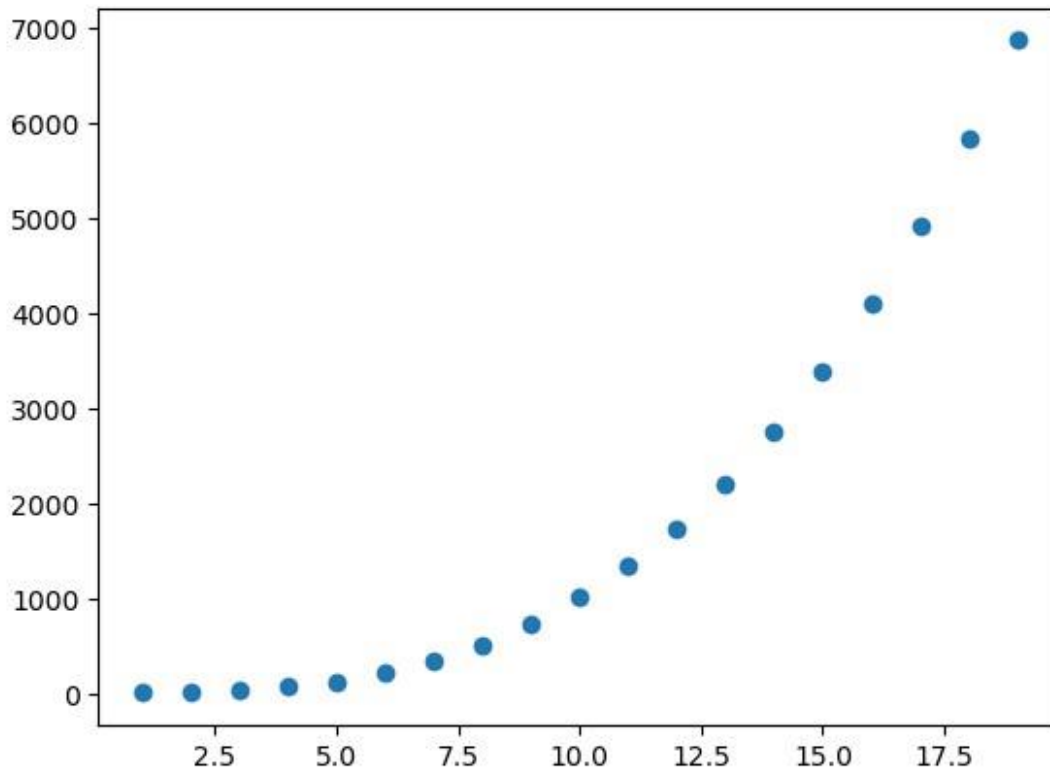
```
X=np.array(np.arange(1,20))
Y=X**3+np.random.rand(19)*10

X
array([ 1,  2,  3,  4,  5,  6,  7,  8,  9, 10, 11, 12, 13, 14, 15, 16,
        17,
        18, 19])

Y
array([ 8.75404301, 16.08069724, 34.36633826, 72.08061354, 125.40025663,
        217.18219078, 345.25871733, 512.47385627,
        736.6442064 , 1009.84971831, 1335.62058038, 1734.53885729,
        2204.96372044, 2748.35439452, 3378.97354101, 4101.88486556,
        4922.48194735, 5834.03732789, 6864.68642583])

plt.scatter(X,Y)

<matplotlib.collections.PathCollection at 0x7aafbe5d01a0>
```



Polynomial Regression Model

```
# Fit a quadratic polynomial to the data coefficients = np.polyfit(X, Y, 2)
# Print the coefficients
print("Coefficients (a, b, c):", coefficients) # Create a polynomial function
from the coefficients polynomial = np.poly1d(coefficients)
```

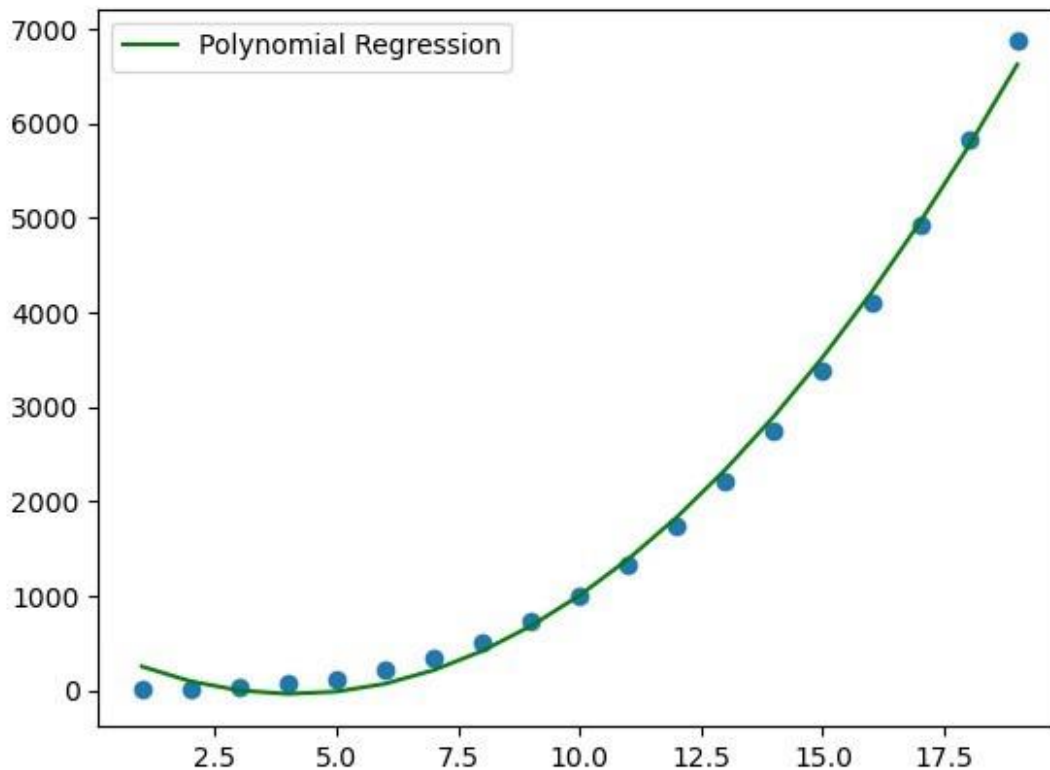
Coefficients (a, b, c): [30.02106986 -246.63638694 469.07911863]

Generate predicted values using the polynomial function model

```
# Generate predicted values using the polynomial function
Y_pred = polynomial(X) Y_pred
array([ 2.52463802e+02,  9.58906242e+01, -6.40413428e-01, -
 3.71293113e+01,
 -1.35760695e+01,  7.00193121e+01,  2.13656833e+02,
 4.17336494e+02,
 6.81058295e+02,  1.00482224e+03,  1.38862832e+03,
 1.83247654e+03,
 2.33636690e+03,  2.90029939e+03,  3.52427403e+03,
 4.20829081e+03,
 4.95234973e+03,  5.75645079e+03,  6.62059399e+03])
```

Visualizing the Model

```
plt.scatter(X,Y)
plt.plot(X,Y_pred,'Green',label='Polynomial Regression') plt.legend() plt.show()
```



Performance Metrics

```
from sklearn.metrics import
mean_absolute_error,mean_squared_error,r2_score
print('Mean Absolute Error :',mean_absolute_error(Y,Y_pred)) print('Mean Squared Error
:',mean_squared_error(Y,Y_pred))
print('R2 Value: ',r2_score(Y,Y_pred))
Mean Absolute Error : 109.40921032360683
Mean Squared Error : 15845.049731937492
R2 Value: 0.996410488067558
```

Part B: Implementation of Polynomial Regression Model on a dataset

Import dataset

```
df=pd.read_csv('archive/Ice_cream selling data.csv') df.head()
```

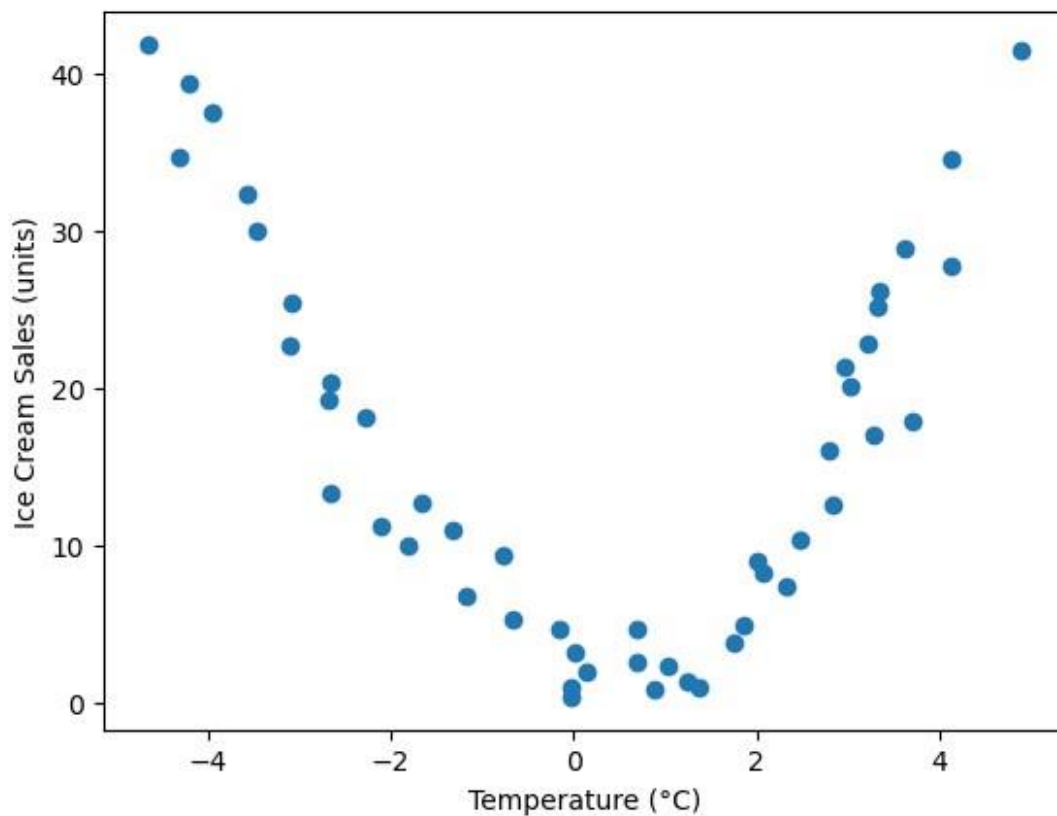
	Temperature (°C)	Ice Cream Sales (units)
0	-4.662263	41.842986
1	-4.316559	34.661120
2	-4.213985	39.383001
3	-3.949661	37.539845
4	-3.578554	32.284531

df.shape (49, 2)

Visualizing the data

```
plt.scatter(df['Temperature (°C)'],df['Ice Cream Sales (units)']) plt.xlabel('Temperature (°C)')
plt.ylabel('Ice Cream Sales (units)')
```

```
Text(0, 0.5, 'Ice Cream Sales (units)')
```



Separating data into Independent/Features Variable and Dependent/Target Variable

```
X=np.array(df['Temperature (°C)'])
```

```
Y=df['Ice Cream Sales (units)']
```

Splitting the data into train and test set

```
from sklearn.model_selection import train_test_split
```

```
X_train,X_test,Y_train,Y_test=train_test_split(X,Y,test_size=0.2,random_state=42)
```

Polynomial Regression model

```
# Fit a cubic polynomial to the data
coefficients = np.polyfit(X_train, Y_train, 3)
# Print the coefficients
print("Coefficients (a, b, c, d):", coefficients) # Create a polynomial function
from the coefficients polynomial = np.poly1d(coefficients)
Coefficients (a, b, c, d): [ 0.05451597  1.87501919 -1.39956426
2.84053099]
```

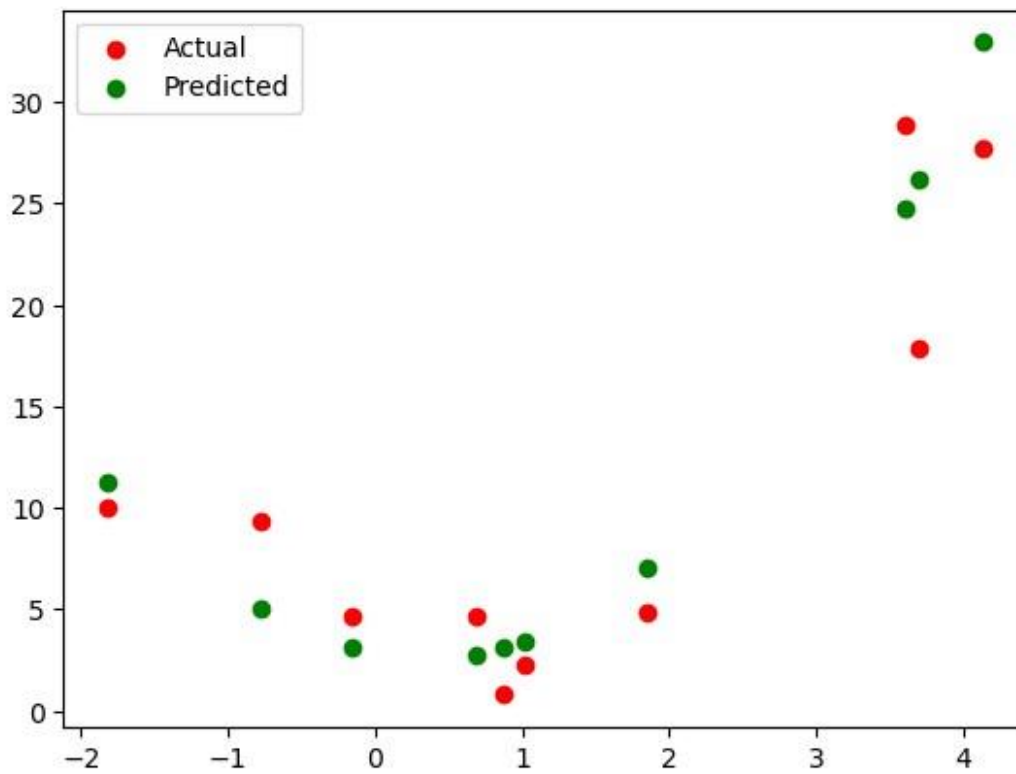
Generate predicted values using the polynomial function

```
# Generate predicted values using the polynomial function
Y_pred = polynomial(X_test) Y_pred

array([11.26173782, 26.15229744, 32.94239669, 24.79940325,
5.0189787 ,
3.43248573, 3.08780469, 2.79001857, 7.01713112,
3.09175474])
```

Visualizing the model on test data

```
#plt.scatter(X_train,Y_train)
plt.scatter(X_test,Y_test,c='red',label='Actual')
plt.scatter(X_test,Y_pred,c='green',label='Predicted') plt.legend() plt.show()
```



Performance Metrics

```
from sklearn.metrics import
mean_absolute_error,mean_squared_error,r2_score
print('Mean Absolute Error :',mean_absolute_error(Y_test,Y_pred)) print('Mean Squared Error
:',mean_squared_error(Y_test,Y_pred))
print('R2 Value: ',r2_score(Y_test,Y_pred))
Mean Absolute Error : 3.2281211297768864
Mean Squared Error : 15.120009451229418
R2 Value: 0.8405107685716922
```

Conclusion:

Merits and Demerits of Polynomial Regression

Merits

1. **Flexibility:**
 - Polynomial regression can model complex relationships between independent and dependent variables by fitting higher-degree polynomials, allowing for curvature in the relationship.
2. **Improved Fit for Non-linear Data:**
 - It performs better than linear regression when the underlying relationship is nonlinear, as it can capture patterns that linear models cannot.
3. **Easy Interpretation of Coefficients:**
 - Each coefficient in the polynomial can still be interpreted in terms of its effect on the dependent variable, similar to linear regression.
4. **Works Well with Small Datasets:**
 - Polynomial regression can provide reasonable fits with smaller datasets when higher-order terms are included, as it can model non-linearity without needing large amounts of data.
5. **Extensive Applicability:**
 - It is applicable in various fields such as finance, biology, and engineering, where relationships between variables are often non-linear.

