```python
# -*- coding: utf-8 -*-
"""ML LAB(BE).ipynb

Automatically generated by Colab.

Original file is located at
    https://colab.research.google.com/drive/1A2vJst99-bU_hKjrI2VdGDEoUaQiUK3p
"""

# Commented out IPython magic to ensure Python compatibility.
# %%html
# <marquee style='width: 55%; color: blue;'><b> Machine Learning Lab by Dr.T.Bhaskar
#  !</b></marquee>

"""# **ASSIGNMENT-1:SIMPLE LINEAR REGRESSION**

"""

#STEP-1: Import Libraries
# Code to read csv file into colaboratory:
!pip install -U -q PyDrive
from pydrive.auth import GoogleAuth
from pydrive.drive import GoogleDrive
from google.colab import auth
from oauth2client.client import GoogleCredentials

#STEP-2: Autheticate E-Mail ID
auth.authenticate_user()
gauth = GoogleAuth()
gauth.credentials = GoogleCredentials.get_application_default()
drive = GoogleDrive(gauth)

## IMPORT CSV FILES FROM DRIVE INTO GOOGLE-COLAB
downloaded = drive.CreateFile({'id':'1ORTNkZgo6uc5VSYopHxl787fPdrs7-LT'}) # replace
the id with id of file you want to access
downloaded.GetContentFile('SLR-Data.csv')
#https://drive.google.com/file/d/1ORTNkZgo6uc5VSYopHxl787fPdrs7-LT/view?usp=sharing

#installation of python libraries
import numpy as np
import matplotlib.pyplot as plt
import pandas as pd

# Load dataset
data = pd.read_csv('SLR-Data.csv')
print(data.shape)
print(data.head())
#print(data.describe())

# Collecting X and Y
```

```python
X = data['No of Hours Spent During(X)'].values
Y = data['Risk Score on a scale of 0-100(Y)'].values

# Calculate Mean X and Y
mean_x = np.mean(X)
mean_y = np.mean(Y)
#print(mean_x)
#print(mean_y)
# Total number of values
m = len(X)
# Using the formula to calculate b1(slope) and b0(intercept)
numer = 0
denom = 0
for i in range(m):
    numer += (X[i] - mean_x) * (Y[i] - mean_y)
    denom += (X[i] - mean_x) ** 2
b1 = numer / denom
b0 = mean_y - (b1 * mean_x)

# Print coefficients:b1,b0
print("Slope,Intercept:",b1,b0)

# Plotting Values and Regression Line
max_x = np.max(X)
min_x = np.min(X)

# Calculating line values x and y
x = np.linspace(min_x, max_x)
y = b0 + b1 * x

# Ploting Line
plt.plot(x, y, color='green', label='Regression Line')
# Ploting Scatter Points
plt.scatter(X, Y, c='blue', label='Scatter Plot')
plt.xlabel('No of Hours Spent During')
plt.ylabel('Risk Score on a scale of 0-100')
plt.legend()
plt.show()

#For  Calculating Root Mean Squares Error
rmse = 0
for i in range(m):
    y_pred = b0 + b1 * X[i]
    rmse += (Y[i] - y_pred) ** 2
rmse = np.sqrt(rmse/m)
print("Root Mean Squares Error:",rmse)
# Calculating Accuracy Score
ss_t = 0
ss_r = 0
for i in range(m):
```

```python
    y_pred = b0 + b1 * X[i]
    ss_t += (Y[i] - mean_y) ** 2
    ss_r += (Y[i] - y_pred) ** 2
r2 = 1 - (ss_r/ss_t)
print("Accuracy:",r2*100)
#predicting a o/p (y) for new value of x
predict_x=int(input('Enter No Hours Spent in Driving:'))
predict_y=(4.58789861*predict_x)+12.584627964022907
plt.scatter(X,Y)
plt.scatter(predict_x,predict_y)
plt.xlabel('No Hours Spent Driving(Predicted_x)')
plt.ylabel('Risk Score on a Scale of 0-100(Predicted_y)')
 #plotting the regression line
plt.scatter(X, Y, c='blue')
plt.plot(x, y, color='green')
# function to show plot
plt.show()


#Home Assignment Height & Weight Problem

downloaded = drive.CreateFile({'id':'1e10Ynfgrc35FtMl2V5qpzTGyuWF4KQsZ'}) # replace
the id with id of file you want to access
downloaded.GetContentFile('hw.csv')

#Read file Height & Weight csv as panda dataframe
import pandas as pd
xyz = pd.read_csv('hw.csv')
print(xyz.head(5))

#Running Height & Weight Program(SLR)
import numpy as np
import matplotlib.pyplot as plt
import pandas as pd

#Load dataset
dataset=pd.read_csv("hw.csv")
# To display dataset
print(dataset)
x=dataset.iloc[:,:-1].values
X=dataset.iloc[:,:-1].values
y=dataset.iloc[:,1].values
print(X)
print(y)
#from sklearn subpackage import linear regression model
from sklearn.linear_model import LinearRegression
regressor=LinearRegression()
regressor.fit(X,y)
#To get the slop
regressor.coef_
#To get the y intercept
```

```python
regressor.intercept_
#To print the equation of line
print("y= "+ str(regressor.coef_) + "X + " + str(regressor.intercept_))

#To get the slop
print("Accuracy:",regressor.score(X,y)*100)
#To plot graph
plt.plot(X,y,'o')
plt.plot(X,regressor.predict(X));
plt.show()
predict_x=int(input('Enter Height:'))
predict_y=(0.67461045*predict_x)-38.45508707607698
plt.scatter(X,y)
plt.scatter(predict_x,predict_y)
plt.xlabel('Enter Height:(Predicted_x)')
plt.ylabel('Enter Weight:(Predicted_y)')
#plotting the Predicted regression line
plt.plot(X,regressor.predict(X),color='green');
plt.show()

"""#**ASSIGNMENT -2:Principal Component Analysis for feature reduction**

"""

# Commented out IPython magic to ensure Python compatibility.
import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
from sklearn.decomposition import PCA
from sklearn.preprocessing import StandardScaler
# %matplotlib inline

url = "https://archive.ics.uci.edu/ml/machine-learning-databases/iris/iris.data"
df = pd.read_csv(url
                 , names=['sepal length','sepal width','petal length','petal
width','target'])
df.head()
#downloaded = drive.CreateFile({'id':'12BY34aCbYLoLjy3gDUMrZEBUf7l5FZsd'}) # replace
the id with id of file you want to access
#downloaded.GetContentFile('iris.csv')
#dataset=pd.read_csv("iris.csv")
#dataset
#https://drive.google.com/file/d/12BY34aCbYLoLjy3gDUMrZEBUf7l5FZsd/view?usp=share_li
nk

"""**Standardize the Data**
```

Since PCA yields a feature subspace that maximizes the variance along the axes, it
makes sense to standardize the data, especially, if it was measured on different
scales. Although, all features in the Iris dataset were measured in centimeters, let

