Submitted by:

Velpucharla Venkata subba reddy

Au72321243054

## Documentation:

**The Design Thinking process is a human-centered approach to problem-solving and innovation. It typically consists of several phases, including problem definition, ideation, prototyping, testing, and implementation. Here, I'll outline the problem statement phase and the subsequent phases of development in the Design Thinking process**.

**Usage**

**Explain how to use your code to perform data preprocessing, exploratory data analysis (EDA), and predictive modeling.**

**Data Preprocessing and EDA**

**Describe the steps to preprocess the data and conduct EDA:**

**Load the Dataset: Replace 'your_dataset.csv' with the path to your dataset file.**

**Data Preprocessing: Explain how to handle missing values and encode categorical variables.**

**Exploratory Data Analysis (EDA): Provide examples of summary statistics and data visualizations.**

**Predictive Modeling**

**Explain how to use the code for predictive modeling:**

**Feature Scaling: Describe the process of feature scaling if required.**

**Initialize and Train the Model: Provide code to initialize and train the predictive model.**

**Make Predictions: Explain how to make predictions on new data.**

**Model Evaluation: Detail the steps for evaluating model performance, including accuracy and classification reports.**

**Contributing**

**Explain how others can contribute to your project, if applicable. Include guidelines for code contributions, bug reports, or feature requests.**

**License**

**Include the license information for your project, if applicable. Specify the open-source license under which your code is released.**

Mention any visualization tools or techniques used to understand the model's behavior.

## 5.Model Testing:

Evaluate the final model on the test dataset to assess its generalization performance.

**6.Model Deployment:**

Explain how the model is deployed in a real-world setting, if applicable. Is it a web application, an API, or something else?

**7.Ethical Considerations:**

Discuss any ethical concerns related to the dataset, model, or predictions, such as privacy, bias, and fairness.

**8.Conclusion:**

Summarize the results and insights obtained from the AI algorithms applied.

Discuss limitations, future work, and potential improvements.

Please provide more specific information about the dataset and problem you have in mind for a more detailed explanation.

**1.** Data Distribution and Summary Statistics: EDA helps you understand the central tendencies and dispersions of the data. It includes statistics like mean, median, variance, and

quartiles. Understanding the distribution of the data can influence the choice of appropriate modeling techniques.

2. Data Visualization: Visualization tools like histograms, box plots, and scatter plots reveal patterns, outliers, and trends in the data. These visualizations can uncover relationships between variables and assist in feature selection and engineering.

3. Data Quality Assessment: EDA exposes data quality issues, such as missing values, duplicate records, and outliers. Dealing with these issues is essential to ensure model accuracy and reliability.

4. Correlations: EDA can reveal the relationships between features, allowing you to identify which features are strongly correlated or exhibit multicollinearity. Understanding these relationships informs feature selection and model building.

5. Categorical Variables: For categorical variables, EDA helps you understand the distribution of categories and the potential impact on the target variable. This information aids in feature engineering and encoding.

6. Outliers: Identifying outliers through EDA is essential, as they can have a significant impact on model performance. Decisions about how to handle outliers depend on their nature and the specific problem.

7. Impact on Performance Predictive Models:

8. Feature Selection and Engineering: Insights from EDA inform decisions about which features to include in the model. You may choose to drop irrelevant or highly correlated features and create new features that capture important patterns.

9. Data Preprocessing: EDA helps you determine how to handle missing values, outliers, and data scaling. Decisions made during preprocessing greatly influence model performance.

10. Model Selection: Understanding the data distribution and relationships between variables guides the choice of an appropriate model. For example, linear regression may be suitable for linear relationships, while decision trees or neural networks may be chosen for more complex patterns.

11. Hyperparameter Tuning: EDA insights can impact hyperparameter tuning. For example, knowing that a dataset

has imbalanced classes may lead to adjusting class weights or using different evaluation metrics.

