# CHATBOT

## *Abstract:*

This code presents a neural network architecture designed for question answering tasks, implemented using the Keras deep learning library. The dataset consists of stories, associated questions, and corresponding answers. Initially, the code preprocesses the data, including tokenization and vectorization using the Tokenizer class provided by Keras. The model architecture involves multiple sequential layers, including embedding layers to convert words into dense vector representations and LSTM layers to process sequences efficiently.

One notable aspect of the architecture is the attention mechanism, implemented through dot product attention, allowing the model to focus on relevant parts of the input story when generating an answer to the question. Additionally, the model utilizes dropout regularization to prevent overfitting during training.

The training process employs the RMSprop optimizer and categorical cross-entropy loss function. To monitor training progress and model performance, the code includes visualization of loss values using matplotlib. Through the provided code, a neural network model capable of effectively answering questions based on provided stories can be trained and evaluated.

## ***Objectives:***

1. Data Preprocessing:

- Tokenize the input text data, including stories, questions, and answers, to convert them into numerical representations.

- Vectorize the tokenized data to prepare it for input into the neural network model.

- Determine the vocabulary size and maximum sequence lengths for stories and questions.

2. Model Architecture Design:

- Design a neural network architecture capable of effectively learning the relationship between stories, questions, and answers.

- Utilize embedding layers to convert words into dense vector representations, allowing the model to capture semantic similarities.

- Implement LSTM layers to process sequential data efficiently, retaining

contextual information from the input stories and questions.

   - Incorporate attention mechanisms to enable the model to focus on relevant parts of the input story when generating an answer to the question.


3. Model Training:

   - Compile the model with appropriate optimization and loss functions, such as RMSprop optimizer and categorical cross-entropy loss, respectively.

   - Train the model using the preprocessed and vectorized training data, including stories, questions, and answers.

   - Monitor training progress and model performance using validation data, visualizing loss values over epochs to ensure convergence and prevent overfitting.

## 4. Model Evaluation:

   - Evaluate the trained model's performance on unseen test data, including stories and questions.

   - Assess the model's accuracy in predicting answers to questions based on input stories.

   - Analyze any potential areas for improvement or further optimization of the model's architecture or training process.

## *Introduction:*

This code presents an implementation of a neural network model designed for question answering tasks using the Keras deep learning library. Question answering (QA) systems aim to develop algorithms capable of understanding and responding to questions posed in natural language, based on provided context or knowledge. The model is trained

on a dataset consisting of stories, associated questions, and corresponding answers, with the objective of learning to generate accurate responses to novel questions given a story context.

The implementation follows a systematic approach, starting with data preprocessing to tokenize and vectorize the input text data, enabling numerical representation for model input. Subsequently, a neural network architecture is designed, comprising embedding layers to convert words into dense vector representations and LSTM layers to process sequential data efficiently, retaining contextual information from the input stories and questions.

An essential aspect of the model architecture is the incorporation of attention mechanisms, facilitating the model's ability to focus on

relevant parts of the input story when generating an answer to the question. This attention mechanism enhances the model's understanding of the relationships between different elements in the input data, leading to more accurate and contextually relevant answers.

During model training, the compiled model is trained using the preprocessed training data, and the training progress is monitored using validation data. The performance of the model is evaluated based on metrics such as accuracy and loss values to assess its effectiveness in the question answering task.

Overall, the code aims to provide a robust framework for training and evaluating neural network models for question answering tasks, with potential applications in various domains

requiring natural language understanding and interaction.

## *Methodology:*

1. **Data Preprocessing:**
   - Normalization: Convert text to lowercase, remove punctuation, and handle contractions.
   - Tokenization: Split sentences into words or subword tokens.
   - Stop Word Removal: Eliminate common words (e.g., "the," "and," "is").
   - Stemming or Lemmatization: Reduce words to their base form.

2. **Model Architecture Design:**
   - Consider using a Transformer architecture with self-attention mechanisms.

- Experiment with different embedding layers (e.g., Word2Vec, GloVe).
- Explore variations of LSTM (Long Short-Term Memory) or GRU (Gated Recurrent Unit) cells.
- Incorporate bidirectional layers for better context understanding.

3. **Model Training:**
- Define an appropriate loss function (e.g., categorical cross-entropy for classification).
- Choose an optimizer (e.g., Adam, SGD) and set learning rate.
- Monitor training using validation data and early stopping.
- Regularize with techniques like dropout or L2 regularization.

4. **Model Evaluation:**

- Calculate metrics such as accuracy, precision, recall, and F1-score.
- Use confusion matrices to understand classification performance.
- Compare against baseline models or previous state-of-the-art results.

