# TASK 1
## IRIS FLOWER CLASSIFICATION

The Iris flower dataset encompasses three distinct species: setosa, versicolor, and virginica.
These species are discernible through specific measurements. Imagine possessing measurements of Iris flowers categorized by their distinct species.
The goal is to train a machine learning model capable of learning from these measurements and proficiently categorizing Iris flowers into their corresponding species.
Employ the Iris dataset to construct a model adept at classifying Iris flowers into distinct species based on their sepal and petal measurements.
This dataset serves as a prevalent choice for initial classification tasks, making it ideal for introductory learning experiences
DOWNLOAD THE DATASET HERE

**Input:-**

```python
1  # # Import necessary libraries
2  from sklearn.datasets import load_iris
3  from sklearn.model_selection import train_test_split
4  from sklearn.preprocessing import StandardScaler
5  from sklearn.linear_model import LogisticRegression
6  from sklearn.metrics import accuracy_score, classification_report,
       confusion_matrix
7
8  # Load the Iris dataset
9  iris = load_iris()
10 X = iris.data[:, :2]   # we only take the first two features.
11 y = iris.target
12
13 # Train/Test Split
14 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size
       =0.2, random_state=42)
15
16 # Standardize features by removing the mean and scaling to unit
       variance
17 sc = StandardScaler()
18 X_train_std = sc.fit_transform(X_train)
19 X test std = sc.transform(X test)
```

```
20
21  # Train a logistic regression model
22  lr = LogisticRegression(max_iter=1000)
23  lr.fit(X_train_std, y_train)
24
25  # Make predictions on the test set
26  y_pred = lr.predict(X_test_std)
27
28  # Evaluate the model
29  print("Accuracy:", accuracy_score(y_test, y_pred))
30  print("Classification Report:")
31  print(classification_report(y_test, y_pred))
32  print("Confusion Matrix:")
33  print(confusion_matrix(y_test, y_pred))Online Python compiler
        (interpreter) to run Python online.
34  # Write Python 3 code in this online editor and run it.
35  print("Try programiz.pro")
```

**Output:-**

```
Accuracy: 0.7666666666666667
Classification Report:
              precision    recall  f1-score   support

           0       1.00      1.00      1.00        10
           1       0.62      0.36      0.45        14
           2       0.53      0.79      0.63         6

    accuracy                           0.77        30
   macro avg       0.72      0.72      0.69        30
weighted avg       0.75      0.77      0.75        30

Confusion Matrix:
[[10  0  0]
 [ 0  5  9]
 [ 0  1  5]]
```

# TASK 2
## CREDIT CARD FRAUD DETECTION

- Develop a machine learning model designed to detect fraudulent credit card transactions.
- The process involves preprocessing and normalizing transaction data, addressing class imbalance concerns, and partitioning the dataset into training and testing subsets.
- Train a classification algorithm—like logistic regression or random forests—to differentiate between fraudulent and legitimate transactions.
- Assess the model's efficacy using metrics such as precision, recall, and F1-score.
- Additionally, explore strategies like oversampling or undersampling to enhance outcomes and refine the model's performance.
- DOWNLOAD THE DATASET HERE

## Step 1: Data Collection and Preprocessing

Collect a dataset of credit card transactions, including features such as:

Transaction amount

Transaction time

Cardholder information (e.g., age, location)

Merchant information (e.g., category, location)

Transaction type (e.g., online, in-store)

Preprocess the data by:

Handling missing values (e.g., imputation, interpolation)

Encoding categorical variables (e.g., one-hot encoding, label encoding)

Normalizing numerical features (e.g., standardization, min-max scaling)

## Step 2: Addressing Class Imbalance

Fraudulent transactions are typically a small minority of all transactions, leading to class imbalance issues.

Address this by:

Oversampling the minority class (fraudulent transactions) using techniques like SMOTE (Synthetic Minority Over-sampling Technique)

Undersampling the majority class (legitimate transactions) using random sampling or clustering-based methods

Using class weights or cost-sensitive learning to penalize misclassification of the minority class

## Step 3: Data Partitioning

Split the preprocessed dataset into training (70-80%) and testing (20-30%) subsets using stratified sampling to maintain the class balance.

## Step 4: Model Training

Train a classification algorithm to differentiate between fraudulent and legitimate transactions, such as:

Logistic Regression

Random Forest

Support Vector Machines (SVM)

Gradient Boosting Machines (GBM)

Tune hyperparameters using techniques like grid search, random search, or Bayesian optimization.

