Neural Machine Translation: a seq2seq implementation to translate English to German In [1]: %%capture !pip install "tensorflow-text"==2.15.0 In [2]: # import necessary libs import numpy as np import re import os import random %matplotlib inline import matplotlib.pyplot as plt import tensorflow\_text as tf\_text import tensorflow as tf from tensorflow.keras import Model from tensorflow.keras.layers import (Layer, Dense, LSTM, Embedding, TextVectorization, Bidirectional, Add, LayerNormalization, Activation) Introduction This notebook is to implement the seq2seq model with Attention proposed by Bahdanau et al., 2015. In [4]: # Let's define some constant variables max\_vocab\_size=16000 DROPOUT=0.3 BUFFER\_SIZE=1024 BATCH\_SIZE=64 embedding\_size=128 hidden\_units=128 data\_file="deu.txt" data\_dir='/content/data/' os.makedirs(data\_dir, exist\_ok=True) **Data Preprocessing** We will use the English to German dataset from Manythings.org. # Download dataset !wget --no-check-certificate 'https://www.manythings.org/anki/deu-eng.zip' -0 deu-eng.zip !unzip deu-eng.zip -d data/ --2024-05-09 17:58:21-- https://www.manythings.org/anki/deu-eng.zip Resolving www.manythings.org (www.manythings.org)... 173.254.30.110 Connecting to www.manythings.org (www.manythings.org)|173.254.30.110|:443... connected. HTTP request sent, awaiting response... 200 OK Length: 10364105 (9.9M) [application/zip] Saving to: 'deu-eng.zip' deu-eng.zip 2024-05-09 17:58:22 (24.7 MB/s) - 'deu-eng.zip' saved [10364105/10364105] Archive: deu-eng.zip inflating: data/deu.txt inflating: data/\_about.txt In [6]: # Take a look at the first lines with open(os.path.join(data\_dir, data\_file)) as f: for n, line in enumerate(f): print(line.strip()) **if** n == 4: break Geh. CC-BY 2.0 (France) Attribution: tatoeba.org #2877272 (CM) & #8597805 (Roujin) Go. Hallo! CC-BY 2.0 (France) Attribution: tatoeba.org #538123 (CM) & #380701 (cburgmer) Hi. Grüß Gott! Hi. CC-BY 2.0 (France) Attribution: tatoeba.org #538123 (CM) & #659813 (Esperantostern) Lauf! CC-BY 2.0 (France) Attribution: tatoeba.org #906328 (papabear) & #941078 (Fingerhut) Run! Lauf! CC-BY 2.0 (France) Attribution: tatoeba.org #4008918 (JSakuragi) & #941078 (Fingerhut) Run. To maintain coherence and prevent loss of word's meaning, contracted terms are expanded. In [7]: en\_contraction\_map = { "let's": "let us", "'d better": " had better", "'s": " is", "'re": " are", "'m": " am", "'11": " will", "'d": " would", "'ve": " have", "'em": " them", "won't": "will not", "n't": " not", "cannot": "can not", ger\_contraction\_map = { "'s": " ist", "ä": "ae", "ö": "oe", "ü": "ue", "ß": "ss", "'ne ": "eine ", "'n ": "ein ", "am ": "an dem " "ans ": "an das " "aufs ": "auf das ", "durchs ": "durch das ", "fuers ": "fuer das ", "hinterm ": "hinter dem ", "im ": "in dem ", "uebers ": "ueber das ", "ums ": "um das ", "unters ": "unter das ", "unterm ": "unter dem ", "vors ": "vor das ", "vorm ": "vor dem ", "zum ": "zu dem ", "ins ": "in das ", "vom ": "von dem" , "beim ": "bei dem ", "zur ": "zu der ", def expand\_contractions(text, mapping): for key, value in mapping.items(): text = re.sub(key, value, text) return text In [8]: # Let's test the function print(expand\_contractions("He definitely didn't do it. He must've been forced to commit crime. He won't do it again.", en\_contraction\_map)) print(expand\_contractions("Hinterm Haus, das am Fluss liegt, steht ein großer Baum.", ger\_contraction\_map)) He definitely did not do it. He must have been forced to commit crime. He will not do it again. Hinterm Haus, das an dem Fluss liegt, steht ein grosser Baum. The dataset are still in unprocessed form. It is necessary to preprocess and store them in appropriate form. In [9]: english = [] german = [] with open(os.path.join(data\_dir, data\_file)) as f: for line in f: line = line.split("CC-BY") if len(line) > 0: sample = line[0]sample = sample.strip().split('\t') english.append(expand\_contractions(sample[0].lower(), en\_contraction\_map)) german.append(expand\_contractions(sample[1].lower(), ger\_contraction\_map)) english = np.array(english) german = np.array(german) In [10]: # Take a look at 5 random examples for i in range(5): rdi = random.randint(0, len(english)) print("{:4} --> {:4}".format(english[rdi], german[rdi])) this road is dangerous. --> der weg ist gefaehrlich. this material is not suitable for a dress. --> dieser stoff ist nicht geeignet fuer ein kleid. tom hates climbing ladders. --> tom steigt nicht gern auf leitern. why do you want to do this? --> warum willst du das tun? he does not get up early. --> er steht nicht frueh auf. In [11]: # english = english[:400] # german = german[:400] In [12]: # Create mask mask = np.full((len(english),), False) train\_mask = np.copy(mask) train\_mask[:int(len(english) \* 0.8)] = True np.random.shuffle(train\_mask) false\_indices = np.where(train\_mask == False)[0] np.random.shuffle(false\_indices) border\_idx = int(len(false\_indices) / 2) val\_mask = np.copy(mask) val\_mask[false\_indices[:border\_idx]] = True test\_mask = np.copy(mask) test\_mask[false\_indices[border\_idx:]] = True In [13]:  $train_raw = ($ tf.data.Dataset .from\_tensor\_slices((english[train\_mask], german[train\_mask])) .shuffle(BUFFER\_SIZE) .batch(BATCH\_SIZE) .prefetch(tf.data.experimental.AUTOTUNE))  $val_raw = ($ tf.data.Dataset .from\_tensor\_slices((english[val\_mask], german[val\_mask])) .shuffle(BUFFER\_SIZE) .batch(BATCH\_SIZE) .prefetch(tf.data.experimental.AUTOTUNE))  $test_raw = ($ tf.data.Dataset .from\_tensor\_slices((english[test\_mask], german[test\_mask])) .shuffle(BUFFER\_SIZE) .batch(BATCH\_SIZE) .prefetch(tf.data.experimental.AUTOTUNE)) It can be inferred from the plots that despite some outliers, sentence length tends to remain stable along both dataset. Therefore, it is not necessary to implement bucketing by length. **Tokenization** Computer obviously cannot handle raw text. Instead, they need to be converted into numerical form for further calculations. Besides, while both removing punctuation and lowercasing all words are common practice in NLP tasks, it is not really the case for Neural Machine Translation. Punctuation is important to mark the start or end of a sentence. Therefore, we may well necessarily tokenize them. In [14]: def text\_standardize(text): # Split accented characters. text = tf\_text.normalize\_utf8(text, 'NFKD') text = tf.strings.lower(text) # Keep space, a to z, and select punctuation. text = tf.strings.regex\_replace(text, '[^ a-z.?!,¿]', '') # Add spaces around punctuation. text = tf.strings.regex\_replace(text, '[.?!,¿]', r' \0 ') # Strip whitespace. text = tf.strings.strip(text) text = tf.strings.join(['[START]', text, '[END]'], separator=' ') return text In [15]: # Vectorizer initial en\_vec = TextVectorization(max\_tokens=max\_vocab\_size, standardize=text\_standardize, ragged=True) ger\_vec = TextVectorization(max\_tokens=max\_vocab\_size, standardize=text\_standardize, ragged=True) In [16]: en\_vec.adapt(train\_raw.map(lambda x, y: x)) ger\_vec.adapt(train\_raw.map(lambda x, y: y)) In [17]: # Vocabulary en\_voc = en\_vec.get\_vocabulary() ger\_voc = ger\_vec.get\_vocabulary() In [18]: # Word to Idx for prediction word\_to\_idx = {} for i in range(len(ger\_voc)): word\_to\_idx[ger\_voc[i]] = i In [19]: # Assign vocab size of each vectorizer input\_vocab\_size = len(en\_vec.get\_vocabulary()) output\_vocab\_size = len(ger\_vec.get\_vocabulary()) print(input\_vocab\_size) print(output\_vocab\_size) 16000 16000 **Data Preparation** Structure the dataset to use the tf.keras.models.Model's fit() method. In [20]: def process\_text(context, target): context = en\_vec(context).to\_tensor() target = ger\_vec(target) targ\_in = target[:, :-1].to\_tensor() targ\_out = target[:, 1:].to\_tensor() return (context, targ\_in), targ\_out train\_ds = train\_raw.map(process\_text, tf.data.AUTOTUNE) val\_ds = val\_raw.map(process\_text, tf.data.AUTOTUNE) test\_ds = test\_raw.map(process\_text, tf.data.AUTOTUNE) Model implementation In seq2seq, we need a RNN block, which is also known as Encoder, to encode the input sequence to a fixed-length vector, then another RNN block called Decoder to decode it. Block generally consists of LSTM cells. Encoder Encoder can be defined multi-layered RNN network. For the sake of simplicity, I will implement it as an one-layer RNN network with 1 cell at each timestep. Each RNN cell receives a source word and previous hidden state as inputs.  $s_i = tanh(Ws_{i-1} + Ux_i)$ According to the formula, the  $i^{th}$  hidden state  $s_i$  is calculated from the  $(i-1)^{th}$  hidden state and the  $i^{th}$  input. # Let's define the encoder class Encoder(Layer): def \_\_init\_\_(self, tokenizer, embedding\_size, hidden\_units, dropout=DROPOUT): Encoder Block in seq2seq :param tokenizer: tokenizer of the source language :param embedding\_size: dimensionality of the embedding layer :param hidden\_units: dimensionality of the output super().\_\_init\_\_() self.tokenizer = tokenizer self.embedding\_size = embedding\_size self.hidden\_units = hidden\_units self.vocab\_size = tokenizer.vocabulary\_size() self.embedding = Embedding(input\_dim=self.vocab\_size, output\_dim=embedding\_size) self.rnn = Bidirectional( merge\_mode="sum", layer=LSTM(units=hidden\_units, dropout=dropout, return\_sequences=True, return\_state=**True**)) def call(self, training=True): :param x: [batch, time\_steps] :param training: is training or not :return: encoder\_hidden\_state: [batch, hidden\_state\_dim] state\_h: [batch, hidden\_state\_dim] state\_c: [batch, hidden\_state\_dim] mask = tf.where(x != 0, True, False) x = self.embedding(x)x, forward\_h, forward\_c, backward\_h, backward\_c = self.rnn(x, mask=mask, training=training) return x, forward\_h + backward\_h, forward\_c + backward\_c **Attention Layer** The Attention Mechanism used in this project is Bahdanau Attention, which is first introduced in the Neural Machine Translation by Jointly Learning to Align and Translate paper by Bahdanau et al., 2015. \ To recap, at each inference step i, Decoder incorporates information from both the previous Decoder timestep  $s_{i-1}$  and all encoder states  $h=(\{h_1,h_2,\ldots,h_{Tx}\})$  to take only the most relevant words to  $y_{i-1}$  through alignment model a.  $e_{ij} = a(s_{i-1},h_j) = v_a^T. tanh(W_as_{i-1} + U_ah_j)$ Therefore, the context vector  $c_i$  can be calculated as in which During the training process, we will implement Teacher Forcing by combining context vector  $c_i$  with Decoder input  $x_i$ . In [22]: class BahdanauAttention(Layer): def \_\_init\_\_(self, hidden\_units): super().\_\_init\_\_() self.Va = Dense(1)self.Wa = Dense(hidden\_units) self.Ua = Dense(hidden\_units) self.norm = LayerNormalization() self.tanh = Activation(tf.keras.activations.tanh) self.add = Add()def call(self, context, x): Calculate the context vector based on all encoder hidden states and previous decoder state. :param: context: tensor, all encoder hidden states :param: x: tensor, previous state from Decoder context\_vector: tensor, the calculated context vector based on the input parameters # Expand dims to ensure scores shape = [batch, Ty, Tx] context = tf.expand\_dims(context, axis=1)  $x = tf.expand_dims(x, axis=2)$ scores = self.Va(self.tanh(self.add([self.Wa(context), self.Ua(x)]))) scores = tf.squeeze(scores, axis=-1) attn\_weights = tf.nn.softmax(scores, axis=-1) # NOTE: context shape = [batch, 1, Tx, feature] so that expand # dim of attention weights context\_vector = tf.expand\_dims(attn\_weights, axis=-1) \* context context\_vector = tf.reduce\_sum(context\_vector, axis=-2) context\_vector = self.norm(context\_vector) context\_vector = self.add([context\_vector, tf.squeeze(x, -2)]) return context\_vector Decoder Encoder and Decoder share the same structure as well as hidden units but the last dense layer at each state which holds for predicting the next word using a softmax. In [23]: # Let's define the decoder class Decoder(Layer): def \_\_init\_\_(self, tokenizer, embedding\_size, hidden\_units, dropout=DROPOUT): Decoder Block in seq2seq :param tokenizer: tokenizer of the source language :param embedding\_size: dimensionality of the embedding layer :param hidden\_units: dimensionality of the output super().