us continue with the transformation of the data onto unit scale (mean=0 and variance=1), which is a requirement for the optimal performance of many machine learning algorithms.
"""

```python
features = ['sepal length', 'sepal width', 'petal length', 'petal width']
x = df.loc[:, features].values

y = df.loc[:,['target']].values

x = StandardScaler().fit_transform(x)

pd.DataFrame(data = x, columns = features).head()

#PCA Projection to 2D
pca = PCA(n_components=2)
principalComponents = pca.fit_transform(x)
principalDf = pd.DataFrame(data = principalComponents
            , columns = ['principal component 1', 'principal component 2'])
principalDf.head(5)

df[['target']].head()

finalDf = pd.concat([principalDf, df[['target']]], axis = 1)
finalDf.head(5)

"""**Visualize 2D Projection**"""

fig = plt.figure(figsize = (8,8))
ax = fig.add_subplot(1,1,1)
ax.set_xlabel('Principal Component 1', fontsize = 15)
ax.set_ylabel('Principal Component 2', fontsize = 15)
ax.set_title('2 Component PCA', fontsize = 20)


targets = ['Iris-setosa', 'Iris-versicolor', 'Iris-virginica']
colors = ['r', 'g', 'b']
for target, color in zip(targets,colors):
    indicesToKeep = finalDf['target'] == target
    ax.scatter(finalDf.loc[indicesToKeep, 'principal component 1']
               , finalDf.loc[indicesToKeep, 'principal component 2']
               , c = color
               , s = 50)
ax.legend(targets)
ax.grid()

#Explained Variance:The explained variance tells us how much information (variance)
#can be attributed to each of the principal components.
pca.explained_variance_ratio_
```

```python
"""Together, the first two principal components contain 95.80% of the information.
The first principal component contains 72.77% of the variance and the second
principal component contains 23.03% of the variance. The third and fourth principal
component contained the rest of the variance of the dataset.

# **ASSIGNMENT-3: Decision Tree**
"""

downloaded = drive.CreateFile({'id':'1jql2mwV15BCFeX52G1PGSCr8Y4jLdn8f'}) # replace
the id with id of file you want to access
downloaded.GetContentFile('DT-Data.csv')

#import packages
import numpy as np
import matplotlib.pyplot as plt
import pandas as pd

#reading Dataset
dataset=pd.read_csv("DT-Data.csv")
X=dataset.iloc[:,:-1]
y=dataset.iloc[:,5].values

#Perform Label encoding
from sklearn.preprocessing import LabelEncoder
labelencoder_X = LabelEncoder()

X = X.apply(LabelEncoder().fit_transform)
print (X)

from sklearn.tree import DecisionTreeClassifier
regressor=DecisionTreeClassifier()
regressor.fit(X.iloc[:,1:5],y)

#Predict value for the given expression
X_in=np.array([0,1,0,1])

y_pred=regressor.predict([X_in])
print ("Prediction:", y_pred)

from six import StringIO
from IPython.display import Image
from sklearn.tree import export_graphviz
import pydotplus
# Create DOT data
dot_data = StringIO()

export_graphviz(regressor, out_file=dot_data, filled=True, rounded=True,
special_characters=True)
# Draw graph
graph = pydotplus.graph_from_dot_data(dot_data.getvalue())
```

```python
graph.write_png('Decision_Tree.png')
# Show graph
Image(graph.create_png())

"""# **ASSIGNMENT-4:Naive Bayes**"""

import numpy as np
import matplotlib.pyplot as plt
import pandas as pd
from sklearn.model_selection import train_test_split
from sklearn.preprocessing import StandardScaler
from sklearn.naive_bayes import GaussianNB
from sklearn.metrics import confusion_matrix

#from sklearn.datasets import load_iris
#dataset = load_iris()
#dataset
downloaded = drive.CreateFile({'id':'12BY34aCbYLoLjy3gDUMrZEBUf7l5FZsd'}) # replace
the id with id of file you want to access
downloaded.GetContentFile('iris.csv')
dataset=pd.read_csv("iris.csv")
dataset
#https://drive.google.com/file/d/12BY34aCbYLoLjy3gDUMrZEBUf7l5FZsd/view?usp=share_li
nk
#from sklearn.datasets import load_iris
#dataset = load_iris()
#dataset
downloaded = drive.CreateFile({'id':'12BY34aCbYLoLjy3gDUMrZEBUf7l5FZsd'}) # replace
the id with id of file you want to access
downloaded.GetContentFile('iris.csv')
dataset=pd.read_csv("iris.csv")
dataset
#https://drive.google.com/file/d/12BY34aCbYLoLjy3gDUMrZEBUf7l5FZsd/view?usp=share_li
nk

#Spliting the dataset in independent and dependent variables
X = dataset.iloc[:,:4].values
y = dataset['variety'].values

# Splitting the dataset into the Training set and Test set
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size = 0.20,
random_state = 82)

# Feature Scaling to bring the variable in a single scale
sc = StandardScaler()
X_train = sc.fit_transform(X_train)
X_test = sc.transform(X_test)

# Fitting Naive Bayes Classification to the Training set with linear kernel
nb = GaussianNB()
```

```python
nb.fit(X_train, y_train)
y_pred = nb.predict(X_test)

# Making the Confusion Matrix
cm = confusion_matrix(y_test, y_pred)
print(cm)

#finding accuracy from the confusion matrix.
a = cm.shape
corrPred = 0
falsePred = 0

for row in range(a[0]):
    for c in range(a[1]):
        if row == c:
            corrPred +=cm[row,c]
        else:
            falsePred += cm[row,c]
print('Correct predictions: ', corrPred)
print('False predictions', falsePred)
print ('\n\nAccuracy of the Naive Bayes Clasification is: ', corrPred/(cm.sum()))

"""# **Assignment-5:SVM**
Predict if cancer is Benign or malignant. Using historical data about patients
diagnosed with cancer enables doctors to differentiate malignant cases and benign
ones are given independent attributes.
"""

# Load the important packages
from sklearn.datasets import load_breast_cancer
import matplotlib.pyplot as plt
from sklearn.inspection import DecisionBoundaryDisplay
from sklearn.svm import SVC

# Load the datasets
cancer = load_breast_cancer()
X = cancer.data[:, :2]
y = cancer.target

#Build the model
svm = SVC(kernel="linear")
# Trained the model
svm.fit(X, y)

# Plot Decision Boundary
DecisionBoundaryDisplay.from_estimator(
                svm,
                X,
                response_method="predict",
                cmap=plt.cm.Spectral,
```

```
                alpha=0.8,
                xlabel=cancer.feature_names[0],
                ylabel=cancer.feature_names[1],
        )

# Scatter plot
plt.scatter(X[:, 0], X[:, 1],
                    c=y,
                    s=20, edgecolors="k")
plt.show()

#Build the model
svm = SVC(kernel="rbf", gamma=0.5, C=1.0)
# Trained the model
svm.fit(X, y)

# Plot Decision Boundary
DecisionBoundaryDisplay.from_estimator(
                svm,
                X,
                response_method="predict",
                cmap=plt.cm.Spectral,
                alpha=0.8,
                xlabel=cancer.feature_names[0],
                ylabel=cancer.feature_names[1],
        )

# Scatter plot
plt.scatter(X[:, 0], X[:, 1],
                    c=y,
                    s=20, edgecolors="k")
plt.show()

"""# **ASSIGNMENT-6:KMeans-Assignment**"""

#import packages
import numpy as np
import matplotlib.pyplot as plt
import pandas as pd

#create dataset using DataFrame
df=pd.DataFrame({'X':[0.1,0.15,0.08,0.16,0.2,0.25,0.24,0.3],
                    'y':[0.6,0.71,0.9,0.85,0.3,0.5,0.1,0.2]})
f1 = df['X'].values
f2 = df['y'].values
X = np.array(list(zip(f1, f2)))
print(X)

#centroid points
C_x=np.array([0.1,0.3])
```

```python
C_y=np.array([0.6,0.2])
centroids=C_x,C_y

#plot the given points
colmap = {1: 'r', 2: 'b'}
plt.scatter(f1, f2, color='k')
plt.show()

#for i in centroids():
plt.scatter(C_x[0],C_y[0], color=colmap[1])
plt.scatter(C_x[1],C_y[1], color=colmap[2])
plt.show()

C = np.array(list((C_x, C_y)), dtype=np.float32)
print (C)

#plot given elements with centroid elements
plt.scatter(f1, f2, c='#050505')
print("point No.6[0.25,0.5] is belongs to blue cluster(cluster no:2)")
plt.scatter(C_x[0], C_y[0], marker='*', s=200, c='r')
plt.scatter(C_x[1], C_y[1], marker='*', s=200, c='b')
plt.show()


#import KMeans class and create object of it
from sklearn.cluster import KMeans
model=KMeans(n_clusters=2,random_state=0)
model.fit(X)
labels=model.labels_
print(labels)

#using labels find population around centroid
count=0
for i in range(len(labels)):
    if (labels[i]==1):
        count=count+1

print('No of population around cluster 2:',count-1)

#Find new centroids
new_centroids = model.cluster_centers_

print('Previous value of m1 and m2 is:')
print('M1==',centroids[0])
print('M1==',centroids[1])

print('Updated value of m1 and m2 is:')
print('M1==',new_centroids[0])
print('M1==',new_centroids[1])
```

```python
#STEP-1: Import Libraries
# Code to read csv file into colaboratory:
!pip install -U -q PyDrive
from pydrive.auth import GoogleAuth
from pydrive.drive import GoogleDrive
from google.colab import auth
from oauth2client.client import GoogleCredentials

#STEP-2: Autheticate E-Mail ID
auth.authenticate_user()
gauth = GoogleAuth()
gauth.credentials = GoogleCredentials.get_application_default()
drive = GoogleDrive(gauth)

#STEP-3: Get File from Drive using file-ID

downloaded = drive.CreateFile({'id':'1e10Ynfgrc35FtMl2V5qpzTGyuWF4KQsZ'}) # replace
the id with id of file you want to access
downloaded.GetContentFile('hw.csv')

"""# **Assignment-7:Gradient Boost Classifier**

# **Implementation Using scikit-learn**
For implementation on a dataset, we will be using the Income Evaluation dataset,
which has information about an individual's personal life and an output of 50K or
<=50. The dataset can be found here
(https://www.kaggle.com/lodetomasi1995/income-classification)

The task here is to classify the income of an individual, when given the required
inputs about his personal life.

First, let's import all required libraries.
"""

# Import all relevant libraries

from sklearn.ensemble import GradientBoostingClassifier

import numpy as np

import pandas as pd

from sklearn.preprocessing import StandardScaler

from sklearn.model_selection import train_test_split

from sklearn.metrics import accuracy_score, confusion_matrix

from sklearn import preprocessing
```

```python
import warnings

warnings.filterwarnings("ignore")

#STEP-1: Import Libraries
# Code to read csv file into colaboratory:
!pip install -U -q PyDrive
from pydrive.auth import GoogleAuth
from pydrive.drive import GoogleDrive
from google.colab import auth
from oauth2client.client import GoogleCredentials

#STEP-2: Autheticate E-Mail ID
auth.authenticate_user()
gauth = GoogleAuth()
gauth.credentials = GoogleCredentials.get_application_default()
drive = GoogleDrive(gauth)

# Get File from Drive using file-ID
downloaded = drive.CreateFile({'id':'1zI-X3zdiuM9u74zQyKIShvAUtPJQ7jUK'}) # replace
the id with id of file you want to access
downloaded.GetContentFile('income_evaluation.csv')
# https://drive.google.com/file/d/1zI-X3zdiuM9u74zQyKIShvAUtPJQ7jUK/view?usp=sharing
 (Dataset Downloads Link)

#Now let's read the dataset and look at the columns to understand the information
better.
#https://drive.google.com/file/d/1zI-X3zdiuM9u74zQyKIShvAUtPJQ7jUK/view?usp=sharing
df = pd.read_csv('income_evaluation.csv')
df.head()

"""**Here my main aim is to tell you how to implement this on python. Now for
training and testing our model, the data has to be divided into train and test
data.**
We will also scale the data to lie between 0 and 1.
"""

df.shape

df.info()

df.isnull().sum()

df.columns

#df.drop(columns=' fnlwgt',inplace=True)
df.columns

X = df.drop(columns=' income')
y = df[' income']
```

```python
from sklearn.preprocessing import LabelEncoder

def label_encoder(a):
    le = LabelEncoder()
    df[a] = le.fit_transform(df[a])

df.columns

label_list = [' workclass', ' education',' marital-status',
        ' occupation', ' relationship', ' race', ' sex',' native-country', ' income']

for i in label_list:
    label_encoder(i)

df.head()

from sklearn.model_selection import train_test_split

X = df.drop([' income'],axis=1).values    # independant features
y = df[' income'].values                                  # dependant variable

# Choose your test size to split between training and testing sets:
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.20,
random_state=42)

print("X_train shape:",X_train.shape)
print("y_test shape:",y_test.shape)
print("X_test shape:",X_test.shape)
print("y_train shape:",y_train.shape)

#Buildimg Gradient Boosting Model
from sklearn.ensemble import GradientBoostingClassifier
from sklearn.model_selection import cross_val_score
from sklearn.model_selection import GridSearchCV

gradient_booster = GradientBoostingClassifier(learning_rate=0.1)
accuracies = cross_val_score(gradient_booster, X_train, y_train, cv=5)
gradient_booster.fit(X_train,y_train)

print("Train Score:",np.mean(accuracies))
print("Test Score:",gradient_booster.score(X_test,y_test))

result_dict_train = {}
result_dict_test = {}
result_dict_train["Gradient-Boost Default Train Score"] = np.mean(accuracies)
result_dict_test["Gradient-Boost Default Test Score"] =
gradient_booster.score(X_test,y_test)

grid = {
```

```python
    'learning_rate':[0.01,0.05,0.1],
    'n_estimators':np.arange(100,500,100),
}

gb = GradientBoostingClassifier()
gb_cv = GridSearchCV(gb, grid, cv = 4)
gb_cv.fit(X_train,y_train)
print("Best Parameters:",gb_cv.best_params_)
print("Train Score:",gb_cv.best_score_)
print("Test Score:",gb_cv.score(X_test,y_test))

result_dict_train
result_dict_test

"""# **Extra Assignment:KNN**"""

downloaded = drive.CreateFile({'id':'1oikTU46hEkvGW_DeFyWos5_6q3cX6h7B'}) # replace
the id with id of file you want to access
downloaded.GetContentFile('knndata.csv')

#Importing Libraries

import numpy as np
import pandas as pd

# To split dataset into its attributes and labels.

dataset=pd.read_csv("knndata.csv")
X=dataset.iloc[:,:-1].values
print(X)
Y=dataset.iloc[:,2].values
print(Y)

# Training of KNN Classification Model using trained data

from sklearn.neighbors import KNeighborsClassifier
classifier = KNeighborsClassifier(n_neighbors=3)
classifier.fit(X,Y)

# Testing  KNN Classification Model using unseen test data

X_test=np.array([6,6])
y_pred = classifier.predict([X_test])
print ('The predicition of classifier is :', y_pred)
classifier = KNeighborsClassifier(n_neighbors=3,weights='distance')
classifier.fit(X,Y)
# predict the class for points(6,6)
X_test=np.array([6,6])
y_pred = classifier.predict([X_test])
print ('The predicition of classifier is :', y_pred)
```

