

1.Attention + Seq2Seq

2.Transformer

## 机器翻译

### 理论部分

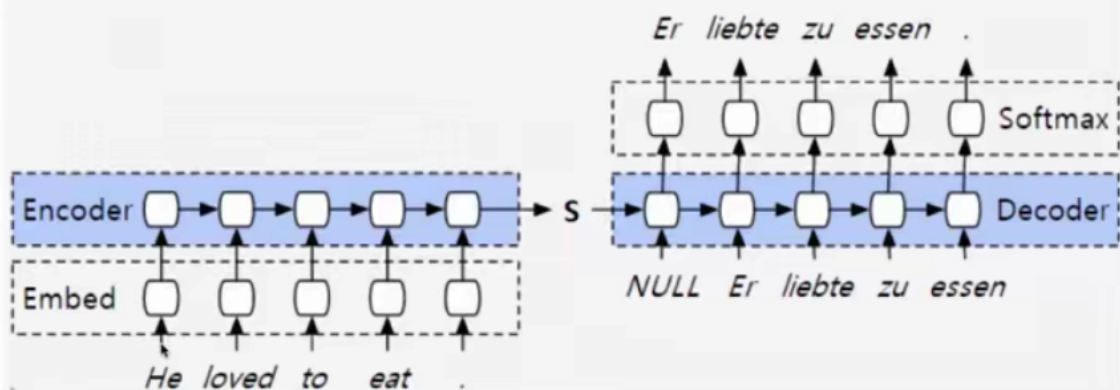
- ◆ Seq2seq模型
- ◆ Attention + Seq2Seq
- ◆ Transformer模型

## 1.Attention + Seq2Seq

## 机器翻译

### 模型思想-Attention

- ◆ 原始的seq2seq



# 机器翻译

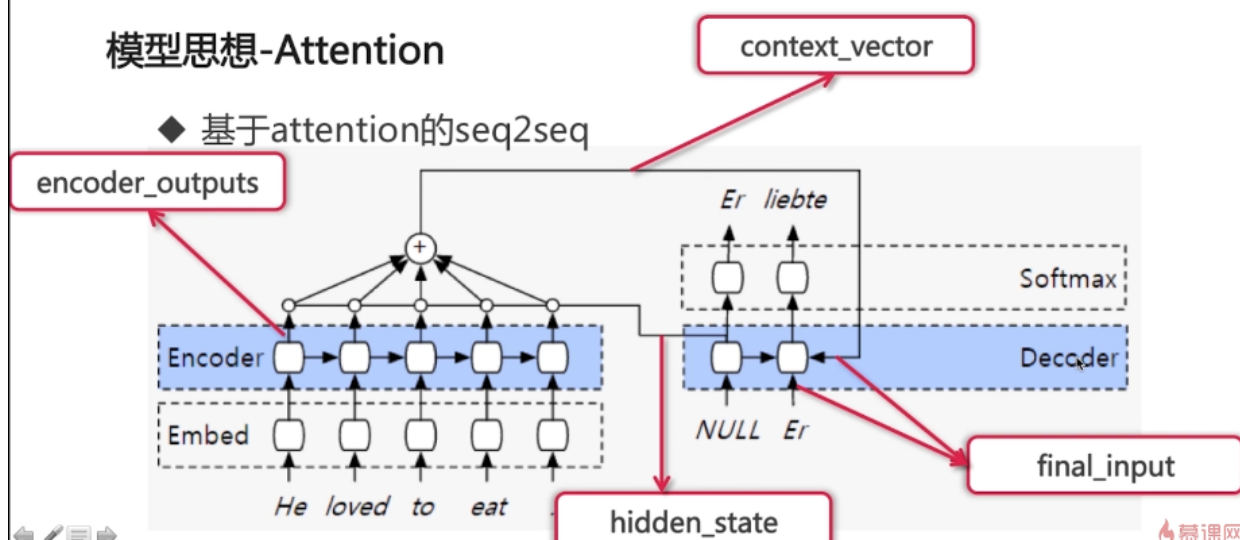
## 模型思想-Attention

- ◆ 原始的seq2seq
  - ◆ Encoder-Decoder结构
  - ◆ 缺点
    - ◆ 定长编码是信息瓶颈
    - ◆ 长度越长，前面输入进RNN的信息就越被稀释