Demerits

1. **Overfitting:**
 - Higher-degree polynomials can fit the training data very well, but they may fail to generalize to unseen data, leading to poor predictive performance. This is especially true in the presence of noise.
2. **Increased Complexity:**
 - The model becomes more complex with higher degrees, making it harder to interpret and understand the influence of individual features.
3. **Extrapolation Issues:**
 - Polynomial regression can yield extreme predictions outside the range of the training data, especially with higher-degree polynomials, leading to unreliable results.
4. **Sensitive to Outliers:**

- Polynomial regression can be heavily influenced by outliers, which may skew the fit and result in misleading interpretations.

5. **Requires Feature Engineering:**

- Polynomial regression necessitates the creation of polynomial features (e.g., x^2 , x^3) from the original variables, increasing the dimensionality of the feature space and potentially leading to multicollinearity issues.

Experiment No.4

Name: Atharv Vijay Deshpande

PRN: 21410044

Batch : EN-3

Title: Implementation of KNN Classifier on appropriate dataset

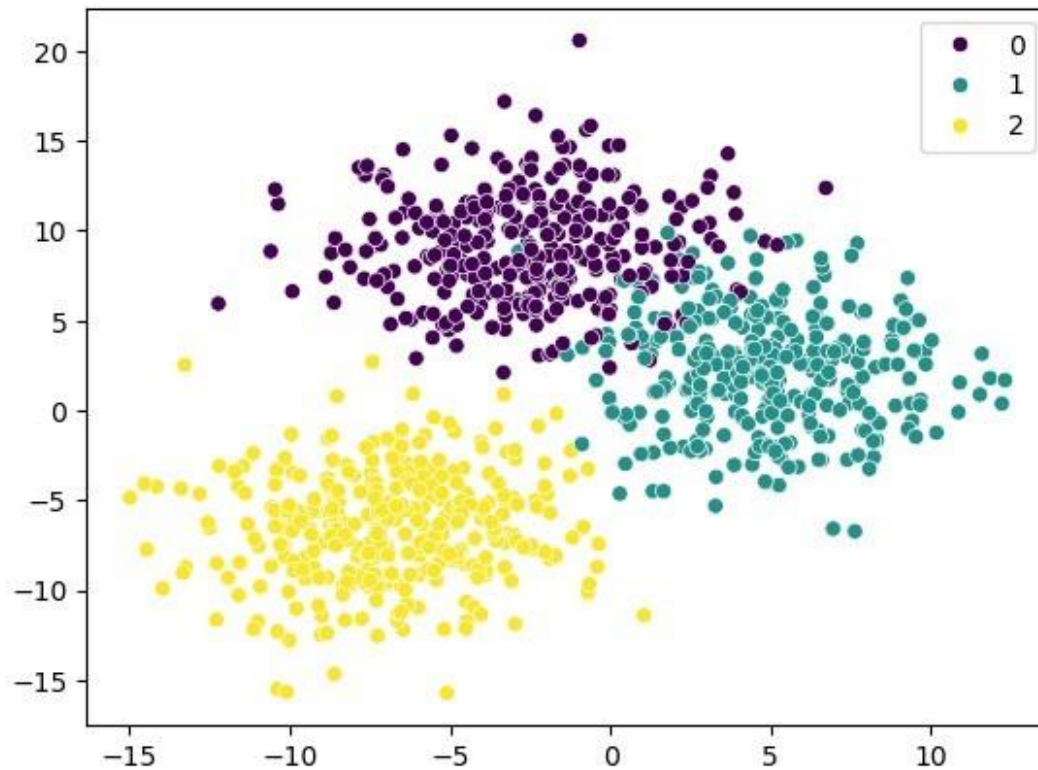
Part A: Implementation of Knn Classifier on Synthetic data

Importing Libraries

```
import numpy as np import pandas  
as pd  
import matplotlib.pyplot as plt import seaborn  
as sns import warnings  
from sklearn.datasets import make_blobs warnings.filterwarnings('ignore')
```

Creating and Visualizing Synthetic data

```
X,Y=make_blobs(n_samples=1000,n_features=2,centers=3,cluster_std=3,random_state=42)  
sns.scatterplot(x=X[:, 0], y=X[:, 1], hue=Y,palette='viridis') plt.show()
```



KNN Classifier Models

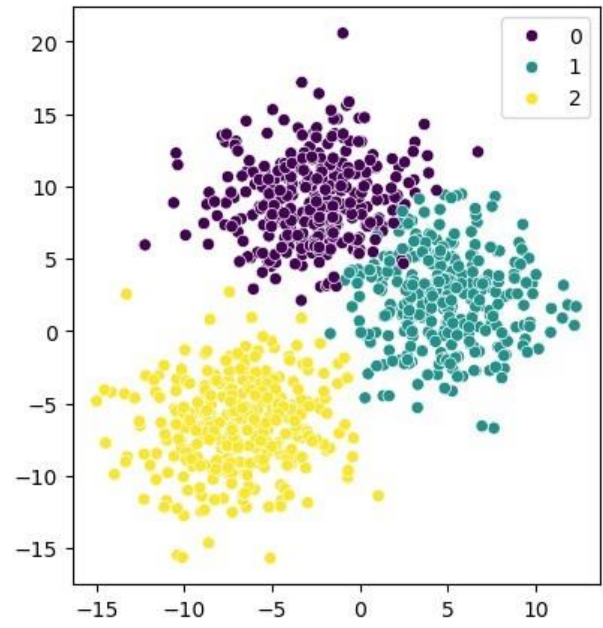
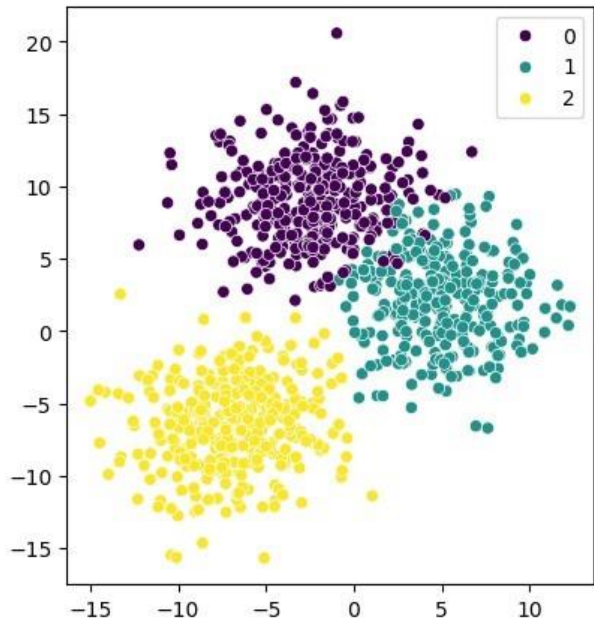
```
from sklearn.neighbors import KNeighborsClassifier
knn_classifier1=KNeighborsClassifier(n_neighbors=3)
knn_classifier2=KNeighborsClassifier(n_neighbors=5) knn_classifier1.fit(X,Y)
KNeighborsClassifier(n_neighbors=3)
knn_classifier2.fit(X,Y) KNeighborsClassifier()
```

Predicting the Labels

```
y_p1=knn_classifier1.predict(X) y_p2=knn_classifier2.predict(X)

plt.figure(figsize=(10,5)) plt.subplot(1,2,1)
sns.scatterplot(x=X[:, 0], y=X[:, 1], hue=y_p1,palette='viridis')

plt.subplot(1,2,2)
sns.scatterplot(x=X[:, 0], y=X[:, 1], hue=y_p2,palette='viridis') <Axes: >
```



Performance Metrics

```

from sklearn.metrics import
accuracy_score,confusion_matrix,classification_report

print('Model with n_neighbors=3')
print('accuracy_score:',accuracy_score(Y,y_p1)) print('confusion_matrix: ')
print(confusion_matrix(Y,y_p1)) print('classification_report: ')
print(classification_report(Y,y_p1))

print('Model with n_neighbors=5')
print('accuracy_score:',accuracy_score(Y,y_p2)) print('confusion_matrix: ')
print(confusion_matrix(Y,y_p2)) print('classification_report: ')
print(classification_report(Y,y_p2))

Model with n_neighbors=3 accuracy_score: 0.979
confusion_matrix:
[[328  6  0]
 [ 13 319  1]
 [  1  0 332]] classification_report:      precision  recall f1-score  support
0      0.96    0.98    0.97    334
1      0.98    0.96    0.97    333
2      1.00    1.00    1.00    333

```

accuracy	0.98	1000	macro avg	0.98	0.98	0.98
1000 weighted avg	0.98	0.98	0.98	1000		

```

Model with n_neighbors=5 accuracy_score: 0.968
confusion_matrix:
[[317 17  0]
 [ 13 319  1]
 [  0  1 332]] classification_report:      precision  recall f1-score  support
0      0.96    0.95    0.95    334
1      0.95    0.96    0.95    333
2      1.00    1.00    1.00    333

```

accuracy	0.97	1000	macro avg	0.97	0.97	0.97
1000 weighted avg	0.97	0.97	0.97	1000		

Part B:Implementation of Knn Classifier on a dataset

Load Penguins dataset from seaborn datasets

```
sns.get_dataset_names()
```

```
['anagrams',  
 'anscombe',  
 'attention',  
 'brain_networks',  
 'car_crashes',  
 'diamonds',  
 'dots',  
 'dowjones',  
 'exercise',  
 'flights',  
 'fmri',  
 'geyser',  
 'glue',  
 'healthexp',  
 'iris',  
 'mpg',  
 'penguins',  
 'planets',  
 'seaice',  
 'taxis',
```

```

'tips',
'titanic']
df=sns.load_dataset('penguins') df.head()

species  island  bill_length_mm  bill_depth_mm  flipper_length_mm
0
1      Adelie  Torgersen        39.5         17.4         186.0
2      Adelie  Torgersen        40.3         18.0         195.0
3      Adelie  Torgersen         NaN          NaN          NaN 4 Adelie  Torgersen        36.7         19.3
193.0
body_mass_g  sex  0
3750.0  Male
1      3800.0  Female
2      3250.0  Female
3         NaN   NaN
4      3450.0  Female
df.info()

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 344 entries, 0 to 343 Data columns (total
7 columns):
#  Column          Non-Null Count  Dtype
---  -
0   species        344 non-null   object
1   island          344 non-null   object
2   bill_length_mm  342 non-null   float64
3   bill_depth_mm   342 non-null   float64
4   flipper_length_mm 342 non-null   float64
5   body_mass_g     342 non-null   float64 6  sex          333 non-null   object
dtypes: float64(4), object(3) memory usage: 18.9+
KB

df.isnull().sum()#Checking for null values
species      0 island      0
bill_length_mm      2
bill_depth_mm      2
flipper_length_mm    2
body_mass_g      2

```

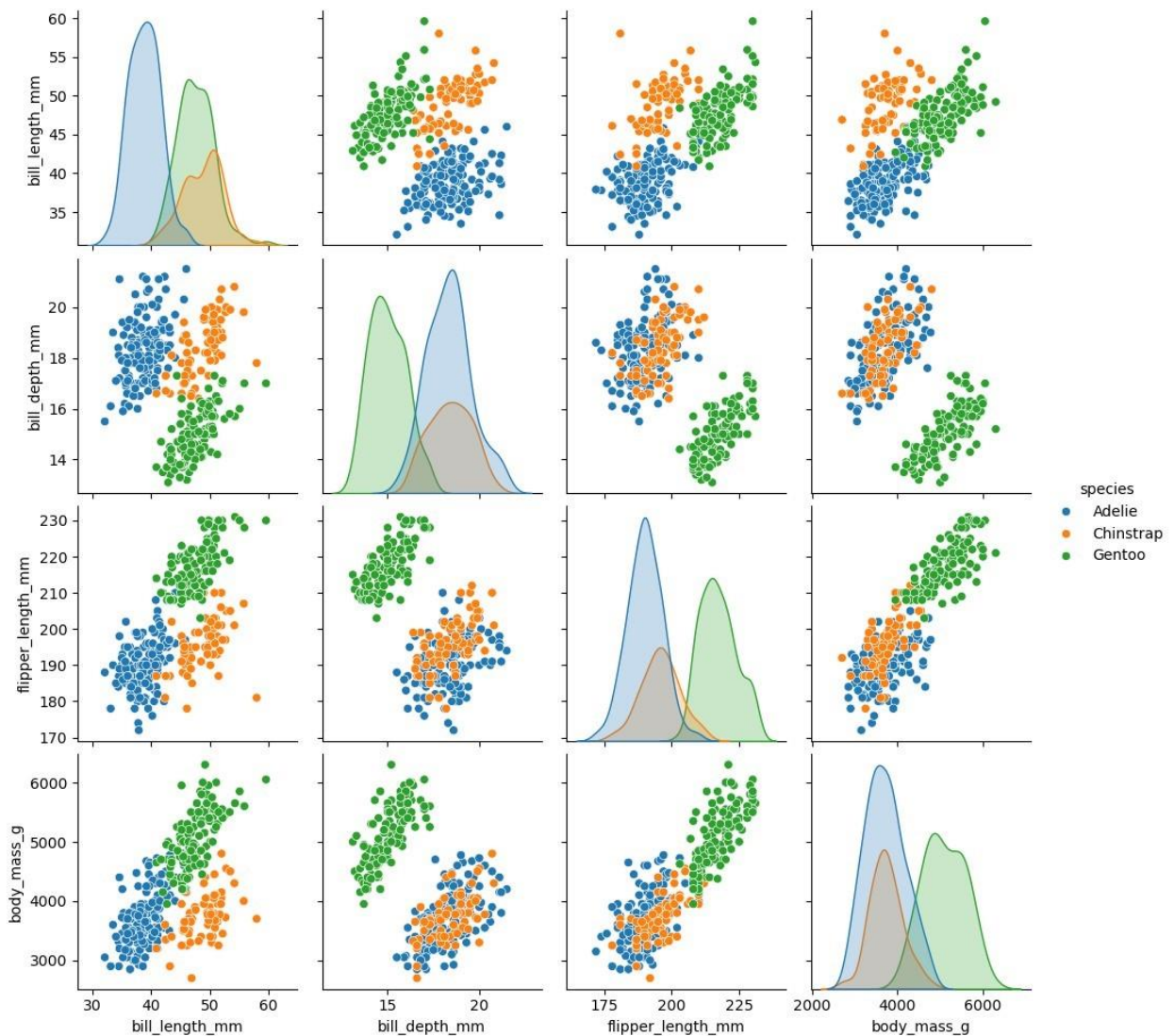
```
sex          11  
dtype: int64
```

```
df.dropna(inplace=True)#removing all null value rows df.shape (333, 7)
```

Visualizing the data for finding best features for classification

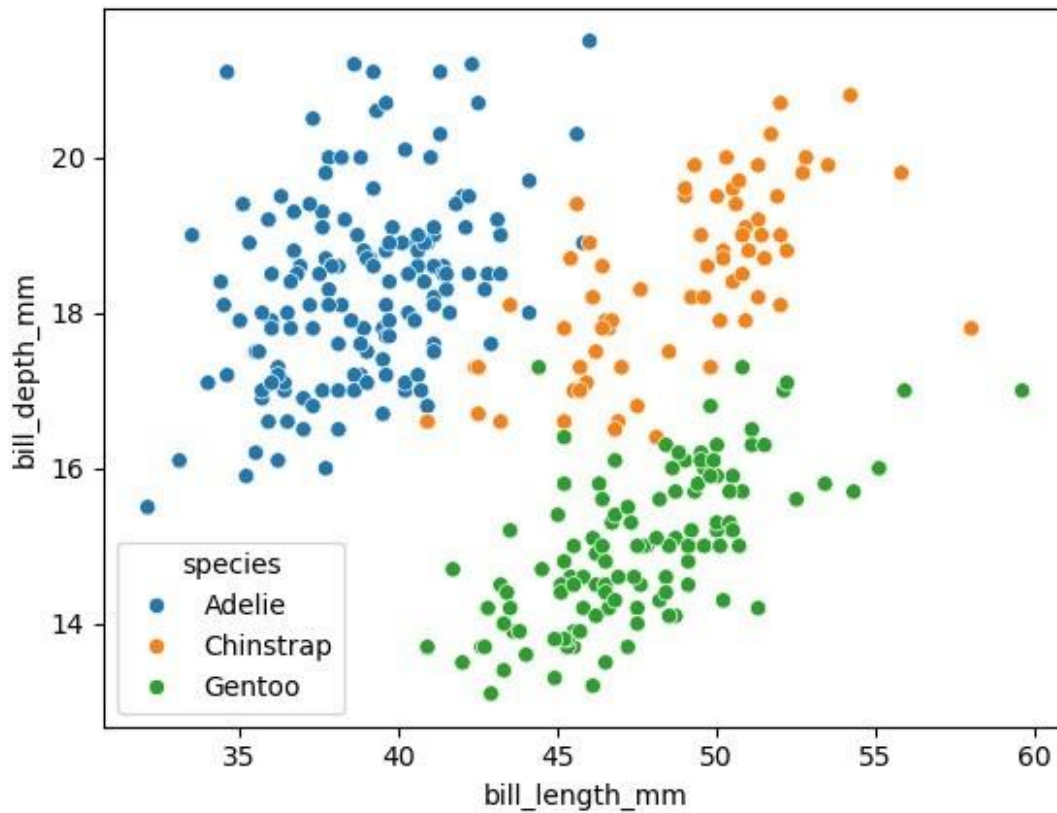
```
sns.pairplot(df,hue='species')
```

```
<seaborn.axisgrid.PairGrid at 0x743e0b04ea80>
```



