```
1 import pandas as pd
2 import re
3 import tensorflow as tf
4 from tensorflow.keras.layers import Embedding, LSTM, Dense
5 from tensorflow.keras.models import Model
6 from tensorflow.keras.preprocessing.text import Tokenizer
7 from tensorflow.keras.preprocessing.sequence import pad_sequences
8 import numpy as np
9 import nltk.translate.bleu_score as bleu
10 import random
11 import string
12 from sklearn.model_selection import train_test_split
13 import os
14 import time
```

```
1 from google.colab import drive
2 drive.mount('/content/drive')
```

```
    Mounted at /content/drive
```

```
1 eng_hin=pd.read_csv('/content/drive/MyDrive/PRL/task/data/synthetic-dataset/train.csv')
2 eng_hin.head()
```

| | English | Hindi | Hinglish | Average rating | Disagreement |
|---|---|---|---|---|---|
| 0 | Program module is a file that contains instruc... | माड्यूल, एक संचिका होती है, जिसमें या तो स्रोत... | module , ek program hoti hai , jismen ya to so... | 7 | 6 |
| 1 | And to Thamud We sent their brother Sali 'h. H... | और (हमने) क़ौमे समूद के पास उनके भाई सालेह को ... | aur hamne aume samood ke pas unke bhaee saleh ... | 6 | 4 |
| 2 | and, when reminded, do not remember\n | और जब उन्हें याद दिलाया जाता है, तो वे याद नही... | aur jab unhen yad dilaya jata hai , to ve yad ... | 10 | 0 |
| 3 | you won the TED Prize 2011.\n | तुम्हें २०११ का टेड प्राइज़ मिल गया है.\n | tumhen २०११ ka ted prize mil gaya hai\n | 9 | 1 |
| 4 | He gone to Kerodemal College of Delhi Universi... | उन्होंने बाद अध्ययन करने के लिए ये दिल्ली विश्... | unhonne bad science karne ke lie ye delhi univ... | 7 | 0 |

```
1 # create a new dataframe of english and hinglish column
2 df = pd.DataFrame()
3 df["english"] = eng_hin["English"]
4 df["hindi"] = eng_hin["Hindi"]
5 df.head()
```

| | english | hindi |
|---|---|---|
| 0 | Program module is a file that contains instruc... | माड्यूल, एक संचिका होती है, जिसमें या तो स्रोत... |
| 1 | And to Thamud We sent their brother Sali 'h. H... | और (हमने) क़ौमे समूद के पास उनके भाई सालेह को ... |
| 2 | and, when reminded, do not remember\n | और जब उन्हें याद दिलाया जाता है, तो वे याद नही... |
| 3 | you won the TED Prize 2011.\n | तुम्हें २०११ का टेड प्राइज़ मिल गया है.\n |
| 4 | He gone to Kerodemal College of Delhi Universi... | उन्होंने बाद अध्ययन करने के लिए ये दिल्ली विश्... |

```
1 eng_hin.dropna(inplace=True)
2 eng_hin.shape
```

```
    (2766, 7)
```

```
1 exclude = set(string.punctuation) # Set of all special characters
2 remove_digits = str.maketrans('', '', string.digits) # Set of all digits
```

```
1 def preprocess(text):
2     '''Function to preprocess English sentence'''
3     text = text.lower() # lower casing
4     text = re.sub("'", '', text) # remove the quotation marks if any
5     text = ''.join(ch for ch in text if ch not in exclude)
```

```
 6    text = text.translate(remove_digits) # remove the digits
 7    text = text.strip()
 8    text = re.sub(" +", " ", text) # remove extra spaces
 9    text = '<start> ' + text + ' <end>'
10    return text
```

```
 1 def preprocess_hin(text):
 2    '''Function to preprocess Marathi sentence'''
 3    text = re.sub("'", '', text) # remove the quotation marks if any
 4    text = ''.join(ch for ch in text if ch not in exclude)
 5    text = re.sub("[२३०८१५७७८९४६]", "", text) # remove the digits
 6    text = text.strip()
 7    text = re.sub(" +", " ", text) # remove extra spaces
 8    text = '<start> ' + text + ' <end>'
