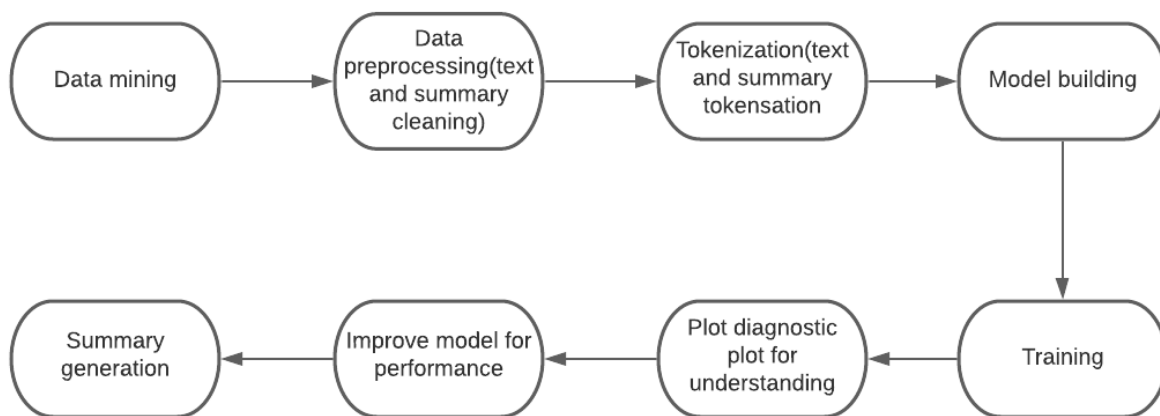


Customer Review Summarizer

Introduction:

Websites such as Amazon and Yelp allow customers to leave reviews for various products. There are usually hundreds of reviews for a single product; each review could be lengthy and repetitive. Therefore automatic review summarization has a huge potential in that it could help customers to make quick decisions on certain products. Summarization is an important challenge of natural language understanding. The aim is to produce a condensed representation of an input text that captures the core meaning of the original.

Conceptual Design:



1. Dataset Mining:

The dataset used consists of reviews of fine foods from Amazon. The data spans a period of more than 10 years, including all ~500,000 reviews up to October 2012. These reviews include product and user information, ratings, plain text review, and summary. It also includes reviews from all other Amazon categories.

2. Data Preprocessing:

Performing basic preprocessing steps is very important before we get to the model building part. Using messy and unclean text data is a potentially disastrous move. So in this step, we will drop all the unwanted symbols, characters, etc. from the text that do not affect the objective of our problem. We'll look at the first 10 rows of the reviews and dataset to get an idea of the preprocessing steps for the summary column.

3. Tokenization:

A tokenizer builds the vocabulary and converts a word sequence to an integer sequence.

4. Model building:

During the model building part, we need to split our dataset into a training and validation set. We'll use 90% of the dataset as the training data and evaluate the performance on the remaining 10% (holdout set). In this step we need to be familiar with the terms such as return sequence, return state, initial state and stacked LSTM.

5. Training:

In the training phase, we will first set up the encoder and decoder. We will then train the model to predict the target sequence offset by one timestep. Let us see in detail on how to set up the encoder and decoder. An Encoder Long Short Term Memory model (LSTM) reads the entire input sequence wherein, at each timestep, one word is fed into the encoder. It then processes the information at every timestep and captures the contextual information present in the input sequence. The decoder is also an LSTM network which reads the entire target sequence word-by-word and predicts the same sequence offset by one timestep. The decoder is trained to predict the next word in the sequence given the previous word.

6. Plot diagnostic:

We will plot a few diagnostic plots to understand the behavior of the model over time. After training, the model is tested on new source sequences for which the target sequence is unknown.

7. Improve model performance:

We will try to increase the training dataset size and build the model. The generalization capability of a deep learning model enhances with an increase in the training dataset size. We can try implementing Bi-Directional LSTM which is capable of capturing the context from both the directions and results in a better context vector

8. Summary generation:

After all the steps are completed we will have the abstractive summary generated. Even though the actual summary and the summary generated by our model do not match in terms of words, both of them are conveying the same meaning. Our model will be able to generate a legible summary based on the context present in the text.