### **1. Linear Regression**
```python
import pandas as pd
import numpy as np
import matplotlib.pyplot as plt

# Load Data
data = pd.read_csv('linear_regression.csv')
X = data.iloc[:, 0].values  # Assuming first column is hours
y = data.iloc[:, 1].values  # Assuming second column is risk

# Linear Regression
coeff = np.polyfit(X, y, 1)
line_eq = f"y = {coeff[0]:.2f}x + {coeff[1]:.2f}"

# Plot
plt.scatter(X, y, color='blue')
plt.plot(X, np.polyval(coeff, X), color='red')
plt.title(f"Best Fit Line: {line_eq}")
plt.xlabel("Feature")
plt.ylabel("Target")
plt.show()
```

---

### **2. PCA for Feature Reduction**
```python
import pandas as pd
from sklearn.decomposition import PCA
import matplotlib.pyplot as plt

# Load Data
data = pd.read_csv('pca_data.csv')
X = data.iloc[:, :-1].values

# PCA
pca = PCA(n_components=2)
X_reduced = pca.fit_transform(X)

# Plot
plt.scatter(X_reduced[:, 0], X_reduced[:, 1], c=data.iloc[:, -1].values,
cmap='viridis')
plt.title("PCA - Reduced to 2 Components")
plt.show()
```

---

### **3. Decision Tree Classifier**
```python
import pandas as pd
from sklearn.tree import DecisionTreeClassifier, export_text

# Load Data
data = pd.read_csv('decision_tree.csv')
X = data.iloc[:, :-1].values
y = data.iloc[:, -1].values

# Train Model
clf = DecisionTreeClassifier().fit(X, y)

# Predict for test data
test = [[20, 30, 0, 1]]  # Example test input
decision = clf.predict(test)[0]

# Decision Tree
tree = export_text(clf, feature_names=data.columns[:-1])
print(tree, "\nDecision:", decision)
```

---

### **4. Naive Bayes Classification**
```python
import pandas as pd
```

```python
from sklearn.naive_bayes import GaussianNB

# Load Data
data = pd.read_csv('naive_bayes.csv')
X = data.iloc[:, :-1].values
y = data.iloc[:, -1].values

# Model
nb = GaussianNB().fit(X, y)
accuracy = nb.score(X, y)

print(f"Accuracy: {accuracy:.2f}")
```

### **5. SVM Classification**
```python
import pandas as pd
from sklearn.svm import SVC

# Load Data
data = pd.read_csv('svm_data.csv')
X = data.iloc[:, :-1].values
y = data.iloc[:, -1].values

# Model
svm = SVC(kernel='linear').fit(X, y)
accuracy = svm.score(X, y)

print(f"Accuracy: {accuracy:.2f}")
```

---

### **6. K-Means Clustering**
```python
import pandas as pd
from sklearn.cluster import KMeans

# Load Data
data = pd.read_csv('kmeans_data.csv').values

# K-Means
kmeans = KMeans(n_clusters=2, init=data[:2, :], n_init=1).fit(data)

# Outputs
print(f"Cluster of P6: {kmeans.labels_[5]}")
print(f"Population of Cluster 2: {sum(kmeans.labels_ == 1)}")
print(f"Updated Centroids: {kmeans.cluster_centers_}")
```

---

### **7. Gradient Boost Classifier**
```python
import pandas as pd
from sklearn.ensemble import GradientBoostingClassifier
from sklearn.model_selection import train_test_split

# Load Data
data = pd.read_csv('gradient_boost.csv')
X = data.iloc[:, :-1].values
y = data.iloc[:, -1].values

# Train-Test Split
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.3,
random_state=42)

# Model
gb = GradientBoostingClassifier().fit(X_train, y_train)
accuracy = gb.score(X_test, y_test)

print(f"Accuracy: {accuracy:.2f}")
```

What is Linear Regression?
Linear Regression is a statistical technique to find a relationship between an
independent variable and a dependent variable by fitting a straight line. It
predicts the value of the dependent variable based on the given input.

What is the goal of Linear Regression?
The goal is to determine the best-fit line that minimizes the error between
predicted and actual values. This line helps explain the relationship between
variables.

What is the equation of Linear Regression?
The equation is $y = b1 * x + b0$, where b1 is the slope (rate of change), and b0 is
the y-intercept (value of y when x is zero).

What does the slope (b1) mean?
The slope indicates how much the dependent variable (y) changes for a one-unit
increase in the independent variable (x). It shows the strength of the relationship.

What does the intercept (b0) mean?
The intercept is the value of y when x is zero. It represents the starting point of
the line on the y-axis.

What is the dependent variable?

The dependent variable is the value we are trying to predict or explain. For example, predicting a house price based on its size.

What is the independent variable?
The independent variable is the input used to predict the dependent variable. For example, house size is an independent variable for predicting house price.

What method is used to find the best-fit line?
The Least Squares Method minimizes the sum of squared differences between actual and predicted values to determine the best-fit line.

How is the slope (b1) calculated?
The formula for slope is:
$b1 = \Sigma((x - mean\_x) * (y - mean\_y)) / \Sigma((x - mean\_x)^2)$. It captures the relationship between x and y.

How is the intercept (b0) calculated?
The intercept is calculated as:
$b0 = mean\_y - (b1 * mean\_x)$. It ensures the line fits the data properly.

What is Residual?
Residual is the error or difference between the actual and predicted values. It is calculated as: Residual = Actual - Predicted.

What is the cost function?
The cost function measures the error in predictions. The Mean Squared Error (MSE) is a common cost function that averages squared residuals.

How is Linear Regression accuracy measured?
Accuracy is measured using metrics like R-squared ($R^2$), which explains the variance, and Root Mean Squared Error (RMSE), which shows prediction error.

What is R-squared ($R^2$)?
$R^2$ measures how well the independent variable explains the dependent variable. A value closer to 1 indicates a better fit.

What is RMSE?
RMSE measures the average prediction error by calculating the square root of the mean squared differences between actual and predicted values.

What is Simple Linear Regression?
Simple Linear Regression uses one independent variable to predict the dependent variable. For example, predicting house price based on its size.

What is Multiple Linear Regression?
Multiple Linear Regression uses two or more independent variables to predict the dependent variable. For example, predicting house price based on size, location, and age.

What is Overfitting?

Overfitting happens when a model performs well on training data but poorly on new data because it memorizes the training data instead of generalizing.

