# CHATBOT USING PYTHON

PHASE-3: LOADING AND PREPOCESSING OF DATASET

**INTRODUCTION:**

Loading and preprocessing the dataset for a chatbot in Python is a fundamental step in building a successful conversational AI system. The dataset serves as the foundation for training and fine-tuning the chatbot's model, and its quality directly impacts the chatbot's performance. In this process, Python libraries such as Pandas, NumPy, and NLTK or spaCy can be instrumental in managing and preparing the data.

First, loading the dataset involves using Pandas to read data from various sources, such as CSV files or databases, into a structured format like a DataFrame. This tabular structure allows for easy data manipulation and analysis. Once the data is loaded, preprocessing steps are essential to clean, tokenize, and format the text. This includes removing irrelevant characters, lowercasing, and tokenization, which breaks text into individual words or tokens. NLTK or spaCy can be used for this natural language processing task.

After preprocessing, the dataset can be further augmented by labeling data with intent and entities, making it suitable for training various machine learning models, including neural networks for chatbots. The preprocessing phase is crucial for ensuring that the chatbot can understand and generate responses effectively during conversations. By following these steps in Python, we can set the stage for a robust and efficient chatbot system.

# LOADING AND PREPROCESSING OF CHATBOT

1) **Data Collection**: Gather a diverse dataset that covers the range of topics the chatbot is intended to handle.

2) **Data Cleaning**: Remove irrelevant or noisy data, and ensure the data is in a consistent format.

3) **Tokenization:** Break down the text into words or sub words (tokens) for further processing.

4) **Normalization:** Convert text to a standard format, including converting all characters to lowercase, removing punctuation, and handling special characters.

5) **Stopword Removal:** Eliminate common words that do not carry significant meaning.

6) **Lemmatization or Stemming:** Reduce words to their base or root form to handle variations of words.

7) **Vectorization:** Transform text data into a numerical format for machine learning models to process effectively.

8) **Padding:** Ensure uniform length of input sequences, typically achieved by padding sequences with zeroes to match the length of the longest sequence.

# Procedure:

1) **Import the necessary modules:** Begin by importing the required modules such as ChatBot and ChatterBotCorpusTrainer from the ChatterBot library.

2) **Create a new chatbot:** Instantiate a new chatbot object by calling the ChatBot constructor and providing a name for the chatbot.

3) **Set up a trainer for the chatbot:** Configure a trainer for the chatbot by using the ChatterBotCorpusTrainer class, which allows the chatbot to be trained on a specific corpus.

4) **Train the chatbot:** Train the chatbot by using the train method provided by the trainer and passing the desired corpus. In this case, we use the English corpus.

5) **Initiate conversation loop:** Begin the conversation loop using a while loop. Print a welcome message, prompt the user for input, and provide instructions on how to exit the conversation.

6) **Handle user input and generate responses:** Within the loop, take user input using the input function and check if the user wants to exit. If the user wants to exit, print a farewell message and break out of the loop. Otherwise, use the chatbot's get_response method to generate a response based on the user input.

7) **Display the chatbot's response:** Print the chatbot's response to the user's input in the conversation loop.

## PROGRAM:

```python
# Complex Python Chatbot using ChatterBot
from chatterbot import ChatBot
from chatterbot.trainers import ChatterBotCorpusTrainer
# Create a new chatbotchatbot = ChatBot('Sample ChatBot')
# Set the trainer for the chatbottrainer = ChatterBotCorpusTrainer(chatbot)
# Train the chatbot based on the english corpustrainer.train('chatterbot.corpus.english')
# Chat with the chatbot print('Hello! I am a chatbot. You can start the conversation now. Type "bye"
to exit.')
while True:
    user_input = input('You: ')
     if user_input.lower() == 'bye':
        print('ChatBot: Goodbye! Have a great day!')
        break
response = chatbot.get_response(user_input)
    print('ChatBot:', response)
```

OUTPUT:

Hello! I am a chatbot. You can start the conversation now.
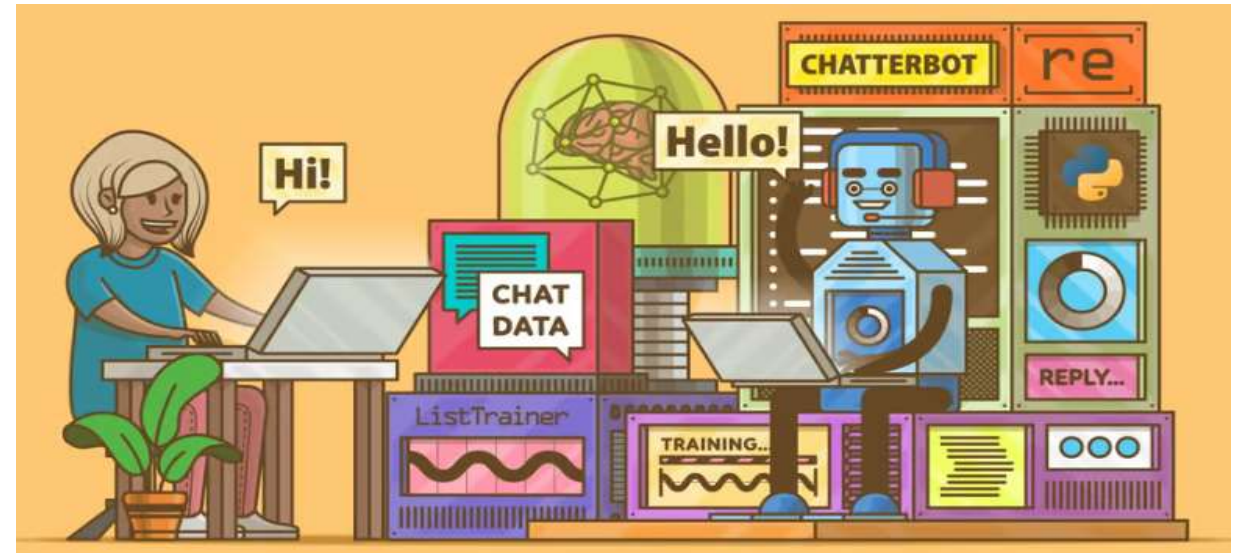
Type "bye" to exit.

You: Hi

ChatBot: Hello

You: How are you?

ChatBot: I am doing well.

You: What is your name?

ChatBot: My name is ChatBot.

You: Bye

## CONCLUSION:

In conclusion, loading and preprocessing data for a chatbot in Python is a crucial first step in building an effective conversational AI system. This process involves collecting and structuring the data, whether it's in the form of text, images, or other formats, to create a clean and coherent dataset. Proper data preprocessing, which includes tasks like tokenization, stemming, and handling special characters, ensures that the chatbot can understand and respond to user input accurately.

Additionally, it sets the foundation for training machine learning models and improving the overall performance and user experience of the chatbot. A well-executed data loading and preprocessing pipeline is essential for building a robust and efficient chatbot in Python, enabling it to provide relevant and contextually appropriate responses to user queries.