Separating data into features and Labels/Target


```
X=df[['bill_length_mm','bill_depth_mm']] Y=df['species']
sns.scatterplot(x='bill_length_mm',y='bill_depth_mm',hue='species',data=df)
<Axes: xlabel='bill_length_mm', ylabel='bill_depth_mm'>
```



Splitting the data into train and test sets

```
from sklearn.model_selection import train_test_split
X_train,X_test,Y_train,Y_test=train_test_split(X,Y,test_size=0.2,random_state=0)
```

KNN Classifier Model

```
from sklearn.neighbors import KNeighborsClassifier
knn_classifier=KNeighborsClassifier(n_neighbors=5) knn_classifier.fit(X_train,Y_train)
KNeighborsClassifier()
```

Predicting Penguin Species on test data

```
y_pred=knn_classifier.predict(X_test) y_pred
array(['Adelie', 'Adelie', 'Gentoo', 'Adelie', 'Adelie', 'Adelie',
       'Gentoo', 'Gentoo', 'Gentoo', 'Chinstrap', 'Gentoo', 'Adelie',
```

```
'Adelie', 'Chinstrap', 'Adelie', 'Adelie', 'Gentoo', 'Adelie',
'Chinstrap', 'Adelie', 'Adelie', 'Adelie', 'Gentoo', 'Gentoo',
'Gentoo', 'Gentoo', 'Adelie', 'Adelie', 'Adelie', 'Adelie',
'Adelie', 'Chinstrap', 'Adelie', 'Adelie', 'Adelie', 'Gentoo',
'Chinstrap', 'Adelie', 'Chinstrap', 'Adelie', 'Gentoo',
'Gentoo',
'Adelie', 'Adelie', 'Adelie', 'Adelie', 'Adelie', 'Adelie',
'Gentoo', 'Adelie', 'Adelie', 'Adelie', 'Gentoo', 'Chinstrap',
'Adelie', 'Adelie', 'Adelie', 'Adelie', 'Adelie', 'Gentoo',
'Adelie', 'Chinstrap', 'Adelie', 'Gentoo', 'Adelie', 'Adelie',
'Gentoo'], dtype=object)
```

Model Performance Evaluation

```
from sklearn.metrics import
accuracy_score, confusion_matrix, classification_report
print('accuracy_score:', accuracy_score(Y_test, y_pred)) print('confusion_matrix: ')
print(confusion_matrix(Y_test, y_pred)) print('classification_report: ')
print(classification_report(Y_test, y_pred))
```

```
accuracy_score: 0.9402985074626866 confusion_matrix:
[[39 0 0]
 [ 2 7 1]
 [ 0 1 17]] classification_report:      precision    recall  f1-score   support
 Adelie      0.95      1.00      0.97       39
 Chinstrap    0.88      0.70      0.78       10
 Gentoo      0.94      0.94      0.94       18
 accuracy                0.94      67  macro avg      0.92      0.88      0.90
 67 weighted avg      0.94      0.94      0.94      67
```

Conclusion:

K-Nearest Neighbors (KNN) Classifier

Merits of KNN Classifier

1. **Simplicity:**
 - KNN is easy to understand and implement, making it a great choice for beginners.
2. **No Training Phase:**

- KNN is a lazy learner, meaning it does not require a training phase. It simply stores the training data and makes predictions based on that data.
- 3. **Flexibility:**
 - It can be used for both classification and regression tasks, providing versatility.
- 4. **Adaptability:**
 - KNN can work with any number of classes and does not make any assumptions about the underlying data distribution.
- 5. **Effectiveness with Large Datasets:**
 - It can perform well with large datasets where decision boundaries are complex.
- 6. **Local Decision Making:**
 - KNN uses local information to make decisions, making it robust to outliers.

Demerits of KNN Classifier

1. **Computational Complexity:**
 - KNN requires computing the distance between the test instance and all training samples, which can be time-consuming, especially with large datasets.
2. **Memory Intensive:**
 - Since KNN stores all the training data, it can consume a significant amount of memory.
3. **Sensitivity to Feature Scaling:**
 - The performance of KNN can be adversely affected by the scale of features. Features need to be normalized or standardized.
4. **Choice of K:**
 - Selecting the optimal number of neighbors (K) can be challenging. A small K can lead to noise, while a large K can smooth out important patterns.
5. **Curse of Dimensionality:**
 - As the number of dimensions increases, the distance between points becomes less meaningful, which can degrade performance.
6. **Imbalanced Data:**
 - KNN can be biased towards the majority class in imbalanced datasets, potentially leading to poor classification performance for minority classes.

Experiment No.5

Name: Atharv Vijay Deshpande

PRN: 21410044

Batch : EN-3

Title: Study of Confusion Matrix with appropriate example

Study of Confusion Matrix

Introduction

In classification problems, it's crucial to evaluate model performance accurately. The **confusion matrix** is one of the most insightful tools for this purpose. It provides a detailed breakdown of correct and incorrect predictions across classes, highlighting the types of errors the model is making, beyond a simple accuracy metric.

Confusion Matrix Definition

A confusion matrix is a table that is often used to describe the performance of a classification model. Each row of the matrix represents the instances in an actual class, while each column represents the instances in a predicted class. For binary classification, it is structured as:

	Predicted Positive	Predicted Negative
Actual Positive	True Positive (TP)	False Negative (FN)
Actual Negative	False Positive (FP)	True Negative (TN)

- **True Positive (TP):** Cases where the model correctly predicted the positive class.
- **True Negative (TN):** Cases where the model correctly predicted the negative class.
- **False Positive (FP):** Cases where the model incorrectly predicted the positive class (Type I error).
- **False Negative (FN):** Cases where the model incorrectly predicted the negative class (Type II error).

Example

Suppose we have a binary classifier predicting whether a patient has a disease (positive class) or not (negative class). The confusion matrix for a test set of 100 patients might look like this:

	Predicted Disease	Predicted No Disease
Actual Disease	45	5
Actual No Disease	10	40

From the above:

- **True Positives (TP) = 45:** Correctly predicted patients with the disease.

- **False Negatives (FN) = 5:** Missed predictions (patients with the disease predicted as no disease).
- **False Positives (FP) = 10:** Incorrectly predicted as having the disease.
- **True Negatives (TN) = 40:** Correctly predicted no disease.

Performance Metrics Derived from the Confusion Matrix

Several performance metrics can be derived to evaluate model performance:

1. **Accuracy:** Proportion of correctly predicted instances (both positive and negative) over the total instances.

$$\text{Accuracy} = \frac{TP + TN}{TP + TN + FP + FN}$$

2. **Precision:** Fraction of correctly predicted positive instances among all predicted positives.

$$\text{Precision} = \frac{TP}{TP + FP}$$

3. **Recall (Sensitivity):** Fraction of correctly predicted positive instances among all actual positives.

$$\text{Recall} = \frac{TP}{TP + FN}$$

4. **F1-Score:** Harmonic mean of precision and recall, giving a balanced measure of performance.

$$\text{F1-Score} = 2 \times \frac{\text{Precision} \times \text{Recall}}{\text{Precision} + \text{Recall}}$$

5. **Specificity:** Proportion of correctly predicted negative instances among all actual negatives.

$$\text{Specificity} = \frac{TN}{TN + FP}$$

Use Case Scenarios

The confusion matrix is particularly useful in scenarios where:

- **Class imbalance** exists, and accuracy alone would be misleading (e.g., fraud detection, medical diagnoses).
- **Error analysis** is critical. By inspecting false positives and false negatives, one can finetune the model to optimize for the most critical metric (e.g., improving recall in cancer detection).

Import Libraries

```
# Import necessary libraries import  
numpy as np import pandas as pd
```

```
from sklearn import datasets  
from sklearn.model_selection import train_test_split from sklearn.svm import SVC  
from sklearn.metrics import classification_report,  
confusion_matrix, ConfusionMatrixDisplay import matplotlib.pyplot as plt  
import seaborn as sns
```

Load the Iris dataset

```
iris = sns.load_dataset('iris') iris.head()
```

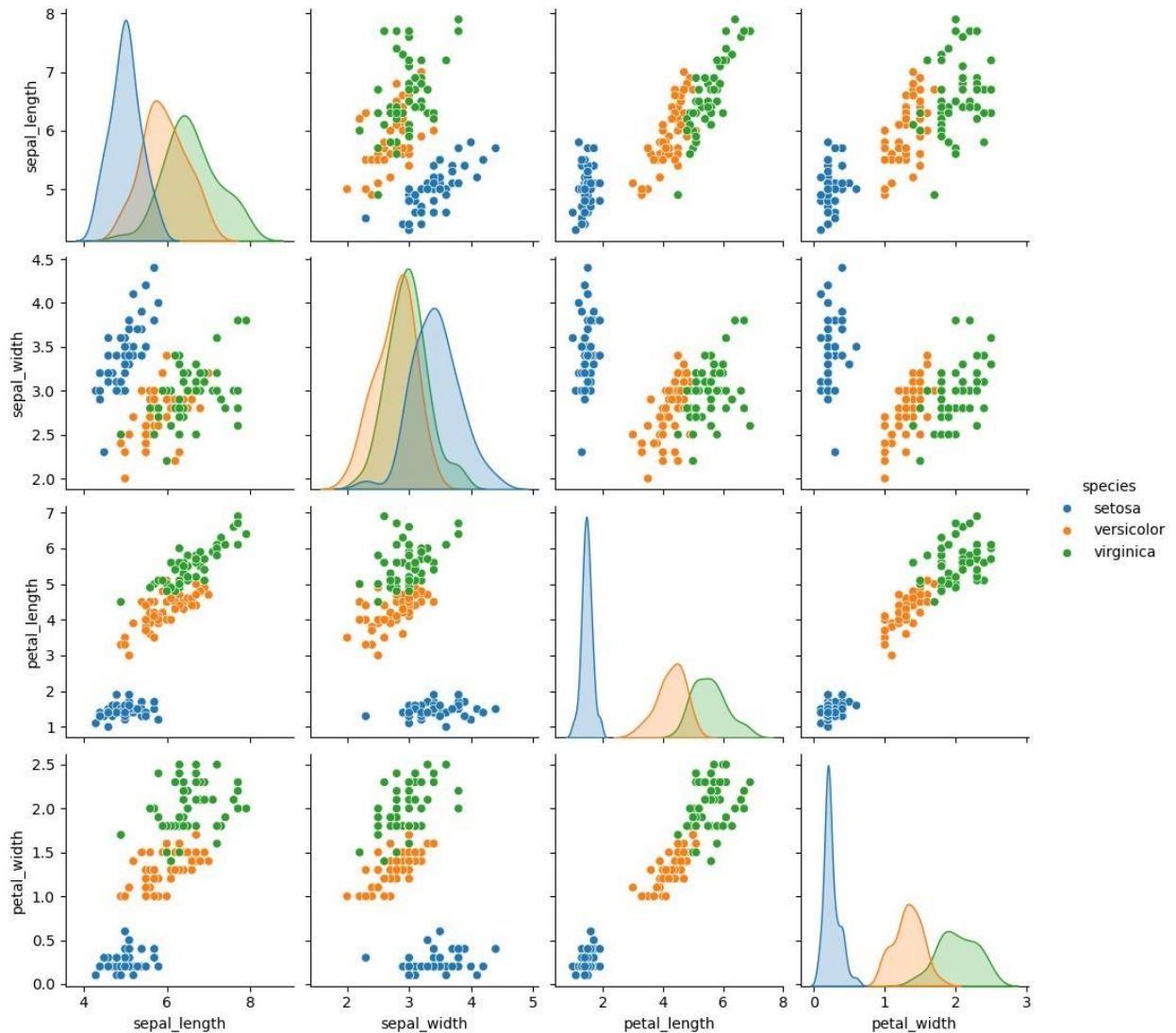
	sepal_length	sepal_width	petal_length	petal_width	species
0	5.1	3.5	1.4	0.2	setosa
1	4.9	3.0	1.4	0.2	setosa
2	4.7	3.2	1.3	0.2	setosa
3	4.6	3.1	1.5	0.2	setosa
4	5.0	3.6	1.4	0.2	setosa

```
iris.info()
```

```
<class 'pandas.core.frame.DataFrame'>  
RangeIndex: 150 entries, 0 to 149 Data columns (total  
5 columns):  
#   Column      Non-Null Count  Dtype  
---  ---  
0    sepal_length  150 non-null    float64  
1    sepal_width   150 non-null    float64  
2    petal_length  150 non-null    float64  
3    petal_width   150 non-null    float64  
4    species       150 non-null    object  
dtypes: float64(4), object(1)  
memory usage: 6.0+ KB
```

Visualize the data

```
sns.pairplot(data=iris, hue='species')  
  
<seaborn.axisgrid.PairGrid at 0x7dd278db8fb0>
```



Separating Feature and Target

```
X=iris[['sepal_width','petal_width']]
Y=iris['species']
```

Split data into train and test set

```
X_train, X_test, Y_train, Y_test = train_test_split(X, Y, test_size=0.2, random_state=42)
```

Create SVM Model

```
svm_model = SVC(kernel='linear', random_state=42) # You can use other kernels like 'rbf', 'poly'
svm_model.fit(X_train, Y_train)
```

```
SVC(kernel='linear', random_state=42)
```

Make predictions on the test data

```
Y_pred = svm_model.predict(X_test)
```

```
array(['versicolor', 'setosa', 'virginica', 'versicolor',
      'versicolor',
      'setosa', 'versicolor', 'virginica', 'versicolor',
      'versicolor',
      'virginica', 'setosa', 'setosa', 'setosa', 'setosa',
      'versicolor',
      'virginica', 'versicolor', 'versicolor', 'virginica', 'setosa',
      'virginica', 'setosa', 'virginica', 'virginica', 'virginica',
      'virginica', 'virginica', 'setosa', 'setosa'], dtype=object)
```