Implementation and Evaluation:

Implementation code :

```
from attention import AttentionLayer
from google.colab import drive
drive.mount('--/content/drive')
```

```

import numpy as np
import pandas as pd
import re
from bs4 import BeautifulSoup
from keras.preprocessing.text import Tokenizer
from keras.preprocessing.sequence import pad_sequences
from nltk.corpus import stopwords
from tensorflow.keras.layers import Input, LSTM, Embedding, Dense, Concatenate,
TimeDistributed
from tensorflow.keras.models import Model
from tensorflow.keras.callbacks import EarlyStopping
import warnings
pd.set_option("display.max_colwidth", 200)
warnings.filterwarnings("ignore")
data=pd.read_csv("/content/drive/MyDrive/NLP//Reviews.csv",nrows=100000)
data.drop_duplicates(subset=['Text'],inplace=True)#dropping duplicates
data.dropna(axis=0,inplace=True)#dropping na
data.info()
<class 'pandas.core.frame.DataFrame'>
Int64Index: 88421 entries, 0 to 99999
Data columns (total 10 columns):
#   Column                Non-Null Count  Dtype
---  ---
0   Id                    88421 non-null  int64
1   ProductId            88421 non-null  object
2   UserId               88421 non-null  object
3   ProfileName          88421 non-null  object
4   HelpfulnessNumerator  88421 non-null  int64
5   HelpfulnessDenominator 88421 non-null  int64
6   Score                88421 non-null  int64
7   Time                 88421 non-null  int64
8   Summary              88421 non-null  object
9   Text                 88421 non-null  object
dtypes: int64(5), object(5)
memory usage: 7.4+ MB
contraction_mapping = {"ain't": "is not", "aren't": "are not", "can't": "cannot", "'cause":
"because", "could've": "could have", "couldn't": "could not",
                        "didn't": "did not", "doesn't": "does not", "don't": "do not", "hadn't": "had not",
                        "hasn't": "has not", "haven't": "have not",

```

"he'd": "he would", "he'll": "he will", "he's": "he is", "how'd": "how did",
"how'd'y": "how do you", "how'll": "how will", "how's": "how is",
"I'd": "I would", "I'd've": "I would have", "I'll": "I will", "I'll've": "I will have",
"I'm": "I am", "I've": "I have", "i'd": "i would",
"i'd've": "i would have", "i'll": "i will", "i'll've": "i will have", "i'm": "i am",
"i've": "i have", "isn't": "is not", "it'd": "it would",
"it'd've": "it would have", "it'll": "it will", "it'll've": "it will have", "it's": "it is",
"let's": "let us", "ma'am": "madam",
"mayn't": "may not", "might've": "might have", "mightn't": "might not",
"mightn't've": "might not have", "must've": "must have",
"mustn't": "must not", "mustn't've": "must not have", "needn't": "need not",
"needn't've": "need not have", "o'clock": "of the clock",
"oughtn't": "ought not", "oughtn't've": "ought not have", "shan't": "shall not",
"sha'n't": "shall not", "shan't've": "shall not have",
"she'd": "she would", "she'd've": "she would have", "she'll": "she will",
"she'll've": "she will have", "she's": "she is",
"should've": "should have", "shouldn't": "should not", "shouldn't've": "should not have",
"so've": "so have", "so's": "so as",
"this's": "this is", "that'd": "that would", "that'd've": "that would have", "that's": "that is",
"there'd": "there would",
"there'd've": "there would have", "there's": "there is", "here's": "here is",
"they'd": "they would", "they'd've": "they would have",
"they'll": "they will", "they'll've": "they will have", "they're": "they are",
"they've": "they have", "to've": "to have",
"wasn't": "was not", "we'd": "we would", "we'd've": "we would have", "we'll": "we will",
"we'll've": "we will have", "we're": "we are",
"we've": "we have", "weren't": "were not", "what'll": "what will", "what'll've": "what will have",
"what're": "what are",
"what's": "what is", "what've": "what have", "when's": "when is", "when've": "when have",
"where'd": "where did", "where's": "where is",
"where've": "where have", "who'll": "who will", "who'll've": "who will have",
"who's": "who is", "who've": "who have",
"why's": "why is", "why've": "why have", "will've": "will have", "won't": "will not",
"won't've": "will not have",
"would've": "would have", "wouldn't": "would not", "wouldn't've": "would not have",
"y'all": "you all",
"y'all'd": "you all would", "y'all'd've": "you all would have", "y'all're": "you all are",
"y'all've": "you all have",
"you'd": "you would", "you'd've": "you would have", "you'll": "you will",
"you'll've": "you will have",

```

        "you're": "you are", "you've": "you have"}

import nltk
nltk.download('stopwords')
stop_words = set(stopwords.words('english'))

def text_cleaner(text,num):
    newString = text.lower()
    newString = BeautifulSoup(newString, "lxml").text
    newString = re.sub(r"([\^])*\\", "", newString)
    newString = re.sub("'", "", newString)
    newString = ''.join([contraction_mapping[t] if t in contraction_mapping else t for t in
newString.split(" ")])
    newString = re.sub(r"\\s\\b", "",newString)
    newString = re.sub("[^a-zA-Z]", " ", newString)
    newString = re.sub('[m]{2,}', 'mm', newString)
    if(num==0):
        tokens = [w for w in newString.split() if not w in stop_words]
    else:
        tokens=newString.split()
    long_words=[]
    for i in tokens:
        if len(i)>1:                                #removing short word
            long_words.append(i)
    return (" ".join(long_words)).strip()

[nltk_data] Downloading package stopwords to /root/nltk_data...
[nltk_data] Unzipping corpora/stopwords.zip.
#call the function
cleaned_text = []
for t in data['Text']:
    cleaned_text.append(text_cleaner(t,0))
#call the function
cleaned_summary = []
for t in data['Summary']:
    cleaned_summary.append(text_cleaner(t,1))
data['cleaned_text']=cleaned_text
data['cleaned_summary']=cleaned_summary
data.replace("", np.nan, inplace=True)
data.dropna(axis=0,inplace=True)
import matplotlib.pyplot as plt

