

Requirement already satisfied: tensorflow in /usr/local/lib/python3.7/dist-packages (2.7.0)
Requirement already satisfied: protobuf<3.9.2 in /usr/local/lib/python3.7/dist-packages (from tensorflow) (3.17.3)
Requirement already satisfied: keras<2.8.1,>=2.7.0rc0 in /usr/local/lib/python3.7/dist-packages (from tensorflow) (2.7.0)
Requirement already satisfied: wrapt<1.11.0 in /usr/local/lib/python3.7/dist-packages (from tensorflow) (1.13.3)
Requirement already satisfied: h5py<2.9.0 in /usr/local/lib/python3.7/dist-packages (from tensorflow) (3.1.0)
Requirement already satisfied: termcolor<=1.1.0 in /usr/local/lib/python3.7/dist-packages (from tensorflow) (1.1.0)
Requirement already satisfied: flatbuffers<3.0,>=1.12 in /usr/local/lib/python3.7/dist-packages (from tensorflow) (2.7.0)
Requirement already satisfied: tensorboard<2.6 in /usr/local/lib/python3.7/dist-packages (from tensorflow) (2.7.0)
Requirement already satisfied: numpy<=1.14.5 in /usr/local/lib/python3.7/dist-packages (from tensorflow) (1.19.5)
Requirement already satisfied: typing-extensions<=3.6.6 in /usr/local/lib/python3.7/dist-packages (from tensorflow) (3.10.0)
Requirement already satisfied: gast<0.5.0,>=0.2.1 in /usr/local/lib/python3.7/dist-packages (from tensorflow) (0.4.0)
Requirement already satisfied: astunparse<=1.6.0 in /usr/local/lib/python3.7/dist-packages (from tensorflow) (1.6.3)
Requirement already satisfied: tensorflow-io-gcs-filesystem<=0.21.0 in /usr/local/lib/python3.7/dist-packages (from tensorflow) (0.22.0)
Requirement already satisfied: libclang<=9.0.1 in /usr/local/lib/python3.7/dist-packages (from tensorflow) (12.0.8)
Requirement already satisfied: tensorflow-estimator<2.8,->2.7.0rc0 in /usr/local/lib/python3.7/dist-packages (from tensorflow) (2.7.0)
Requirement already satisfied: google-auth<=1.6.3 in /usr/local/lib/python3.7/dist-packages (from tensorflow) (3.3.0)
Requirement already satisfied: six<=1.12.0 in /usr/local/lib/python3.7/dist-packages (from tensorflow) (1.15.0)
Requirement already satisfied: keras-preprocessing<=1.1.1 in /usr/local/lib/python3.7/dist-packages (from tensorflow) (1.1.2)
Requirement already satisfied: absl-py<0.4.0,>=0.2.1 in /usr/local/lib/python3.7/dist-packages (from tensorflow) (0.12.0)
Requirement already satisfied: wheel<1.0,>=0.32.0 in /usr/local/lib/python3.7/dist-packages (from tensorflow) (0.2.0)
Requirement already satisfied: grpcio<2.0.0,>=1.24.3 in /usr/local/lib/python3.7/dist-packages (from tensorflow) (1.42.0)
Requirement already satisfied: cached-property in /usr/local/lib/python3.7/dist-packages (from h5py<=2.9.0) (1.5.2)
Requirement already satisfied: google-auth<0.5,>=0.4.1 in /usr/local/lib/python3.7/dist-packages (from tensorflow) (2.7.0)
Requirement already satisfied: tensorboard-data-server<0.7.0,>=0.6.0 in /usr/local/lib/python3.7/dist-packages (from tensorflow) (0.6.1)
Requirement already satisfied: google-auth-oauthlib<0.6.0,>=0.4.6 in /usr/local/lib/python3.7/dist-packages (from tensorflow) (0.4.6)
