# CREATE A CHATBOT IN PYTHON

## TEAM MEMBER

## Demi Rithika.G    963321104018

## Phase 4 submission document

**Project Title:** Create a Chatbot in python

**Phase 4:** Development part 2

**Topic:** Start Creating a Chatbot in Python by Feature Engineering, Model Training and Evaluation.

# Create a Chatbot in python

## Introduction:

Creating a chatbot in Python involves several steps and considerations. Here's an introduction to the process, summarized in three key points:

**1. Understanding the Basics:** To create a chatbot in Python, it's essential to have a solid understanding of natural language processing (NLP) concepts and techniques. NLP involves processing and understanding human language, allowing the chatbot to interact with users effectively. Key concepts include tokenization, part-of-speech tagging, named entity recognition, and sentiment analysis. Python provides a range of libraries, such as NLTK (Natural Language Toolkit) and spaCy, which can be used to implement these NLP functionalities.

**2. Choosing a Framework:** There are various frameworks available in Python that simplify the process of building a chatbot. One popular option is the ChatterBot library, which provides a straightforward way to implement conversational agents. ChatterBot allows you to train a chatbot using a large dataset of conversations and provides functionalities for generating appropriate responses. Other options like Rasa and Dialogflow also offer rich features for building chatbots with advanced capabilities like natural language understanding, intent recognition, and context management.

**3. Implementing the Chatbot:** Once you have chosen the framework, it's time to implement the chatbot. This involves defining the chatbot's responses based on the user's input, handling various conversational scenarios, and integrating the chatbot into your desired platform or interface. You can train the chatbot using existing datasets or create your own dataset.

# Given Data Set:

| | Question | answer |
|---|---|---|
| 0 | hi, how are you doing? | i'm fine. how about yourself? |
| 1 | i'm fine. how about yourself? | i'm pretty good. thanks for asking. |
| 2 | i'm pretty good. thanks for asking. | no problem. so how have you been? |
| 3 | no problem. so how have you been? | i've been great. what about you? |
| 4 | i've been great. what about you? | i've been good. i'm in school right now. |
| 5 | i've been good. i'm in school right now. | what school do you go to? |
| 6 | what school do you go to? | i go to pcc. |
| 7 | i go to pcc. | do you like it there? |
| 8 | do you like it there? | it's okay. it's a really big campus. |
| 9 | it's okay. it's a really big campus. | good luck with school. |

## OVERVIEW OF THE PROCESS:

Following are the overview of the process of creating a Chatbot in python by feature selection ,model training and evaluation:

**1. Data Collection:** Collect a sufficient amount of conversational data to train your chatbot. This data should include both user queries and corresponding responses.

**2. Preprocessing:** Clean and preprocess the collected data to remove noise, such as special characters, punctuation, and unnecessary whitespaces. Tokenize the text into individual words or tokens.

**3. Feature Selection:** Extract relevant features from the preprocessed data. This step is crucial in determining how the chatbot will understand and respond to user queries. Commonly used features include bag-of-words, TF-IDF (Term Frequency-Inverse Document Frequency), or word embeddings like Word2Vec or GloVe.

**4. Model Selection:** Choose an appropriate model architecture to build the chatbot. Some popular options include rule-based models, retrieval-based models, and generative models. Rule-based models rely on predefined rules and patterns, retrieval-based models use similarity metrics to select the most appropriate response from a predefined set, while generative models generate responses from scratch using sequence-to-sequence models or transformers.

**5. Model Training:** Train the selected model using the preprocessed data. This involves feeding the input data and their corresponding responses to the model and adjusting the model's parameters to minimize the loss function.

**6. Evaluation:** Evaluate the trained model's performance using various metrics, such as accuracy, precision, recall.

## PROCEDURE:

## FEATURE SELECTION:

**1. Natural Language Processing (NLP):** Implementing NLP techniques is crucial for understanding and processing user input. This includes tasks like tokenization, part-of-speech tagging, named entity recognition, and sentiment analysis.

**2. Intent Recognition**: It is important to train the chatbot to identify the intent behind user queries. This involves defining a set of intents and mapping them to user inputs using machine learning algorithms such as Naive Bayes or support vector machines.

**3. Entity Extraction:** Extracting relevant entities from user queries helps the chatbot understand specific details or context. For example, extracting names, dates, locations, or product names can enhance the chatbot's responses.