```

```

text_word_count = []
summary_word_count = []

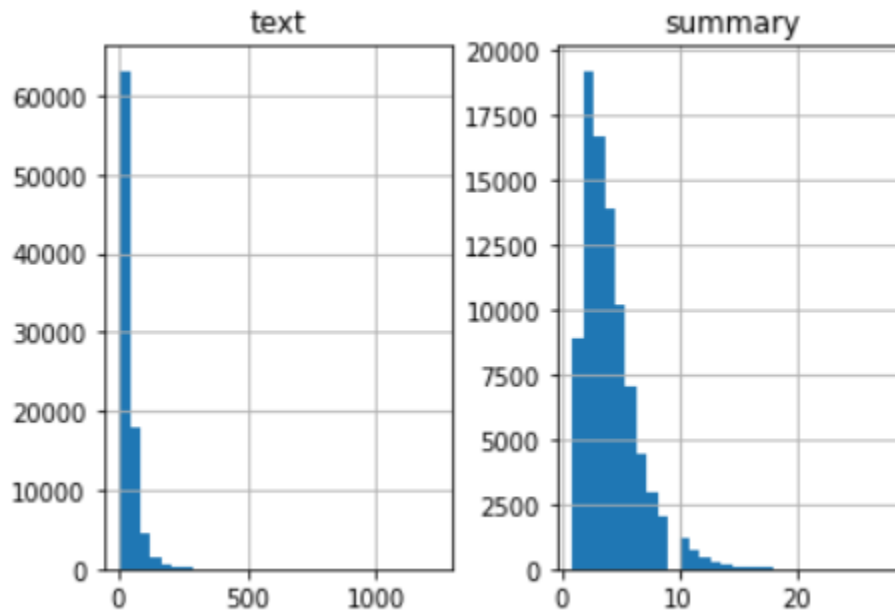
# populate the lists with sentence lengths
for i in data['cleaned_text']:
    text_word_count.append(len(i.split()))

for i in data['cleaned_summary']:
    summary_word_count.append(len(i.split()))

length_df = pd.DataFrame({'text':text_word_count, 'summary':summary_word_count})

length_df.hist(bins = 30)
plt.show()

```



```

cnt=0
for i in data['cleaned_summary']:
    if(len(i.split())<=8):
        cnt=cnt+1
print(cnt/len(data['cleaned_summary']))
0.9424907471335922
max_text_len=30
max_summary_len=8
cleaned_text=np.array(data['cleaned_text'])
cleaned_summary=np.array(data['cleaned_summary'])