Requirement already satisfied: requests<3,>=2.21.0 in /usr/local/lib/python3.7/dist-packages (from tensorflow) (2.25.0)
Requirement already satisfied: tensorflow-plugin-api<=1.6.0 in /usr/local/lib/python3.7/dist-packages (from tensorflow) (1.0.1)
Requirement already satisfied: cachedools<5.0,>=2.0.0 in /usr/local/lib/python3.7/dist-packages (from tensorflow) (2.6.6-tensorflow) (3.3.6)
Requirement already satisfied: setuptools<=41.0.8 in /usr/local/lib/python3.7/dist-packages (from tensorflow) (52.0.6-tensorflow) (1.0.1)
Requirement already satisfied: rx<=0.2.1 in /usr/local/lib/python3.7/dist-packages (from tensorflow) (2.6.3-tensorboard) (0.2.8)
Requirement already satisfied: pyasn1<0.6.0,>=0.4.6 in /usr/local/lib/python3.7/dist-packages (from google-auth) (0.5.1)
Requirement already satisfied: rsa<4.5,>=3.1.4 in /usr/local/lib/python3.7/dist-packages (from google-auth) (4.7.2)
Requirement already satisfied: cachetools<5.0,>=2.0.0 in /usr/local/lib/python3.7/dist-packages (from google-auth) (4.2.4)
Requirement already satisfied: requests<0.5,>=0.7.0 in /usr/local/lib/python3.7/dist-packages (from google-auth-oauthlib) (2.25.0)
Requirement already satisfied: importlib-metadata<=4.4 in /usr/local/lib/python3.7/dist-packages (from google-auth-oauthlib) (4.8.2)
Requirement already satisfied: importlib-metadata<=4.4 in /usr/local/lib/python3.7/dist-packages (from google-auth-oauthlib) (4.8.2)
Requirement already satisfied: pyasn1<0.6.0,>=0.4.6 in /usr/local/lib/python3.7/dist-packages (from pyasn1-modules) (0.5.1)
Requirement already satisfied: certifi<=2021.4.17 in /usr/local/lib/python3.7/dist-packages (from requests) (2021.10.8)
Requirement already satisfied: idna<3.0,>=2.5 in /usr/local/lib/python3.7/dist-packages (from requests) (3.0.2)
Requirement already satisfied: urllib3<1.25.0,>=1.25.1,!=1.25.1,!=1.26,!=1.26.1 in /usr/local/lib/python3.7/dist-packages (from requests) (1.24.3)
Requirement already satisfied: chardet<4.0,>=3.0.2 in /usr/local/lib/python3.7/dist-packages (from requests) (3.0.4)
Requirement already satisfied: oauthlib<3.0.0 in /usr/local/lib/python3.7/dist-packages (from requests-oauthlib) (3.1.1)
Requirement already satisfied: oauthlib<3.0.0,>=3.0.1 in /usr/local/lib/python3.7/dist-packages (from requests-oauthlib) (3.1.1)

```
In [ ]:

#title
#Importing libraries

import pandas as pd
from tensorflow.keras.preprocessing.text import Tokenizer
from tensorflow.keras.preprocessing.sequence import pad_sequences
import numpy as np
from tensorflow.keras.models import Model, Sequential
from tensorflow.keras.layers import GRU, Input, Dense, TimeDistributed, Activation, RepeatVector, Bidirectional
from tensorflow.keras.losses import sparse_categorical_crossentropy
from tensorflow.keras.optimizers import Adam
from sklearn.model_selection import train_test_split
from tensorflow.keras.layers import Dense, LSTM, Embedding, RepeatVector