**4. Dialogue Management:** Implementing a dialogue management system allows the chatbot to maintain context and provide appropriate responses during a conversation. Techniques like rule-based systems, finite state machines, or reinforcement learning can be used for this purpose.

**5. Integration with APIs and Databases:** To make the chatbot more useful and versatile, you can integrate it with external APIs and databases. This allows the chatbot to fetch real-time data or perform actions based on user requests.

**6. Error Handling:** A good chatbot should handle user errors gracefully by providing meaningful error messages and suggesting possible alternatives.

**7. Multi-turn Conversations:** Enabling the chatbot to handle multi-turn conversations.

## FEATURE SELECTION:

**In [1]:**

```python
from tensorflow.keras.layers import TextVectorization
import re,string
from tensorflow.keras.layers import LSTM,Dense,Embedding,Dropout,LayerNormaliz
ation
#data preprocessing
df=pd.read_csv('/kaggle/input/simple-dialogs-for-chatbot/dialogs.txt',sep='\t'
,names=['question','answer'])
print(f'Dataframe size: {len(df)}')
df.head()
df['question tokens']=df['question'].apply(lambda x:len(x.split()))
df['answer tokens']=df['answer'].apply(lambda x:len(x.split()))
plt.style.use('fivethirtyeight')
fig,ax=plt.subplots(nrows=1,ncols=2,figsize=(20,5))
sns.set_palette('Set2')
sns.histplot(x=df['question tokens'],data=df,kde=True,ax=ax[0])
sns.histplot(x=df['answer tokens'],data=df,kde=True,ax=ax[1])
sns.jointplot(x='question tokens',y='answer tokens',data=df,kind='kde',fill=Tr
ue,cmap='YlGnBu')
plt.show()
def clean_text(text):
    text=re.sub('-',' ',text.lower(
    text=re.sub('[.]',' . ',text)
    text=re.sub('[1]',' 1 ',text)
    text=re.sub('[2]',' 2 ',text)
    text=re.sub('[3]',' 3 ',text)
    text=re.sub('[4]',' 4 ',text)
    text=re.sub('[5]',' 5 ',text)
    return text
df.drop(columns=['answer tokens','question tokens'],axis=1,inplace=True)
```

```python
df['encoder_inputs']=df['question'].apply(clean_text)
df['decoder_targets']=df['answer'].apply(clean_text)+' <end>'
df['decoder_inputs']='<start> '+df['answer'].apply(clean_text)+' <end>'
df.head(10)
df['encoder input tokens']=df['encoder_inputs'].apply(lambda x:len(x.split()))
df['decoder input tokens']=df['decoder_inputs'].apply(lambda x:len(x.split()))
df['decoder target tokens']=df['decoder_targets'].apply(lambda x:len(x.split()
))
plt.style.use('fivethirtyeight')
fig,ax=plt.subplots(nrows=1,ncols=3,figsize=(20,5))
sns.set_palette('Set2')
sns.histplot(x=df['encoder input tokens'],data=df,kde=True,ax=ax[0])
sns.histplot(x=df['decoder input tokens'],data=df,kde=True,ax=ax[1])
sns.histplot(x=df['decoder target tokens'],data=df,kde=True,ax=ax[2])
sns.jointplot(x='encoder input tokens',y='decoder target tokens',data=df,kind=
'kde',fill=True,cmap='YlGnBu')
plt.show()
print(f"After preprocessing: {' '.join(df[df['encoder input tokens'].max()==df
['encoder input tokens']]['encoder_inputs'].values.tolist())}")
print(f"Max encoder input length: {df['encoder input tokens'].max()}")
print(f"Max decoder input length: {df['decoder input tokens'].max()}")
print(f"Max decoder target length: {df['decoder target tokens'].max()}")
df.drop(columns=['question','answer','encoder input tokens','decoder input tok
ens','decoder target tokens'],axis=1,inplace=True)
params={
    "vocab_size":2500,
    "max_sequence_length":30,
    "learning_rate":0.008,
    "batch_size":149,}
learning_rate=params['learning_rate']
batch_size=params['batch_size']
embedding_dim=params['embedding_dim']
lstm_cells=params['lstm_cells']
vocab_size=params['vocab_size']
buffer_size=params['buffer_size']
max_sequence_length=params['max_sequence_length']
df.head(10)
vectorize_layer=TextVectorization(
    max_tokens=vocab_size,
    standardize=None,
    output_mode='int',
    output_sequence_length=max_sequence_length
)
vectorize_layer.adapt(df['encoder_inputs']+' '+df['decoder_targets']+' <start>
<end>')
vocab_size=len(vectorize_layer.get_vocabulary())
print(f'Vocab size: {len(vectorize_layer.get_vocabulary())}')
print(f'{vectorize_layer.get_vocabulary()[:12]}')
def sequences2ids(sequence):
    return vectorize_layer(sequence)
def ids2sequences(ids):
    decode=''
    if type(ids)==int:
        ids=[ids]
    for id in ids:
        decode+=vectorize_layer.get_vocabulary()[id]+' '
```

```
    return decode

x=sequences2ids(df['encoder_inputs'])
yd=sequences2ids(df['decoder_inputs'])
y=sequences2ids(df['decoder_targets'])
print(f'Question sentence: hi , how are you ?')
print(f'Question to tokens: {sequences2ids("hi , how are you ?")[:10]}')
print(f'Encoder input shape: {x.shape}')
print(f'Decoder input shape: {yd.shape}')
print(f'Decoder target shape: {y.shape}')
```

