

Assignment 03

Implementing English-to-French Machine Translation Using RNN or Transformer

Batool Ali Akbar
ba07612@st.habib.edu.pk

I. CHOSEN MODEL

A translation task is a common use case of a sequence-to-sequence model. An encoder-decoder model implemented via RNN is one of the earliest sequence-to-sequence models. However, the model is quite slow and cannot preserve long dependencies. Therefore, I have chosen a transformer based model for this assignment. Unlike RNN, a transformer can consider all parts of a sentence simultaneously, capturing both past and future context which improves the overall accuracy of the translation.

II. IMPLEMENTATION

This entire section focuses on how the model was implemented and evaluated. The provided dataset was quite large with about 175,622 data entries. This kind of data processing needed huge computation power which my system could not support and hence threw an error entailing resource exhaustion. Therefore, I had to cut down the file to 60,000 entries.

A. Data Preprocessing

The following block of code preprocesses the data read from the file using the built-in pandas library. Line 4 reads only the first 60,000 entries from the entire data file. Lines 6 and 7 turn the entire dataset into lower case. Lines 9 to 16 tokenize each word, assigning them an integer and converting each sentence into an array of integers, each denoting the corresponding word of the sentence. Lines 18 to 22 use padding to ensure equal length for each sentence. Lines 24 and 25 split the data into training sample (80%), validation sample (10%) and testing sample (10%).

```

1
2 def preprocess_data(file_path):
3
4     data = pd.read_csv(file_path)[:60000]
5
6     eng_text = data['English words/sentences'].apply
7     (lambda x: x.lower().strip())
8     french_text = data['French words/sentences'].
9     apply(lambda x: x.lower().strip())
10
11     eng_tokenizer = Tokenizer()
12     french_tokenizer = Tokenizer()
13
14     eng_tokenizer.fit_on_texts(eng_text)
15     french_tokenizer.fit_on_texts(french_text)
16
17     eng_sequences = eng_tokenizer.texts_to_sequences
18     (eng_text)
19     french_sequences = french_tokenizer.
20     texts_to_sequences(french_text)

```

```

18 max_eng_len = max(len(seq) for seq in
    eng_sequences)
19 max_french_len = max(len(seq) for seq in
    french_sequences)
20
21 eng_padded = pad_sequences(eng_sequences, maxlen=
    max_eng_len, padding='post')
22 french_padded = pad_sequences(french_sequences,
    maxlen=max_french_len, padding='post')
23
24 eng_train, eng_temp, french_train, french_temp =
    train_test_split(eng_padded, french_padded,
    test_size=0.2, random_state=42)
25 eng_val, eng_test, french_val, french_test =
    train_test_split(eng_temp, french_temp,
    test_size=0.5, random_state=42)
26
27 return eng_train, eng_val, eng_test,
    french_train, french_val, french_test,
    eng_tokenizer, french_tokenizer, max_eng_len,
    max_french_len

```

Listing 1: Data Preprocessing

B. Model Implementation

Code listing 2 provides an implementation of the transformer model. The model has been built on the architecture illustrated in Figure 1.