What is Underfitting?
Underfitting occurs when a model is too simple and fails to capture the relationship in the data. This results in poor predictions.

What are the assumptions of Linear Regression?
Key assumptions include linearity between variables, constant variance of residuals, independence of errors, no multicollinearity, and normal distribution of residuals.

What is Linearity?
Linearity assumes that the relationship between the independent and dependent variables is a straight line.

What is Homoscedasticity?
Homoscedasticity means the variance of residuals is constant across all levels of the independent variable, ensuring reliability in predictions.

What is Multicollinearity?
Multicollinearity occurs when independent variables are highly correlated, making it hard to determine their individual effects on the dependent variable.

How is data split for Linear Regression?
Data is typically split into training (e.g., 80%) and testing (e.g., 20%) sets to train the model and evaluate its accuracy.

What Python libraries are used for Linear Regression?
Popular libraries include scikit-learn for modeling, Pandas and NumPy for data handling, and Matplotlib for visualization.

Which scikit-learn function is used for Linear Regression?
The LinearRegression() function from scikit-learn is used to create and train Linear Regression models.

How is a model trained in scikit-learn?
The model is trained using the fit() method, where input features (X) and output labels (y) are passed as arguments.

How are predictions made in scikit-learn?
Predictions are made using the predict() method, which uses the trained model to generate outputs for new inputs.

What is Gradient Descent?
Gradient Descent is an optimization method to minimize the cost function by iteratively updating parameters in the direction of least error.

Why is Standardization important?
Standardization scales all variables to have a mean of 0 and a standard deviation of 1, ensuring fair comparison and faster model convergence.

What is Cross-Validation?
Cross-Validation splits the data into multiple sets to train and test the model, helping assess its performance on unseen data.

What is the difference between Training and Testing Error?
Training error measures the model's performance on the training dataset, while testing error measures how well it generalizes to new data.

What are the effects of outliers?
Outliers can skew the regression line, leading to inaccurate predictions. Removing or addressing outliers improves model performance.

How to handle missing data?
Missing data can be handled by imputing values (e.g., mean, median) or dropping rows/columns with too many missing values.

What is Polynomial Regression?
Polynomial Regression is an extension of Linear Regression that fits a curve by including higher-order terms of the independent variable.

What is the Correlation Coefficient?
The correlation coefficient measures the strength and direction of the linear relationship between two variables, ranging from -1 to 1.

What is a p-value?
A p-value indicates the significance of coefficients. A small p-value (e.g., < 0.05) suggests a strong relationship between variables.

Can Linear Regression handle non-linear relationships?
No, Linear Regression assumes linear relationships. For non-linear data, methods like Polynomial Regression or decision trees are better.

What is the importance of the intercept in predictions?
The intercept ensures the regression line correctly predicts values when the independent variable is zero.

What are common applications of Linear Regression?
Linear Regression is used in price prediction, sales forecasting, demand analysis, and trend identification in various fields.

What is PCA?
Principal Component Analysis (PCA) is a dimensionality reduction technique used to reduce the number of features in a dataset while retaining as much variance as possible. It transforms the original features into a new set of orthogonal features called principal components.

Why is PCA important for feature reduction?

PCA reduces the complexity of a dataset by projecting it onto a smaller number of components that capture the most variance. This helps improve computational efficiency and prevent overfitting in machine learning models.

What is the role of standardization in PCA?
Standardization ensures that each feature has a mean of 0 and variance of 1. This is important because PCA is sensitive to the scale of the data, and features with larger scales could disproportionately affect the results.

What is the first step in applying PCA to a dataset?
The first step is to standardize the data, which involves transforming the features to have zero mean and unit variance using techniques like StandardScaler in scikit-learn.

Why do we use PCA on datasets with features of different scales?
PCA is sensitive to the scales of features, so if features are on different scales, they could dominate the principal components. Standardization equalizes the influence of each feature on the analysis.

What is the output of PCA when transforming the data?
The output of PCA is a set of new features, called principal components, which are linear combinations of the original features. These components capture the most variance in the data.

How does PCA work mathematically?
PCA finds the eigenvectors (principal components) of the covariance matrix of the dataset. It projects the data onto these eigenvectors to reduce dimensionality while preserving the most important information.

What is the covariance matrix in PCA?
The covariance matrix is a square matrix that represents the covariance between different features in the dataset. PCA uses this matrix to identify the principal components that maximize variance.

What does "explained variance" mean in PCA?
Explained variance indicates how much information (or variance) each principal component captures. It helps to understand how well the PCA transformation has preserved the original data's structure.

What is the significance of the eigenvalues in PCA?
Eigenvalues correspond to the amount of variance explained by each principal component. Larger eigenvalues indicate that the component captures more information from the dataset.

What is the purpose of the explained_variance_ratio_ attribute in PCA?
The explained_variance_ratio_ attribute provides the proportion of variance explained by each principal component. It helps determine how many components are needed to explain most of the variance in the data.

What is the role of PCA(n_components=2) in the given code?

In the code, PCA(n_components=2) reduces the dataset to two principal components, allowing for a 2D visualization of the data and preserving the most significant variance.

How do you visualize PCA results in 2D?
To visualize PCA results, you can plot the first two principal components using a scatter plot. This shows how the data is distributed along these components.

What is the function of StandardScaler() in PCA?
StandardScaler() standardizes the dataset by removing the mean and scaling the data to unit variance. This ensures that PCA operates on a uniform scale across all features.

What is the advantage of using PCA for visualization?
PCA helps visualize high-dimensional data by reducing it to 2 or 3 dimensions, making it easier to identify patterns, clusters, or outliers in the data.

What does the term "Principal Component" mean?
A Principal Component is a new feature created by PCA that is a linear combination of the original features. These components are ordered by the amount of variance they capture from the data.

What does the term "Variance" mean in PCA?
Variance measures the spread of data points. In PCA, variance is crucial because it helps determine how much of the data's information is retained by each principal component.

How can you reduce dimensions using PCA?
You reduce dimensions by selecting the top principal components that explain most of the variance. For example, using PCA(n_components=2) reduces the data to 2 components.

Why is PCA commonly used in machine learning?
PCA is used in machine learning to reduce the feature space, speeding up model training, decreasing overfitting, and improving model generalization by removing noise.

What does a 2D PCA plot show?
A 2D PCA plot shows how the data is distributed along the two principal components, helping visualize relationships and clusters in the dataset.

What do the colors represent in the PCA visualization?
In the PCA visualization, the colors represent different categories or classes of data points. For example, in the Iris dataset, each color represents a different flower species.

Why do we plot the target variable in PCA?
Plotting the target variable in PCA helps understand how well the data points are separated by the principal components, giving insights into the effectiveness of the dimensionality reduction.

What is the target variable in the Iris dataset?
The target variable in the Iris dataset is the species of the Iris flower, which has three categories: Iris-setosa, Iris-versicolor, and Iris-virginica.

What is the benefit of using PCA for feature selection?
PCA helps identify the most important features by creating components that explain the most variance, allowing you to keep only those components and reduce redundant or less useful features.

What is the cumulative explained variance in PCA?
Cumulative explained variance is the total variance explained by the first n components. It helps decide how many components to retain for preserving most of the dataset's information.

How does PCA handle correlation between features?
PCA eliminates correlation between features by transforming them into orthogonal principal components, where each component is uncorrelated with others.

What is the impact of retaining too many components in PCA?
Retaining too many components may lead to overfitting and computational inefficiency, as the additional components may capture noise rather than useful variance.

What is the impact of retaining too few components in PCA?
Retaining too few components may result in losing important information, leading to underfitting and poor model performance due to insufficient variance capture.

What does the term "orthogonal components" mean in PCA?
Orthogonal components are uncorrelated and independent. In PCA, the principal components are orthogonal to each other, ensuring each captures unique variance in the data.

How do you interpret the explained variance ratio?
The explained variance ratio tells you how much variance each principal component explains in relation to the total variance. It helps assess the importance of each component.

What is the role of pd.concat() in the given code?
pd.concat() is used to merge the principal components with the target variable (target) into a single DataFrame for easier analysis and visualization.

What is the function of ax.scatter() in the PCA visualization?
ax.scatter() creates a scatter plot by plotting the data points based on the first two principal components. Each point is colored according to its species in the Iris dataset.

How do you interpret the PCA visualization?
In the PCA visualization, points that are closer together represent similar data, while points that are farther apart represent different species or categories.

Why does PCA focus on maximizing variance?
Maximizing variance ensures that the new components capture the most informative aspects of the data, which are important for making accurate predictions.

What does the "Principal Component 1" represent?
Principal Component 1 is the direction in the data that captures the most variance. It is the most significant linear combination of the original features.

What does the "Principal Component 2" represent?
Principal Component 2 is the second most significant direction in the data, orthogonal to Principal Component 1, capturing the next largest variance.

What are the limitations of PCA?
PCA assumes linear relationships between features and may not perform well with non-linear data. It also requires the data to be standardized and can be sensitive to outliers.

What other dimensionality reduction techniques are similar to PCA?
Other dimensionality reduction techniques include Linear Discriminant Analysis (LDA), t-SNE, and Autoencoders, which can handle different types of data and relationships.

What is the main advantage of using PCA over raw features?
PCA reduces noise and redundancy, making models faster and more efficient, while improving generalization by focusing on the most informative components.

Can PCA be applied to non-numeric data?
PCA requires numeric data, as it relies on mathematical operations like covariance and eigenvalues. For non-numeric data, other techniques like One-Hot Encoding or embedding methods are used first.