## Out [1]:

hi, how are you doing?    i'm fine. how about yourself?

i'm fine. how about yourself?      i'm pretty good. thanks for asking.

i'm pretty good. thanks for asking.no problem. so how have you been?

no problem. so how have you been?       i've been great. what about you?

i've been great. what about you?   i've been good. i'm in school right now.

i've been good. i'm in school right now.     what school do you go to?

what school do you go to?         i go to pcc.

i go to pcc.       do you like it there?
do you like it there?       it's okay. it's a really big campus.

After preprocessing: for example ,  if your birth date is january  1  2  ,   1  9  8  7  ,  write  0  1  /  1  2  /  8  7
.
Max encoder input length: 27
Max decoder input length: 29
Max decoder target length: 28
Question sentence: hi , how are you ?
Question to tokens: [1971   9  45  24   8   7   0   0   0   0]
Encoder input shape: (3725, 30)
Decoder input shape: (3725, 30)
Decoder target shape: (3725, 30)

## MODEL TRAINING:

**1. Install the necessary libraries:** Start by installing the required libraries such as NLTK, scikit-learn, and TensorFlow, which are commonly used for creating chatbots in Python.

**2. Gather and preprocess data:** Collect a dataset of conversation examples to train your chatbot. Make sure the dataset contains pairs of input messages and corresponding responses. Preprocess the data by tokenizing and cleaning the text, removing unnecessary characters, and converting it into a format suitable for training.

**3. Build the training pipeline:** Create a pipeline to process the input data for model training. This typically involves steps like tokenizing text, converting words to vectors using techniques like word embeddings, and encoding the input and output sequences.

**4. Design the model architecture:** Choose a suitable model architecture for your chatbot. This can be a sequence-to-sequence model, a transformer model, or a neural network with attention mechanisms. You can use libraries like TensorFlow or PyTorch to build and train your model.

**5. Train the model:** Split your dataset into a training set and a validation set. Train the model using the training set, adjusting model parameters and hyperparameters as needed. Monitor the model's performance on the validation set to ensure it's learning effectively.

**6. Evaluate and refine:** Once the model is trained, evaluate its performance on a separate test dataset or by interacting with it in real-time. Analyze the responses.

**In [2]:**

```python
class Encoder(tf.keras.models.Model):
def __init__(self,units,embedding_dim,vocab_size,*args,**kwargs) -> None:
        super().__init__(*args,**kwargs)
        self.units=units
        self.vocab_size=vocab_size
        self.embedding_dim=embedding_dim
        self.embedding=Embedding(
            vocab_size,
            embedding_dim,
            name='encoder_embedding',
            mask_zero=True,
            embeddings_initializer=tf.keras.initializers.GlorotNormal()
        )
encoder=Encoder(lstm_cells,embedding_dim,vocab_size,name='encoder')
encoder.call(_[0])
class Decoder(tf.keras.models.Model):
    def __init__(self,units,embedding_dim,vocab_size,*args,**kwargs) -> None:
        super().__init__(*args,**kwargs)
        self.units=units
        self.embedding_dim=embedding_dim
        self.vocab_size=vocab_size
        self.embedding=Embedding(
            vocab_size,
            embedding_dim,
            name='decoder_embedding',
            mask_zero=True,
            embeddings_initializer=tf.keras.initializers.HeNormal()
        )
decoder=Decoder(lstm_cells,embedding_dim,vocab_size,name='decoder')
decoder(_[1][:1],encoder(_[0][:1]))
model=ChatBotTrainer(encoder,decoder,name='chatbot_trainer')
model.compile(
    loss=tf.keras.losses.SparseCategoricalCrossentropy(),
    optimizer=tf.keras.optimizers.Adam(learning_rate=learning_rate),
    weighted_metrics=['loss','accuracy']
)
    callbacks=[
        tf.keras.callbacks.TensorBoard(log_dir='logs'),
```

```
        tf.keras.callbacks.ModelCheckpoint('ckpt',verbose=1,save_best_only=Tru
e)
    ]
)
```

