# PHASE 4 DEVELOPMENT PART 2
# PROJECT SUBMISSION

## PHASE 4 INSTRUCTIONS:

*In this section continue building the project by performing different activities like feature engineering, model training, evaluation etc as per the instructions in the project.*

*Building a project that involves loading and preprocessing a dataset, performing feature engineering, model training, and evaluation is a common workflow in data science and machine learning. I'll provide you with a step-by-step guide on how to go about it:*

## *DEVELOPMENT:*

### 1. DEFINE YOUR PROBLEM:

*- Start by clearly defining your project's objectives and the problem you want to solve. Understanding the problem is crucial for selecting the right dataset and modeling approach.*

### 2. DATA COLLECTION:

- Gather the data relevant to your problem. This might involve web scraping, data acquisition from databases, or using publicly available datasets.

### 3. DATA PREPROCESSING:

- Clean the data by handling missing values, removing duplicates, and dealing with outliers. Preprocessing might also include data transformation, such as normalization or scaling.

### 4. EXPLORATORY DATA ANALYSIS (EDA):

   - *Perform EDA to gain insights into the dataset. This involves visualizing the data, identifying patterns, and understanding the relationships between different variables.*

### 5. FEATURE ENGINEERING:

   - *Create or modify features to improve the performance of your model. Feature engineering can involve techniques like one-hot encoding, feature scaling, and creating new features from existing ones.*

### 6. DATA SPLITTING:

   - *Split your dataset into training, validation, and testing sets. The training set is used to train the model, the validation set for hyperparameter tuning, and the testing set for evaluating the final model's performance.*

### 7. MODEL SELECTION:

   - *Choose a machine learning or deep learning model that's suitable for your problem. The selection depends on your data and objectives. Common choices include linear regression, decision trees, random forests, neural networks, etc.*

### 8. MODEL TRAINING:

   - *Train your selected model on the training dataset. This involves optimizing the model's parameters to fit the data. The process may include cross-validation for hyperparameter tuning.*

### 9. MODEL EVALUATION:

   - Evaluate the model's performance using appropriate metrics (e.g., accuracy, F1-score, RMSE). For classification problems, you can use a confusion matrix. For regression problems, you can plot the predicted vs. actual values.

### 10. MODEL VALIDATION:

   - *Validate your model's performance on the validation dataset. Make necessary adjustments to improve its performance.*

### 11. HYPERPARAMETER TUNING:

   - Fine-tune the model's hyperparameters to optimize its performance. This may involve grid search, random search, or Bayesian optimization.

### 12. TESTING AND FINAL EVALUATION:

   - Finally, evaluate your model on the testing dataset, which it has never seen before. This will provide an unbiased assessment of its performance.

### 13. DEPLOYMENT:

   - If your model performs well, consider deploying it to make predictions in a real-world application.

### 14. DOCUMENTATION:

   - Document the entire process, including data sources, preprocessing steps, model selection, hyperparameters, and evaluation results.

*Remember that this is a high-level overview, and the specific steps and tools you use will depend on the project, data, and your objectives. Additionally, it's important to iterate and refine your approach as you learn more about the problem and the data during the project.*

*A chatbot from scratch involves multiple steps, including dataset loading, preprocessing, feature engineering, model training, and evaluation. I'll provide a basic example using Python and libraries like NLTK for natural language processing and scikit-learn for feature engineering. In practice, you might use more advanced techniques and larger datasets for a production-ready chatbot.*

## PROGRAM:

```python
import nltk

import numpy as np

from sklearn.feature_extraction.text import TfidfVectorizer

from sklearn.metrics.pairwise import cosine_similarity


# Download NLTK data

nltk.download('punkt')

nltk.download('stopwords')


# Sample dataset

corpus = [

    "Hello, how are you?",

    "I'm doing well, thank you.",

    "What's the weather today?",
```

```python
    "It's sunny outside.",

    "Tell me a joke.",

    "Why was the math book sad? Because it had too many problems!",

    "Goodbye!",

    "See you later!"

]


# Preprocessing the dataset

from nltk.corpus import stopwords

from nltk.tokenize import word_tokenize


stop_words = set(stopwords.words('english'))


def preprocess_text(text):

    # Tokenization

    words = word_tokenize(text)


    # Removing stopwords and punctuation

    words = [word.lower() for word in words if word.isalnum() and word not in stop_words]


    return " ".join(words)


corpus = [preprocess_text(text) for text in corpus]
```

```python
# Feature Engineering

tfidf_vectorizer = TfidfVectorizer()

tfidf_matrix = tfidf_vectorizer.fit_transform(corpus)


# Chatbot function

def chatbot_response(user_input, corpus, tfidf_vectorizer, tfidf_matrix):

    user_input = preprocess_text(user_input)

    user_tfidf = tfidf_vectorizer.transform([user_input])


    # Calculate cosine similarities

    similarities = cosine_similarity(user_tfidf, tfidf_matrix)


    # Find the index of the most similar response

    most_similar_index = np.argmax(similarities)


    return corpus[most_similar_index]


# Chat with the bot

while True:

    user_input = input("You: ")

    if user_input.lower() == 'quit':

        print("Chatbot: Goodbye!")

        break

    response = chatbot_response(user_input, corpus, tfidf_vectorizer, tfidf_matrix)
```

```
print("Chatbot:", response)
```

*This code creates a basic chatbot that uses TF-IDF (Term Frequency-Inverse Document Frequency) to calculate similarity between user input and predefined responses. You can expand upon this example by using more extensive datasets, more advanced NLP techniques, and incorporating a broader range of responses for a more sophisticated chatbot.*