Evaluate the model using a confusion matrix and classification report

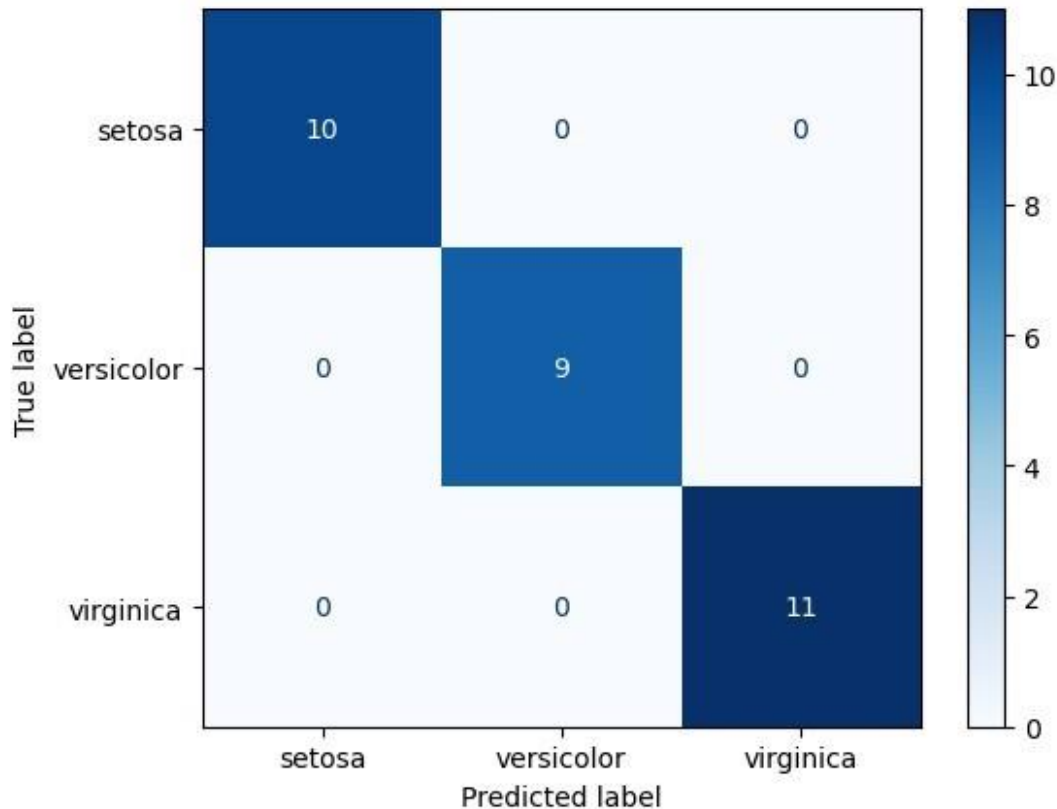
```
labels = iris['species'].unique() # The unique class labels for the target
cm = confusion_matrix(Y_test, Y_pred) print("Confusion Matrix:\n", cm)
disp = ConfusionMatrixDisplay(confusion_matrix=cm,
display_labels=labels) disp.plot(cmap='Blues')
print("\nClassification Report:\n", classification_report(Y_test, Y_pred))
```

Confusion Matrix:

```
[[10  0  0]
 [ 0  9  0]
 [ 0  0 11]]
```

Classification Report:

	precision	recall	f1-score	support					
setosa	1.00	1.00	1.00	10	versicolor	1.00	1.00	1.00	
			9	9	virginica	1.00	1.00	1.00	11
accuracy			1.00	30	macro avg	1.00	1.00	1.00	
30 weighted avg	1.00	1.00	1.00	30					



Conclusion:

The **study of the confusion matrix** provides a deeper understanding of a classification model's performance, beyond what a single metric like accuracy can offer. By examining the true positives, true negatives, false positives, and false negatives, the confusion matrix allows us to assess not only how well the model predicts each class, but also the types of errors it makes. This is particularly important when dealing with imbalanced datasets or when the cost of false positives and false negatives differs.

For instance, in the **Iris dataset** example, using a Support Vector Machine (SVM) model, we could observe the confusion matrix detailing how the model classifies each species of iris flower (setosa, versicolor, and virginica). The matrix helps us understand where the model makes errors, such as misclassifying virginica as versicolor or vice versa. Using metrics derived from the confusion matrix, like **precision**, **recall**, and **F1-score**, we can evaluate how the model balances between false positives and false negatives across each class.

Through this experiment, we learned:

- **Accuracy** alone may not be sufficient to evaluate model performance, especially in cases where class imbalance or uneven error costs exist.
- The confusion matrix provides insights into specific types of errors (false positives and false negatives), which can guide model tuning and help improve model performance for critical applications, such as medical diagnoses or fraud detection.
- By visualizing the confusion matrix, it becomes easier to interpret the performance of the classifier and take steps to mitigate the model's weaknesses, like increasing recall for one class or improving precision for another.

Overall, the confusion matrix serves as a vital tool in understanding and improving classification models by offering a granular view of predictions and errors.

Experiment No.6

Name: Atharv Vijay Deshpande

PRN: 21410044

Batch : EN-3

Title: Implementation of Support Vector Machine

Part A: Implementation of SVM Model on Synthetic data

Importing Libraries

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

Creating and Visualizing synthetic dataset

```
# Create a dataframe for apples
apples = pd.DataFrame({
    'weight': np.random.uniform(200.0, 300.0, size=100),
    'circumference': np.random.uniform(40.0, 50.0, size=100),
    'fruit': 'Apple'
})

# Create a dataframe for oranges
oranges = pd.DataFrame({
    'weight': np.random.uniform(100.0, 200.0, size=100),
    'circumference': np.random.uniform(20.0, 40.0, size=100),
    'fruit': 'Orange'
})

# Concatenate the two dataframes
df = pd.concat([apples, oranges])
df.head()
  weight  circumference  fruit
0  223.253891         43.690499  Apple
1   298.310891         48.903060  Apple
2   244.214135         44.880593  Apple
3   210.796607         48.338364  Apple
4   263.650477         49.161789  Apple

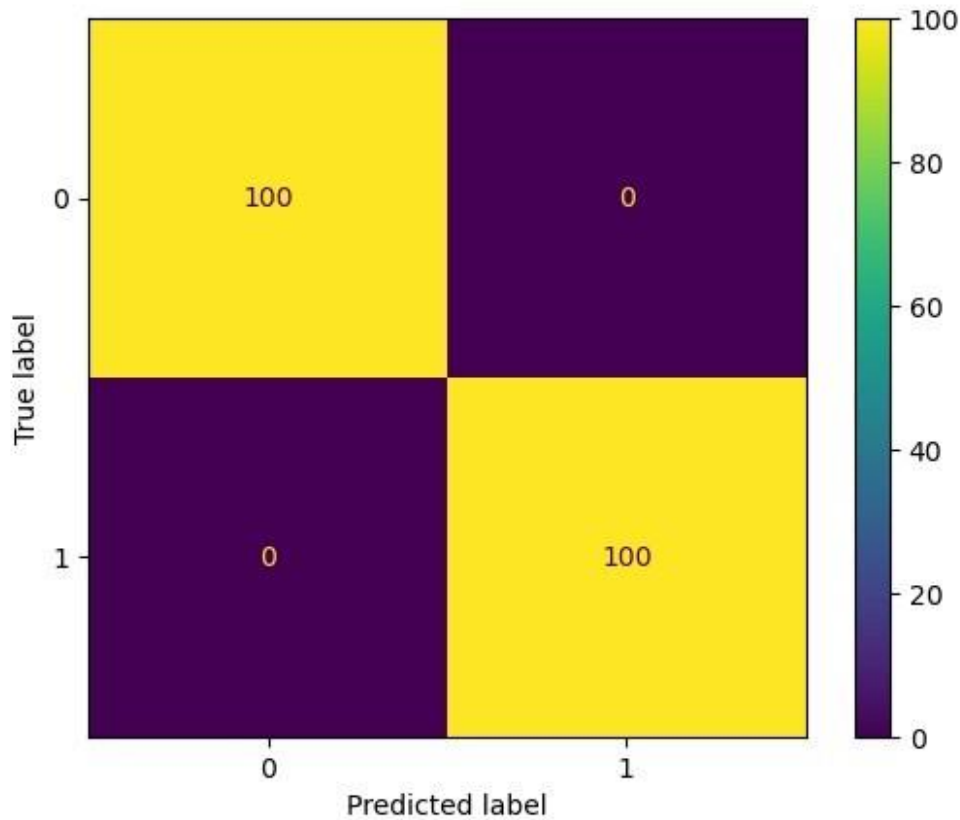
sns.scatterplot(x='weight', y='circumference', hue='fruit', data=df)
plt.show()
```


'Orange',

Performance Metrics

```
from sklearn.metrics import accuracy_score, confusion_matrix, classification_report,
ConfusionMatrixDisplay
print(accuracy_score(y, y_pred))
print(confusion_matrix(y, y_pred))
print(classification_report(y, y_pred))
ConfusionMatrixDisplay(confusion_matrix(y, y_pred)).plot()
plt.show()
```

1.0								
[[100 0]								
[0 100]]								
	precision	recall	f1-score	support				
Apple	1.00	1.00	1.00	100				
Orange	1.00	1.00	1.00	100				
accuracy			1.00	200	macro avg	1.00	1.00	1.00
200 weighted avg		1.00	1.00	1.00	200			

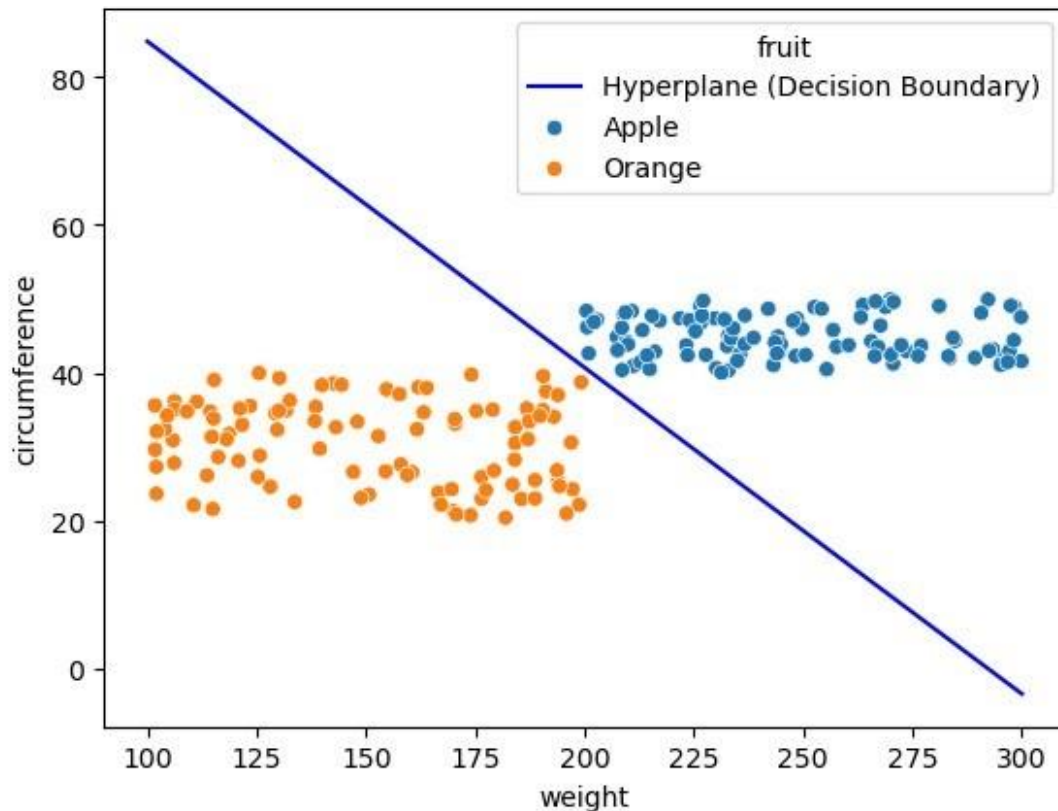


Get the coefficients (w) and intercept (b) of the hyperplane

```
w = model.coef_[0] b =
model.intercept_[0] print('w:',w)
print('b:',b)
w: [-0.18971444 -0.42958037]
b: 55.41741104047813
```

Visualizing the SVM Model Hyperplane

```
x_vals = np.linspace(100, 300, 100) y_vals = -
(w[0]/w[1])*x_vals - b/w[1]
plt.plot(x_vals, y_vals, color='blue', label='Hyperplane (Decision Boundary)')
sns.scatterplot(x='weight', y='circumference', hue='fruit', data=df) plt.show()
```



Part B: Implementation of SVM Model on a dataset

```
# Import necessary libraries
import numpy as np
import pandas as pd
from sklearn import datasets
from sklearn.model_selection import train_test_split
from sklearn.svm import SVC
from sklearn.metrics import classification_report, confusion_matrix
import matplotlib.pyplot as plt
import seaborn as sns

# Import the dataset
iris = sns.load_dataset('iris')
iris.head()
```

	sepal_length	sepal_width	petal_length	petal_width	species
0	5.1	3.5	1.4	0.2	setosa
1	4.9	3.0	1.4	0.2	setosa
2	4.7	3.2	1.3	0.2	setosa
3	4.6	3.1	1.5	0.2	setosa
4	5.0	3.6	1.4	0.2	setosa

```
# Separating feature and target variables
```

```

X=iris[['sepal_width','petal_width']] Y=iris['species']

#Split data into train and test set
X_train, X_test, Y_train, Y_test = train_test_split(X, Y, test_size=0.2, random_state=42)

X_train.shape, X_test.shape, Y_train.shape, Y_test.shape ((120, 2), (30, 2), (120,), (30,))

#Create SVM Model
svm_model = SVC(kernel='linear', random_state=42) # You can use other kernels like 'rbf', 'poly'
svm_model.fit(X_train, Y_train) SVC(kernel='linear', random_state=42)

#Make Predictions
Y_pred = svm_model.predict(X_test) Y_pred

array(['versicolor', 'setosa', 'virginica', 'versicolor',
       'versicolor',
       'setosa', 'versicolor', 'virginica', 'versicolor',
       'versicolor',
       'virginica', 'setosa', 'setosa', 'setosa', 'setosa',
       'versicolor',
       'virginica', 'versicolor', 'versicolor', 'virginica', 'setosa',
       'virginica', 'setosa', 'virginica', 'virginica', 'virginica',
       'virginica', 'virginica', 'setosa', 'setosa'], dtype=object)

#Evaluate the model performance
from sklearn.metrics import accuracy_score, confusion_matrix, classification_report,
ConfusionMatrixDisplay
labels = iris['species'].unique() # The unique class labels for the target
print(accuracy_score(Y_test, Y_pred)) print(confusion_matrix(Y_test, Y_pred))
print(classification_report(Y_test, Y_pred))
ConfusionMatrixDisplay(confusion_matrix(Y_test,
Y_pred),display_labels=labels).plot() plt.show()

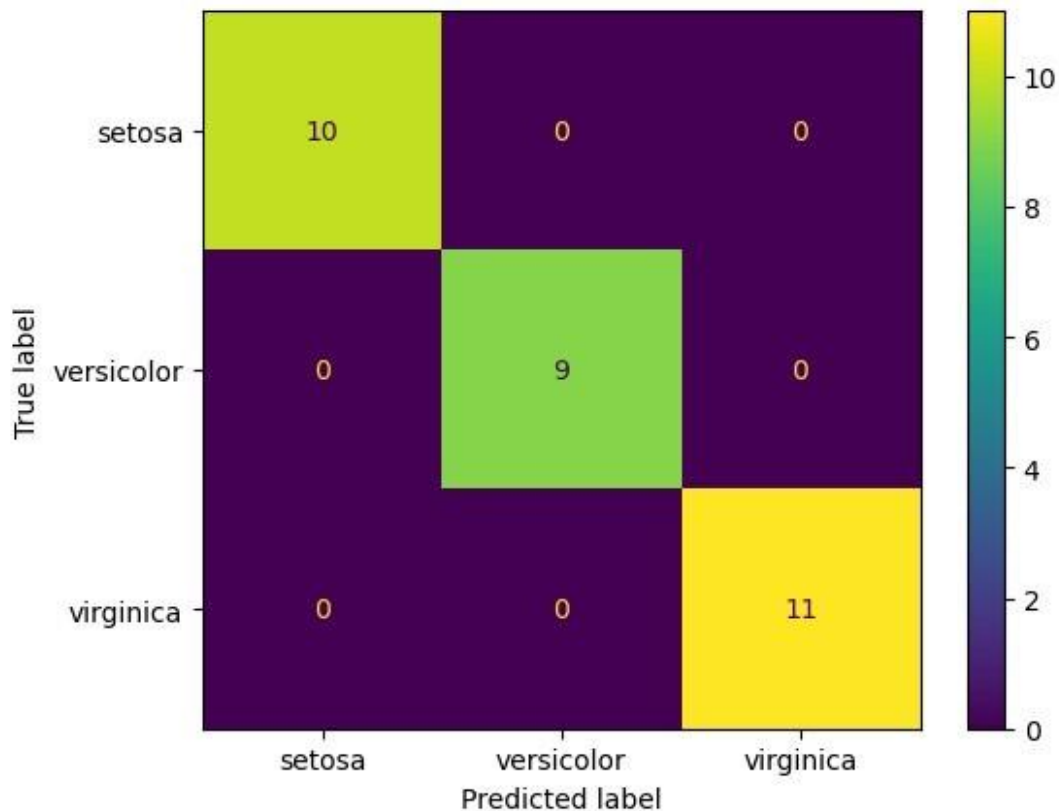
```

```

1.0
[[10 0 0]
 [ 0 9 0]
 [ 0 0 11]]
precision recall f1-score support

```

setosa	1.00	1.00	1.00	10	versicolor	1.00	1.00	1.00
			9	virginica	1.00	1.00	1.00	11
accuracy			1.00	30	macro avg	1.00	1.00	1.00
30 weighted avg	1.00	1.00	1.00	30				



Get the coefficients (w) and intercept (b) of the hyperplanes

```
w0 = svm_model.coef_[0]
w1 = svm_model.coef_[1]
w2 = svm_model.coef_[2]
b = svm_model.intercept_[0]
print('w0:',w[0])
print('w1:',w[1])
print('w2:',w[2])
print('b:',b)

w0: [ 0.88988726 -1.94381961]
w1: [ 0.45917166 -1.85798801]
w2: [ 0.90852965 -3.63658848]
b: -1.1919105788982298
```

Conclusion:

Merits and Demerits of Support Vector Machine (SVM)

Merits of SVM

1. Effective in High-Dimensional Spaces:

- SVM is particularly effective in high-dimensional spaces and is still effective when the number of dimensions exceeds the number of samples.
- 2. **Versatile Kernel Trick:**
 - SVM can be adapted to various classification tasks using different kernel functions (linear, polynomial, radial basis function, etc.), allowing it to capture complex relationships.
- 3. **Robust to Overfitting:**
 - Due to the concept of maximizing the margin between classes, SVM is less prone to overfitting, especially in high-dimensional spaces.
- 4. **Clear Margin of Separation:**
 - SVM provides a clear margin of separation between classes, making the classification process intuitive and interpretable.
- 5. **Memory Efficiency:**
 - SVM uses a subset of training points in the decision function (support vectors), making it memory efficient compared to other algorithms that might require all training data.
- 6. **Works Well with Unbalanced Data:**
 - SVM can handle unbalanced datasets by adjusting the penalty parameter, which helps in improving model performance.