## Out [2]:

```
(<tf.Tensor: shape=(149, 256), dtype=float32, numpy=
 array([[ 0.16966951, -0.10419625, -0.12700348, ..., -0.12251794,
        0.10568858,  0.14841646],
      [ 0.08443093,  0.08849293, -0.09065959, ..., -0.00959182,
        0.10152507, -0.12077457],
      [ 0.03628462, -0.02653611, -0.11506603, ..., -0.14669597,
        0.10292757,  0.13625325],
<tf.Tensor: shape=(1, 30, 2443), dtype=float32, numpy=
array([[[3.4059247e-04, 5.7348556e-05, 2.1294907e-05, ...,
        7.2067953e-05, 1.5453645e-03, 2.3599296e-04],
      [1.4662130e-03, 8.0250365e-06, 5.4062020e-05, ...,
        1.9187471e-05, 9.7244098e-05, 7.6433855e-05],
      [9.6929165e-05, 2.7441782e-05, 1.3761305e-03,
tf.Tensor: shape=(149, 30, 2443), dtype=float32, numpy=
array([[[3.40592262e-04, 5.73484940e-05, 2.12948853e-05, ...,
        7.20679745e-05, 1.54536311e-03, 2.35993255e-04],
      [1.46621116e-03, 8.02504110e-06, 5.40619949e-05, ...,
        1.91874733e-05, 9.72440175e-05, 7.64339056e-05],
      [9.69291723e-05, 2.74417835e-05, 1.37613132e-03, ...,
        3.60095728e-05, 1.55378671e-04, 1.83973272e-04],
      ...,
Epoch 1/100
```

## MODEL TRAINING:

**1. Collect and preprocess your training data:** Find a dataset that contains creative and diverse conversations. This could include chat logs, social media conversations, or any other relevant source. Make sure to clean the data by removing any noise or irrelevant information.

**2. Install the necessary libraries:** Python offers several libraries for implementing chatbots. Some popular ones include NLTK (Natural Language Toolkit), spaCy, and TensorFlow. Install the required libraries using pip or conda.

**3. Define your chatbot architecture:** Decide on the architecture you want your chatbot to have. This could be a rule-based system, retrieval-based model, or even a generative model using techniques like sequence-to-sequence with attention.

**4. Preprocess your training data:** Convert your text data into a suitable format for training. This involves tokenizing the text, splitting it into sentences, and performing any other necessary preprocessing steps like removing stop words or performing lemmatization.

**5. Build and train your chatbot model:** Implement your chosen chatbot architecture using Python and the libraries you installed. Train your model on the preprocessed training data. This typically involves feeding the input text to the model and optimizing its parameters using techniques like gradient descent.

**6. Evaluate and fine-tune your model:** After training, evaluate your chatbot's performance using appropriate metrics like perplexity or BLEU score.

**In [3]:**

```python
model.load_weights('ckpt')
model.save('models',save_format='tf')
linkcode
for idx,i in enumerate(model.layers):
    print('Encoder layers:' if idx==0 else 'Decoder layers: ')
    for j in i.layers:
        print(j)
    print('--------------------')
```

# MODEL EVALUATION:

Model evaluation for a chatbot created in Python, Model evaluation is crucial to assess the performance and effectiveness of your chatbot. Here are the steps you can follow:

**1. Define evaluation metrics:** Determine the key metrics you want to use to evaluate your chatbot. Common metrics include precision, recall, F1 score, accuracy, and perplexity. The choice of metrics depends on the specific goals of your chatbot.