from tensorflow.keras.callbacks import ModelCheckpoint

In [ ]:

#title
#reading the train dataset

#english corpus
eng_train=pd.read_csv('/content/drive/MyDrive/Colab Notebooks/PARALLEL-CORPUS (TRAIN)/IITB-en-hi-en.csv')

#hindi corpus
hi_train=pd.read_csv('/content/drive/MyDrive/Colab Notebooks/PARALLEL-CORPUS (TRAIN)/IITB-en-hi-hi.csv')

#reading the validation dataset

#english corpus
eng_val=pd.read_csv('/content/drive/MyDrive/Colab Notebooks/PARALLEL-CORPUS (VALIDATION)/validation-en.csv')

#hindi corpus
hi_val=pd.read_csv('/content/drive/MyDrive/Colab Notebooks/PARALLEL-CORPUS (VALIDATION)/validation-hi.csv')

#reading the test dataset

#english corpus
eng_test=pd.read_csv('/content/drive/MyDrive/Colab Notebooks/PARALLEL-CORPUS (TEST)/test-en.csv')

#hindi corpus
hi_test=pd.read_csv('/content/drive/MyDrive/Colab Notebooks/PARALLEL-CORPUS (TEST)/test-hi.csv')
```

```
COMBINING BOTH THE LANGUAGE PAIRS TO FORM PARALLEL CORPUS OF TRAIN, TEST AND VALIDATION

In [ ]:

#title
#TRAINING DATA

df_train=pd.concat([eng_train,hi_train],axis=1)

print("*****13,\"NTRAINING DATA\"n\",sep='')

df_train.head()

=====
TRAINING DATA

Out[ ]:

ENGLISH CORPUS

HINDI CORPUS

0 Give your application an accessibility workout. अपने अनुप्रयोग को पहुंचनीयता सुधारें ता आप इसे
1 Accersicer Accessibility Explorer. एक्सेसिबिलिटी एक्सेप्लोरर
2 The default plugin layout for the bottom panel. निम्नलेखन पैनल के लिए डिफ़ॉल्ट प्लगइन
3 The default plugin layout for the top panel. ऊपरी पैनल के लिए डिफ़ॉल्ट प्लगइन
4 A list of plugins that are disabled by default. उन प्लगइन्स की सूची जिन्हें डिफ़ॉल्ट रूप से नि...
```

```
In [ ]:

#title
#VALIDATION DATA

df_val=pd.concat([eng_val,hi_val],axis=1)

print("*****13,\"NVALIDATION DATA\"n\",sep='')

df_val.head()

=====
VALIDATION DATA

Out[ ]:

ENGLISH CORPUS

HINDI CORPUS

0 Students of the Dattatreya City Municipal corp... दत्तत्रेया शहर के छात्रों का स्थानीय नगरपालिका
1 With encouragement from Principal Sandhya Medp... प्रधानाचार्य के प्रोत्साहन से
2 Rajesh Gaurav the President of the MNPA teacher... राजेश गौरव संघ के अध्यक्ष के रूप में चुने गए
3 Ramesh Saapute examined the fort. रामेश सापुटे किला की जाँच
4 Students like Nikhil Kavle Darshan Godekar Sah... निखिल कावले जैसे विद्यार्थी किला की जाँच करते हैं।
```

```
In [ ]:

#title
#TESTING DATA

df_test=pd.concat([eng_test,hi_test],axis=1)

print("*****13,\"NTESTING DATA\"n\",sep='')

df_test.head()

=====
TESTING DATA

Out[ ]:

ENGLISH CORPUS

HINDI CORPUS

0 A black box in your car? कार में काला बॉक्स?
1 As America's road planners struggle to find th... क्योंकि अमेरिका के राज्दाय योजनाकारों को यह खोजने में
2 The devices which track every mile a motorist ... वह डिवाइस हैं जो मोटरवाहन द्वारा चले गए प्रत्येक मील
3 The usually dull arena of highway planning has... आमतौर पर खाली सड़क नियोजन का क्षेत्र अब रोमांचक
4 Libertarians have joined environmental groups ... स्वतंत्रतावादी लोग पर्यावरण संरक्षण समूहों में शामिल हुए...
```