 9    return text
```

```
 1 eng_hin['english'] = df['english'].apply(preprocess)
 2 eng_hin['hindi'] = df['hindi'].apply(preprocess_hin)
 3
 4 eng_hin.rename(columns={"english": "english", "hindi": "hindi"},inplace=True)
 5
 6 eng_hin.head()
```

| | English | Hindi | Hinglish | Average rating | Disagreement | english | hindi |
|---|---|---|---|---|---|---|---|
| 0 | Program module is a file that contains instruc... | माड्यूल, एक संचिका होती है, जिसमें या तो स्रोत... | module , ek program hoti hai , jismen ya to so... | 7 | 6 | <start> program module is a file that contains... | <start> माड्यूल एक संचिका होती है जिसमें या तो... |
| 1 | And to Thamud We sent their brother Sali 'h. H... | और (हमने) क़ौमे समूद के पास उनके भाई सालेह को ... | aur hamne aume samood ke pas unke bhaee saleh ... | 6 | 4 | <start> and to thamud we sent their brother sa... | <start> और हमने क़ौमे समूद के पास उनके भाई साल... |
| 2 | and, when reminded, do not remember\n | और जब उन्हें याद दिलाया जाता है, तो वे याद नही... | aur jab unhen yad dilaya jata hai , to ve yad ... | 10 | 0 | <start> and when reminded do not remember <end> | <start> और जब उन्हें याद दिलाया जाता है तो वे ... |
| 3 | you won the TED Prize 2011.\n | तुम्हें २०११ का टेड प्राइज़ मिल गया है.\n | tumhen २०११ ka ted prize mil gaya hai\n | 9 | 1 | <start> you won the ted prize <end> | <start> तुम्हें का टेड प्राइज़ मिल गया है <end> |
| | He gone to | उन्होंने बाद | unhonne had | | | <start> he gone to | <start> उन्होंने |

```
 1 def tokenize(lang):
 2
 3    lang_tokenizer = tf.keras.preprocessing.text.Tokenizer(filters='')
 4    lang_tokenizer.fit_on_texts(lang)
 5
 6    tensor = lang_tokenizer.texts_to_sequences(lang)
 7
 8    tensor = tf.keras.preprocessing.sequence.pad_sequences(tensor,padding='post',maxlen=20,dtype='int32')
 9
10    return tensor, lang_tokenizer
```

```
 1 def load_dataset():
 2
 3    input_tensor, inp_lang_tokenizer = tokenize(eng_hin['english'].values)
 4    target_tensor, targ_lang_tokenizer = tokenize(eng_hin['hindi'].values)
 5
 6    return input_tensor, target_tensor, inp_lang_tokenizer, targ_lang_tokenizer
```

```
 1 input_tensor, target_tensor, inp_lang, targ_lang = load_dataset()
```

```
 1 max_length_targ, max_length_inp = target_tensor.shape[1], input_tensor.shape[1]
```

```
 1 input_tensor_train, input_tensor_val, target_tensor_train, target_tensor_val = train_test_split(input_tensor, target_ten
 2
 3 print(len(input_tensor_train), len(target_tensor_train), len(input_tensor_val), len(target_tensor_val))
```

```
 2212 2212 554 554
```

```
1 BUFFER_SIZE = len(input_tensor_train)
2 BATCH_SIZE = 32
3 N_BATCH = BUFFER_SIZE//BATCH_SIZE
4 embedding_dim = 256
5 units = 1024
6 steps_per_epoch = len(input_tensor_train)//BATCH_SIZE
7
8 vocab_inp_size =len(inp_lang.word_index.keys())
9 vocab_tar_size =len(targ_lang.word_index.keys())
10
11 dataset = tf.data.Dataset.from_tensor_slices((input_tensor_train, target_tensor_train)).shuffle(BUFFER_SIZE)
12 dataset = dataset.batch(BATCH_SIZE, drop_remainder=True)
```

```
1 embeddings_index = dict()
2 f = open('/content/drive/MyDrive/PRL/task/data/synthetic-dataset/glove.6B.300d.txt')
3 for line in f:
4     values = line.split()
5     word = values[0]
6     coefs = np.asarray(values[1:], dtype='float32')
7     embeddings_index[word] = coefs
8 f.close()
9
10 embedding_matrix = np.zeros((vocab_inp_size+1, 300))
11 for word, i in inp_lang.word_index.items():