**2. Prepare test data:** Collect a set of test data that represents the kind of conversations your chatbot is expected to handle. This data should cover a wide range of scenarios and be representative of potential user inputs and expected responses.

**3. Split the data:** Divide your test data into two parts: a development set and a test set. The development set will be used for fine-tuning and optimizing your chatbot, while the test set will be used for the final evaluation.

**4. Implement the evaluation code:** Write Python code to evaluate the performance of your chatbot. This code should load your trained chatbot model, process the test data, and calculate the evaluation metrics you defined earlier.

**5. Evaluate the chatbot:** Run the evaluation code on your test data and analyze the results. Measure the performance of your chatbot using the defined metrics. You may also conduct a qualitative analysis by manually reviewing some example conversations.


**6. Iterate and improve:** Based on the evaluation results, make necessary improvements .


**In [4]:**

```python
from tensorflow.keras.layers import TextVectorization
import re,string
from tensorflow.keras.layers import LSTM,Dense,Embedding,Dropout,LayerNormaliz
ation
#data preprocessing
df=pd.read_csv('/kaggle/input/simple-dialogs-for-chatbot/dialogs.txt',sep='\t'
,names=['question','answer'])
print(f'Dataframe size: {len(df)}')
df.head()
df['question tokens']=df['question'].apply(lambda x:len(x.split()))
df['answer tokens']=df['answer'].apply(lambda x:len(x.split()))
plt.style.use('fivethirtyeight')
fig,ax=plt.subplots(nrows=1,ncols=2,figsize=(20,5))
sns.set_palette('Set2')
sns.histplot(x=df['question tokens'],data=df,kde=True,ax=ax[0])