Demerits of SVM

1. **Computationally Intensive:**
 - Training SVMs can be time-consuming, especially with large datasets, as the algorithm's complexity grows with the number of samples.
2. **Choice of Kernel and Parameters:**
 - Selecting the appropriate kernel and tuning parameters can be challenging. Poor choices can lead to underfitting or overfitting.
3. **Sensitivity to Noisy Data:**
 - SVM can be sensitive to noise in the data, especially when there are overlapping classes. This can affect the placement of the decision boundary.
4. **Difficulties with Large Datasets:**
 - While SVMs are effective in high-dimensional spaces, they struggle with very large datasets, where algorithms like logistic regression or decision trees might perform better.
5. **Binary Classification:**
 - SVM is inherently a binary classifier. While it can be extended to multi-class classification (using techniques like one-vs-one or one-vs-all), this adds complexity and can reduce performance.
6. **Interpretability:**
 - While SVM provides a clear margin of separation, understanding the model's decisions can be less interpretable compared to simpler models like linear regression or decision trees.

Experiment No.7

Name: Atharv Vijay Deshpande

PRN: 21410044

Batch : EN-3

Title: Implementation of K-Means Clustering Algorithm on appropriate dataset

Part A: Implementation of K-Means Clustering Algorithm on Synthetic data

Import Libraries

```
import pandas as pd
import matplotlib.pyplot as plt
import seaborn as sns
from sklearn.datasets import make_blobs
from sklearn.cluster import KMeans
```

Creating Synthetic data

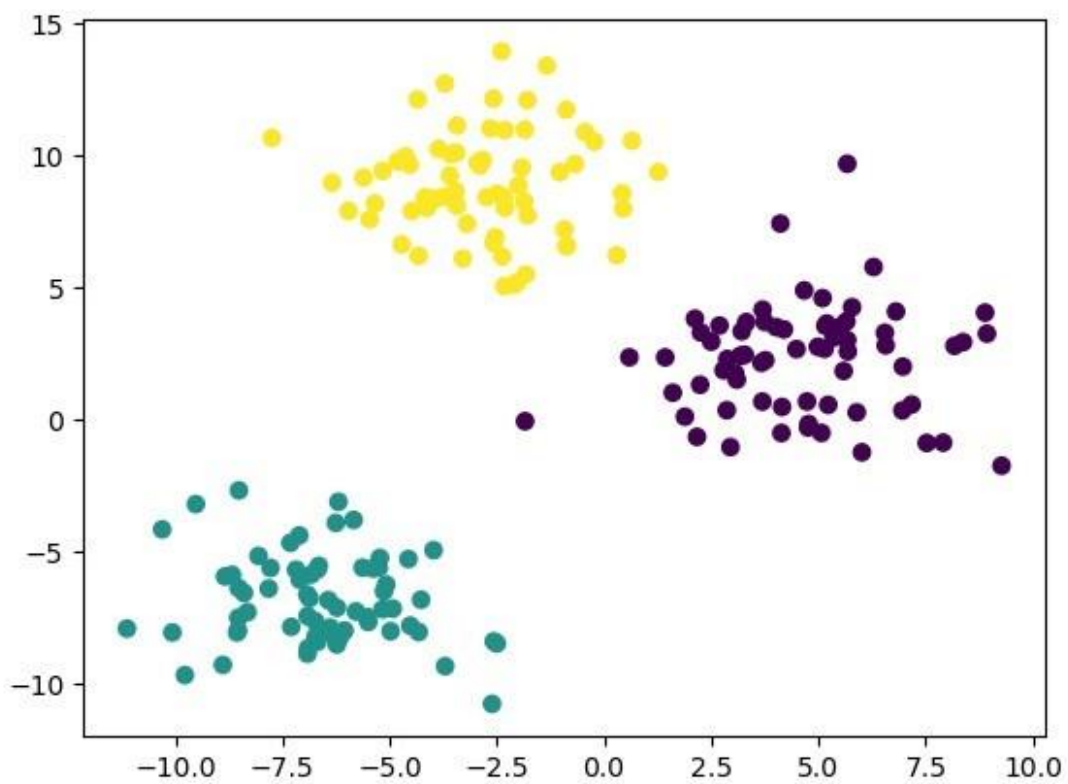
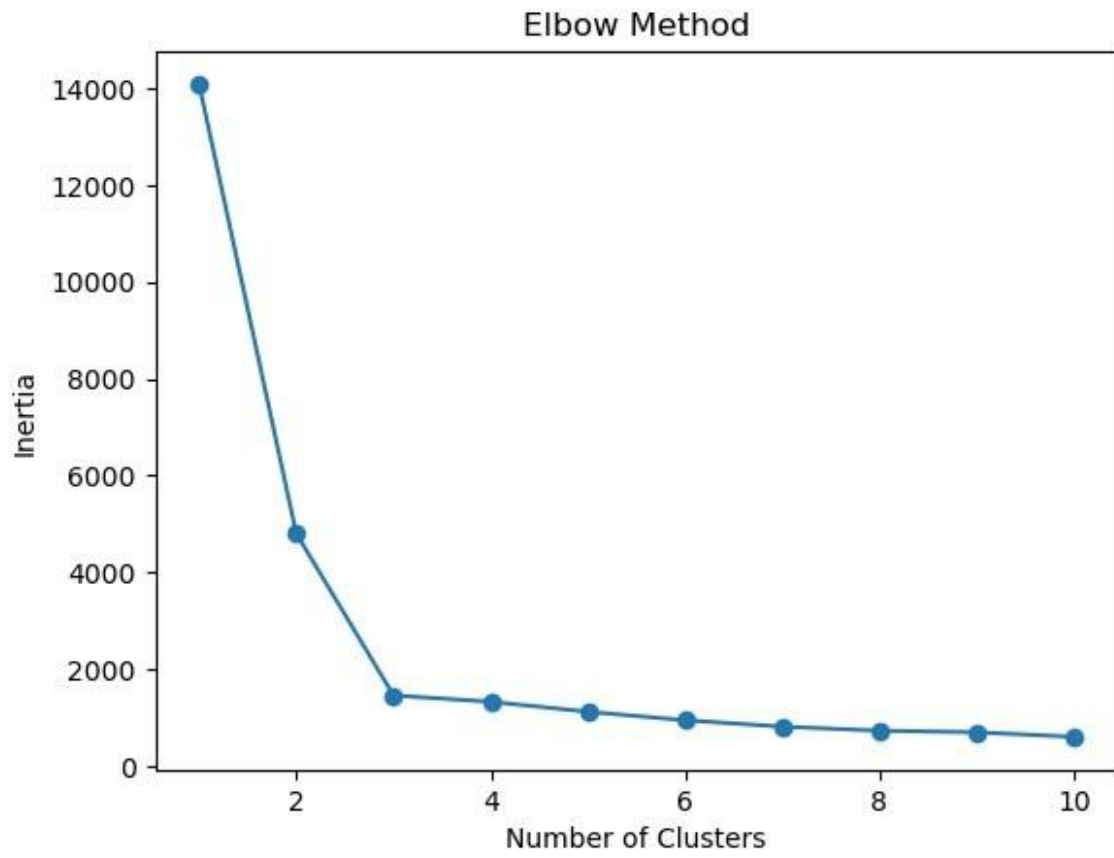
```
df=make_blobs(n_samples=200,n_features=2,centers=3,cluster_std=2,random_state=42)
data=list(zip(df[0][:,0],df[0][:,1]))
```

Finding the optimal no. of clusters for given data

```
inertias=[]
for i in range(1,11):
    kmeans=KMeans(n_clusters=i)
    kmeans.fit(data)
    inertias.append(kmeans.inertia_)
plt.plot(range(1,11),inertias,'-o')
plt.title('Elbow Method')
plt.xlabel('Number of Clusters')
plt.ylabel('Inertia')
plt.show()

kmeans=KMeans(n_clusters=3)
kmeans.fit(data)

plt.scatter(df[0][:,0],df[0][:,1],c=kmeans.labels_)
plt.show()
```



Part B: Implementation of K-Means Clustering Algorithm on a dataset

#Importing dataset

```
df=pd.read_csv('archive/Mall_Customers.csv') df.head()
```

	CustomerID	Gender	Age	Annual Income (k\$)	Spending Score (1-100)
0	1	Male	19	15	39
1	2	Male	21	15	81
2	3	Female	20	16	6
3	4	Female	23	16	77
4	5	Female	31	17	40

```
x = df.iloc[:, [3, 4]].values
```

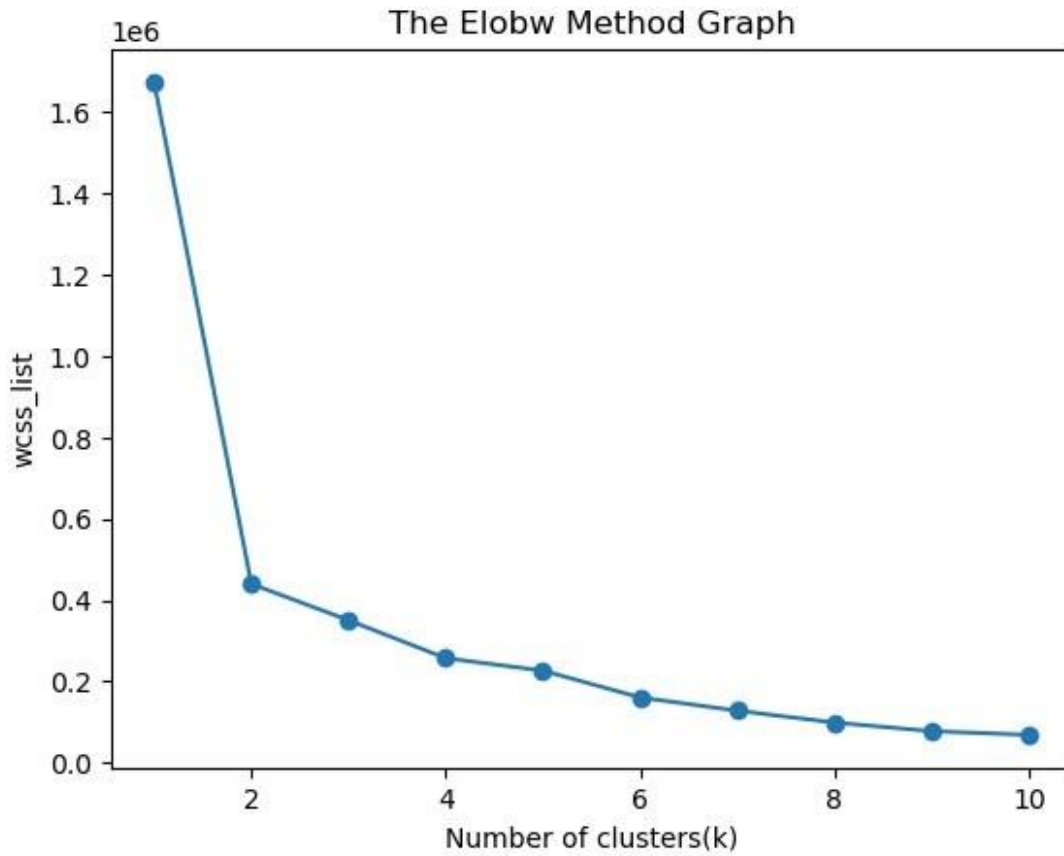
#finding optimal number of clusters using the elbow method `from sklearn.cluster import`

`KMeans`

```
wcss_list= [] #Initializing the list for the values of WCSS
```

```
#Using for loop for iterations from 1 to 10. for i in range(1, 11): kmeans = KMeans(n_clusters=i,  
init='k-means++', random_state= 42)
```

```
    kmeans.fit(x)  
    wcss_list.append(kmeans.inertia_) plt.plot(range(1,  
11), wcss_list, '-o') plt.title('The Elbow Method Graph')  
plt.xlabel('Number of clusters(k)') plt.ylabel('wcss_list')  
plt.show()
```

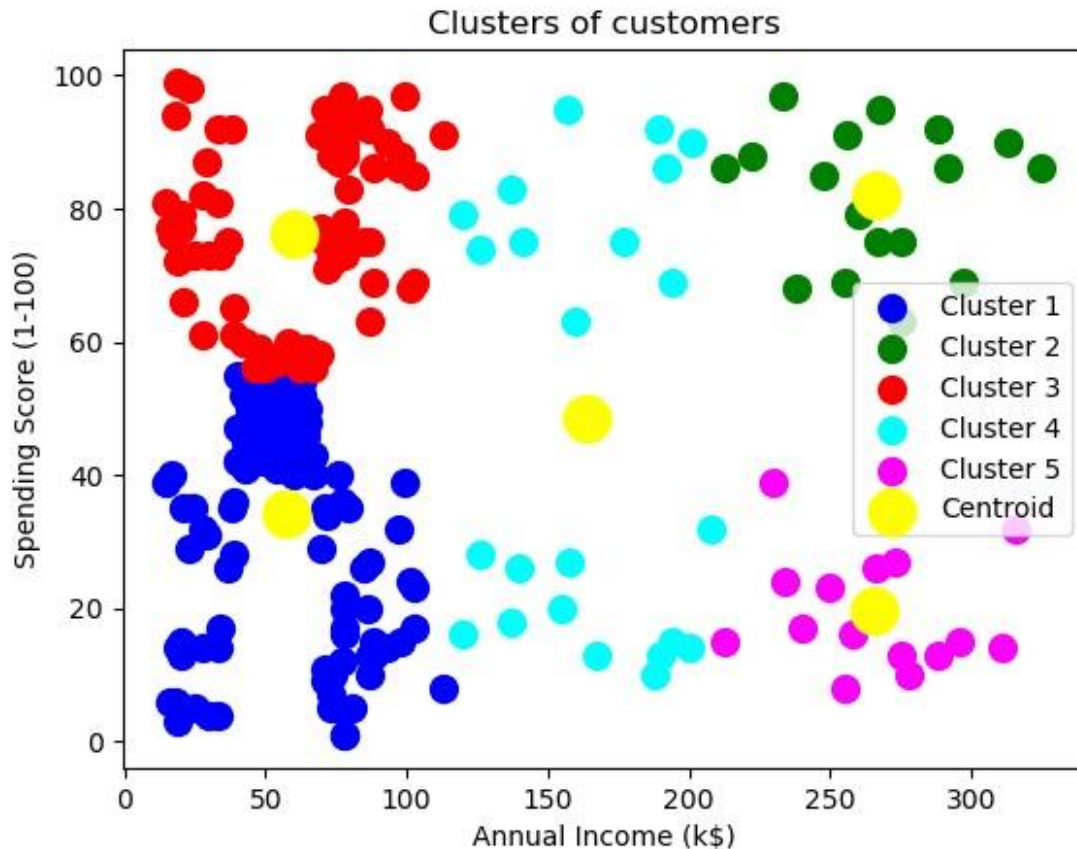


#training the K-means model on a dataset

```
kmeans = KMeans(n_clusters=5, init='k-means++', random_state=42) y_predict= kmeans.fit_predict(x)
```