12     embedding_vector = embeddings_index.get(word)
13     if embedding_vector is not None:
14         embedding_matrix[i] = embedding_vector
```

```
1 class Encoder(tf.keras.Model):
2     def __init__(self, vocab_size, embedding_dim, enc_units, batch_sz):
3         super(Encoder, self).__init__()
4         self.batch_sz = batch_sz
5         self.enc_units = enc_units
6         self.embedding = tf.keras.layers.Embedding(input_dim=vocab_size, output_dim=embedding_dim, name="embedding_layer
7         self.gru = tf.keras.layers.GRU(units, return_sequences=True, return_state=True, recurrent_activation='sigmoid',
8
9     def call(self, x, hidden):
10        x = self.embedding(x)
11        output, state = self.gru(x, initial_state = hidden)
12        return output, state
13
14    def initialize_hidden_state(self):
15        return tf.zeros((self.batch_sz, self.enc_units))
```

```
1 class Decoder(tf.keras.Model):
2     def __init__(self, vocab_size, embedding_dim, dec_units, batch_sz):
3         super(Decoder, self).__init__()
4         self.batch_sz = batch_sz
5         self.dec_units = dec_units
6         self.embedding = tf.keras.layers.Embedding(vocab_size, embedding_dim)
7         self.gru = tf.keras.layers.GRU(units, return_sequences=True, return_state=True, recurrent_activation='sigmoid',
8         self.fc = tf.keras.layers.Dense(vocab_size)
9
10                # used for attention
11        self.W1 = tf.keras.layers.Dense(self.dec_units)
12        self.W2 = tf.keras.layers.Dense(self.dec_units)
13        self.V = tf.keras.layers.Dense(1)
14
15    def call(self, x, hidden, enc_output):
16
17        hidden_with_time_axis = tf.expand_dims(hidden, 1)
18
19        score = self.V(tf.nn.tanh(self.W1(enc_output) + self.W2(hidden_with_time_axis)))
20
21        attention_weights = tf.nn.softmax(score, axis=1)
22
23        context_vector = attention_weights * enc_output
24        context_vector = tf.reduce_sum(context_vector, axis=1)
25
26        x = self.embedding(x)
27
```

```
28          x = tf.concat([tf.expand_dims(context_vector, 1), x], axis=-1)
29
30          output, state = self.gru(x)
31
32          output = tf.reshape(output, (-1, output.shape[2]))
33
34          x = self.fc(output)
35
36          return x, state, attention_weights
37
38      def initialize_hidden_state(self):
39          return tf.zeros((self.batch_sz, self.dec_units))
```

```
1 tf.keras.backend.clear_session()
2
3 encoder = Encoder(vocab_inp_size+1, 300, units, BATCH_SIZE)
4 decoder = Decoder(vocab_tar_size+1, embedding_dim, units, BATCH_SIZE)
```

```
1 optimizer = tf.keras.optimizers.Adam()
2 loss_object = tf.keras.losses.SparseCategoricalCrossentropy(from_logits=True,
3                                                              reduction='none')
4
5
6 def loss_function(real, pred):
7   mask = tf.math.logical_not(tf.math.equal(real, 0))
8   loss_ = loss_object(real, pred)
9
10  mask = tf.cast(mask, dtype=loss_.dtype)
11  loss_ *= mask
12
13  return tf.reduce_mean(loss_)
```

```
1 checkpoint_dir = './training_checkpoints'
2 checkpoint_prefix = os.path.join(checkpoint_dir, "ckpt")
3 checkpoint = tf.train.Checkpoint(optimizer=optimizer,
4                                  encoder=encoder,
5                                  decoder=decoder)
```

```
1 @tf.function
2 def train_step(inp, targ, enc_hidden):
3   loss = 0
4
5   with tf.GradientTape() as tape:
6     enc_output, enc_hidden = encoder(inp, enc_hidden)