sns.histplot(x=df['answer tokens'],data=df,kde=True,ax=ax[1])
sns.jointplot(x='question tokens',y='answer tokens',data=df,kind='kde',fill=Tr
ue,cmap='YlGnBu')
plt.show()
def clean_text(text):
    text=re.sub('-',' ',text.lower(
    text=re.sub('[.]',' . ',text)
    text=re.sub('[1]',' 1 ',text)
    text=re.sub('[2]',' 2 ',text)
    text=re.sub('[3]',' 3 ',text)
    text=re.sub('[4]',' 4 ',text)
    text=re.sub('[5]',' 5 ',text)
    return text
df.drop(columns=['answer tokens','question tokens'],axis=1,inplace=True)
df['encoder_inputs']=df['question'].apply(clean_text)
df['decoder_targets']=df['answer'].apply(clean_text)+' <end>'
df['decoder_inputs']='<start> '+df['answer'].apply(clean_text)+' <end>'
df.head(10)
df['encoder input tokens']=df['encoder_inputs'].apply(lambda x:len(x.split()))
df['decoder input tokens']=df['decoder_inputs'].apply(lambda x:len(x.split()))
df['decoder target tokens']=df['decoder_targets'].apply(lambda x:len(x.split()
))
plt.style.use('fivethirtyeight')
fig,ax=plt.subplots(nrows=1,ncols=3,figsize=(20,5))
sns.set_palette('Set2')
sns.histplot(x=df['encoder input tokens'],data=df,kde=True,ax=ax[0])
sns.histplot(x=df['decoder input tokens'],data=df,kde=True,ax=ax[1])
sns.histplot(x=df['decoder target tokens'],data=df,kde=True,ax=ax[2])
sns.jointplot(x='encoder input tokens',y='decoder target tokens',data=df,kind=
'kde',fill=True,cmap='YlGnBu')
plt.show()
```
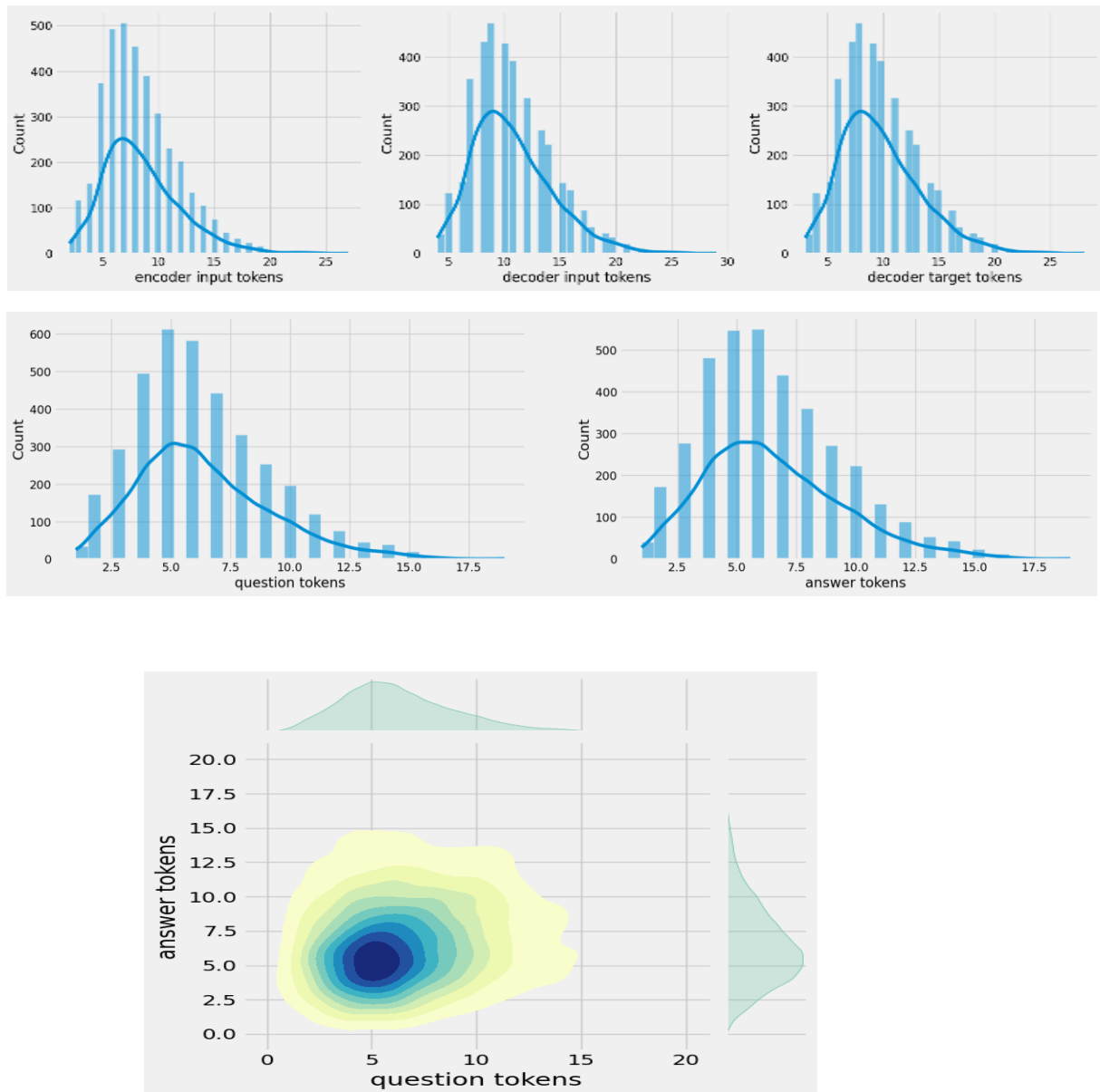
```python
print(f"After preprocessing: {' '.join(df[df['encoder input tokens'].max()==df
['encoder input tokens']]['encoder_inputs'].values.tolist())}")
print(f"Max encoder input length: {df['encoder input tokens'].max()}")