#visualizing the clusters

```
plt.scatter(x[y_predict == 0, 0], x[y_predict == 0, 1], s = 100, c =
'blue', label = 'Cluster 1') #for first cluster
plt.scatter(x[y_predict == 1, 0], x[y_predict == 1, 1], s = 100, c =
'green', label = 'Cluster 2') #for second cluster
plt.scatter(x[y_predict == 2, 0], x[y_predict == 2, 1], s = 100, c =
'red', label = 'Cluster 3') #for third cluster
plt.scatter(x[y_predict == 3, 0], x[y_predict == 3, 1], s = 100, c =
'cyan', label = 'Cluster 4') #for fourth cluster
plt.scatter(x[y_predict == 4, 0], x[y_predict == 4, 1], s = 100, c =
'magenta', label = 'Cluster 5') #for fifth cluster
plt.scatter(kmeans.cluster_centers[:, 0], kmeans.cluster_centers[:,
1], s = 300, c = 'yellow', label = 'Centroid') plt.title('Clusters of
customers') plt.xlabel('Annual Income (k$)') plt.ylabel('Spending
Score (1-100)') plt.legend() plt.show()
```



Conclusion

Merits and Demerits of K-Means Clustering Algorithm

Merits of K-Means

1. **Simplicity and Ease of Implementation:**
 - K-Means is straightforward to understand and implement, making it a popular choice for clustering tasks.
2. **Efficiency:**
 - The algorithm is computationally efficient, especially with a small number of clusters, as its time complexity is generally $O(n * k * i)$, where n is the number of data points, k is the number of clusters, and i is the number of iterations.
3. **Scalability:**
 - K-Means can handle large datasets efficiently and can be easily scaled to a large number of samples.
4. **Works Well with Spherical Clusters:**
 - K-Means performs well when the clusters are spherical and equally sized, making it effective for certain types of data distributions.
5. **Easy to Interpret:**

- The results of K-Means are easy to interpret, as the algorithm assigns cluster labels that can be directly analyzed.
6. **Online Clustering:**
- Variants of K-Means allow for online clustering, which can update cluster centers incrementally as new data arrives.

Demerits of K-Means

1. **Choosing the Number of Clusters (k):**
 - The requirement to specify the number of clusters in advance can be challenging, and the wrong choice can lead to poor clustering results.
2. **Sensitivity to Initial Conditions:**
 - K-Means can converge to different solutions based on the initial placement of centroids, leading to potential inconsistencies in clustering results. This issue can be mitigated using techniques like K-Means++ for better initialization.
3. **Assumption of Spherical Clusters:**
 - The algorithm assumes that clusters are spherical and of similar size, which can lead to poor performance on datasets with non-spherical or varying-sized clusters.
4. **Outliers and Noisy Data:**
 - K-Means is sensitive to outliers and noise, which can distort the placement of centroids and negatively impact clustering performance.
5. **Limited to Numerical Data:**
 - K-Means works primarily with numerical data and cannot handle categorical features directly, requiring additional preprocessing steps.
6. **Local Optima:**
 - The algorithm can converge to a local minimum, meaning it might not find the best overall solution. Running K-Means multiple times with different initializations can help address this issue.

K-Means clustering is a widely used algorithm due to its simplicity and efficiency. However, users should be aware of its limitations, particularly regarding the choice of the number of clusters, sensitivity to initial conditions, and its assumptions about data distribution. Proper preprocessing and experimentation can help mitigate some of these issues, making K-Means a valuable tool for exploratory data analysis and clustering tasks.

Experiment No.8

Name: Atharv Vijay Deshpande

PRN: 21410044

Batch : EN-3

Title: Implementation of Anomaly Detection model to identify unusual pattern

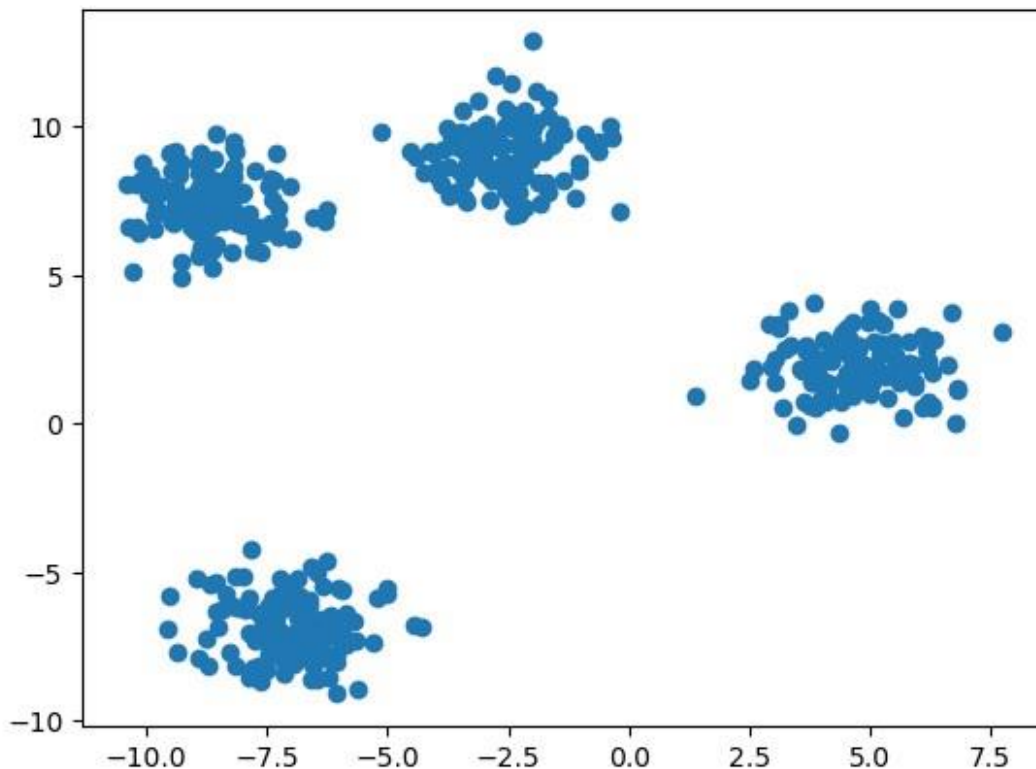
Import Libraries

```
import numpy as np
import matplotlib.pyplot as plt
from sklearn.cluster import DBSCAN
from sklearn.datasets import make_blobs
```

Make a Dataset for Clustering Experimentation

```
X,y=make_blobs(n_samples=500,centers=4,random_state=42,cluster_std=1)
plt.scatter(X[:,0],X[:,1])
```

<matplotlib.collections.PathCollection at 0x7ab975ad77d0>



Add Some Anomalies/Outliers in the Dataset

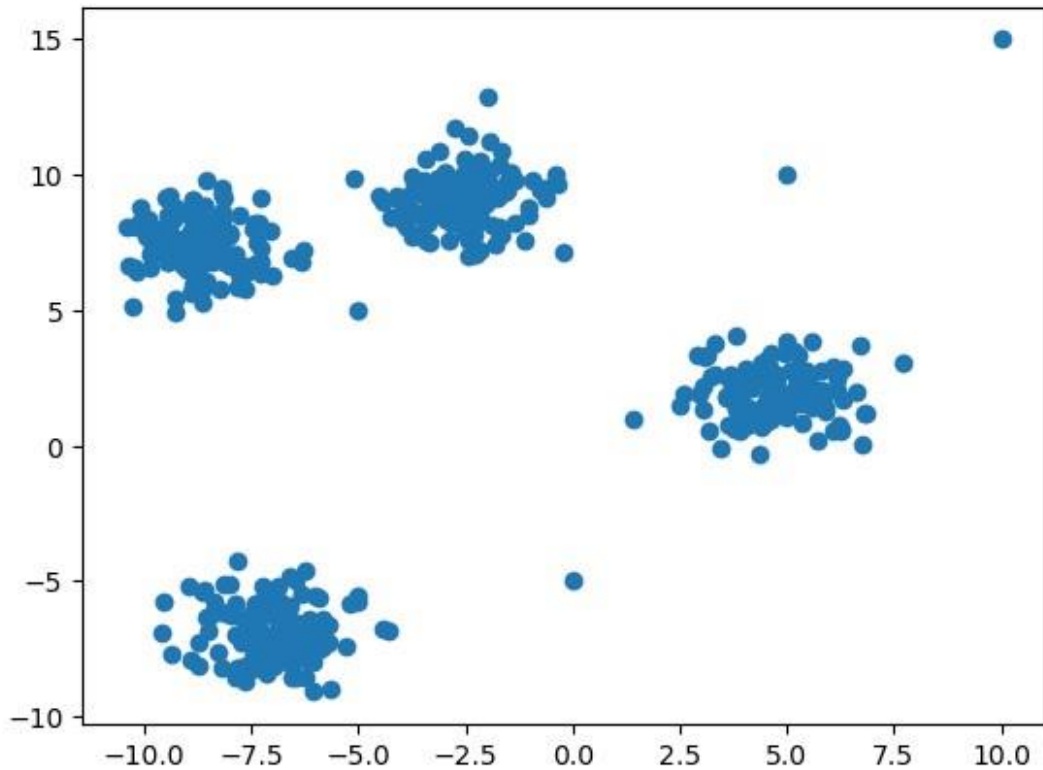
```
import matplotlib.pyplot as plt
import numpy as np
```

```

anomalies = np.array([[10, 15], [-5, 5], [0, -5], [5, 10]]) X = np.concatenate((X,
anomalies))
y = np.concatenate((y, np.array([4, 4, 4, 4]))) #Assign a new cluster label for outliers

plt.scatter(X[:, 0], X[:, 1]) plt.show()

```



Make a DBSCAN Clustering model and find the Clusters

```

dbsacn=DBSCAN(eps=1.5,min_samples=10)
labels=dbsacn.fit_predict(X) labels
array([ 0, 1, 2, 3, 1, 1, 0, 1, 2, 1, 2, 3, 2, 3, 1, 2,
3,
      0, 0, 3, 2, 3, 2, 0, 0, 1, 1, 0, 0, 3, 1, 3, 3,
3,
      1, 1, 2, 2, 0, 0, 1, 2, 3, 3, 3, 2, 2, 2, 1, 0,
1,
      3, 0, 1, 2, 3, 3, 0, 1, 0, 0, 3, 1, 0, 2, 1, 1,
0,
      2, 1, 2, 1, 1, 0, 3, 0, 3, 1, 2, 3, 1, 2, 1, 3,
0,

```

0, 0, 0, 2, 3, 0, 1, 2, 1, 2, 0, 3, 2, 3, 0, 2,

2,

2, 0, 0, 3, 3, 0, 3, 0, 1, 0, 0, 0, 0, 1, 2, 0,

1,

1, 3, 2, 1, 2, 0, 2, 2, 1, 1, 0, 0, 2, 2, 1, 2,

0,

0, 0, 2, 2, 1, 0, 2, 2, 1, 0, 0, 3, 3, 3, 1, 1,

2,

2, 3, 0, 3, 0, 1, 1, 0, 0, 2, 2, 1, 3, 2, 1, 0,

0,

1, 3, 3, 0, 0, 3, 3, 1, 1, 1, 3, 0, 3, 3, 0, 0,

3,

2, 3, 1, 1, 0, 0, 1, 3, 1, 3, 3, 0, 1, 3, 3, 1,

0,

2, 0, 1, 2, 2, 0, 0, 1, 0, 3, 3, 1, 3, 0, 3, 2,

2,

3, 0, 3, 2, 3, 3, 1, 1, 2, 1, 2, 3, 1, 0, 1, 3,

1,

2, 2, 2, 1, 3, 2, 0, 0, 2, 3, 3, 1, 3, 3, 3, 3,

2,

2, 0, 3, 1, 2, 3, 3, 0, 3, 1, 1, 2, 3, 2, 1, 3,

3,

2, 0, 3, 3, 3, 2, 0, 1, 3, 2, 2, 3, 1, 3, 2, 0,

2,

1, 1, 3, 3, 2, 0, 2, 0, 2, 1, 2, 1, 0, 2, 2, 3,

2,

1, 1, 2, 2, 1, 0, 2, 2, 2, 1, 0, 2, 2, 2, 0, 2,

2,

0, 0, 0, 1, 0, 0, 0, 0, 0, 3, 0, 2, 0, 3, 3, 0,

2,

2, 2, 3, 3, 1, 3, 1, 3, 2, 1, 1, 2, 1, 0, 0, 1,

0,

```

3, 1, 3, 2, 3, 1, 1, 0, 2, 1, 2, 2, 0, 3, 1, 2,
1,
0, 2, 0, 2, 0, 0, 1, 0, 3, 1, 0, 1, 3, 1, 0, 1,
0,
2, 3, 3, 1, 1, 2, 3, 3, 3, 2, 0, 1, 3, 2, 2, 1,
0,
1, 1, 1, 3, 2, 1, 3, 2, 3, 3, 2, 3, 2, 3, 0, 1,
2,
3, 1, 0, 3, 2, 2, 1, 1, 1, 0, 1, 2, 0, 3, 3, 0,
2,
0, 2, 3, 1, 1, 1, 3, 2, 2, 0, 3, 1, 1, 1, 2, 1,
3,
2, 0, 2, 1, 2, 3, 2, 3, 2, 1, 3, 0, 0, 3, 3, 0,
0,
3, 1, 2, 1, 2, 3, 0, 0, 3, 1, 0, 0, 2, 0, 0, 3,
0,
1, 3, 0, 3, 1, 1, 2, -1, -1, -1, -1])

```

Find the Anomalies

Note: The Anomalies are with label as -1

```

pred_anomalies=X[labels==-1]
pred_anomalies

array([[10., 15.],
       [-5.,  5.],
       [ 0., -5.],
       [ 5., 10.]])