PREPROCESSING THE DATASET

TOKENIZE

For a neural network to predict on text data, it first has to be turned into data it can understand. Text data like "dog" is a sequence of character encodings. Since a neural network is a series of multiplication and addition operations, the input data needs to be numbers.

We can turn each character into a number or each word into a number. These are called character and word ids, respectively. Character ids are used for character level models that generate text predictions for each character. A word level model uses word ids that generate text predictions for each word. Word level models tend to learn better, since they are lower in complexity, so we'll use those.

Turn each sentence into a sequence of word ids using Keras's Tokenizer function. Use this function to tokenize english_sentences and french_sentences in the cell below.

```
In [ ]:

#title
def tokenization(lines):
    tokenizer = Tokenizer()
    tokenizer.fit_on_texts(lines)
    return tokenizer

In [ ]:

#title
# train data

# prepare english train tokenizer
eng_train_tokenizer = tokenization(df_train["ENGLISH CORPUS"])
eng_train_vocab_size = len(eng_train_tokenizer.word_index) + 1

eng_train_length = 8
print('English Train Vocabulary Size: %d' % eng_train_vocab_size)

# validation data
eng_val_tokenizer = tokenization(df_val["ENGLISH CORPUS"])
eng_val_vocab_size = len(eng_val_tokenizer.word_index) + 1

eng_val_length = 8
print('English Validation Vocabulary Size: %d' % eng_val_vocab_size)

English Train Vocabulary Size: 21729
English Validation Vocabulary Size: 2379

In [ ]:

#title
# train data

# prepare hindi train tokenizer
hi_train_tokenizer = tokenization(df_train["HINDI CORPUS"])
hi_train_vocab_size = len(hi_train_tokenizer.word_index) + 1

hi_train_length = 8
print('Hindi Train Vocabulary Size: %d' % hi_train_vocab_size)

# validation data
hi_val_tokenizer = tokenization(df_val["HINDI CORPUS"])
hi_val_vocab_size = len(hi_val_tokenizer.word_index) + 1

hi_val_length = 8
print('Hindi Validation Vocabulary Size: %d' % hi_val_vocab_size)

Hindi Train Vocabulary Size: 33823
Hindi Validation Vocabulary Size: 2672

In [ ]:

#title
df['HINDI_TOKENS'] = df.swifter.apply(lambda row: hi_tokenizer(row['HINDI CORPUS']), axis=1)
```

Padding

When batching the sequence of word ids together, each sequence needs to be the same length. Since sentences are dynamic in length, we can add padding to the end of the sequences to make them the same length.