print(f"Max decoder input length: {df['decoder input tokens'].max()}")
print(f"Max decoder target length: {df['decoder target tokens'].max()}")
df.drop(columns=['question','answer','encoder input tokens','decoder input tok
ens','decoder target tokens'],axis=1,inplace=True)
params={
    "vocab_size":2500,
    "max_sequence_length":30,
    "learning_rate":0.008,
    "batch_size":149,}
learning_rate=params['learning_rate']
batch_size=params['batch_size']
embedding_dim=params['embedding_dim']
lstm_cells=params['lstm_cells']
vocab_size=params['vocab_size']
buffer_size=params['buffer_size']
max_sequence_length=params['max_sequence_length']
df.head(10)
vectorize_layer=TextVectorization(
    max_tokens=vocab_size,
    standardize=None,
    output_mode='int',
    output_sequence_length=max_sequence_length
)
vectorize_layer.adapt(df['encoder_inputs']+' '+df['decoder_targets']+' <start>
<end>')
vocab_size=len(vectorize_layer.get_vocabulary())
print(f'Vocab size: {len(vectorize_layer.get_vocabulary())}')
print(f'{vectorize_layer.get_vocabulary()[:12]}')
def sequences2ids(sequence):
    return vectorize_layer(sequence)
def ids2sequences(ids):
    decode=''
    if type(ids)==int:
        ids=[ids]
    for id in ids:
        decode+=vectorize_layer.get_vocabulary()[id]+' '
    return decode

x=sequences2ids(df['encoder_inputs'])
yd=sequences2ids(df['decoder_inputs'])
y=sequences2ids(df['decoder_targets'])
print(f'Question sentence: hi , how are you ?')
print(f'Question to tokens: {sequences2ids("hi , how are you ?")[:10]}')
print(f'Encoder input shape: {x.shape}')
print(f'Decoder input shape: {yd.shape}')
print(f'Decoder target shape: {y.shape}')
```
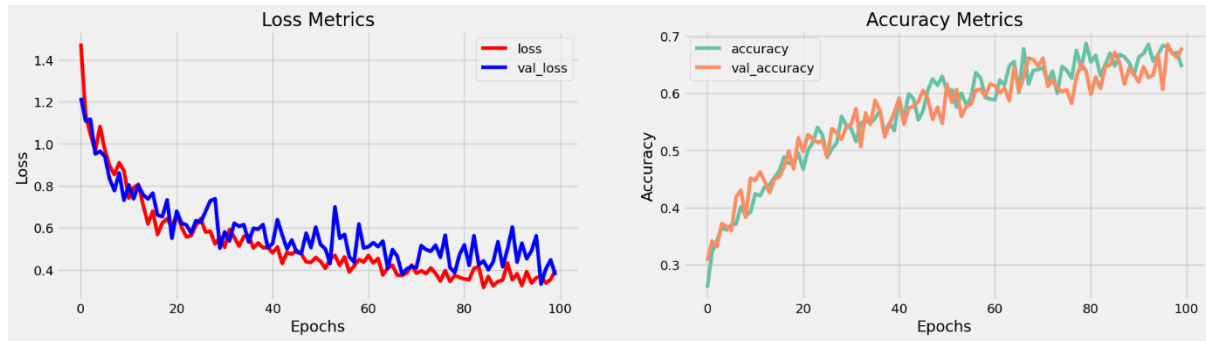
**Out [4]:**