1. What is a Decision Tree Classifier?
A Decision Tree Classifier is a machine learning algorithm used for classification tasks. It splits the data into branches to make decisions based on the values of input features. The tree-like structure helps visualize decisions and outcomes.

2. What is the purpose of the target variable in the Decision Tree?
The target variable represents the output or the decision the model is trying to predict based on the input features. In this case, it is whether the customer buys the product ("Buys").

3. What is label encoding?
Label encoding converts categorical values into numerical values. This transformation is necessary for algorithms that can only work with numerical input.

4. What is the importance of the root node in a decision tree?
The root node is the starting point of a decision tree. It is the feature that best

splits the data to maximize the classification accuracy.

5. What is the role of a Decision Tree in predictive modeling?
A Decision Tree is used for making predictions by learning from past data. It splits the data based on feature values and predicts the target variable for new inputs.

6. What is overfitting in Decision Trees?
Overfitting occurs when the model learns too much from the training data, capturing noise and details that do not generalize well to new data. This results in a high training accuracy but poor test performance.

7. What is pruning in Decision Trees?
Pruning is the process of removing parts of the tree that do not provide additional predictive power. It helps prevent overfitting.

8. How is the split point chosen in a Decision Tree?
The algorithm chooses the feature and its threshold that results in the best split, typically maximizing information gain or minimizing impurity (Gini or Entropy).

9. What does the Gini Index measure in a Decision Tree?
The Gini Index measures the impurity of a split. A lower Gini value indicates a purer node, meaning the data points in that node are more similar.

10. What is entropy in Decision Trees?
Entropy measures the disorder or uncertainty in the dataset. A node with high entropy means the data points are mixed, and a node with low entropy means the data points are more homogeneous.

11. How does a Decision Tree handle continuous features?
Continuous features are split based on a threshold value. For example, an "age" feature might be split at 25 years old, creating two branches: age < 25 and age >= 25.

12. What is the importance of data preprocessing in Decision Trees?
Data preprocessing ensures that the features are in a suitable format for the model. For example, categorical features need to be label encoded, and features should be cleaned for missing values.

13. What are the advantages of using Decision Trees?
Decision Trees are easy to understand, visualize, and interpret. They can handle both numerical and categorical data and are non-parametric, meaning they do not require assumptions about the underlying data distribution.

14. What are the disadvantages of using Decision Trees?
Decision Trees can easily overfit, especially with deep trees. They can also be sensitive to small changes in data, and their performance can degrade with noisy data.

15. How do you determine the depth of a Decision Tree?
The depth of a tree is determined by the number of splits (or levels) from the root

node to the leaf nodes. Limiting the depth helps prevent overfitting.

16. What is the Decision Tree's criterion for selecting the best split?
The criterion is typically either Gini Impurity or Information Gain (Entropy). These criteria evaluate the quality of a split.

17. What are the steps in building a Decision Tree model?
The steps include loading the dataset, preprocessing the data, splitting the dataset into features and target variables, training the model, and making predictions.

18. How do you predict values using a trained Decision Tree?
A trained Decision Tree predicts by traversing the tree from the root to a leaf, following the splits based on feature values to arrive at a decision.

19. What is the purpose of the export_graphviz function in Decision Trees?
The export_graphviz function generates a textual representation of the decision tree, which can be visualized as an image using tools like pydotplus.

20. What does the predict function in a Decision Tree return?
The predict function returns the predicted class or target variable for the input data based on the learned decision rules.

21. What is the significance of the leaf nodes in a Decision Tree?
Leaf nodes represent the final decision or outcome in the tree. Each leaf node corresponds to a class label or target variable value.

22. How do you prevent a Decision Tree from overfitting?
Overfitting can be controlled by setting parameters such as max depth, min samples per split, or by pruning the tree.

23. What is cross-validation, and why is it important in Decision Trees?
Cross-validation is a technique to evaluate the model's performance on unseen data. It helps prevent overfitting by ensuring the model generalizes well to new data.

24. What is feature selection, and why is it important in Decision Trees?
Feature selection is the process of choosing the most relevant features for the model. This improves the model's efficiency and reduces the risk of overfitting.

25. How does a Decision Tree handle missing data?
Decision Trees can handle missing data by using strategies like assigning the most frequent value or predicting missing values based on the majority of other features.

26. What is a Decision Boundary in a Decision Tree?
A Decision Boundary is the line or region that separates different classes based on the Decision Tree model. It represents the decision rule used by the tree.

27. What is the purpose of using LabelEncoder in preprocessing?
LabelEncoder is used to convert categorical variables into numerical labels, making them suitable for machine learning algorithms.

28. What is a split in a Decision Tree?
A split is a decision point in the tree where the data is divided into two or more branches based on a feature's value.

29. What does the fit function in a Decision Tree do?
The fit function trains the Decision Tree model by learning the relationships between the features and the target variable.

30. Why is feature scaling not required in Decision Trees?
Feature scaling is not required because Decision Trees do not rely on distances between points, unlike algorithms such as k-nearest neighbors or support vector machines.

31. What is a categorical feature in Decision Trees?
A categorical feature is a feature that takes on discrete values, such as "Gender" (Male/Female) or "Marital Status" (Married/Single).

32. What is the entropy formula used in Decision Trees?
The entropy formula is:

$$H(S) = -\sum_{i=1}^{n} p_i \log_2(p_i)$$

i

 )
Where
$p$
$i$
p
i

  is the probability of class
$i$
i in the dataset
$S$
S.

33. What is Information Gain in Decision Trees?
Information Gain measures how much uncertainty is reduced after splitting a dataset based on a specific feature. It is calculated as the difference between entropy before and after the split.

34. How do Decision Trees handle numerical features?
Numerical features are split by thresholds. For example, a feature "Age" might be split at age > 30, creating two branches for ages greater than or less than 30.

35. What is the max_depth parameter in a Decision Tree?
The max_depth parameter controls the maximum depth of the tree. Limiting depth helps to avoid overfitting by preventing the tree from becoming too complex.

36. How do Decision Trees handle multi-class classification?
Decision Trees can handle multi-class classification by using splits that separate multiple classes into distinct branches.

37. What does the "root node" represent in a Decision Tree?
The root node is the topmost node and represents the first feature to split on. It is the feature that best separates the data according to the chosen criterion.

38. What is the min_samples_split parameter in a Decision Tree?
The min_samples_split parameter defines the minimum number of samples required to split an internal node. It helps prevent overfitting by ensuring nodes have enough data before splitting.

39. What is the role of min_samples_leaf in a Decision Tree?
The min_samples_leaf parameter sets the minimum number of samples that a leaf node must have. This helps avoid creating leaves that are too specific to the training data.

40. How do you visualize the trained Decision Tree model?
The trained Decision Tree model can be visualized using tools like export_graphviz and pydotplus, which generate graphical representations of the tree.

These questions and answers cover the key concepts related to Decision Trees, including the process of building, evaluating, and interpreting them.


1. What is Naive Bayes Classification?
Answer:
Naive Bayes is a probabilistic machine learning algorithm based on Bayes' Theorem, assuming that the features are conditionally independent given the class. It is used for classification tasks.

2. What is Bayes' Theorem?
Answer:
Bayes' Theorem calculates the probability of a class given the features using the formula:

$P$
$($
$C$
$|$
$X$
$)$
$=$
$P$
$($
$X$
$|$
$C$
$)$
$\cdot$
$P$
$($
$C$
$)$
$P$
$($
$X$
$)$

$$P(C|X) = \frac{P(X|C) \cdot P(C)}{P(X)}$$


Where:

$P$
$($
$C$
$|$
$X$
$)$

$P(C|X)$ is the probability of the class given the features.

$P$
$($
$X$
$|$
$C$
$)$
P(X|C) is the probability of the features given the class.
$P$
$($
$C$
$)$
P(C) is the prior probability of the class.
$P$
$($
$X$
$)$
P(X) is the total probability of the features.
3. What types of Naive Bayes classifiers exist?
Answer:
There are three main types of Naive Bayes classifiers:

Gaussian Naive Bayes (used for continuous data).
Multinomial Naive Bayes (used for discrete count data).
Bernoulli Naive Bayes (used for binary/boolean features).
4. Why is Naive Bayes considered "naive"?
Answer:
It is called "naive" because it assumes that the features are conditionally
independent of each other, which is rarely true in real-world data. Despite this
assumption, Naive Bayes often performs surprisingly well.

5. What is the Gaussian Naive Bayes classifier?
Answer:
The Gaussian Naive Bayes classifier is used when the features are continuous and
assumes that the data follows a normal distribution (Gaussian distribution). It
calculates the probability based on the mean and variance of the features for each
class.

6. What are the advantages of Naive Bayes?
Answer:

Simple and easy to implement.
Works well with high-dimensional data.
Computationally efficient, especially for large datasets.
Performs well with categorical and continuous features.
7. What are the limitations of Naive Bayes?
Answer:

Assumes independence between features, which is rarely true in real-world data.
Performance can degrade when the independence assumption is strongly violated.
Can struggle with zero probabilities in categorical data (Laplace smoothing can

mitigate this).
8. How is Naive Bayes used for classification?
Answer:
Naive Bayes calculates the posterior probability for each class given the input
features and predicts the class with the highest posterior probability.

9. What is a confusion matrix?
Answer:
A confusion matrix is a table that is used to evaluate the performance of a
classification algorithm. It shows the actual vs predicted classifications,
including true positives, false positives, true negatives, and false negatives.

10. How do you compute accuracy from a confusion matrix?
Answer:
Accuracy is calculated as:
Accuracy
=
Correct Predictions
Total Predictions
Accuracy=
Total Predictions
Correct Predictions


Where correct predictions are the sum of true positives and true negatives.