7     encoder.get_layer('embedding_layer_encoder').set_weights([embedding_matrix])
8     dec_hidden = enc_hidden
9
10    dec_input = tf.expand_dims([targ_lang.word_index['<start>']] * BATCH_SIZE, 1)
11
12    for t in range(1, targ.shape[1]):
13      predictions, dec_hidden, _ = decoder(dec_input, dec_hidden, enc_output)
14
15      loss += loss_function(targ[:, t], predictions)
16
17      dec_input = tf.expand_dims(targ[:, t], 1)
18
19  batch_loss = (loss / int(targ.shape[1]))
20
21  variables = encoder.trainable_variables + decoder.trainable_variables
22
23  gradients = tape.gradient(loss, variables)
24
25  optimizer.apply_gradients(zip(gradients, variables))
26
27  return batch_loss
```

```
1 EPOCHS = 100
2
3 for epoch in range(EPOCHS):
4   start = time.time()
5
6   enc_hidden = encoder.initialize_hidden_state()
```

```
7    total_loss = 0
8
9    for (batch, (inp, targ)) in enumerate(dataset.take(steps_per_epoch)):
10     batch_loss = train_step(inp, targ, enc_hidden)
11     total_loss += batch_loss
12
13     if batch % 100 == 0:
14       print(f'Epoch {epoch+1} Batch {batch} Loss {batch_loss.numpy():.4f}')
15   if (epoch + 1) % 2 == 0:
16     checkpoint.save(file_prefix=checkpoint_prefix)
17
18   print(f'Epoch {epoch+1} Loss {total_loss/steps_per_epoch:.4f}')
19   print(f'Time taken for 1 epoch {time.time()-start:.2f} sec\n')
```

```
Time taken for 1 epoch 8.04 sec

Epoch 87 Batch 0 Loss 0.0049
Epoch 87 Loss 0.0027
Time taken for 1 epoch 7.15 sec

Epoch 88 Batch 0 Loss 0.0048
Epoch 88 Loss 0.0031
Time taken for 1 epoch 8.04 sec

Epoch 89 Batch 0 Loss 0.0007
Epoch 89 Loss 0.0033
Time taken for 1 epoch 7.16 sec

Epoch 90 Batch 0 Loss 0.0033
Epoch 90 Loss 0.0027
Time taken for 1 epoch 7.95 sec

Epoch 91 Batch 0 Loss 0.0133
Epoch 91 Loss 0.0029
Time taken for 1 epoch 7.17 sec

Epoch 92 Batch 0 Loss 0.0006
Epoch 92 Loss 0.0030
Time taken for 1 epoch 8.25 sec

Epoch 93 Batch 0 Loss 0.0007
Epoch 93 Loss 0.0040
Time taken for 1 epoch 7.18 sec

Epoch 94 Batch 0 Loss 0.0142
Epoch 94 Loss 0.0036
Time taken for 1 epoch 7.92 sec

Epoch 95 Batch 0 Loss 0.0327
Epoch 95 Loss 0.0062
Time taken for 1 epoch 7.18 sec

Epoch 96 Batch 0 Loss 0.0008
Epoch 96 Loss 0.0111
Time taken for 1 epoch 7.93 sec

Epoch 97 Batch 0 Loss 0.0160
Epoch 97 Loss 0.0155
Time taken for 1 epoch 7.20 sec

Epoch 98 Batch 0 Loss 0.0259
Epoch 98 Loss 0.0537
Time taken for 1 epoch 7.94 sec

Epoch 99 Batch 0 Loss 0.0447
Epoch 99 Loss 0.1171
Time taken for 1 epoch 7.19 sec

Epoch 100 Batch 0 Loss 0.1225
Epoch 100 Loss 0.1011
Time taken for 1 epoch 7.96 sec
```

```
1 def evaluate(sentence):
2   attention_plot = np.zeros((max_length_targ, max_length_inp))
3
4   sentence = preprocess(sentence)
5
6   inputs = [inp_lang.word_index[i] for i in sentence.split(' ')]
7   inputs = tf.keras.preprocessing.sequence.pad_sequences([inputs],maxlen=20, padding='post')
```

```
 8    inputs = tf.convert_to_tensor(inputs)
 9
10    result = ''
11
12    hidden = [tf.zeros((1, units))]
13    enc_out, enc_hidden = encoder(inputs, hidden)
14
15    dec_hidden = enc_hidden
16    dec_input = tf.expand_dims([targ_lang.word_index['<start>']], 0)
17
18    for t in range(max_length_targ):
19      predictions, dec_hidden, attention_weights = decoder(dec_input,
20                                                             dec_hidden,
21                                                             enc_out)
22      # storing the attention weights to plot later on