```

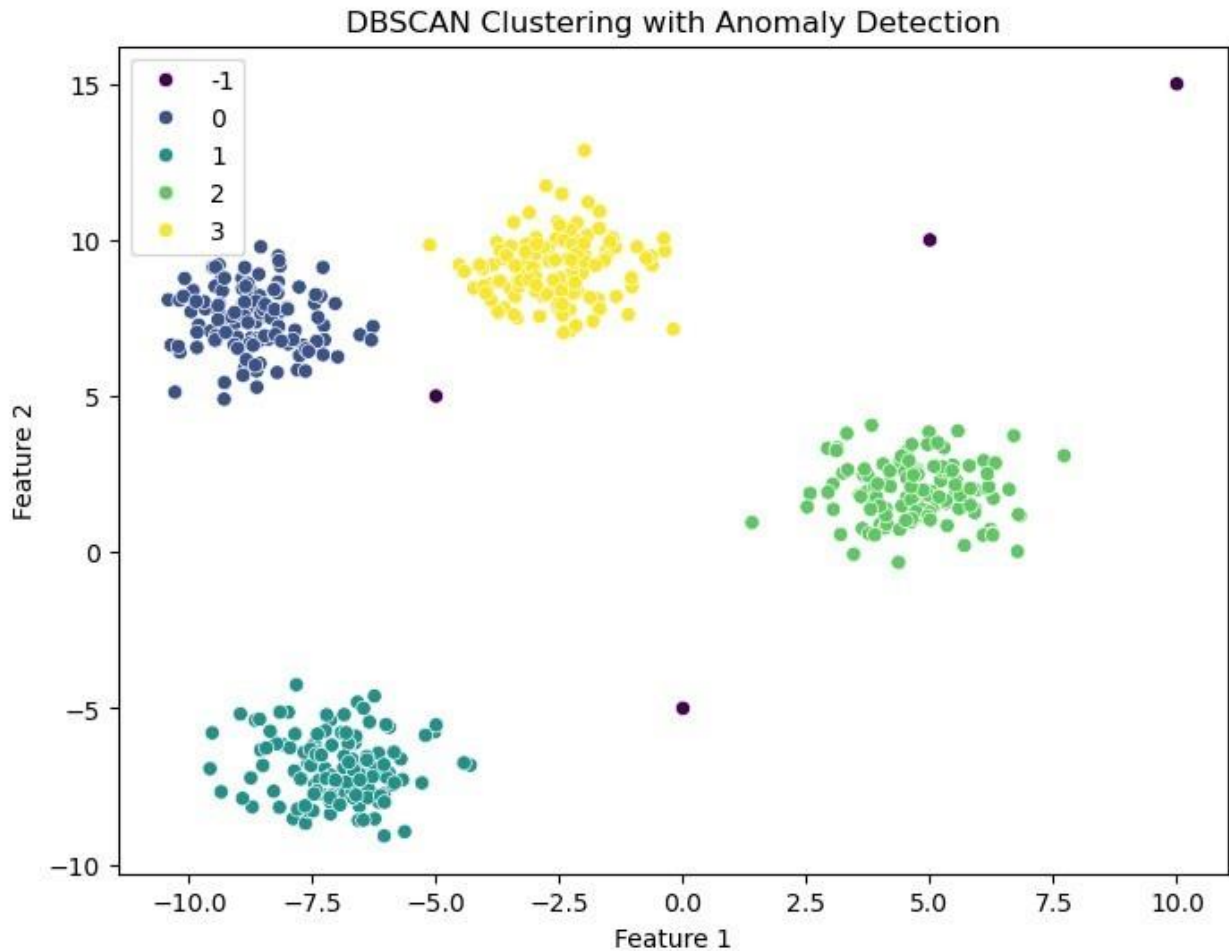
visualization of anomalies found by DBSCAN model

```

import matplotlib.pyplot as plt
import seaborn as sns

plt.figure(figsize=(8, 6))
sns.scatterplot(x=X[:, 0], y=X[:, 1], hue=labels, palette="viridis", legend="full")
plt.title("DBSCAN Clustering with Anomaly Detection")
plt.xlabel("Feature 1")
plt.ylabel("Feature 2")
plt.show()

```



Merits of Anomaly Detection Using DBSCAN

1. **Density-Based Approach:** DBSCAN (Density-Based Spatial Clustering of Applications with Noise) is highly effective for identifying anomalies as outliers in regions with low data density, making it well-suited for detecting unusual patterns in datasets with irregular distributions.
2. **No Need for Predefined Number of Clusters:** Unlike methods like K-means, DBSCAN does not require the number of clusters to be specified beforehand, allowing for flexible anomaly detection without prior knowledge of the dataset.
3. **Detection of Arbitrarily Shaped Clusters:** DBSCAN can detect clusters of arbitrary shape, which is an advantage in datasets where anomalies occur in irregular regions that are not easily captured by traditional clustering algorithms.
4. **Robust to Noise:** The algorithm naturally classifies points in sparse regions as noise, effectively identifying outliers (anomalies) without requiring separate training or labeling of outliers.

Demerits of Anomaly Detection Using DBSCAN

1. **Sensitive to Hyperparameters:** DBSCAN relies on two important hyperparameters: `eps` (the maximum distance between two samples for them to be considered part of the same cluster) and `min_samples` (the minimum number of points required to form a dense region). Choosing these values incorrectly can lead to poor detection of anomalies or false positives.
2. **Scalability Issues with High Dimensional Data:** DBSCAN struggles with highdimensional data as the density concept breaks down in high dimensions (the curse of dimensionality). This limits its performance on complex, large-scale datasets.
3. **Difficulty Handling Varying Densities:** DBSCAN performs less effectively when the dataset contains clusters of varying densities. It may fail to separate dense clusters from noise or anomalies in such cases.
4. **Non-Deterministic in Borderline Cases:** In cases where points lie on the border between dense and sparse regions, DBSCAN's classification may be nondeterministic, leading to instability in detecting certain anomalies.

Conclusion

The implementation of an anomaly detection model using DBSCAN proved effective for identifying unusual patterns in low-dimensional datasets where density-based separation of normal points and outliers is feasible. DBSCAN's ability to detect arbitrary-shaped clusters without specifying the number of clusters provides flexibility and robustness in certain anomaly detection tasks. However, the sensitivity of DBSCAN to the choice of hyperparameters and its challenges in handling high-dimensional data or clusters of varying densities limit its effectiveness in more complex scenarios.

In conclusion, DBSCAN is a powerful tool for anomaly detection in specific contexts where data is well-structured with distinct densities, but its limitations must be carefully considered in large-scale or high-dimensional applications.

Experiment No.9

Name: Atharv Vijay Deshpande

PRN: 21410044

Batch : EN-3

Title: Implementation of Single Layer and Multilayer perceptron

Single Layer Perceptron (SLP)

A **Single Layer Perceptron (SLP)** is the simplest form of a neural network. It consists of an input layer that directly connects to an output layer, with no hidden layers in between. The SLP is mainly used as a linear classifier, which means it can only solve problems that are linearly separable.

Key Components of a Single Layer Perceptron

1. **Inputs (Features):**
 - The input data that is used for prediction, represented as a vector. For instance, in a binary classification problem, the input might be $[x_1, x_2, \dots, x_n]$.
2. **Weights:**
 - Each input is associated with a weight, which signifies how important that feature is for determining the output. The weights are adjusted during the training process.
3. **Bias:**
 - An additional term added to the weighted sum of inputs to shift the decision boundary. It allows the perceptron to classify inputs that are not necessarily centered around the origin.
4. **Weighted Sum:**
 - The perceptron computes the weighted sum of inputs:
$$z = w_1x_1 + w_2x_2 + \dots + w_nx_n + b$$
where:
 - (w) are the weights, •
 - (x) are the inputs, •
 - (b) is the bias.
5. **Activation Function:**
 - An activation function decides whether a neuron should be activated or not. For a single layer perceptron, a **step function** or **sign function** is commonly used:

$$f(z) = \begin{cases} 1 & \text{if } z \geq 0 \\ 0 & \text{if } z < 0 \end{cases}$$

This function converts the weighted sum into the output (either 0 or 1 in the case of binary classification).

6. Output:

- The perceptron produces a binary output, which can be used for classification tasks (e.g., distinguishing between two classes).

Training a Single Layer Perceptron

Training involves adjusting the weights and bias to minimize the classification error. This is typically done using the **Perceptron Learning Rule**, which updates the weights based on the error between predicted and actual outputs:

$$w_i \leftarrow w_i + \alpha (y - \hat{y}) x_i$$

where:

- w_i : weight for input x_i
- α : learning rate
- y : actual output
- \hat{y} : predicted output

Limitations of Single Layer Perceptron

- **Linearly Separable Data:** A single layer perceptron can only classify data that is linearly separable. This means the classes must be separable by a straight line (in 2D) or a hyperplane (in higher dimensions).
- **No Hidden Layers:** Since there are no hidden layers, the model lacks the complexity to learn more intricate patterns in the data.

Summary

A Single Layer Perceptron is a simple neural network that classifies inputs based on a weighted sum of features. It is an effective model for linearly separable problems but cannot handle more complex, non-linear data. To solve non-linear problems, more advanced networks like **MultiLayer Perceptrons (MLP)** with hidden layers are needed.

```
import numpy as np

#Activation function def sigmoid(x):
return 1 / (1 + np.exp(-x))
```



```
#Another Activation function def step(x):  
return np.where(x>=0, 1, 0)
```

```
#Single Perceptron with aggregation and activation def perceptron(x,w,b):  
return step(np.dot(x,w)+b)  
  
x=np.array([[0,0],[0,1],[1,0],[1,1]])  
y=np.array([0,0,0,1]) w=np.array([1,1]) b=-1.5  
  
perceptron(x,w,b)#and gate  
array([0, 0, 0, 1])  
def perceptron_learn(X, y, lr, epochs): w =  
np.zeros(X.shape[1]) b = 0  
  
for epoch in range(epochs): for i in range(X.shape[0]):  
y_hat = step(np.dot(X[i], w) + b) w = w + lr * (y[i] -  
y_hat) * X[i] b = b + lr * (y[i] - y_hat) return w, b  
  
perceptron_learn(x,y,1,100)
```

Multi-Layer Perceptron (MLP)

A **Multi-Layer Perceptron (MLP)** is a type of artificial neural network that consists of multiple layers: an input layer, one or more hidden layers, and an output layer. Unlike a Single Layer Perceptron (SLP), MLPs can solve non-linearly separable problems due to the presence of hidden layers and non-linear activation functions.

Key Components of a Multi-Layer Perceptron

1. **Input Layer:**
 - This layer receives the input features from the dataset. Each feature is represented as an input node in this layer.
 - Example: For a dataset with 3 features, the input layer would have 3 nodes.
2. **Hidden Layers:**
 - MLPs have one or more hidden layers between the input and output layers. Each hidden layer contains neurons that process inputs using weights, biases, and activation functions.

- The number of neurons and layers can be adjusted to increase the model's capacity to learn complex patterns.
- **Activation Functions:** Non-linear functions like ReLU (Rectified Linear Unit), Sigmoid, or Tanh are applied to the neurons in hidden layers.

$$z^{(l)} = w^{(l)} x^{(l)} + b^{(l)}$$

where:

- $w^{(l)}$ are the weights for layer (l)
- $x^{(l)}$ are the inputs to layer (l)
- $b^{(l)}$ is the bias term
- The activation function (f) is applied as:

$$a^{(l)} = f(z^{(l)})$$

3. Output Layer:

- The output layer produces the final predictions. The number of neurons in the output layer depends on the task:
 - For **binary classification**, there is typically 1 output neuron with a Sigmoid activation function.
 - For **multi-class classification**, the output layer uses the Softmax function with one neuron for each class.
- **Sigmoid Activation** for binary classification:

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

- **Softmax Activation** for multi-class classification:

$$\text{Softmax}(z_j) = \frac{e^{z_j}}{\sum_{k=1}^n e^{z_k}}$$

How MLP Works

1. Forward Propagation:

- Input data passes through the network layer by layer. Each neuron computes the weighted sum of its inputs, adds a bias, and applies a non-linear activation function to produce an output. This process is repeated for each layer until the final output is produced in the output layer.

2. Loss Function:

- The model uses a loss function to measure how far the predicted outputs are from the actual targets. Common loss functions are:
 - **Binary Cross-Entropy** for binary classification.
 - **Categorical Cross-Entropy** for multi-class classification.

$$L(y, \hat{y}) = -\frac{1}{N} \sum_{i=1}^N (y_i \log(\hat{y}_i) + (1 - y_i) \log(1 - \hat{y}_i))$$

where (y_i) is the actual label, and (\hat{y}_i) is the predicted probability.

3. Backpropagation:

- After computing the loss, the MLP uses backpropagation to adjust the weights and biases. The gradients of the loss function with respect to each weight and bias are computed using the chain rule, and the weights are updated using gradient descent:

$$w^{(l)} \leftarrow w^{(l)} - \alpha \frac{\partial L}{\partial w^{(l)}}$$

where alpha is the learning rate.

4. Optimization:

- An optimization algorithm like **Stochastic Gradient Descent (SGD)** or **Adam** is used to minimize the loss function by adjusting the weights iteratively.

Advantages of MLP

- **Non-Linear Problems:** MLPs can learn complex, non-linear relationships between input features and outputs.
- **Multiple Layers:** The presence of hidden layers allows the network to learn hierarchical patterns in data.
- **Universal Approximation:** MLPs with enough neurons and layers can approximate any continuous function.

Limitations of MLP

- **Computational Cost:** Training deep networks with many layers and neurons can be computationally expensive.
- **Risk of Overfitting:** MLPs with too many parameters can overfit to the training data, especially if the dataset is small.
- **Hyperparameter Tuning:** Finding the optimal architecture (number of layers, neurons, learning rate) can be challenging and requires careful tuning.

Summary

A Multi-Layer Perceptron is a powerful neural network that can model complex relationships between inputs and outputs. It works through forward propagation and backpropagation, using non-linear activation functions and optimization algorithms to minimize the loss and improve predictions. While MLPs are versatile, they require careful tuning and computational resources to be effective.

```
import numpy as np
from sklearn.neural_network import MLPClassifier from sklearn.metrics
import accuracy_score
```

```
# XOR inputs and outputs
```

```
X = np.array([[0, 0], [0, 1], [1, 0], [1, 1]]) y = np.array([0, 1, 1, 0])
```

```
# Initialize the MLPClassifier with a hidden layer of 2 neurons mlp =
MLPClassifier(hidden_layer_sizes=(4,4), activation='relu', solver='adam', max_iter=5000,
random_state=42)
```

```
# Train the model on all the data mlp.fit(X, y)
```

```
# Predict on the training set y_pred =
mlp.predict(X)
```

```
# Print the results print("Predictions:", y_pred)
print("Accuracy:", accuracy_score(y, y_pred))
Predictions: [0 1 1 0]
Accuracy: 1.0
```

MLP Implementation on a Dataset

```
# Import necessary libraries
from sklearn.neural_network import MLPClassifier from sklearn.datasets
import load_breast_cancer from sklearn.model_selection import
train_test_split from sklearn.preprocessing import StandardScaler
from sklearn.metrics import accuracy_score, classification_report,
confusion_matrix, ConfusionMatrixDisplay

# Load the Breast Cancer dataset cancer_data =
load_breast_cancer()
X, y = cancer_data.data, cancer_data.target

# Split the dataset into training and testing sets
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=42)

# Standardize features by removing the mean and scaling to unit
variance
scaler = StandardScaler()
X_train = scaler.fit_transform(X_train)
X_test = scaler.transform(X_test)

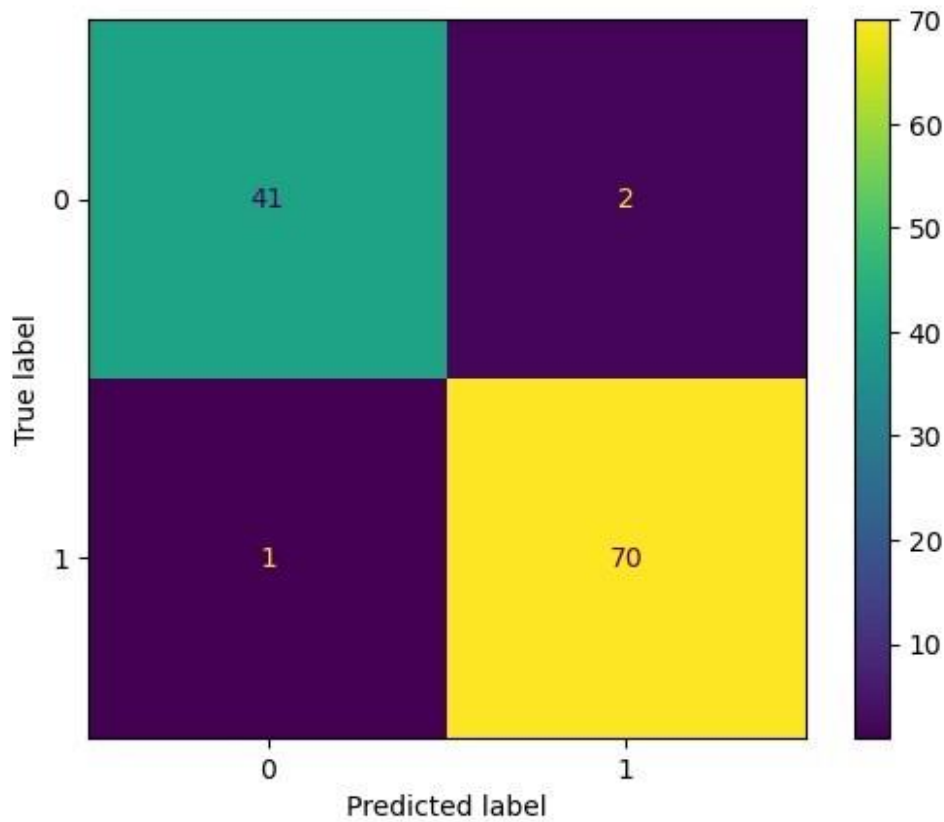
# Create an MLPClassifier model
mlp = MLPClassifier(hidden_layer_sizes=(64, 32), max_iter=1000, random_state=42)

# Train the model on the training data
mlp.fit(X_train, y_train)
```

```
# Make predictions on the test data y_pred =
mlp.predict(X_test)

# Calculate the accuracy of the model accuracy =
accuracy_score(y_test, y_pred) print(f'Accuracy: {accuracy:.2f} ")
cm=confusion_matrix(y_test, y_pred)
ConfusionMatrixDisplay(cm).plot()
print(classification_report(y_test, y_pred))
Accuracy: 0.97
```

	precision	recall	f1-score	support
0	0.98	0.95	0.96	43
1	0.97	0.99	0.98	71
accuracy			0.97	114
macro avg	0.97	0.97	0.97	
weighted avg	0.97	0.97	0.97	114



Conclusion:

In this experiment, we explored the implementation of both Single Layer Perceptron (SLP) and Multilayer Perceptron (MLP) to solve a classification task, specifically focusing on the XOR problem,

which is a classic non-linearly separable problem. The experiment demonstrated the following key conclusions:

1. **Single Layer Perceptron Limitations:** The Single Layer Perceptron, despite being a powerful linear classifier, was unable to solve the XOR problem. This failure highlights the limitation of the SLP in solving problems that require non-linear decision boundaries, as XOR cannot be separated by a single linear hyperplane.
2. **Multilayer Perceptron (MLP) Capability:** The Multilayer Perceptron, equipped with hidden layers and non-linear activation functions like ReLU, successfully solved the XOR problem. This showcases the strength of MLPs in learning non-linear relationships between inputs and outputs by leveraging the depth of the network.
3. **Effect of Hidden Layers:** The presence of hidden layers in MLP allows for hierarchical learning, where the network can learn more complex patterns and transformations. A single hidden layer with a small number of neurons (in our case, two neurons) was sufficient to model the XOR problem.
4. **Training and Convergence:** Increasing the number of training iterations and using an optimizer like Adam ensured that the MLP converged to a solution with high accuracy. The experiment demonstrated that while SLPs are simple and fast, MLPs require more computational resources and training time, but provide significantly better performance for non-linear tasks.
5. **Practical Implications:** This experiment underlines the importance of selecting appropriate model architectures based on the nature of the problem. For linearly separable tasks, SLP may be sufficient, but for more complex problems involving non-linear boundaries, MLPs with hidden layers are essential.