Make sure all the English sequences have the same length and all the Hindi sequences have the same length by adding padding to the end of each sequence using Keras function.

```
In [ ]:

#title
# encode and pad sequences
def encode_sequences(tokenizer, length, lines):
    seq = tokenizer.texts_to_sequences(lines)
    # pad sequences with 0 values
    seq = pad_sequences(seq, maxlen=length, padding='post')
    return seq

In [ ]:

#title
# prepare training data
train_x = encode_sequences(eng_train_tokenizer, eng_train_length, df_train["ENGLISH CORPUS"])
train_y = encode_sequences(hi_train_tokenizer, hi_train_length, df_train["HINDI CORPUS"])

# prepare validation data
val_x = encode_sequences(eng_val_tokenizer, eng_val_length, df_val["ENGLISH CORPUS"])
val_y = encode_sequences(hi_val_tokenizer, hi_val_length, df_val["HINDI CORPUS"])

In [ ]:

#title
def count(series):
    l = []
    for lis in series:
        for words in lis:
            if words not in l:
                l.append(words)
    return l

In [ ]:

#title
def generate_batch(X = df_train["ENGLISH CORPUS"], y = df_train["HINDI CORPUS"], batch_size = 128):
    """Generate a batch of data"""
    while True:
        X=df_train["ENGLISH CORPUS"]
        y=df_train["HINDI CORPUS"]
        for j in range(0, len(X), batch_size):
            encoder_input_data = np.zeros((batch_size, max_length_src),dtype='float32')
            decoder_input_data = np.zeros((batch_size, max_length_tar),dtype='float32')
            decoder_target_data = np.zeros((batch_size, max_length_tar, num_decoder_tokens),dtype='float32')
            for i, (input_text, target_text) in enumerate(zip(X[j:batch_size], y[j:batch_size])):
                for t, word in enumerate(input_text.split()):
                    encoder_input_data[i][t] = input_token_index[word]
                if word in input_token_index.keys():
                    encoder_input_data[i, t] = input_token_index[word] # encoder input seq
            for t, word in enumerate(target_text.split()):
                if word in target_token_index.keys():
                    decoder_input_data[i, t-1] = target_token_index[word] # decoder input seq
                    if t>0:
                        decoder_target_data[i, t] = target_token_index[word] # decoder target seq
                    # does not include the START token
                    # Offset by one timestep
                    decoder_target_data[i, t - 1, target_token_index[word]-1] = 1.
            yield((encoder_input_data, decoder_input_data), decoder_target_data)
```

```
In [ ]:

#title
latent_dim=300
```

```
In [ ]:

#title
# Encoder
encoder_inputs = Input(shape=(None,))
enc_emb = Embedding(count(df['ENGLISH TOKENS']), latent_dim, mask_zero = True)(encoder_inputs)
encoder_lstm = LSTM(latent_dim, return_state=True)
encoder_outputs, state_h, state_c = encoder_lstm(enc_emb, return_sequences=True, return_state=True)
# We discard encoder_outputs, as we only keep the states.
encoder_states = [state_h, state_c]
```

```
In [ ]:

#title
# Set up the decoder, using 'encoder_states' as initial state.
decoder_inputs = Input(shape=(None,))
dec_emb_layer = Embedding(count(df['HINDI_TOKENS']), latent_dim, mask_zero = True)
dec_emb = dec_emb_layer(decoder_inputs)
# We set up our decoder to return full output sequences,
# and to return internal states as well. We don't use the
# return states in the training model, but we will use them in inference.
decoder_lstm = LSTM(latent_dim, return_sequences=True, return_state=True)
decoder_outputs, _, _ = decoder_lstm(dec_emb, initial_state=encoder_states)
decoder_output2, _, _ = decoder_lstm(decoder_outputs, activation='softmax')
decoder_outputs = decoder_output2

# Define the model that will turn
# 'encoder_input_data' & 'decoder_input_data' into 'decoder_output_data'
model = Model([encoder_inputs, decoder_input_data], decoder_outputs)
```

```
In [ ]:

#title
model.compile(optimizer='rmsprop', loss='categorical_crossentropy')
```

```
In [ ]:

#title
model.summary()
```

Layer (type)	Output Shape	Param #	Connected to
input_1 (InputLayer)	[(None, None)]	0	[]
input_2 (InputLayer)	[(None, None)]	0	[]
embedding (Embedding)	(None, None, 300)	875400	['input_1[0][0]']
embedding_1 (Embedding)	(None, None, 300)	11239200	['input_2[0][0]']
lstm (LSTM)	[(None, 300), (None, 300), (None, 300)]	721200	['embedding[0][0]']
lstm_1 (LSTM)	[(None, None, 300), (None, 300), (None, 300)]	721200	['embedding_1[0][0]', 'lstm[0][1]', 'lstm[0][2]']
dense (Dense)	(None, None, 37464)	11276664	['lstm_1[0][0]']

=====

Total params: 32,093,664
Trainable params: 32,093,664
Non-trainable params: 0

```
In [ ]:

#title
train_samples = len(df_train["ENGLISH CORPUS"])
val_samples = len(df_val["ENGLISH CORPUS"])
batch_size = 128
epochs = 100

In [ ]:

#title
max_length_src=max(df["ENGLISH SENTENCE LENGTH"])
max_length_tar=max(df["HINDI SENTENCE LENGTH"])
num_decoder_tokens=count(df["HINDI_TOKENS"])

In [ ]:

#title
def count(series):
    l = []
    for lis in series:
        for words in lis:
            if words not in l:
                l.append(words)
    return l

In [ ]:

#title
input_token_index = dict([(word, i+1) for i, word in enumerate(sorted(count(df["ENGLISH TOKENS"])))])
target_token_index = dict([(word, i+1) for i, word in enumerate(sorted(count(df["HINDI_TOKENS"])))])
```

```
In [ ]:

#title
model.fit_generator(generator = generate_batch(df_train["ENGLISH CORPUS"], df_train["HINDI CORPUS"], batch_size=batch_size,
epochs=epochs, validation_data = train_samples/batch_size, validation_data = generate_batch(df_test["ENGLISH CORPUS"], df_test["HINDI CORPUS"], batch_size=batch_size, validation_data = val_samples/batch_size))

/usr/local/lib/python3.7/dist-packages/ipykernel_launcher.py:6: UserWarning: 'Model.fit_generator' is deprecated and will be removed in a future version. Please use 'Model.fit', which supports generators.
```

