phase 4 project

BY:

 NAME : H. Kishore kannan

NM ID : au922121106034

REGISTER NO: 922121106034

# 1. Exploratory Data Analysis (EDA):

**EDA is a crucial first step to understand your data. You can use Python libraries like Pandas, Matplotlib, and Seaborn to perform the following tasks:**

  **- Load your dataset.**

  **- Examine basic statistics like mean, median, standard deviation, etc.**

  **- Visualize data distributions, relationships, and outliers using histograms, scatter plots, and box plots.**

  **- Identify missing data and decide on handling strategies (imputation or removal).**

  **- Perform correlation analysis to understand feature relationships.**

# 2. Feature Engineering:

Feature engineering involves creating new features or transforming existing ones to improve model performance. Some common techniques include:

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

- Scaling and normalizing numerical features.

- Creating interaction features, aggregations, or statistical features.

- Handling time-related data if applicable (e.g., extracting day of the week, month, etc.).

- Feature selection to choose the most relevant features.

# 3. Predictive Modeling:

Building predictive models depends on the nature of your problem (classification, regression, etc.). You can use machine learning libraries like scikit-learn or deep learning frameworks like TensorFlow
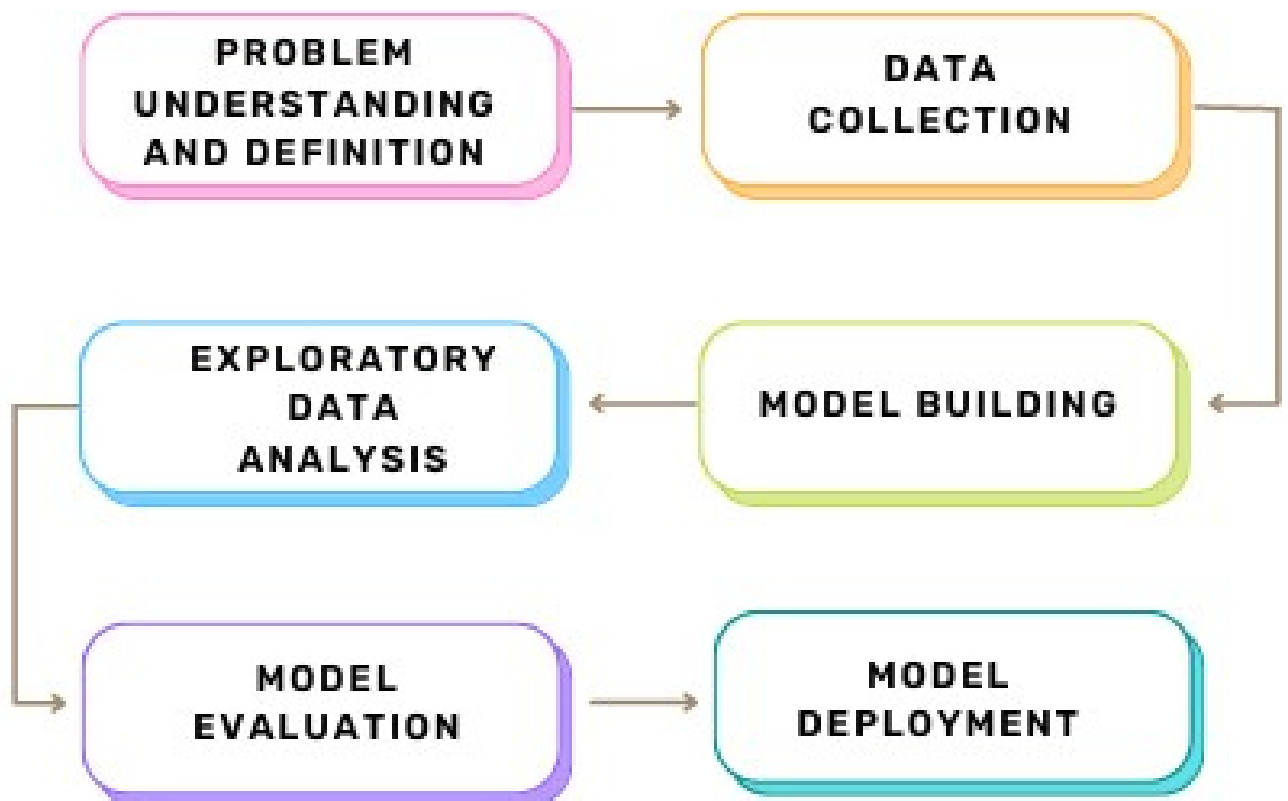
or PyTorch. Here are the steps to follow:

   - Split your dataset into training and testing sets for model evaluation.