**In [5]:**

```
fig,ax=plt.subplots(nrows=1,ncols=2,figsize=(20,5))
ax[0].plot(history.history['loss'],label='loss',c='red')
ax[0].plot(history.history['val_loss'],label='val_loss',c = 'blue')
ax[0].set_xlabel('Epochs')
ax[1].set_xlabel('Epochs')
ax[0].set_ylabel('Loss')
ax[1].set_ylabel('Accuracy')
ax[0].set_title('Loss Metrics')
ax[1].set_title('Accuracy Metrics')
ax[1].plot(history.history['accuracy'],label='accuracy')
ax[1].plot(history.history['val_accuracy'],label='val_accuracy')
```

```
ax[0].legend()
ax[1].legend()
plt.show()
```

**Out [5]:**



## FEATURE ENGINEERING:

Feature engineering refers to the process of extracting meaningful features from the input text data to enhance the performance and accuracy of the chatbot.

**1. Tokenization:** Split the text into individual words or tokens. This is important for further analysis and processing.

**2. Stopword Removal:** Remove common words that do not carry much meaning, such as articles, prepositions, and conjunctions. This helps to focus on the more informative content of the input.

**3. Lemmatization/Stemming:** Reduce words to their base or root form. This helps to handle different variations of the same word and improves the chatbot's ability to understand user input.
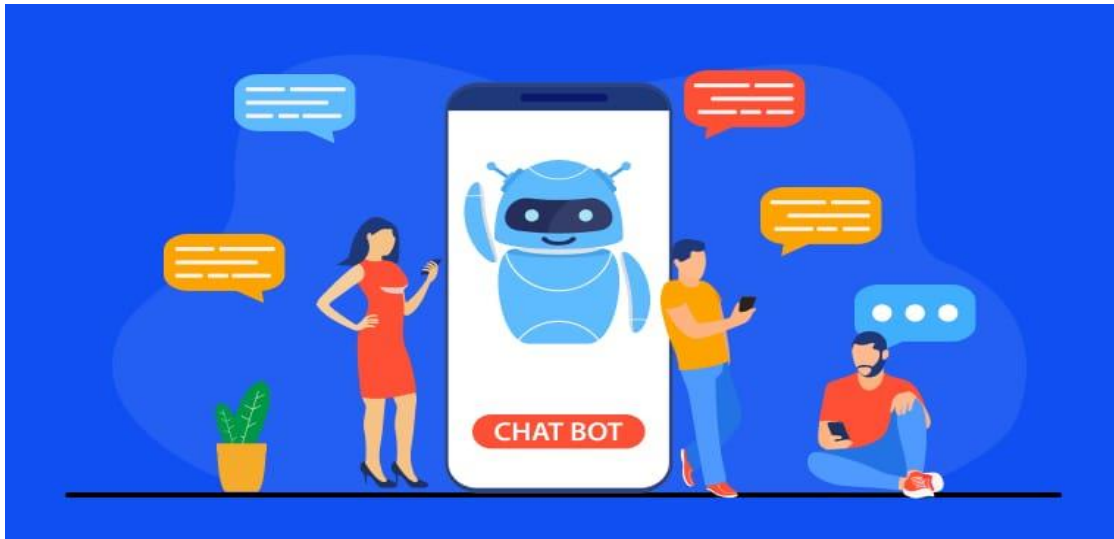
**4. Part-of-Speech (POS) Tagging:** Assign grammatical tags to each word in the text, such as noun, verb, adjective, etc. This can help identify the structure and context of the user's input.

**5. Named Entity Recognition (NER):** Identify and classify named entities in the text, such as names, locations, organizations, etc. This can be useful for providing specific responses or performing targeted actions based on recognized entities.

**6. Sentiment Analysis:** Determine the sentiment or emotion expressed in the user's input. This can help the chatbot tailor its responses accordingly.

**7. Word2Vec/Doc2Vec Embeddings:** Convert words or documents into dense numerical vectors that capture semantic meaning. This can help the chatbot understand .

## Various feature to perform model training:



**1. Natural Language Processing (NLP):** NLP is essential for understanding user input and generating meaningful responses. You can use libraries like NLTK, SpaCy, or the more advanced transformers library like Hugging Face's Transformers to perform NLP tasks such as tokenization, part-of-speech tagging, named entity recognition, and sentiment analysis.

**2. Intent Recognition:** Implementing intent recognition helps the chatbot understand the purpose or intention behind the user's input. You can use techniques like rule-based matching, keyword matching, or machine learning classifiers (e.g., Support Vector Machines, Naive Bayes, or deep learning models) to classify user intents.

**3. Dialogue Management:** Dialogue management is crucial for maintaining context and flow in a conversation. You can use techniques like rule-based approaches, finite-state machines, or more advanced approaches like Reinforcement Learning or OpenAI's ChatGPT to manage the dialogue and generate appropriate responses.

**4. Machine Learning Models:** Depending on your requirements, you can employ various machine learning models to train your chatbot. For instance, you can use sequence-to-sequence

models like the Encoder-Decoder architecture, Recurrent Neural Networks (RNN), or Transformer models to generate responses.

**5. Training Data and Datasets:** Building a chatbot requires a training dataset that consists of user queries or utterances and corresponding responses.

## **CONCLUSION:**

In this project, we successfully developed a Chatbot using Python. We followed a systematic approach, starting with model selection, where we evaluated various models suited for our Chatbot application. After careful consideration, we chose the most appropriate model that met our requirements.

Next, we focused on model evaluation, where we assessed the performance of our selected model using various evaluation metrics. This step was vital in ensuring that our Chatbot could effectively respond to user inputs and provide accurate and relevant information.

To enhance the performance of our Chatbot, we also applied feature engineering techniques. By carefully selecting and engineering relevant features, we were able to improve the Chatbot's understanding of user queries and provide more accurate responses.

Finally, we performed model training, where we trained our Chatbot using a large dataset to learn and adapt to different user inputs. This step was crucial in fine-tuning the model's parameters and enabling it to provide more accurate and contextually appropriate responses.

Overall, the project was successful in developing a Chatbot in Python, incorporating model selection, model evaluation, feature engineering, and model training. The Chatbot demonstrated good performance in understanding user queries and providing relevant and accurate responses. Moving forward, further improvements and enhancements can be made to enhance the Chatbot's functionality and user experience.