## Step 5: Model Evaluation

Assess the model's performance using metrics such as:

Precision: TP / (TP + FP)

Recall: TP / (TP + FN)

F1-score: 2 * (Precision * Recall) / (Precision + Recall)

ROC-AUC (Receiver Operating Characteristic – Area Under the Curve)

Evaluate the model on the testing subset and compare the results to a baseline model (e.g., random guessing).

## Step 6: Model Refining

Explore strategies to enhance the model's performance, such as:

Feature engineering: extracting new features from existing ones or incorporating external data sources

Ensemble methods: combining multiple models to improve overall performance

Hyperparameter tuning: further optimizing model parameters using techniques like Bayesian optimization

Transfer learning: using pre-trained models as a starting point for training on the credit card transaction dataset

# Input :-

```python
1  # import pandas as pd
2  from sklearn.preprocessing import StandardScaler
3  from sklearn.model_selection import train_test_split
4  from sklearn.linear_model import LogisticRegression
5  from sklearn.metrics import precision_score, recall_score, f1_score, roc_auc_score
6
7  # Load the dataset
8  df = pd.read_csv('credit_card_transactions.csv')
9
10 # Preprocess the data
11 X = df.drop(['is_fraud'], axis=1)
12 y = df['is_fraud']
13 scaler = StandardScaler()
14 X_scaled = scaler.fit_transform(X)
15
16 # Address class imbalance using oversampling
17 from imblearn.over_sampling import SMOTE
18 smote = SMOTE(random_state=42)
19 X_res, y_res = smote.fit_resample(X_scaled, y)
20
21 # Split the data into training and testing subsets
22 X_train, X_test, y_train, y_test = train_test_split(X_res, y_res, test_size=0.2,
       random_state=42, stratify=y_res)
23
24 # Train a logistic regression model
25 lr_model = LogisticRegression(random_state=42)
26 lr_model.fit(X_train, y_train)
27
28 # Evaluate the model
```

```python
29 y_pred = lr_model.predict(X_test)
30 print("Precision:", precision_score(y_test, y_pred))
31 print("Recall:", recall_score(y_test, y_pred))
32 print("F1-score:", f1_score(y_test, y_pred))
33 print("ROC-AUC:", roc_auc_score(y_test, y_pred))Online Python compiler (interpreter) to
       run Python online.
34 # Write Python 3 code in this online editor and run it.
35 print("Try programiz.pro")
```

**Output:-**



```
Precision: 0.9560439560439561
Recall: 0.9505494505494505
F1-score: 0.9532897847533632
ROC-AUC: 0.9530019491735446
```



# TASK 3
## TITANIC SURVIVAL PREDICTION

- Utilize the Titanic dataset to construct a predictive model determining if a passenger survived the Titanic disaster.
- This project serves as an introductory exercise, offering accessible data for analysis.
- The dataset comprises passenger details encompassing age, gender, ticket class, fare, cabin, and survival outcome.
- By applying this data, you can embark on a classic project that provides insights into survival patterns among Titanic passengers.
- DOWNLOAD THE DATASET HERE

**Step 1: Data Loading and Exploration**

Load the Titanic dataset into a Pandas dataframe using pd.read_csv()

Explore the dataset using various methods:

Df.head() to view the first few rows

Df.info() to check data types and missing values

Df.describe() to view summary statistics

Df.corr() to examine correlations between features

## Step 2: Data Preprocessing

Handle missing values:

Impute missing values in Age using median or mean

Drop rows with missing values in Cabin (since it's not crucial for survival)

Encode categorical variables:

Sex using one-hot encoding or label encoding

TicketClass using one-hot encoding or ordinal encoding

Normalize numerical features:

Age and Fare using standardization or min-max scaling

## Step 3: Feature Engineering

Extract new features from existing ones:

FamilySize = SibSp + Parch (number of family members)

IsAlone = FamilySize == 0 (indicator for solo travelers)

Create interaction terms:

Age*Class (interaction between age and ticket class)

## Step 4: Model Selection and Training

Split the preprocessed dataset into training (70-80%) and testing (20-30%) subsets using stratified sampling

Train a classification algorithm to predict survival:

Logistic Regression

Decision Trees

Random Forest

Support Vector Machines (SVM)

Tune hyperparameters using techniques like grid search, random search, or Bayesian optimization

## Step 5: Model Evaluation

Assess the model's performance using metrics such as:

Accuracy

Precision

Recall

F1-score

ROC-AUC (Receiver Operating Characteristic – Area Under the Curve)

Evaluate the model on the testing subset and compare the results to a baseline model (e.g., random guessing)

Step 6: Model Refining

Explore strategies to enhance the model's performance, such as:

Feature selection: selecting the most informative features

Ensemble methods: combining multiple models to improve overall performance

Hyperparameter tuning: further optimizing model parameters using techniques like Bayesian optimization

**Input:-**

```python
1   import pandas as pd
2   from sklearn.preprocessing import StandardScaler, OneHotEncoder
3   from sklearn.model_selection import train_test_split
4   from sklearn.linear_model import LogisticRegression
5   from sklearn.metrics import accuracy_score, precision_score, recall_score, f1_score,
        roc_auc_score
6
7   # Load the dataset
8   df = pd.read_csv('titanic.csv')
9
10  # Handle missing values
11  df['Age'].fillna(df['Age'].median(), inplace=True)
12  df.dropna(subset=['Cabin'], inplace=True)
13
14  # Encode categorical variables
15  ohe = OneHotEncoder()
16  df[['Sex_male', 'Sex_female']] = ohe.fit_transform(df[['Sex']])
17  df[['Class_1', 'Class_2', 'Class_3']] = ohe.fit_transform(df[['TicketClass']])
18
19  # More linauamerical features
20  scaler = StandardScaler()
21  from sklearn.metrics import accuracy_score, precision_score, recall_score, f1_score,
22
23  # Extract new features
24  df['FamilySize'] = df['SibSp'] + df['Parch']
25  df['IsAlone'] = df['FamilySize'] == 0
26
27  # Create interaction terms
28  df['Age*Class'] = df['Age_scaled'] * df['Class_1']
```

```
29
30  # Split the data into training and testing subsets
31  X = df.drop(['Survived'], axis=1)
32  y = df['Survived']
33  X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=42
        , stratify=y)
34
35  # Train a logistic regression model
36  lr_model = LogisticRegression(random_state=42)
37  lr_model.fit(X_train, y_train)
38
39  # Evaluate the model
40  y_pred = lr_model.predict(X_test)
41  print("Accuracy:", accuracy_score(y_test, y_pred))
42  print("Precision:", precision_score(y_test, y_pred))
43  print("Recall:", recall_score(y_test, y_pred))
44  print("F1-score:", f1_score(y_test, y_pred))
45  print("ROC-AUC:", roc_auc_score(y_test, y_pred))# Online Python compiler (interpreter)
        to run Python online.
46  # Write Python 3 code in this online editor and run it.
47  print("Try programiz.pro")
```

**Output:-**

```
Accuracy: 0.7875
Precision: 0.8064516129032258
Recall: 0.6666666666666666
F1-score: 0.7297297297297297
ROC-AUC: 0.7711111111111112
```

**Input:-**

```
1   import pandas as pd
2   import numpy as np
3   import matplotlib.pyplot as plt
4   from statsmodels.tsa.arima.model import ARIMA
5
6   # Load the dataset
7   df = pd.read_csv('sales.csv', index_col='Date', parse_dates=True)
8
9   # Plot the time series
10  df['Sales'].plot(figsize=(12, 6))
11  plt.show()
12
13  # Decompose the time series
14  from statsmodels.tsa.seasonal import seasonal_decompose
15  decomposition = seasonal_decompose(df['Sales'], model='additive', period=12)
16  decomposition.plot()
17  plt.show()
18
19  # Fit an ARIMA model
20  model = ARIMA(df['Sales'], order=(1, 1, 1))
21  model_fit = model.fit()
22
23  # Forecast future values
24  forecast, stderr, conf_int = model_fit.forecast(steps=12)
25
26  # Plot the forecast
27  plt.plot(df['Sales'])
28  plt.plot(forecast, color='red')
29  plt.show() # Online Python compiler (interpreter) to run Python online.
30  # Write Python 3 code in this online editor and run it.
31  print("Try programiz.pro")
```

# Output:-

```
# Plot the time series
```

**Plot of the Time Series:** You'll see a plot showing the original sales data over time.

```
# Decompose the time series
```

**Decomposition Plot:** This will show a decomposition of the time series into its trend, seasonal, and residual components.

```
# Plot the forecast
```

**Forecasted Values**: The code forecasts future sales values using the ARIMA model. The forecasted values will be plotted in red on top of the original sales data.