11. What is the purpose of feature scaling in Naive Bayes?
Answer:
Feature scaling helps normalize the range of independent variables or features,
ensuring that no feature dominates due to its scale, particularly important when
using Gaussian Naive Bayes with continuous data.

12. What is StandardScaler in scikit-learn?
Answer:
StandardScaler standardizes the features by removing the mean and scaling to unit
variance, ensuring that each feature contributes equally to the model.

13. What does train_test_split do in scikit-learn?
Answer:
The train_test_split function splits the dataset into two parts: one for training
the model and the other for testing the model's performance.

14. How do you interpret the output of a confusion matrix?
Answer:
The confusion matrix consists of four values:

True Positives (TP): Correctly predicted positive values.
True Negatives (TN): Correctly predicted negative values.
False Positives (FP): Incorrectly predicted positive values.
False Negatives (FN): Incorrectly predicted negative values.

15. What is the role of Laplace smoothing in Naive Bayes?
Answer:
Laplace smoothing (also known as additive smoothing) is used to handle zero probabilities in the dataset by adding a small constant (usually 1) to all feature probabilities.

16. What is the difference between Gaussian Naive Bayes and Multinomial Naive Bayes?
Answer:

Gaussian Naive Bayes is used for continuous features, assuming a normal distribution.
Multinomial Naive Bayes is used for discrete count data, where features represent the frequency of occurrence.
17. What is a probability density function (PDF)?
Answer:
A probability density function (PDF) is a function that describes the likelihood of a random variable to take on a particular value. In Gaussian Naive Bayes, it's used to model continuous features.

18. How does Naive Bayes handle categorical data?
Answer:
Naive Bayes handles categorical data by calculating the probability of each category of a feature for each class and uses these probabilities for classification.

19. What is the fit function used for in scikit-learn's Naive Bayes?
Answer:
The fit function is used to train the Naive Bayes model on the training data. It learns the parameters (mean, variance, or probability) required to classify new data.

20. What is the predict function used for in scikit-learn's Naive Bayes?
Answer:
The predict function is used to classify new, unseen data based on the trained Naive Bayes model.

21. What is the role of the score function in Naive Bayes?
Answer:
The score function returns the accuracy of the Naive Bayes model on the test data. It measures how well the model performs.

22. What is the effect of incorrect feature scaling on Naive Bayes?
Answer:
Incorrect scaling may distort the features and affect the model's assumptions, leading to inaccurate predictions, especially in Gaussian Naive Bayes.

23. What is the purpose of train_test_split's random_state parameter?
Answer:
The random_state parameter ensures reproducibility by controlling the randomness of the data splitting process.

24. What does the fit_transform method do in StandardScaler?
Answer:
The fit_transform method first fits the scaler to the training data (computes the mean and standard deviation) and then transforms the data by standardizing it.

25. What is the transform method in StandardScaler used for?
Answer:
The transform method is used to apply the scaling learned from the training data to the test data.

26. What is the primary advantage of Naive Bayes over other classifiers?
Answer:
Naive Bayes is computationally efficient, especially when dealing with high-dimensional data, and it performs well even with small datasets.

27. How do you handle missing values in Naive Bayes?
Answer:
You can handle missing values by either removing the rows with missing values, imputing the missing values using mean/median/mode, or using models that can handle missing data.

28. What happens when Naive Bayes assumes that the features are not independent?
Answer:
The model will perform poorly because the conditional independence assumption is violated, leading to incorrect probability estimations.

29. What is the difference between precision and recall?
Answer:

Precision is the ratio of correctly predicted positive observations to the total predicted positives.
Recall is the ratio of correctly predicted positive observations to all actual positives.
30. What is an ROC curve?
Answer:
An ROC (Receiver Operating Characteristic) curve is a graphical representation of the model's performance at various classification thresholds, showing the trade-off between true positive rate and false positive rate.

31. How do you compute F1-score?
Answer:
F1-score is the harmonic mean of precision and recall:
F1-score
=
2
×
Precision
×
Recall
Precision

+
Recall
F1-score=2×
Precision+Recall
Precision×Recall


32. Can Naive Bayes handle multi-class classification?
Answer:
Yes, Naive Bayes can handle multi-class classification by calculating the posterior
probabilities for each class and selecting the class with the highest probability.

33. How do you interpret the probabilities outputted by Naive Bayes?
Answer:
The probabilities represent the likelihood of each class given the input features,
and the class with the highest probability is the predicted class.

34. What is cross-validation?
Answer:
Cross-validation is a technique used to evaluate a model's performance by splitting
the data into multiple folds and training/testing on different combinations of these
folds.

35. How can Naive Bayes be improved if it doesn't perform well?
Answer:

Feature engineering to remove irrelevant features or add new ones.
Combining Naive Bayes with other models (ensemble methods).
Applying smoothing techniques or handling missing data more effectively.
36. What is the significance of the Iris dataset?
Answer:
The Iris dataset is a famous dataset used for classification, containing features
like sepal length, sepal width, petal length, and petal width to classify flowers
into three species.

37. How do you visualize the performance of Naive Bayes?
Answer:
You can visualize performance using confusion matrices, ROC curves, and
precision-recall curves.

38. How do you deal with highly imbalanced classes in Naive Bayes?
Answer:

You can balance the classes using techniques like oversampling or undersampling.
Adjust the class priors to reflect the imbalance.
39. What is the purpose of class_prior in Naive Bayes?
Answer:
The class_prior parameter sets the prior probabilities of the classes, which is
important in imbalanced datasets.

40. What performance metrics are most useful in evaluating Naive Bayes?
Answer:
Common metrics include accuracy, precision, recall, F1-score, and AUC-ROC.


Here are 40 potential viva questions with answers related to the SVM Classification Assignment you provided:

1. What is the SVM algorithm?
Answer: SVM (Support Vector Machine) is a supervised machine learning algorithm used for classification and regression tasks. It works by finding the hyperplane that best separates the data points of different classes.
2. What is the role of the hyperplane in SVM?
Answer: The hyperplane is a decision boundary that separates different classes in the feature space. In the case of binary classification, it separates the positive and negative classes. SVM aims to find the hyperplane that maximizes the margin between the two classes.
3. What are support vectors?
Answer: Support vectors are the data points that are closest to the hyperplane and influence its position and orientation. These points are critical in determining the optimal hyperplane.
4. What does the kernel function do in SVM?
Answer: The kernel function maps data into higher-dimensional space to make it easier to find a linear separating hyperplane. Common kernels include linear, polynomial, and radial basis function (RBF) kernels.
5. What is the difference between a linear kernel and a radial basis function (RBF) kernel?
Answer: A linear kernel is used when data can be separated by a straight line (or hyperplane in higher dimensions). An RBF kernel is used for non-linear data and transforms the data into a higher-dimensional space where a linear separator can be used.
6. What is the importance of the 'C' parameter in SVM?
Answer: The 'C' parameter controls the trade-off between achieving a low error on the training data and maintaining a large margin. A higher value of 'C' leads to a smaller margin but fewer misclassifications, while a lower value of 'C' leads to a larger margin but may allow some misclassifications.
7. What does the 'gamma' parameter control in an RBF kernel?
Answer: The 'gamma' parameter defines the influence of a single training example. A small value of gamma means that the influence of each point is far-reaching, whereas a large gamma means that the influence is localized and more sensitive to individual data points.
8. What is the 'decision boundary' in SVM?
Answer: The decision boundary is the line (in 2D) or hyperplane (in higher dimensions) that divides the feature space into different classes. SVM creates a decision boundary that separates the classes with the largest possible margin.
9. How does SVM handle non-linearly separable data?
Answer: SVM uses kernel functions to map non-linearly separable data into

higher-dimensional space where a linear hyperplane can be used to separate the classes.

10. What is meant by a 'soft margin' in SVM?

Answer: A soft margin allows some misclassifications in the data to achieve better generalization. It is controlled by the 'C' parameter, where a higher 'C' reduces misclassifications but increases the risk of overfitting.

11. What is the 'hard margin' in SVM?

Answer: A hard margin is used when the data is linearly separable, and no misclassifications are allowed. It is the strictest form of SVM, where all data points must be correctly classified.

12. What does the fit method in SVM do?

Answer: The fit method trains the SVM model by finding the optimal hyperplane that maximizes the margin between the classes using the training data.

13. What is the role of the score method in SVM?

Answer: The score method computes the accuracy of the model on the provided test data by comparing the predicted labels with the true labels.

14. Explain the importance of feature scaling in SVM.

Answer: Feature scaling is important in SVM because the algorithm is sensitive to the distance between data points. If the features are on different scales, the model may give disproportionate importance to features with larger scales, affecting the performance.

15. What is the kernel='linear' argument used for in SVM?

Answer: The kernel='linear' argument specifies that a linear kernel should be used, meaning SVM will try to find a linear decision boundary to separate the classes.

16. What is the role of the gamma parameter in an RBF kernel?

Answer: The gamma parameter controls the shape of the decision boundary in an RBF kernel. A small gamma creates a smoother decision boundary, while a large gamma creates a more complex, nonlinear decision boundary.

17. How is the model evaluated in your SVM classification?

Answer: The model is evaluated using accuracy, which is the percentage of correctly classified instances. Confusion matrices, precision, recall, and F1 score can also be used to evaluate the model's performance.

18. What is the purpose of the DecisionBoundaryDisplay function?

Answer: The DecisionBoundaryDisplay function visualizes the decision boundaries of the trained model in the feature space, showing how the model separates the classes.

