In this lab, you will explore basic Python concepts and machine learning using Scikit-learn. You will learn how to use NumPy for numerical operations, Pandas for data manipulation, and Matplotlib for plotting. You will also learn to create functions, use global and local variables, and apply control flow structures such as `if-else` and loops. Finally, we will train a simple machine learning model using Scikit-learn.

After completing this lab, you will be able to:

1.  Understand Python basics, including NumPy, Pandas, and plotting.
2.  Create functions and work with global and local variables.
3.  Use flow control structures like `if-else` and `for` loops.
4.  Use Google Colab to run Python scripts.
5.  Train a basic machine learning model using Scikit-learn.


To get started, make sure you:

-   Have a Google account to access Google Colab.

You can open Google Colab by visiting https://colab.research.google.com/ .


# Task 1: Introduction to Python Basics

### NumPy Basics:

NumPy is a powerful library for numerical operations in Python.

1.  Import NumPy and create an array.
2.  Perform operations on arrays.

```
1. import numpy as np
2.
3. # Create a NumPy array
4. arr = np.array([1, 2, 3, 4, 5])
5.
6. # Access elements of the array
7. print(arr[0])  # First element
8. print(arr[1:3])  # Elements from index 1 to 2
9.
```

```
1. # Array operations
2. arr = arr + 5
3. print(arr)
4.
5. # Create a 2D array
6. arr_2d = np.array([[1, 2, 3], [4, 5, 6]])
7. print(arr_2d)
8.
```

### Pandas Basics:

Pandas is a powerful library for data manipulation.

1.  Import Pandas and create a DataFrame.

2. Perform basic data manipulation.

```python
1. import pandas as pd
2.
3. # Create a DataFrame
4. data = {'Name': ['Alice', 'Bob', 'Charlie'], 'Age': [25, 30, 35]}
5. df = pd.DataFrame(data)
6.
7. # Access elements in DataFrame
8. print(df.head())  # First few rows
9. print(df['Name'])  # Accessing a specific column
10.
```

```python
1. # Add a new column
2. df['City'] = ['New York', 'Los Angeles', 'Chicago']
3. print(df)
4.
5. # Filtering data
6. df_filtered = df[df['Age'] > 28]
7. print(df_filtered)
8.
```

## Creating Functions:

Functions in Python allow you to organize your code into reusable blocks.

Define a simple function.

```python
1. # Create a function that adds two numbers
2. def add_numbers(a, b):
3.     return a + b
4.
5. # Call the function
6. result = add_numbers(10, 5)
7. print(result)
8.
```

## Global and Local Variables

Understanding variable scope is crucial in programming.

Define global and local variables.

```python
1. # Global variable
2. x = 10
3.
4. def change_x():
5.     # Local variable (within function)
6.     x = 5
7.     print(f"Inside function: x = {x}")
8.
9. # Call the function and print global variable
10. change_x()
11. print(f"Outside function: x = {x}")
12.
```

## Flow Control (If-Else Statements):

Flow control allows you to execute different code based on conditions.

Use `if-else` statements.

```
1. age = 20
2.
3. # If-else statement
4. if age >= 18:
5.     print("You are an adult.")
6. else:
7.     print("You are a minor.")
8.
```

**Loops (For Loops):**

Loops help you iterate over a sequence of values.

Use a `for` loop to iterate over a list.

```
1. # List of names
2. names = ['Alice', 'Bob', 'Charlie']
3.
4. # For loop
5. for name in names:
6.     print(name)
7.
```

**Plotting with Matplotlib:**

Matplotlib is a popular library for creating visualizations.

Plot a simple graph.

```
 1. import matplotlib.pyplot as plt
 2.
 3. # Create data
 4. x = [1, 2, 3, 4, 5]
 5. y = [10, 20, 25, 30, 35]
 6.
 7. # Create a plot
 8. plt.plot(x, y)
 9. plt.title("Simple Line Plot")
10. plt.xlabel("X-axis")
11. plt.ylabel("Y-axis")
12. plt.show()
13.
```

# Task 2: Introduction to Scikit-learn

Scikit-learn is a popular library for machine learning in Python. You will use it to load datasets, split data, and train machine learning models. First, Install Scikit-learn (available in Colab by default, but you can install it if needed).

**Key Components of Scikit-learn**

1. **Datasets**: Scikit-learn provides easy access to several standard datasets, including the Iris dataset, Boston housing data, and handwritten digits. We will use the Iris dataset, a simple dataset of 150 iris flowers, classified into three species.