**12.** Model Evaluation: EDA informs the selection of appropriate evaluation metrics based on the nature of the problem and the dataset. For classification tasks, metrics like accuracy, precision, recall, and F1-score may be selected based on the distribution of classes and the costs associated with false positives and false negatives.

**13.** Overfitting Mitigation: Insights from EDA can help identify potential overfitting issues. Adjusting model complexity, regularization techniques, and cross-validation strategies can be influenced by EDA findings.

# Submission:

**# Import necessary libraries**

**import pandas as pd**

**import numpy as np**

**import matplotlib.pyplot as plt**

**import seaborn as sns**

**from sklearn.model_selection import train_test_split**

**from sklearn.preprocessing import StandardScaler**

```python
from sklearn.ensemble import RandomForestClassifier

from sklearn.metrics import accuracy_score, classification_report

import xgboost as xgb


# Load the dataset

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


# Data Preprocessing

# Handle missing values

data.dropna(inplace=True)


# Encode categorical variables

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


# Split the dataset into features and target

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

y = data['target']


# Split the dataset into training and testing sets

X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2,
random_state=42)


# Exploratory Data Analysis (EDA)
```

```python
# Summary statistics
print(data.describe())


# Data visualization
plt.figure(figsize=(10, 6))
sns.heatmap(data.corr(), annot=True, cmap='coolwarm')
plt.show()


# Predictive Modeling
# Feature scaling
scaler = StandardScaler()
X_train = scaler.fit_transform(X_train)
X_test = scaler.transform(X_test)
# Initialize and train a predictive model
model = RandomForestClassifier(n_estimators=100,
random_state=42)
model.fit(X_train, y_train)
# Make predictions on the test set
y_pred = model.predict(X_test)
# Model evaluation
accuracy = accuracy_score(y_test, y_pred)
print(f'Accuracy: {accuracy:.2f}')


# Generate a classification report
```

```python
    print(classification_report(y_test, y_pred))


# Optionally, you can also use XGBoost as a different predictive model
# model = xgb.XGBClassifier(n_estimators=100, random_state=42)
# model.fit(X_train, y_train)
# y_pred = model.predict(X_test)
# ...


# Further hyperparameter tuning and cross-validation can be applied to improve model performance.
```


**READ.me File:**

# Project Title

A brief description of your project and its purpose.

## Table of Contents

- [Getting Started](#getting-started)
  - [Prerequisites](#prerequisites)
  - [Installing Dependencies](#installing-dependencies)
- [Usage](#usage)

## Getting Started

Provide instructions for running the code on a local machine or a specific environment. Include information about prerequisites and dependencies.

### Prerequisites

List any software, libraries, or tools that need to be installed before running the code. Include versions if necessary.

- Python (>=3.7)

- Pandas

- Numpy

- Matplotlib

- Seaborn

- Scikit-learn

- XGBoost (if applicable)

### Installing Dependencies

You can install the required dependencies using `pip`. Run the following command:

```bash
pip install -r requirements.txt
```

## Usage

Explain how to use your code to perform data preprocessing, exploratory data analysis (EDA), and predictive modeling.

## Data Preprocessing and EDA

Describe the steps to preprocess the data and conduct EDA:

Load the Dataset: Replace 'your_dataset.csv' with the path to your dataset file.

Data Preprocessing: Explain how to handle missing values and encode categorical variables.

Exploratory Data Analysis (EDA): Provide examples of summary statistics and data visualizations.

## Predictive Modeling

Explain how to use the code for predictive modeling:

Feature Scaling: Describe the process of feature scaling if required.

**Initialize and Train the Model: Provide code to initialize and train the predictive model.**

**Make Predictions: Explain how to make predictions on new data.**

**Model Evaluation: Detail the steps for evaluating model performance, including accuracy and classification reports.**

**Contributing**

**Explain how others can contribute to your project, if applicable. Include guidelines for code contributions, bug reports, or feature requests.**

**License**

**Include the license information for your project, if applicable. Specify the open-source license under which your code is released.**