19. Why do you use a scatter plot along with the decision boundary?

Answer: The scatter plot shows the actual data points, while the decision boundary helps to visualize how the SVM model classifies those data points. Together, they provide a clear picture of how well the model separates the classes.

20. Why did you use only the first two features of the dataset?

Answer: We used only the first two features to make the problem visually understandable. In higher-dimensional spaces, visualization becomes difficult, so a 2D plot is created using just two features.

21. Can SVM be used for regression problems?

Answer: Yes, SVM can also be used for regression (SVR). In SVR, the objective is to find a function that approximates the target values, unlike classification where the goal is to separate different classes.

22. What is overfitting in SVM?

Answer: Overfitting occurs when the model is too complex and captures noise in the data rather than general patterns. This often happens when the 'C' parameter is too

large, leading to a very narrow margin that overfits the training data.
23. What is underfitting in SVM?
Answer: Underfitting occurs when the model is too simple to capture the underlying patterns in the data. It happens when the 'C' parameter is too small, leading to a larger margin and more misclassifications.
24. What is the difference between SVM and Logistic Regression?
Answer: SVM is a margin-based classifier that tries to find the optimal hyperplane, whereas logistic regression is a probability-based classifier that models the probability of class membership using a sigmoid function.
25. What is the bias-variance tradeoff in SVM?
Answer: The bias-variance tradeoff refers to the balance between underfitting and overfitting. A simpler model (higher bias) may have higher errors, while a more complex model (higher variance) may overfit the training data.
26. What happens if you use a polynomial kernel?
Answer: A polynomial kernel allows SVM to learn non-linear decision boundaries by mapping the input data into a higher-dimensional space, creating a more flexible model compared to the linear kernel.
27. How do you handle multi-class classification in SVM?
Answer: SVM handles multi-class classification using strategies like "one-vs-one" or "one-vs-rest." In one-vs-one, a classifier is trained for each pair of classes, while in one-vs-rest, a classifier is trained for each class against all others.
28. What is the cost of using an RBF kernel compared to a linear kernel?
Answer: The RBF kernel can capture more complex patterns and provides more flexibility, but it is computationally more expensive than a linear kernel.
29. What is the significance of the 'C' parameter in terms of the margin?
Answer: A larger value of 'C' results in a smaller margin but fewer misclassifications, while a smaller 'C' results in a larger margin but more misclassifications.
30. What would happen if you did not perform feature scaling?
Answer: Without feature scaling, the SVM algorithm may give more importance to features with larger numerical ranges, leading to biased results.
31. What are the advantages of SVM over other classification algorithms?
Answer: SVM is effective in high-dimensional spaces, works well with non-linear data using kernels, and is robust against overfitting, especially in high-dimensional space.
32. How do you choose the best kernel for your problem?
Answer: The choice of kernel depends on the data. If the data is linearly separable, a linear kernel works well. For non-linear data, the RBF kernel is commonly used.
33. How do you select the right value for the 'gamma' parameter?
Answer: The value of gamma can be tuned using techniques like cross-validation. A grid search over possible values of gamma helps to find the best one.
34. What is the role of cross-validation in SVM?
Answer: Cross-validation helps assess the generalization ability of the model by dividing the data into multiple subsets and training/testing the model on different combinations of subsets.
35. What does the term 'dual form' in SVM refer to?
Answer: The dual form of SVM optimization problem is derived from the primal form. It involves maximizing the Lagrangian multipliers instead of the original problem, and is useful when using kernel functions.
36. How does SVM compare to k-NN?

Answer: SVM is more powerful when it comes to handling high-dimensional data and finding a clear margin, while k-NN is a non-parametric method that works by measuring the distance between points.

37. Why is it important to visualize the decision boundary in SVM?
Answer: Visualizing the decision boundary helps in understanding how well the model separates different classes and whether it overfits or underfits.

38. Can SVM handle imbalanced data?
Answer: SVM can handle imbalanced data by adjusting the class weights, which penalize the misclassification of the minority class more heavily.

39. What is the effect of using a high-dimensional feature space in SVM?
Answer: Using high-dimensional feature space can increase the complexity and computational cost but can also lead to better separation of classes if the data is not linearly separable.

40. What is the significance of support vectors in the decision-making process of SVM?
Answer: Support vectors are the critical data points that influence the position of the decision boundary. These points are the most difficult to classify and define the optimal hyperplane.

1. What is K-Means Clustering?
Answer: K-Means is a partitioning method in unsupervised learning that divides a dataset into K clusters based on their similarity, minimizing intra-cluster variance.

2. How does K-Means clustering work?
Answer: K-Means initializes K centroids, assigns each data point to the nearest centroid, and then updates the centroids by calculating the mean of the points in each cluster. This process repeats until convergence.

3. What is the objective of K-Means clustering?
Answer: The objective of K-Means is to partition the data into K clusters such that the sum of squared distances from each point to its assigned centroid is minimized.

4. What is the significance of the centroid in K-Means?
Answer: The centroid represents the center of each cluster and is used to define the cluster by minimizing the variance within the cluster.

5. Why do we need to choose the number of clusters (K) in K-Means?
Answer: The number of clusters (K) must be predefined, and it influences the quality of the clustering results. Various methods like the Elbow Method are used to choose K.

6. How do you determine the value of K in K-Means clustering?
Answer: The value of K can be determined using methods like the Elbow Method, Silhouette Analysis, or by domain knowledge.

7. What is the Elbow Method?
Answer: The Elbow Method involves plotting the within-cluster sum of squares (WCSS)

against different values of K and selecting the K where the decrease in WCSS slows down, forming an "elbow."

8. What is the role of distance in K-Means clustering?

Answer: Distance, typically Euclidean, is used to measure the similarity between data points and centroids. Points closer to a centroid are assigned to that cluster.

9. How does K-Means handle outliers?

Answer: K-Means is sensitive to outliers because they can significantly affect the position of the centroid. Outliers can cause the centroids to shift, leading to poor clustering.

10. What happens if K is set too high or too low?

Answer: If K is set too high, the algorithm may overfit, creating many small clusters. If K is too low, it may underfit, merging distinct clusters into one.

11. What is the importance of initializing centroids in K-Means?

Answer: The initialization of centroids affects the convergence and the final clustering result. Poor initialization may lead to suboptimal solutions.

12. What methods can be used to initialize centroids?

Answer: The common methods for initializing centroids include random initialization, K-Means++, or choosing centroids based on domain knowledge.

13. How does K-Means handle new data points after clustering?

Answer: New data points are assigned to the nearest centroid based on the Euclidean distance.

14. What is the role of the K-Means algorithm's iteration?

Answer: The iteration ensures that centroids are updated and data points are reassigned to clusters until the centroids no longer change, indicating convergence.

15. Can K-Means be used for non-numerical data?

Answer: No, K-Means requires numerical data because it relies on calculating distances between points and centroids, which requires numerical representation.

16. What is the stopping criterion in K-Means clustering?

Answer: The algorithm stops when the centroids do not change significantly between iterations, or a maximum number of iterations is reached.

17. Why is the K-Means algorithm sensitive to initial centroids?

Answer: K-Means can converge to a local minimum depending on the initial centroids, which might result in suboptimal clustering.

18. What is the impact of large K on K-Means performance?

Answer: A large K can increase computation time and lead to overly specific clusters, resulting in overfitting.

19. What is the relationship between K-Means and the concept of variance?

Answer: K-Means minimizes the variance within clusters by updating centroids to reduce the squared Euclidean distance of points from the centroid.

20. How do you evaluate the quality of clustering in K-Means?

Answer: The quality can be evaluated using metrics like within-cluster sum of squares (WCSS), silhouette score, or external validation indices like Adjusted Rand Index.

21. What are the limitations of K-Means?

Answer: K-Means assumes clusters to be spherical and equally sized, is sensitive to outliers, and requires specifying K in advance.

22. What would happen if you set K=1 in K-Means?

Answer: If K=1, all points would be assigned to a single cluster, which is generally not useful for clustering purposes.

23. How does K-Means handle non-linear separability?

Answer: K-Means may struggle with non-linearly separable data, as it assumes a linear relationship between points and centroids.

24. What are the differences between K-Means and hierarchical clustering?

Answer: K-Means is a partitional method that requires specifying K, while hierarchical clustering does not require K and produces a tree-like structure (dendrogram).

25. What is the computational complexity of the K-Means algorithm?

Answer: The computational complexity is $O(n * K * t)$, where n is the number of data points, K is the number of clusters, and t is the number of iterations.

26. What is the effect of the choice of metric in K-Means clustering?

Answer: K-Means uses Euclidean distance by default, but using a different metric may lead to different cluster assignments and centroid locations.

27. What is the role of the n_init parameter in K-Means?

Answer: The n_init parameter controls the number of times the algorithm is run with different centroid initializations. The best result from these runs is chosen.

28. Can K-Means clustering be used for anomaly detection?

Answer: Yes, K-Means can be used for anomaly detection by identifying points that do not belong to any cluster or are far from any centroid.

29. What is the significance of the fit method in K-Means?

Answer: The fit method trains the model by assigning data points to clusters and updating centroids until convergence.

30. How would you apply K-Means clustering to image compression?

Answer: In image compression, K-Means can be used to quantize the pixel colors by grouping similar colors into clusters, reducing the number of unique colors.

31. What is the effect of scaling on K-Means clustering?

Answer: Feature scaling (normalization or standardization) is important for K-Means because it is sensitive to the scale of data and may give biased results if the features have different scales.

32. What is the significance of KMeans.labels_?

Answer: KMeans.labels_ contains the cluster assignments for each data point after the model is fitted.

33. What does the KMeans.cluster_centers_ attribute represent?

Answer: KMeans.cluster_centers_ contains the coordinates of the centroids of the clusters after the model has converged.