   - Select appropriate algorithms (e.g., linear regression, decision trees, neural networks) and train them on the training data.

   - Tune hyperparameters to optimize model performance using techniques like grid search or random search.

   - Evaluate models using appropriate metrics (e.g., accuracy, mean squared error) and choose the best-performing one.

   - Validate the model on the testing data to assess its generalization performance.

   - Interpret the model results to gain insights into the problem.

Remember to iterate on these steps, refine your model, and potentially consider more advanced techniques like cross-validation, ensemble methods, and deep learning architectures if needed. The effectiveness of your project greatly depends on the quality of your EDA and feature engineering, so

invest time in those stages●

```
┌────────────────────┐        ┌────────────────────┐
│     PROBLEM        │        │      DATA          │
│  UNDERSTANDING     │  ───>  │   COLLECTION       │
│  AND DEFINITION    │        │                    │
└────────────────────┘        └────────────────────┘
                                        │
┌────────────────────┐        ┌────────────────────┐
│   EXPLORATORY      │        │                    │
│      DATA          │  <───  │  MODEL BUILDING    │ <──
│    ANALYSIS        │        │                    │
└────────────────────┘        └────────────────────┘
   │
┌────────────────────┐        ┌────────────────────┐
│     MODEL          │        │     MODEL          │
│   EVALUATION       │  ───>  │  DEPLOYMENT        │
└────────────────────┘        └────────────────────┘
```

# EXAMPLES WITH PROGRAM:..

Certainly, I can provide some code examples for each stage of building your AI-driven exploration and prediction project:

**1. Exploratory Data Analysis (EDA):**

Here's an example using Python and the Pandas library for EDA:

```python
import pandas as pd

import matplotlib.pyplot as plt

import seaborn as sns


# Load your dataset

data = pd.read_csv('your_dataset.csv')


# Basic statistics

print(data.describe())


# Data visualization

plt.figure(figsize=(10, 6))

sns.histplot(data['feature1'], bins=20)

plt.title('Distribution of Feature 1')

plt.show()


# Identify missing data
```

```
missing_data = data.isnull().sum()

print(missing_data)


# Correlation analysis

correlation_matrix = data.corr()

sns.heatmap(correlation_matrix, annot=True)

plt.show()
```

**2. Feature Engineering:**

Feature engineering depends on your dataset and problem. Here's a general example:

```python
# Encoding categorical variables

data = pd.get_dummies(data, columns=['categorical_feature'])


# Scaling numerical features

from sklearn.preprocessing import StandardScaler

scaler = StandardScaler()

data['numerical_feature'] = scaler.fit_transform(data['numerical_feature'])


# Creating interaction features

data['interaction_feature'] = data['feature1'] * data['feature2']


# Feature selection

from sklearn.feature_selection import SelectKBest, f_regression

X = data.drop('target', axis=1)

y = data['target']
```

```
X_new = SelectKBest(f_regression, k=5).fit_transform(X, y)
```

**3. Predictive Modeling:**

Let's use scikit-learn to create a simple linear regression model as an example:

```python
from sklearn.model_selection import train_test_split

from sklearn.linear_model import LinearRegression

from sklearn.metrics import mean_squared_error


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


# Create and train the model
model = LinearRegression()

model.fit(X_train, y_train)


# Make predictions
y_pred = model.predict(X_test)


# Evaluate the model
mse = mean_squared_error(y_test, y_pred)

print(f"Mean Squared Error: {mse}")
```

These are simplified examples. In a real project, you would need to adapt the code to your specific dataset and problem. Additionally, you may explore more advanced models and techniques based on the characteristics of your data and the performance of your initial models.