Epoch 1/100

```
In [ ]:

#title
# Encode the input sequence to get the "thought vectors"
encoder_model = Model(encoder_inputs, encoder_states)

# Decoder setup
keep_track = False
# We keep track of the states of the previous time step
decoder_state_input_h = Input(shape=(latent_dim,))
decoder_state_input_c = Input(shape=(latent_dim,))
decoder_states_inputs = [decoder_state_input_h, decoder_state_input_c]

dec_emb2 = dec_emb_layer(decoder_inputs) # Get the embeddings of the decoder sequence

# To predict the next word in the sequence, set the initial states to the states from the previous time step
decoder_output2, state_h2, state_c2 = decoder_lstm(dec_emb2, initial_state=decoder_states_inputs)
decoder_output2, state_h2, state_c2 = decoder_lstm(decoder_output2, state_h2, state_c2)
decoder_output2 = decoder_output2 # A dense softmax layer to generate prob dist. over the
```

```
In [ ]:

#title
decoder_model = Model([decoder_inputs] + decoder_states_inputs, [decoder_output2, decoder_state2])

In [ ]:

#title
def decode_sequence(input_seq):
    # Encode the input as state vectors.
    states_value = encoder_model.predict(input_seq)
    # Generate empty target sequence of length 1.
    target_seq = np.zeros((1,1))
    # Populate the first character of target sequence with the start character.
    target_seq[0, 0] = target_token_index['START']

    # Sampling loop for a batch of sequences
    # (to simplify, here we assume a batch of size 1).
    stop_condition = False
    while not stop_condition:
        output_tokens, h, c = decoder_model.predict([target_seq] + states_value)

        # Sample a token
        sampled_token_index = np.argmax(output_tokens[0, -1, :])
        sampled_char = reverse_target_char_index[sampled_token_index]
        decoder_sentence += sampled_char

        # Exit condition: either hit max length
        # or find stop character.
        if (sampled_char == 'END' or len(decoded_sentence) > 50):
            stop_condition = True

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

```
In [ ]:

#title
train_gen = generate_batch(X_train, y_train, batch_size = 1)
k=1

In [ ]:

#title
k+=1
(input_seq, actual_output), _ = next(train_gen)
decoded_sentence = decode_sequence(input_seq)
print('Input English sentence:', X_train[k:k+1].values[0])
print('Actual Hindi Translation:', y_train[k:k+1].values[0][6:-4])
print('Predicted Hindi Translation:', decoded_sentence[:-4])

In [ ]:

#title
k+=1
(input_seq, actual_output), _ = next(train_gen)
decoded_sentence = decode_sequence(input_seq)
print('Input English sentence:', X_train[k:k+1].values[0])
print('Actual Hindi Translation:', y_train[k:k+1].values[0][6:-4])
print('Predicted Hindi Translation:', decoded_sentence[:-4])

In [ ]:

#title
k+=1
(input_seq, actual_output), _ = next(train_gen)
decoded_sentence = decode_sequence(input_seq)
print('Input English sentence:', X_train[k:k+1].values[0])
print('Actual Hindi Translation:', y_train[k:k+1].values[0][6:-4])
print('Predicted Hindi Translation:', decoded_sentence[:-4])
```

Blue Score calculation

```
In [ ]:

#title
a = y_train[k:k+1].values[0][6:-4]
b = decoded_sentence[:-4]

In [ ]:

#title
from nltk.translate.bleu_score import sentence_bleu
score = sentence_bleu(a, b)
print('Bleu score:', '%3f'%score)
```

=====THE END=====