```

```
short_text=[]
short_summary=[]
```

```
for i in range(len(cleaned_text)):
    if(len(cleaned_summary[i].split())<=max_summary_len and
len(cleaned_text[i].split())<=max_text_len):
        short_text.append(cleaned_text[i])
        short_summary.append(cleaned_summary[i])
```

```
df=pd.DataFrame({'text':short_text,'summary':short_summary})
df['summary']=df['summary'].apply(lambda x : 'sostok '+ x + ' eostok')
from sklearn.model_selection import train_test_split
x_tr,x_val,y_tr,y_val=train_test_split(np.array(df['text']),np.array(df['summary']),test_size=0.1,ra
ndom_state=0,shuffle=True)
from keras.preprocessing.text import Tokenizer
from keras.preprocessing.sequence import pad_sequences
```

```
#prepare a tokenizer for reviews on training data
x_tokenizer = Tokenizer()
x_tokenizer.fit_on_texts(list(x_tr))
thresh=4
```

```
cnt=0
tot_cnt=0
freq=0
tot_freq=0
```

```
for key,value in x_tokenizer.word_counts.items():
    tot_cnt=tot_cnt+1
    tot_freq=tot_freq+value
    if(value<thresh):
        cnt=cnt+1
        freq=freq+value
```

```
print("% of rare words in vocabulary:",(cnt/tot_cnt)*100)
print("Total Coverage of rare words:",(freq/tot_freq)*100)
% of rare words in vocabulary: 66.12339930151339
Total Coverage of rare words: 2.953684513790566
#prepare a tokenizer for reviews on training data
x_tokenizer = Tokenizer(num_words=tot_cnt-cnt)
```

```
x_tokenizer.fit_on_texts(list(x_tr))
```

```
#convert text sequences into integer sequences
```

```
x_tr_seq = x_tokenizer.texts_to_sequences(x_tr)
```

```
x_val_seq = x_tokenizer.texts_to_sequences(x_val)
```

```
#padding zero upto maximum length
```

```
x_tr = pad_sequences(x_tr_seq, maxlen=max_text_len, padding='post')
```

```
x_val = pad_sequences(x_val_seq, maxlen=max_text_len, padding='post')
```

```
#size of vocabulary ( +1 for padding token)
```

```
x_voc = x_tokenizer.num_words + 1
```

```
x_voc
```

```
8440
```

```
#prepare a tokenizer for reviews on training data
```

```
y_tokenizer = Tokenizer()
```

```
y_tokenizer.fit_on_texts(list(y_tr))
```

```
thresh=6
```

```
cnt=0
```

```
tot_cnt=0
```

```
freq=0
```

```
tot_freq=0
```

```
for key,value in y_tokenizer.word_counts.items():
```

```
    tot_cnt=tot_cnt+1
```

```
    tot_freq=tot_freq+value
```

```
    if(value<thresh):
```

```
        cnt=cnt+1
```

```
        freq=freq+value
```

```
print("% of rare words in vocabulary:",(cnt/tot_cnt)*100)
```

```
print("Total Coverage of rare words:",(freq/tot_freq)*100)
```

```
% of rare words in vocabulary: 78.12740675541863
```

```
Total Coverage of rare words: 5.3921899389571895
```

```
#prepare a tokenizer for reviews on training data
```

```
y_tokenizer = Tokenizer(num_words=tot_cnt-cnt)
```

```
y_tokenizer.fit_on_texts(list(y_tr))
```

```
#convert text sequences into integer sequences
```



```
y_tr_seq = y_tokenizer.texts_to_sequences(y_tr)
y_val_seq = y_tokenizer.texts_to_sequences(y_val)
```

```
#padding zero upto maximum length
```

```
y_tr = pad_sequences(y_tr_seq, maxlen=max_summary_len, padding='post')
y_val = pad_sequences(y_val_seq, maxlen=max_summary_len, padding='post')
```

```
#size of vocabulary
```

```
y_voc = y_tokenizer.num_words + 1
ind=[]
```

```
for i in range(len(y_tr)):
```

```
    cnt=0
```

```
    for j in y_tr[i]:
```

```
        if j!=0:
```

```
            cnt=cnt+1
```

```
    if(cnt==2):
```

```
        ind.append(i)
```

```
y_tr=np.delete(y_tr,ind, axis=0)
```

```
x_tr=np.delete(x_tr,ind, axis=0)
```

```
ind=[]
```

```
for i in range(len(y_val)):
```

```
    cnt=0
```

```
    for j in y_val[i]:
```

```
        if j!=0:
```

```
            cnt=cnt+1
```

```
    if(cnt==2):
```

```
        ind.append(i)
```

```
y_val=np.delete(y_val,ind, axis=0)
```

```
x_val=np.delete(x_val,ind, axis=0)
```

```
from keras import backend as K
```

```
K.clear_session()
```

```
latent_dim = 300
```

```
embedding_dim=100
```

```
# Encoder
```

```
encoder_inputs = Input(shape=(max_text_len,))
```

#embedding layer

```
enc_emb = Embedding(x_voc, embedding_dim, trainable=True)(encoder_inputs)
```