23      attention_weights = tf.reshape(attention_weights, (-1, ))
24      attention_plot[t] = attention_weights.numpy()
25      predicted_id = tf.argmax(predictions[0]).numpy()
26
27      result += targ_lang.index_word[predicted_id] + ' '
28
29      if targ_lang.index_word[predicted_id] == '<end>':
30        return result,attention_plot
31
32      # the predicted ID is fed back into the model
33      dec_input = tf.expand_dims([predicted_id], 0)
34
35    return result,attention_plot
```

```
1 input_sentence= 'please ensure that you use the appropriate form '
2 print('Input sentence in english : ',input_sentence)
3 predicted_output,attention_plot=evaluate(input_sentence)
4 print('Predicted sentence in hindi : ',predicted_output)
```

```
    Input sentence in english :  please ensure that you use the appropriate form
    Predicted sentence in hindi :  आपकी मांग को बदलने का गाना मंजूर किया हुआ <end>
```

```
1 input_sentence='and do something with it to change the world '
2 print('Input sentence in english : ',input_sentence)
3 predicted_output,attention_plot=evaluate(input_sentence)
4 print('Predicted sentence in hindi : ',predicted_output)
```

```
    Input sentence in english :  and do something with it to change the world
    Predicted sentence in hindi :  और अगर मालूम हो <end>
```

## ▾ OpenInAPP sentences output

```
1 input_sentence='So even if its a big video I will clearly mention all the products '
2 print('Input sentence in english : ',input_sentence)
3 predicted_output,attention_plot=evaluate(input_sentence)
4 print('Predicted sentence in hindi : ',predicted_output)
```

```
    Input sentence in english :  So even if its a big video I will clearly mention all the products
    Predicted sentence in hindi :  तो सच है <end>
```

```
1 input_sentence='I was waiting for my bag '
2 print('Input sentence in english : ',input_sentence)
3 predicted_output,attention_plot=evaluate(input_sentence)
4 print('Predicted sentence in hindi : ',predicted_output)
```

```
    Input sentence in english :  I was waiting for my bag
    Predicted sentence in hindi :  मैंने मैं धैर्य न छूट जाएं। <end>
```

```
1
```

```
1    input_sentence='definitely share your feedback in the comment section '
2    print('Input sentence in english : ',input_sentence)
3    predicted_output,attention_plot=evaluate(input_sentence)
4    print('Predicted sentence in hindi : ',predicted_output)
```

```
Input sentence in english :  definitely share your feedback in the comment section
---------------------------------------------------------------------------
KeyError                                  Traceback (most recent call last)
<ipython-input-66-854b97a97912> in <cell line: 3>()
      1 input_sentence='definitely share your feedback in the comment section '
      2 print('Input sentence in english : ',input_sentence)
----> 3 predicted_output,attention_plot=evaluate(input_sentence)
      4 print('Predicted sentence in hindi : ',predicted_output)

                               ↕ 1 frames
<ipython-input-61-4ac5775c6a63> in <listcomp>(.0)
      4     sentence = preprocess(sentence)
      5
----> 6     inputs = [inp_lang.word_index[i] for i in sentence.split(' ')]
      7     inputs =
tf.keras.preprocessing.sequence.pad_sequences([inputs],maxlen=20, padding='post')
      8     inputs = tf.convert_to_tensor(inputs)

KeyError: 'definitely'
```

> SEARCH STACK OVERFLOW

1

● 0s    completed at 9:36 PM                                    ● ✕