\_\_init\_\_() self.tokenizer = tokenizer self.embedding\_size = embedding\_size self.hidden\_units = hidden\_units self.vocab = tokenizer.get\_vocabulary() self.vocab\_size = tokenizer.vocabulary\_size() self.embedding = Embedding(input\_dim=self.vocab\_size, output\_dim=embedding\_size) self.rnn = LSTM(units=hidden\_units, dropout=dropout, return\_sequences=True, return\_state=True) self.attention = BahdanauAttention(hidden\_units) self.dense = Dense(self.vocab\_size) def call(self, context, x, encoder\_state, training=True, return\_state**=False**): :param context: all encoder states :param x: all initial decoder states :param encoder\_state: last state from encoder :param training: :param return\_state: :return: logits: state\_h: hidden state state\_c: cell state mask = tf.where(x != 0, True, False)x = self.embedding(x)decoder\_outputs, state\_h, state\_c = self.rnn(x, initial\_state=encoder\_state, mask=mask, training=training) dense\_inputs = self.attention(context, decoder\_outputs) logits = self.dense(dense\_inputs) if return\_state: return logits, state\_h, state\_c return logits seq2seq Now we have got the Encoder and Decoder. Let's combine them into seq2seq model. In [24]: class NMT(Model): @classmethod def add\_method(cls, fun): setattr(cls, fun.\_\_name\_\_, fun) **return** fun def \_\_init\_\_(self, input\_tokenizer, output\_tokenizer, embedding\_size, hidden\_units): Initialize an instance for Neural Machine Translation Task :param input\_tokenizer: tokenizer of the input language :param output\_tokenizer: tokenizer of the output language :param embedding\_size: dimensionality of embedding layer :param hidden\_units: dimensionality of the output super().\_\_init\_\_() self.input\_tokenizer = input\_tokenizer self.output\_tokenizer = output\_tokenizer self.embedding\_size = embedding\_size self.hidden\_units = hidden\_units self.encoder = Encoder(input\_tokenizer, embedding\_size, hidden\_units) self.decoder = Decoder(output\_tokenizer, embedding\_size, hidden\_units) def call(self, inputs): encoder\_inputs, decoder\_inputs = inputs encoder\_outputs, state\_h, state\_c = self.encoder(encoder\_inputs) logits = self.decoder(encoder\_outputs, decoder\_inputs, [state\_h, state\_c]) return logits def get\_config(self): config = super().get\_config() config.update({ "input\_tokenizer": tf.keras.utils.serialize\_keras\_object(self.input\_tokenizer), "output\_tokenizer": tf.keras.utils.serialize\_keras\_object(self.output\_tokenizer), "embedding\_size": self.embedding\_size, "hidden\_units": self.hidden\_units }) return {\*\*config} In [25]: @NMT.add\_method def translate(self, next\_inputs, maxlen=40): def sampling(logits): probs = tf.nn.softmax(logits) dist = probs.numpy().squeeze() idx = np.random.choice(range(self.decoder.vocab\_size), p=dist) return idx translation = []next\_inputs = expand\_contractions(next\_inputs.lower(), en\_contraction\_map) next\_idx = np.asarray(self.encoder.tokenizer(next\_inputs)) while next\_idx.ndim != 2: next\_idx = tf.expand\_dims(next\_idx, axis=0) encoder\_outputs, state\_h, state\_c = self.encoder(next\_idx, training=False) next\_inputs = "[START]" next\_idx = np.asarray(word\_to\_idx[next\_inputs]) for i in range(maxlen): while next\_idx.ndim != 2: next\_idx = tf.expand\_dims(next\_idx, axis=0) logits, state\_h, state\_c = self.decoder(encoder\_outputs, next\_idx, [state\_h, state\_c], training=False, return\_state=True) next\_idx = sampling(logits) next\_inputs = self.decoder.vocab[next\_idx] if next\_inputs == "[END]": break elif next\_inputs == "[UNK]": continue translation.append(next\_inputs) return " ".join(translation) In [26]: model = NMT(en\_vec, ger\_vec, embedding\_size, hidden\_units) model.compile(optimizer=tf.keras.optimizers.Adam(0.005), loss=tf.keras.losses.SparseCategoricalCrossentropy(from\_logits=True), metrics=["accuracy"]) In [30]: history = model.fit(train\_ds, validation\_data=val\_ds) Epoch 1/5 Epoch 2/5 Epoch 3/5 Epoch 4/5 Epoch 5/5 In [28]: model.evaluate(test\_ds) [0.9508467316627502, 0.7981686592102051] Out[28]:

model.save\_weights("model\_v8.weights.h5")

In [39]: model.translate("Come on!")

'komm demnaechst !'