34. In the assignment, which cluster does P6 belong to?

Answer: P6 [0.25, 0.5] belongs to cluster 2, as per the final labels from the K-Means algorithm.

35. What is the population of the cluster around m2 (P8)?

Answer: The population of cluster 2 (around m2/P8) is 4, based on the K-Means results.

36. What is the updated value of m1 and m2 after clustering?

Answer: The updated values of m1 and m2 after the clustering are:
m1: [0.15, 0.725]
m2: [0.25, 0.2]

37. What is the role of KMeans.fit_predict() method?

Answer: fit_predict() fits the K-Means model and simultaneously returns the predicted cluster labels for each data point.

38. What does the KMeans.inertia_ attribute represent?

Answer: KMeans.inertia_ represents the sum of squared distances from each point to its assigned cluster centroid, indicating the compactness of the clusters.

39. What would you do if the K-Means algorithm doesn't converge?
Answer: If K-Means doesn't converge, consider increasing the number of iterations, changing the initialization method, or checking for issues in the data.
40. How does K-Means clustering relate to vector quantization?
Answer: K-Means clustering is essentially a vector quantization technique, where each cluster center represents a quantized vector used to represent the data.
These questions cover various aspects of K-Means clustering, from the basic algorithm to advanced topics like initialization and evaluation techniques.

Here are 40 potential viva questions related to the **Gradient Boost Classifier** assignment, with answers:

### 1. **What is a Gradient Boosting Classifier?**
   - **Answer**: Gradient Boosting Classifier is an ensemble machine learning algorithm that combines multiple weak learners (typically decision trees) to form a strong learner. It works by iteratively fitting trees to residuals (errors) of the previous trees.

### 2. **What is the purpose of using Gradient Boosting?**
   - **Answer**: The purpose is to create a predictive model by combining multiple weak models (usually decision trees) into one strong model, improving prediction accuracy.

### 3. **What are the key parameters of a GradientBoostingClassifier?**
   - **Answer**: Key parameters include `learning_rate`, `n_estimators`, `max_depth`, `subsample`, and `min_samples_split`.

### 4. **What is the role of the `learning_rate` parameter in Gradient Boosting?**
   - **Answer**: The learning rate controls the contribution of each tree to the final model. A lower learning rate requires more trees to achieve a similar result but can lead to a more accurate model.

### 5. **How does Gradient Boosting handle overfitting?**
   - **Answer**: Gradient Boosting mitigates overfitting through early stopping, limiting the depth of trees, and using regularization techniques such as shrinkage (learning rate).

### 6. **What is the `n_estimators` parameter in Gradient Boosting?**
   - **Answer**: It specifies the number of boosting stages (trees) to be added. More trees can lead to a more complex model but also may increase the risk of overfitting.

### 7. **What is cross-validation, and why is it used in the Gradient Boosting implementation?**
   - **Answer**: Cross-validation is a technique for assessing the performance of a model by splitting the dataset into multiple parts and training and testing on different splits. It helps prevent overfitting and gives a better estimate of the model's performance.

### 8. **Why is the `train_test_split` method used in the code?**
   - **Answer**: The `train_test_split` method splits the dataset into training and testing sets, enabling the model to be trained on one subset of data and evaluated on another.

### 9. **How do you handle missing data in the dataset?**
   - **Answer**: Missing data can be handled by either removing missing rows or filling them with the mean, median, or mode values, or by using algorithms that can handle missing data.

### 10. **What is Label Encoding, and why is it used in this assignment?**
   - **Answer**: Label Encoding converts categorical data into numeric labels. It is used here to encode categorical features, which allows them to be used in machine learning models.

### 11. **What is the difference between training accuracy and test accuracy?**
   - **Answer**: Training accuracy measures how well the model performs on the training data, while test accuracy measures how well it generalizes to unseen data.

### 12. **What does the `.score()` method return?**
   - **Answer**: The `.score()` method returns the accuracy of the model, which is the percentage of correctly predicted instances in the test dataset.

### 13. **What is the significance of using `GridSearchCV`?**
   - **Answer**: `GridSearchCV` performs an exhaustive search over a specified parameter grid, helping find the best combination of hyperparameters to optimize model performance.

### 14. **Explain the purpose of the `cv` parameter in `GridSearchCV`.**
   - **Answer**: The `cv` parameter specifies the number of cross-validation folds to use. It ensures that the hyperparameters are evaluated more robustly.

### 15. **What is the `learning_rate` and how does it affect the model's performance?**
   - **Answer**: The `learning_rate` controls how much each tree contributes to the final model. A smaller value leads to slower convergence, but can potentially increase model performance.

### 16. **Why is feature scaling important in machine learning?**
   - **Answer**: Feature scaling ensures that all features contribute equally to the model by normalizing their values, which is especially important for models that rely on distance or gradient-based optimization.

### 17. **How does Gradient Boosting improve over individual decision trees?**
   - **Answer**: Gradient Boosting improves by combining weak models (trees) into a strong one, focusing on correcting errors from previous trees, leading to better overall performance.

### 18. **What is the `max_depth` parameter, and how does it affect the model?**
   - **Answer**: The `max_depth` parameter controls the maximum depth of each tree.

Limiting the depth prevents overfitting by keeping trees from becoming too complex.

### 19. **What is the meaning of `n_estimators=100` in the code?**
   - **Answer**: It means that 100 boosting stages (trees) will be used to train the model.

### 20. **What do you understand by ensemble methods?**
   - **Answer**: Ensemble methods combine multiple models to improve performance. Gradient Boosting is an ensemble method where weak learners (trees) are combined to create a strong learner.

### 21. **How do you evaluate the performance of a classification model?**
   - **Answer**: Performance can be evaluated using metrics such as accuracy, precision, recall, F1-score, confusion matrix, and ROC-AUC.

### 22. **What is a confusion matrix?**
   - **Answer**: A confusion matrix is a table used to describe the performance of a classification model by showing the true positives, true negatives, false positives, and false negatives.

### 23. **What is overfitting, and how can it be avoided in Gradient Boosting?**
   - **Answer**: Overfitting occurs when the model learns noise or random fluctuations in the training data. It can be avoided by using regularization, tuning hyperparameters, and performing cross-validation.

### 24. **What does the term `shrinkage` mean in Gradient Boosting?**
   - **Answer**: Shrinkage refers to scaling the contribution of each tree by the learning rate to prevent overfitting and allow for a better generalization.

### 25. **What is the difference between boosting and bagging?**
   - **Answer**: Boosting sequentially trains weak models, correcting errors from previous models, while bagging trains models in parallel on different subsets of the data.

### 26. **What are the benefits of using Gradient Boosting?**
   - **Answer**: Gradient Boosting can achieve high accuracy, handle various types of data, and provide a reliable model even with noisy datasets.

### 27. **What is the impact of the `subsample` parameter?**
   - **Answer**: The `subsample` parameter controls the fraction of samples used to fit each tree. Lower values lead to more regularization but can reduce the model's performance.

### 28. **How can we visualize the performance of the Gradient Boosting model?**
   - **Answer**: Performance can be visualized using confusion matrices, ROC curves, or learning curves.

### 29. **What are residuals in the context of Gradient Boosting?**
   - **Answer**: Residuals are the differences between the true values and the predictions made by the model. In Gradient Boosting, new trees are trained to

predict the residuals of previous trees.

### 30. **What does `GridSearchCV.best_params_` return?**
   - **Answer**: It returns the best combination of hyperparameters found by `GridSearchCV` that maximized model performance.

### 31. **What happens if you use a very high number of trees in Gradient Boosting?**
   - **Answer**: Using too many trees can lead to overfitting, especially if the learning rate is too high or if trees are too deep.

### 32. **What is a hyperparameter, and why is tuning them important?**
   - **Answer**: Hyperparameters are parameters set before training a model, like learning rate and number of trees. Tuning them helps optimize model performance.

### 33. **Explain the role of `min_samples_split` in controlling the complexity of the model.**
   - **Answer**: The `min_samples_split` parameter controls the minimum number of samples required to split an internal node. Higher values prevent creating overly complex trees.

### 34. **Why did you choose a `test_size=0.2` for splitting the dataset?**
   - **Answer**: A test size of 0.2 ensures that 20% of the data is used for testing while the rest is used for training, providing a good balance between training and evaluation data.

### 35. **What is the impact of the `n_estimators` hyperparameter in tuning the model?**
   - **Answer**: Increasing the number of estimators (trees) can improve model accuracy but might also increase overfitting if not controlled by parameters like `learning_rate`.

### 36. **What are the advantages of using Gradient Boosting over Random Forest?**
   - **Answer**: Gradient Boosting typically performs better on tasks with complex relationships between features, whereas Random Forest is generally better at handling large, noisy datasets.

### 37. **What is the main objective of the `train_test_split` function?**
   - **Answer**: The `train_test_split` function splits the data into separate training and testing sets to evaluate the model's ability to generalize to new, unseen data.

### 38. **What are out

liers, and how does Gradient Boosting handle them?**
   - **Answer**: Outliers are extreme values that differ significantly from most of the data. Gradient Boosting can be less sensitive to outliers due to its iterative nature.

### 39. **What does boosting aim to minimize in Gradient Boosting?**

- **Answer**: Boosting aims to minimize the loss function, typically mean squared error for regression tasks and log-loss for classification tasks.

### 40. **Why is the evaluation of model performance important in machine learning?**
   - **Answer**: Evaluating model performance is crucial for understanding its effectiveness, ensuring it generalizes well to new data, and selecting the best model for deployment.