Example: Loading the Iris dataset.

```
1. from sklearn.datasets import load_iris
2. iris = load_iris()
3. X = iris.data
4. y = iris.target
5.
```

2.  **Model Selection**: Scikit-learn offers functions to split data into training and testing sets, and it provides tools like cross-validation to validate models effectively.

```
1. from sklearn.model_selection import train_test_split
2.
3. # Split dataset into training and testing sets
4. X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.3, random_state=42)
5.
```

3.  **Preprocessing**: Preprocessing ensures your data is in the right format for machine learning models. Common preprocessing tasks include scaling data, normalizing features, and encoding categorical variables.

    - **Standardization**: Scaling features to have zero mean and unit variance.
    - **Encoding categorical features**: Convert text labels into numerical labels.

```
1. from sklearn.preprocessing import StandardScaler
2.
3. scaler = StandardScaler()
4. X_train_scaled = scaler.fit_transform(X_train)
5. X_test_scaled = scaler.transform(X_test)
6.
```

```
1. from sklearn.preprocessing import LabelEncoder
2.
3. encoder = LabelEncoder()
4. y_train_encoded = encoder.fit_transform(y_train)
5.
```

4.  **Model Training and Evaluation**: Scikit-learn provides many algorithms for classification, regression, and clustering. These algorithms follow a simple workflow:

    - **Fit the model**: Train the model on the training data.
    - **Predict**: Make predictions on new data.
    - **Evaluate**: Measure the model's performance using metrics like accuracy for classification or mean squared error for regression.

**Example: Classification with Scikit-learn**

We will train an SVM classifier on the Iris dataset using the following steps.

1.  Load and split the data.
2.  Preprocess the data.
3.  Train the model.
4.  Evaluate the model.

```
1. from sklearn.model_selection import train_test_split
2. from sklearn.datasets import load_iris
3.
4. iris = load_iris()
5. X = iris.data
6. y = iris.target
7.
8. # Split dataset
9. X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.3,
random_state=42)
```

```
1. from sklearn.preprocessing import StandardScaler
2.
3. scaler = StandardScaler()
4. X_train_scaled = scaler.fit_transform(X_train)
5. X_test_scaled = scaler.transform(X_test)
6.
```

```
1. from sklearn.svm import SVC
2.
3. svm_model = SVC()
4. svm_model.fit(X_train_scaled, y_train)
5.
```

```
1. from sklearn.metrics import accuracy_score
2.
3. y_pred = svm_model.predict(X_test_scaled)
4. accuracy = accuracy_score(y_test, y_pred)
5. print(f"Model Accuracy: {accuracy * 100:.2f}%")
6.
```

## Cross-Validation and Model Tuning

Use **Cross-validation** and **GridSearchCV** to optimize model performance.

```
1. from sklearn.model_selection import cross_val_score, GridSearchCV
2.
3. # Perform 5-fold cross-validation
4. cv_scores = cross_val_score(SVC(), X_train_scaled, y_train, cv=5)
5. print(f"Cross-Validation Accuracy: {cv_scores.mean() * 100:.2f}%")
6.
7. # Grid search for hyperparameter tuning
8. param_grid = {'C': [0.1, 1, 10], 'kernel': ['linear', 'rbf']}
9. grid_search = GridSearchCV(SVC(), param_grid, cv=5)
10. grid_search.fit(X_train_scaled, y_train)
11. best_model = grid_search.best_estimator_
12.
```

## Pipelines for Streamlined Workflow

Use **Pipelines** to streamline preprocessing and training.

```
1. from sklearn.pipeline import Pipeline
2.
3. pipeline = Pipeline([
4.     ('scaler', StandardScaler()),
5.     ('svm', SVC()) ])
7.
```

```
 8.  pipeline.fit(X_train, y_train)
 9.  y_pred = pipeline.predict(X_test)
10.  print(f"Accuracy: {accuracy_score(y_test, y_pred) * 100:.2f}%")
11.
```

## Practice Tasks

To reinforce your learning, complete the following practice tasks. These tasks build on the concepts you've learned today.

1.  Array Operations (NumPy):
    I.  Create a 3x3 matrix using NumPy, and perform addition, subtraction, and element-wise multiplication with another matrix of the same size.
    II.  Extract the second row and the last column of the matrix using slicing.

2.  DataFrame Manipulation (Pandas):
    I.  Create a DataFrame with three columns: 'Name', 'Age', and 'Score'. Add a new column that categorizes the score into "Pass" or "Fail" based on a threshold.
    II.  Filter the DataFrame to show only the rows where the score is above 70.

3.  Function Creation:
    Write a Python function that calculates the factorial of a given number.

4.  Machine Learning Challenge:
    I.  Train a Decision Tree model on the Iris dataset instead of the SVM. Use *sklearn.tree.DecisionTreeClassifier* and compare its accuracy to the SVM model.
    II.  Try visualizing the decision tree using Matplotlib and the *plot_tree* function from Scikit-learn.