#encoder lstm 1

```
encoder_lstm1 =
```

```
LSTM(latent_dim, return_sequences=True, return_state=True, dropout=0.4, recurrent_dropout=0.4  
)
```

```
encoder_output1, state_h1, state_c1 = encoder_lstm1(enc_emb)
```

#encoder lstm 2

```
encoder_lstm2 =
```

```
LSTM(latent_dim, return_sequences=True, return_state=True, dropout=0.4, recurrent_dropout=0.4  
)
```

```
encoder_output2, state_h2, state_c2 = encoder_lstm2(encoder_output1)
```

#encoder lstm 3

```
encoder_lstm3=LSTM(latent_dim, return_state=True,
```

```
return_sequences=True, dropout=0.4, recurrent_dropout=0.4)
```

```
encoder_outputs, state_h, state_c= encoder_lstm3(encoder_output2)
```

Set up the decoder, using `encoder_states` as initial state.

```
decoder_inputs = Input(shape=(None,))
```

#embedding layer

```
dec_emb_layer = Embedding(y_voc, embedding_dim, trainable=True)
```

```
dec_emb = dec_emb_layer(decoder_inputs)
```

```
decoder_lstm = LSTM(latent_dim, return_sequences=True,
```

```
return_state=True, dropout=0.4, recurrent_dropout=0.2)
```

```
decoder_outputs, decoder_fwd_state, decoder_back_state =
```

```
decoder_lstm(dec_emb, initial_state=[state_h, state_c])
```

Attention layer

```
attn_layer = AttentionLayer(name='attention_layer')
```

```
attn_out, attn_states = attn_layer([encoder_outputs, decoder_outputs])
```

Concat attention input and decoder LSTM output

```
decoder_concat_input = Concatenate(axis=-1, name='concat_layer')([decoder_outputs, attn_out])
```

#dense layer

```
decoder_dense = TimeDistributed(Dense(y_voc, activation='softmax'))
decoder_outputs = decoder_dense(decoder_concat_input)
```

```
# Define the model
```

```
model = Model([encoder_inputs, decoder_inputs], decoder_outputs)
```

```
model.summary()
```

```
WARNING:tensorflow:Layer lstm will not use cuDNN kernel since it doesn't meet the cuDNN
kernel criteria. It will use generic GPU kernel as fallback when running on GPU
```

```
WARNING:tensorflow:Layer lstm_1 will not use cuDNN kernel since it doesn't meet the
cuDNN kernel criteria. It will use generic GPU kernel as fallback when running on GPU
```

```
WARNING:tensorflow:Layer lstm_2 will not use cuDNN kernel since it doesn't meet the
cuDNN kernel criteria. It will use generic GPU kernel as fallback when running on GPU
```

```
WARNING:tensorflow:Layer lstm_3 will not use cuDNN kernel since it doesn't meet the
cuDNN kernel criteria. It will use generic GPU kernel as fallback when running on GPU
```

```
Model: "model"
```

Layer (type)	Output Shape	Param #	Connected to
input_1 (InputLayer)	[(None, 30)]	0	
embedding (Embedding)	(None, 30, 100)	844000	input_1[0][0]
lstm (LSTM)	[(None, 30, 300), (N 481200		embedding[0][0]
input_2 (InputLayer)	[(None, None)]	0	
lstm_1 (LSTM)	[(None, 30, 300), (N 721200		lstm[0][0]
embedding_1 (Embedding)	(None, None, 100)	198900	input_2[0][0]
lstm_2 (LSTM)	[(None, 30, 300), (N 721200		lstm_1[0][0]

lstm_3 (LSTM)	[(None, None, 300), 481200	embedding_1[0][0] lstm_2[0][1] lstm_2[0][2]
---------------	----------------------------	---

attention_layer (AttentionLayer ((None, None, 300), 180300	lstm_2[0][0] lstm_3[0][0]
--	------------------------------

concat_layer (Concatenate)	(None, None, 600) 0	lstm_3[0][0] attention_layer[0][0]
----------------------------	---------------------	---------------------------------------

time_distributed (TimeDistribut (None, None, 1989) 1195389	concat_layer[0][0]
--	--------------------