In conclusion, while SLPs are limited to solving linearly separable problems, MLPs with hidden layers significantly expand the range of solvable problems by learning non-linear mappings, making them more suitable for real-world tasks like XOR and beyond.

Experiment No.10

Name: Atharv Vijay Deshpande

PRN: 21410044

Batch : EN-3

Title: Use Convolutional Neural Network (CNN) on suitable dataset

```
import tensorflow as tf
from tensorflow.keras import datasets, layers, models
import matplotlib.pyplot
as plt
import numpy as np
```

1. Loading CIFAR-10 Dataset

```
#Load the dataset
(X_train, y_train), (X_test, y_test) = datasets.cifar10.load_data()
X_train.shape
(50000, 32, 32, 3)
```

Explanation: The CIFAR-10 dataset is loaded. It contains 60,000 32x32 color images in 10 classes, with 50,000 training images and 10,000 test images.


```

X_test.shape
(10000, 32, 32, 3)
y_train.shape
(50000, 1)
y_train[:5]
array([[6],
       [9],
       [9],
       [4],
       [1]], dtype=uint8)
y_train = y_train.reshape(-1,) y_train[:5]
array([6, 9, 9, 4, 1], dtype=uint8)
y_test = y_test.reshape(-1,)

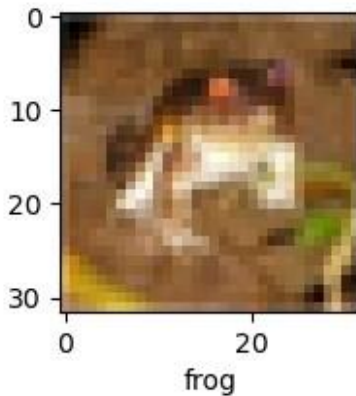
classes =
["airplane", "automobile", "bird", "cat", "deer", "dog", "frog", "horse", "ship", "truck"]

```

```

def plot_sample(X, y, index): plt.figure(figsize = (15,2))
plt.imshow(X[index]) plt.xlabel(classes[y[index]])
plot_sample(X_train, y_train, 0)

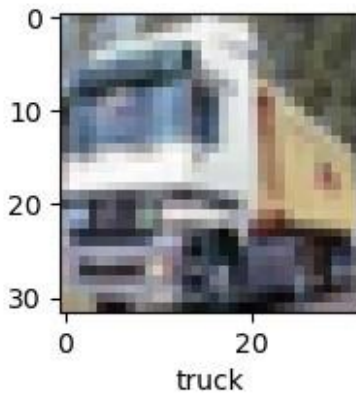
```



```

plot_sample(X_train, y_train, 1)

```



2. Preprocessing Data

#Normalizing the training data

`X_train = X_train / 255.0`

`X_test = X_test / 255.0`

Explanation: The pixel values of images are normalized by dividing by 255 to scale them between 0 and 1.

```
#Build simple artificial neural network for image classification ann = models.Sequential([
    layers.Flatten(input_shape=(32,32,3)),    layers.Dense(3000, activation='relu'),
    layers.Dense(1000, activation='relu'),    layers.Dense(10, activation='softmax')
])
```

```
ann.compile(optimizer='SGD',
            loss='sparse_categorical_crossentropy',    metrics=['accuracy'])
```

```
ann.fit(X_train, y_train, epochs=5)
```

/home/aspatil/anaconda3/lib/python3.12/site-packages/keras/src/layers/resizing/flatten.py:37:

UserWarning: Do not pass an

`input_shape`/`input_dim` argument to a layer. When using Sequential models, prefer using an `Input(shape)` object as the first layer in the model instead.

super().__init__(**kwargs) 2024-10-14

23:03:53.105745: I

external/local_xla/xla/stream_executor/cuda/cuda_executor.cc:998] successful NUMA node read from SysFS had negative value (-1), but there must be at least one NUMA node, so returning NUMA node zero. See more at

<https://github.com/torvalds/linux/blob/v6.0/Documentation/ABI/testing/>

sysfs-bus-pci#L344-L355

2024-10-14 23:03:53.106673: W

tensorflow/core/common_runtime/gpu/gpu_device.cc:2251] Cannot dlopen some GPU libraries. Please make sure the missing libraries mentioned above are installed properly if you would like to use GPU.

Follow the guide at <https://www.tensorflow.org/install/gpu> for how to download and setup the required libraries for your platform. Skipping registering GPU devices...

Epoch 1/5

2024-10-14 23:03:53.804697: W

external/local_tsl/tsl/framework/cpu_allocator_impl.cc:83] Allocation of 614400000 exceeds 10% of free system memory.

1563/1563 ————— 73s 46ms/step - accuracy: 0.3042 -
loss: 1.9276

Epoch 2/5

1563/1563 ————— 72s 46ms/step - accuracy: 0.4162 -
loss: 1.6495

Epoch 3/5

1563/1563 ————— 68s 44ms/step - accuracy: 0.4509 -
loss: 1.5561

Epoch 4/5

1563/1563 ————— 68s 43ms/step - accuracy: 0.4777 -
loss: 1.4892

Epoch 5/5

1563/1563 ————— 68s 43ms/step - accuracy: 0.4953 -
loss: 1.4309

<keras.src.callbacks.history.History at 0x72a657877bc0>

You can see that at the end of 5 epochs, accuracy is at around 49%

```
from sklearn.metrics import confusion_matrix, classification_report
import numpy as np
y_pred = ann.predict(X_test)
y_pred_classes = [np.argmax(element) for element in y_pred]
```

```
print("Classification Report: \n", classification_report(y_test, y_pred_classes))
```

313/313

3s 9ms/step

Classification Report

	precision	recall	f1-score	support				
0	0.55	0.57	0.56	1000				
1	0.62	0.52	0.57	1000				
2	0.42	0.24	0.31	1000				
3	0.44	0.20	0.27	1000				
4	0.45	0.26	0.33	1000				
5	0.51	0.23	0.32	1000				
6	0.38	0.74	0.50	1000				
7	0.51	0.53	0.52	1000				
8	0.66	0.58	0.62	1000				
9	0.35	0.78	0.48	1000				
accuracy			0.47	10000	macro avg	0.49	0.47	0.45
10000 weighted avg		0.49	0.47	0.45	10000			

3. Building the CNN Model

```
#Now let us build a convolutional neural network to train our images
cnn = models.Sequential([
    layers.Conv2D(filters=32, kernel_size=(3, 3), activation='relu', input_shape=(32, 32, 3)),
    layers.MaxPooling2D((2, 2)),

    layers.Conv2D(filters=64, kernel_size=(3, 3), activation='relu'),
    layers.MaxPooling2D((2, 2)),

    layers.Flatten(),
    layers.Dense(64, activation='relu'),
    layers.Dense(10, activation='softmax')
])
```

/home/aspatil/anaconda3/lib/python3.12/site-packages/keras/src/layers/convolutional/base_conv.py:107: UserWarning: Do not pass an `input_shape`/`input_dim` argument to a layer. When using Sequential models, prefer using an `Input(shape)` object as the first layer in the model instead.

```
super().__init__(activity_regularizer=activity_regularizer, **kwargs)
```

Explanation: A sequential CNN model is built:

The first convolutional layer has 32 filters with a 3x3 kernel and ReLU activation, followed by a 2x2 max pooling layer.

The second convolutional layer has 64 filters, followed by another 2x2 max pooling. The model is flattened and then connected to a dense layer with 64 units and ReLU activation. The output layer uses softmax activation for multi-class classification (10 classes).

4. Compiling the Model

```
cnn.compile(optimizer='adam',  
            loss='sparse_categorical_crossentropy', metrics=['accuracy'])
```

Explanation: The model is compiled using the Adam optimizer, sparse categorical crossentropy as the loss function (since the labels are integers), and accuracy as a performance metric.

5. Training the Model `cnn.fit(X_train, y_train,`

```
epochs=10)
```

Epoch 1/10

2024-10-14 23:12:34.726665: W

external/local_tsl/tsl/framework/cpu_allocator_impl.cc:83] Allocation of 614400000 exceeds 10% of free system memory.

1563/1563 ————— 24s 15ms/step - accuracy: 0.3864 -

loss: 1.6857 Epoch 2/10

1563/1563 ————— 24s 15ms/step - accuracy: 0.5996 -

loss: 1.1473 Epoch 3/10

1563/1563 ————— 24s 16ms/step - accuracy: 0.6542 -

loss: 0.9885 Epoch 4/10

1563/1563 ————— 25s 16ms/step - accuracy: 0.6896 -

loss: 0.8905 Epoch 5/10

1563/1563 —————

—————

————— 26s

17ms/step - accuracy:

0.7229 -

loss: 0.8042 Epoch 6/10

1563/1563 ————— 26s 16ms/step - accuracy: 0.7376 -

loss: 0.7552 Epoch 7/10

1563/1563 ————— 25s 16ms/step - accuracy: 0.7562 -

loss: 0.7063 Epoch 8/10

1563/1563 ————— 25s 16ms/step - accuracy: 0.7679 -

loss: 0.6659 Epoch 9/10

1563/1563 ————— 25s 16ms/step - accuracy: 0.7817 -

loss: 0.6268 Epoch 10/10

```
1563/1563 ————— 25s 16ms/step - accuracy: 0.7972 -  
loss: 0.5844 <keras.src.callbacks.history.History at 0x72a5f40ef560>
```

Explanation: The model is trained for 10 epochs using the training data, and the test data is used for validation after each epoch.

With CNN, at the end 5 epochs, accuracy was at around 70% which is a significant improvement over ANN. CNN's are best for image classification and gives superb accuracy. Also computation is much less compared to simple ANN as maxpooling reduces the image dimensions while still preserving the features

6. Evaluating the Model

```
cnn.evaluate(X_test,y_test)  
313/313 ————— 2s 5ms/step - accuracy: 0.7031 - loss:  
0.9155  
[0.9262094497680664, 0.7027000188827515]
```

Explanation: The model is evaluated on the test set, and the test accuracy is printed.

7. Making Predictions

```
y_pred = cnn.predict(X_test) y_pred[:5]  
313/313 —————  
————— 2s 5ms/step  
array([[3.55052273e-03, 1.55897200e-04, 1.39184936e-03,  
9.35529172e-  
01,  
1.67694248e-04, 2.22749654e-02, 3.54748257e-02,  
1.96549998e-  
04,  
1.14560092e-03, 1.13041526e-  
04],
```

```

[2.01558592e-04, 4.26799478e-03, 4.41421116e-06, 7.37760502e-
07,
1.49627688e-08, 2.92013453e-08, 6.21419858e-08, 4.23154418e-
08,
9.95396197e-01, 1.28925662e-04],
[6.84590042e-02, 5.66818342e-02, 2.95170187e-03, 1.44834840e-
03,
1.16028148e-03, 5.02392824e-04, 4.70734631e-05, 1.27827586e-
03,
8.63696575e-01, 3.77456215e-03],
[9.77042556e-01, 5.57158282e-03, 1.75640616e-03, 2.78189266e-
03,
1.18206302e-03, 1.82632648e-05, 1.14431069e-03, 3.25264082e-
05,
1.01074576e-02, 3.62811814e-04],
[1.11192406e-07, 1.79055169e-05, 1.47151144e-03, 1.60174351e-
03,
3.74911875e-01, 6.63788989e-04, 6.21324658e-01, 5.06102026e-
07,
4.57737906e-06, 3.22937558e-06]], dtype=float32)

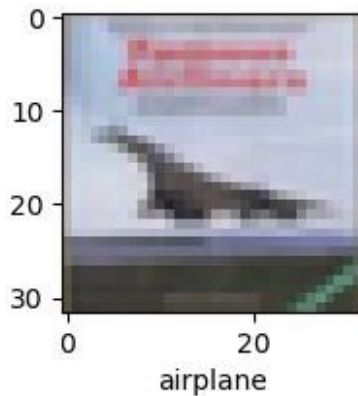
```

Explanation: Predictions are made on the test images. The highest probability class is selected as the predicted class.

```

y_classes = [np.argmax(element) for element in y_pred] y_classes[:5]
[3, 8, 8, 0, 6]
y_test[:5]
array([3, 8, 8, 0, 6], dtype=uint8)
plot_sample(X_test, y_test,3)

```



```
classes[y_classes[3]]
```

```
'airplane'
```

```
classes[y_classes[3]]
```

```
'airplane'
```

Conclusion

In this experiment, we used a **Convolutional Neural Network (CNN)** to classify images from the CIFAR-10 dataset, which contains 60,000 images categorized into 10 classes. CNNs are especially effective for image classification tasks because they automatically detect important features like edges, textures, and shapes through convolutional layers. We compared the CNN model's performance with a basic Artificial Neural Network (ANN) to highlight the advantages of CNNs in such tasks.

1. Model Architecture:

- CNNs consist of convolutional layers that can detect spatial hierarchies in images by learning local patterns. Max-pooling layers reduce the spatial dimensions, making the model computationally efficient.
- In contrast, ANNs process each pixel individually and do not consider the spatial relationships between pixels, making them less effective for image data.

2. Performance:

- The CNN model achieved a significantly higher accuracy on the CIFAR-10 dataset compared to an ANN.
- The local feature extraction capability of CNNs makes them far superior for image classification tasks, as they recognize shapes and patterns more effectively than ANNs, which tend to overfit or struggle to generalize for high-dimensional data like images.

3. Training Time and Complexity:

- While CNNs are more complex and take longer to train than ANNs due to additional operations like convolution and pooling, the increased accuracy and ability to handle image data justify this added complexity.
- ANNs, being simpler, train faster but are not well-suited for image data and struggle to capture the hierarchical patterns necessary for tasks like image classification.

Final Comparison:

- **CNN:** Achieved higher accuracy, performed better in handling image data, and showed better generalization for unseen test data.
- **ANN:** Faster training but underperformed on the same dataset, primarily because it doesn't capture spatial relationships in images effectively.

Thus, CNNs are a more suitable and effective choice for image-based tasks compared to traditional ANNs, particularly when dealing with datasets like CIFAR-10 that require spatial pattern recognition.