## *Code:*

```
import pickle

import numpy as np

from keras.preprocessing.sequence import pad_sequences

from keras.preprocessing.text import Tokenizer

from keras.models import Sequential, Model
```

```python
from keras.layers import Embedding, Input,
Activation, Dense, Permute, Dropout,
Concatenate, add, dot, LSTM

import matplotlib.pyplot as plt


with
open("/kaggle/input/dataset/test_qa2201201
45430-220818-175426.txt", "rb") as fp:
    test_data = pickle.load(fp)


with open("/content/train_qa220120145526-
220818-175522.txt", "rb") as fp1:
    train_data = pickle.load(fp1)


vocab = set()
all_data = test_data + train_data
for story, question, ans in all_data:
    vocab = vocab.union(set(story))
```

```python
    vocab = vocab.union(set(question))

vocab.add('yes')
vocab.add('no')

vocab_len = len(vocab) + 1
max_story_len = max([len(data[0]) for data in
all_data])
max_ques_len = max([len(data[1]) for data in
all_data])

tokenizer = Tokenizer(filters=[])
tokenizer.fit_on_texts(vocab)

train_story_text = []
train_ques_text = []
train_ans = []
for story, ques, ans in train_data:
```

```python
        train_story_text.append(story)

        train_ques_text.append(ques)

        train_ans.append(ans)


train_story_seq =
tokenizer.texts_to_sequences(train_story_tex
t)




def vectorize_data(data,
word_index=tokenizer.word_index,
max_story_len=max_story_len,
max_ques_len=max_ques_len):

    X = []

    Xq = []

    Y = []

    for story, ques, ans in data:

        x = [word_index[word.lower()] for word
in story]
```

```python
        xq = [word_index[word.lower()] for word
in ques]

        y = np.zeros(len(word_index) + 1)

        y[word_index[ans]] = 1

        X.append(x)

        Xq.append(xq)

        Y.append(y)

        return (pad_sequences(X,
maxlen=max_story_len), pad_sequences(Xq,
maxlen=max_ques_len), np.array(Y))


inputs_train, queries_train, answers_train =
vectorize_data(train_data)

inputs_test, queries_test, answers_test =
vectorize_data(test_data)


input_sequence = Input((max_story_len,))
```

```python
question = Input((max_ques_len,))

input_encoder_m = Sequential()

input_encoder_m.add(Embedding(input_dim=vocab_len, output_dim=64))

input_encoder_m.add(Dropout(0.2))


input_encoder_c = Sequential()

input_encoder_c.add(Embedding(input_dim=vocab_len, output_dim=max_ques_len))

input_encoder_c.add(Dropout(0.2))


ques_encoder = Sequential()

ques_encoder.add(Embedding(input_dim=vocab_len, input_length=max_ques_len, output_dim=64))

ques_encoder.add(Dropout(0.2))
```

```
input_encoded_m =
input_encoder_m(input_sequence)

input_encoded_c =
input_encoder_c(input_sequence)

ques_encoded = ques_encoder(question)


match = dot([input_encoded_m,
ques_encoded], axes=(2, 2))

match = Activation('softmax')(match)


response = add([match, input_encoded_c])

response = Permute((2, 1))(response)


answer = Concatenate()([response,
ques_encoded])

answer = LSTM(32)(answer)

answer = Dropout(0.5)(answer)

answer = Dense(vocab_len)(answer)
```

```python
answer = Activation('softmax')(answer)

model = Model([input_sequence, question],
answer)

model.compile(optimizer="rmsprop",
loss="categorical_crossentropy",
metrics=['accuracy'])

history = model.fit([inputs_train,
queries_train], answers_train,
validation_data=([inputs_test, queries_test],
answers_test), batch_size=30, epochs=15)

plt.plot(history.history['loss'])

plt.plot(history.history['val_loss'])

plt.title('Model loss')

plt.ylabel('Loss')

plt.xlabel('Epoch')
```

```
plt.legend(['Train', 'Validation'], loc='upper
left')
```

```
plt.show()
```

## Conclusion:

In conclusion, the provided code implements a question answering model using a neural network architecture. The methodology encompasses data preprocessing, model architecture design, model training, evaluation, and documentation.

Throughout the implementation, the code tokenizes and vectorizes the input data, designs a neural network architecture with appropriate layers such as embedding and LSTM, trains the model with optimization techniques, evaluates its performance on test

data, and documents the entire process for transparency and reproducibility.

The trained model demonstrates the ability to generate accurate answers to questions based on the provided context. By systematically following the methodology outlined in the code, valuable insights can be gained into the effectiveness of the question answering approach and potential avenues for future research and improvement.

Thanks and regards,

Yashas H G