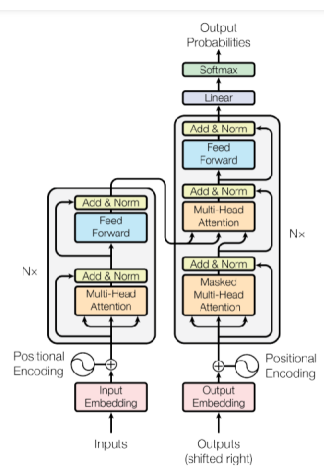


Fig. 1: Transformers Architecture

```

1 def transformer_model(input_vocab_size,
2   target_vocab_size, input_max_len, target_max_len
3   , embedding_dim=256, num_heads=4,
4   feed_forward_dim=512):
5
6   # Encoder
7   inputs = Input(shape=(input_max_len,))
8   enc_emb = Embedding(input_vocab_size,
9   embedding_dim)(inputs)
10  positional_encoding_input = positional_encoding(
11  input_max_len, embedding_dim)
12  enc_emb = enc_emb + positional_encoding_input
13
14  for _ in range(2):
15      # Multi-Head Attention
16      attn_output = MultiHeadAttention(num_heads=
17  num_heads, key_dim=embedding_dim)(enc_emb,
18  enc_emb, enc_emb)
19      # Add & Norm
20      norm_1 = LayerNormalization()(attn_output)
21      add_1 = norm_1 + enc_emb
22      # Feed Forward Network
23      ffn_output = Dense(feed_forward_dim,
24  activation='relu')(add_1)
25      ffn_output = Dense(embedding_dim)(ffn_output)
26
27      # Add & Norm
28      norm_2 = LayerNormalization()(ffn_output)
29      add_2 = norm_2 + add_1
30      enc_emb = add_2
31  enc_output = enc_emb
32
33  # Decoder
34  dec_inputs = Input(shape=(target_max_len,))
35  dec_emb = Embedding(target_vocab_size,
36  embedding_dim)(dec_inputs)
37  positional_encoding_target = positional_encoding(
38  target_max_len, embedding_dim)
39  dec_emb = dec_emb + positional_encoding_target
40
41  for _ in range(2):
42      # Masked Multi-Head Attention (Self-
43  Attention)
44      attn_output = MultiHeadAttention(num_heads=
45  num_heads, key_dim=embedding_dim)(dec_emb,
46  dec_emb, dec_emb, attention_mask=
47  create_causal_mask(target_max_len))
48      # Add & Norm
49      norm_1 = LayerNormalization()(attn_output)
50      add_1 = dec_emb + norm_1
51      # Encoder-Decoder Attention
52      enc_dec_attn_output = MultiHeadAttention(
53  num_heads=num_heads, key_dim=embedding_dim)(
54  add_1, enc_output, enc_output)
55      # Add & Norm
56      norm_2 = LayerNormalization()(
57  enc_dec_attn_output)
58      add_2 = norm_1 + add_1
59      # Feed Forward Network
60      ffn_output = Dense(feed_forward_dim,
61  activation='relu')(add_2)
62      ffn_output = Dense(embedding_dim)(ffn_output)
63
64      # Add & Norm
65      norm_3 = LayerNormalization()(ffn_output)
66      add_3 = ffn_output + norm_3
67      dec_emb = add_3
68  dec_output = dec_emb
69
70  # Final Dense Layer
71  output = Dense(target_vocab_size, activation='
72  softmax')(dec_output)

```

```

53 model = Model(inputs=[inputs, dec_inputs],
54   outputs=output)
55 return model

```

Listing 2: Model Implementation

The current implementation uses 2 layers of encoder and 2 layers of decoder with an embedding dimension of 256, 4 attention heads for each attention mechanism and RELU as an activation function. It then uses built-in functions provided by the tensorflow library to create a transformer model with an architecture similar to the one mentioned in Figure 1. Lines 4 to 23 are responsible for encoding the input sequence while the lines beyond 25 are the decoder part of the model. Line 52 applies the softmax function to the output of the decoder to normalize the values and line 53 then compiles the entire model.

C. Model Training

Python's builtin library 'tensorflow' is used to train the model. The code below illustrates how the model has been trained.

```

1 model = transformer_model(input_vocab_size,
2   target_vocab_size, max_eng_len, max_french_len)
3 model.compile(optimizer='adam', loss='
4   sparse_categorical_crossentropy', metrics=['
5   accuracy'])
6
7 history = model.fit(
8   [eng_train, french_train],
9   np.expand_dims(french_train, -1),
10  epochs=10,
11  batch_size=64,
12  validation_data=([eng_val, french_val], np.
13  expand_dims(french_val, -1))

```

Listing 3: Model Training

The model is compiled using the Adam optimizer along with sparse categorical cross-entropy loss. The model is then trained over 10 epochs with a batch size of 64. During training, the input English sentences are paired with their corresponding French translations. The training also includes validation data to monitor the model's performance on unseen data. The model.fit() function performs the training, with both training and validation accuracy being tracked throughout the process.

D. Evaluation

The model has been evaluated using the BLEU score. The following code depicts how the BLEU score was calculated.

```

1 def evaluate_model(model, eng_test, french_test,
2   eng_tokenizer, french_tokenizer):
3
4   # Predicting the French sentences
5   french_pred = model.predict([eng_test,
6   french_test])
7   french_pred_text = [
8       french_tokenizer.sequences_to_texts([np.
9       argmax(pred, axis=-1)])[0]
10      for pred in french_pred
11  ]

```

```

10
11 # Decoding the French sentences
12 french_test_decoded = [
13     ' '.join([french_tokenizer.index_word[token]
14         for token in sentence if token != 0])
15     for sentence in french_test
16 ]
17
18 # Prepare references and hypotheses for BLEU
19 # score calculation
20 references = [[sent.split()] for sent in
21     french_test_decoded]
22 hypotheses = [sent.split() for sent in
23     french_pred_text]
24
25 # BLEU score calculation
26 smoothing = SmoothingFunction().method1
27 bleu_score = corpus_bleu(references, hypotheses,
28     smoothing_function=smoothing)
29 print(f'BLEU Score: {bleu_score}')

```

Listing 4: BLEU Score Calculation

The code uses the transformer model defined above to make predictions for the test set. It then converts the predictions into proper texts and then compares the predicted output with the expected output to calculate the BLEU score using the 'nltk' library.