=====
Total params: 4,823,389

Trainable params: 4,823,389

Non-trainable params: 0

```
model.compile(optimizer='rmsprop', loss='sparse_categorical_crossentropy')
es = EarlyStopping(monitor='val_loss', mode='min', verbose=1, patience=2)
history=model.fit([x_tr,y_tr[:, :-1]], y_tr.reshape(y_tr.shape[0],y_tr.shape[1], 1)[:,:1:],
,epochs=50,callbacks=[es],batch_size=128, validation_data=([x_val,y_val[:, :-1]],
y_val.reshape(y_val.shape[0],y_val.shape[1], 1)[:,:1:])))
```

Encode the input sequence to get the feature vector

```
encoder_model = Model(inputs=encoder_inputs,outputs=[encoder_outputs, state_h, state_c])
```

Decoder setup

Below tensors will hold the states of the previous time step

```
decoder_state_input_h = Input(shape=(latent_dim,))
```

```
decoder_state_input_c = Input(shape=(latent_dim,))
```

```
decoder_hidden_state_input = Input(shape=(max_text_len,latent_dim))
```

Get the embeddings of the decoder sequence

```
dec_emb2= dec_emb_layer(decoder_inputs)
```

To predict the next word in the sequence, set the initial states to the states from the previous time step

```
decoder_outputs2, state_h2, state_c2 = decoder_lstm(dec_emb2,  
initial_state=[decoder_state_input_h, decoder_state_input_c])
```

```
#attention inference
```

```
attn_out_inf, attn_states_inf = attn_layer([decoder_hidden_state_input, decoder_outputs2])  
decoder_inf_concat = Concatenate(axis=-1, name='concat')([decoder_outputs2, attn_out_inf])
```

```
# A dense softmax layer to generate prob dist. over the target vocabulary
```

```
decoder_outputs2 = decoder_dense(decoder_inf_concat)
```

```
# Final decoder model
```

```
decoder_model = Model(  
    [decoder_inputs] + [decoder_hidden_state_input, decoder_state_input_h,  
    decoder_state_input_c],  
    [decoder_outputs2] + [state_h2, state_c2])
```

```
def decode_sequence(input_seq):
```

```
    # Encode the input as state vectors.
```

```
    e_out, e_h, e_c = encoder_model.predict(input_seq)
```

```
    # Generate empty target sequence of length 1.
```

```
    target_seq = np.zeros((1,1))
```

```
    # Populate the first word of target sequence with the start word.
```

```
    target_seq[0, 0] = target_word_index['sostok']
```

```
    stop_condition = False
```

```
    decoded_sentence = "
```

```
    while not stop_condition:
```

```
        output_tokens, h, c = decoder_model.predict([target_seq] + [e_out, e_h, e_c])
```

```
    # Sample a token
```

```
    sampled_token_index = np.argmax(output_tokens[0, -1, :])
```

```
    sampled_token = reverse_target_word_index[sampled_token_index]
```

```
    if(sampled_token!='eostok'):
```

```
        decoded_sentence += ' '+sampled_token
```

```
    # Exit condition: either hit max length or find stop word.
```

```
    if (sampled_token == 'eostok' or len(decoded_sentence.split()) >= (max_summary_len-1)):
```

```

        stop_condition = True

    # Update the target sequence (of length 1).
    target_seq = np.zeros((1,1))
    target_seq[0, 0] = sampled_token_index

    # Update internal states
    e_h, e_c = h, c

    return decoded_sentence
def seq2summary(input_seq):
    newString=""
    for i in input_seq:
        if((i!=0 and i!=target_word_index['sostok']) and i!=target_word_index['eostok']):
            newString=newString+reverse_target_word_index[i]+' '
    return newString

def seq2text(input_seq):
    newString=""
    for i in input_seq:
        if(i!=0):
            newString=newString+reverse_source_word_index[i]+' '
    return newString
for i in range(0,5):
    print("Review:",seq2text(x_tr[i]))
    print("Original summary:",seq2summary(y_tr[i]))
    print("Predicted summary:",decode_sequence(x_tr[i].reshape(1,max_text_len)))
    print("\n")

```

Review: gave caffeine shakes heart anxiety attack plus tastes unbelievably bad stick coffee tea
soda thanks

Original summary: hour

Predicted summary: green tea

Review: got great course good belgian chocolates better

Original summary: would like to give it stars but

Predicted summary: delicious

Review: one best flavored coffees tried usually like flavored coffees one great serve company love

Original summary: delicious

Predicted summary: great coffee

Review: salt separate area pain makes hard regulate salt putting like salt go ahead get product

Original summary: tastes ok packaging

Predicted summary: salt

Review: really like product super easy order online delivered much cheaper buying gas station stocking good long drives

Original summary: turkey jerky is great

Predicted summary: great