III. PERFORMANCE

The BLEU score for the model was calculated to be 0.87. It suggests that the model's translations are close to the reference translations. This high score shows that the model's output is generally fluent and accurate in terms of word choice, syntax, and meaning. This can also be seen in the following 3 samples that were printed out while testing.

Expected Translation	Predicted Translation
je vais en ville	je vais en ville
les chaussures se vendent par paires	les chaussures se dames par réussirait
nous étions chanceux	nous étions chanceux

Table 1: Translation Samples in Test Set

Furthermore, the accuracy and loss of the model for the training and the validation set over the 10 epochs was also recorded. The values are mentioned in the table below.

Epoch	Training Accuracy	Training Loss	Validation Accuracy	Validation Loss
1	0.7822	1.6869	0.9533	0.3502
2	0.9334	0.4752	0.9647	0.2900
3	0.9608	0.2602	0.9720	0.2403
4	0.9763	0.1465	0.9739	0.2199
5	0.9771	0.1233	0.9757	0.2097
6	0.9815	0.0930	0.9773	0.2100
7	0.9815	0.0852	0.9662	0.2408
8	0.9827	0.0744	0.9812	0.1919
9	0.9885	0.0476	0.9756	0.2161
10	0.9836	0.0713	0.9841	0.1842

Table 2: Training and Validation Accuracy and Loss per Epoch

This trend can also be seen in the following graphs.

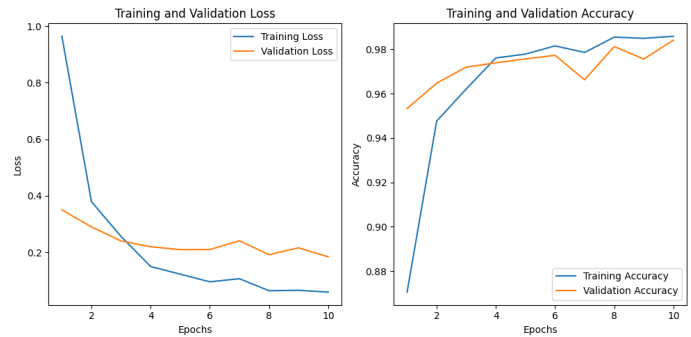


Fig. 2: Loss and Accuracy over 10 Epochs

The graph on the left depicts the training loss (blue line) and the validation loss (orange line). Both of them decrease initially, indicating that the model is learning and improving its predictions. The training is initially quite high but drops rapidly in the first few epochs, then gradually levels out, suggesting that the model is fitting well on the training data. The validation loss starts from quite a small value and continues to decrease slowly. However, the loss for both training and validation set fluctuates slightly in later epochs, suggesting some generalization challenges but overall the trend is consistent.

The graph on the right depicts the training accuracy (blue line) and the validation accuracy (orange line). The training accuracy starts off with a low value but then increases quickly in the first few epochs and then stabilizes near 0.98, indicating that the model performs well on the training data. The validation accuracy starts at around 0.95 and follows a similar trend, achieving a high accuracy close to that of the training set. The slight fluctuation in the validation accuracy suggests some variation in performance on the validation set, but the overall trend shows good model generalization with high accuracy on both the training and validation sets.

Overall, the model is performing well, with a high BLEU score, low loss and high accuracy on both training and validation datasets. The small gap between training and validation metrics suggests minimal overfitting, indicating that the model is likely generalizing well to unseen data.

IV. REFLECTION

Although the implementation of the model has been done entirely using the libraries, I still faced some challenges. I initially planned to implement an RNN model but the model was taking much longer than expected to predict an output. Therefore, I had to switch to a transformers based model. Moreover, the size of the dataset was quite large. The normal VS code compiler could not process the entire data. I had to move to Google Colab and switch from CPU to GPU to run the model in a reasonable amount of time. However, this had its own challenges. Google Colab had a limited time GPU

access which once again slowed down the entire pace of the assignment.