# 机器翻译

## 模型思想-Attention

- ◆ 基于attention的seq2seq



## 机器翻译

### 模型思想-Attention

- ◆ EO: encoder各个位置的输出
- ◆ H: decoder某一步的隐含状态
- ◆ FC: 全连接层
- ◆ X: decoder的一个输入

- ◆  $\text{score} = \text{FC}(\tanh(\text{FC}(\text{EO}) + \text{FC}(\text{H})))$  [Bahdanau注意力]
- ◆ 另一选项:  $\text{score} = \text{EO} * \text{W} * \text{H}$  [luong注意力]
- ◆  $\text{attention\_weights} = \text{softmax}(\text{score}, \text{axis} = 1)$
- ◆  $\text{context} = \text{sum}(\text{attention\_weights} * \text{EO}, \text{axis} = 1)$
- ◆  $\text{final\_input} = \text{concat}(\text{context}, \text{embed}(\text{x}))$

```
1 %matplotlib inline
2 import matplotlib as mpl
3 import matplotlib.pyplot as plt
4 import numpy as np
5 import os
6 import pandas as pd
7 import sklearn
8 import sys
9 import tensorflow as tf
10 import time
11 from tensorflow import keras
12
13 # 1. preprocessing data
14 # 2. build model
15 # 2.1 encoder
16 # 2.2 attention
17 # 2.3 decoder
18 # 2.4 loss & optimizer
19 # 2.5 train
20 # 3. evaluation
21 # 3.1 give sentence, return translated results
22 # 3.2 visualize results (attention)
23
```

```

24 "1. preprocessing data"
25 import unicodedata
26 import re
27 from sklearn.model_selection import train_test_split
28
29 # 全为unicode, 词表过大, 转为ascii码, 则只有255个字符
30 # NFD 如果某个unicode由多个ascii码组成, 则将其拆开
31 # Mn为注音
32 def unicode_to_ascii(s):
33     return ''.join(c for c in unicodedata.normalize('NFD', s) if u
34         nicodeata.category(c) != 'Mn')
35
36 def preprocess_sentence(w):
37     w = unicode_to_ascii(w.lower().strip())
38     w = re.sub(r"([?!.,:])", r" \1 ", w)
39     # 去除重复空格
40     w = re.sub(r'" "+', " ", w)
41     # replacing everything with space except (a-z, A-Z, ".", "?",
42     "!", ",")
43     w = re.sub(r"[^a-zA-Z?!.,:]+", " ", w)
44     # 去除前后空格
45     w = w.rstrip().strip()
46     # adding a start and an end token to the sentence
47     # so that the model know when to start and stop predicting.
48     w = '<start> ' + w + ' <end>'
49     return w
50
51 data_path = './spa.txt'
52 # 1. Remove the accents
53 # 2. Clean the sentences
54 # 3. Return word pairs in the format: [ENGLISH, SPANISH]
55 def create_dataset(path, num_examples):
56     lines = open(path, encoding='UTF-
57         8').read().strip().split('\n')
58     word_pairs = [[preprocess_sentence(w) for w in
59         line.split('\t')] for line in lines[:num_examples]]

```

```

56 # 解包 [(1,2),(3,4)] -> [(1,3),(2,4)]
57 return zip(*word_pairs)
58
59 en, sp = create_dataset(data_path, None)
60
61 def max_length(tensor):
62     return max(len(t) for t in tensor)
63
64 def tokenize(lang):
65     lang_tokenizer =
66     tf.keras.preprocessing.text.Tokenizer(filters='', split=' ')
67     # 统计词频, 生成词表
68     lang_tokenizer.fit_on_texts(lang)
69     # word -> id
70     tensor = lang_tokenizer.texts_to_sequences(lang)
71     tensor = tf.keras.preprocessing.sequence.pad_sequences(tensor,
72     padding='post')
73     return tensor, lang_tokenizer
74
75 def load_dataset(path, num_examples=None):
76     # creating cleaned input, output pairs
77     targ_lang, inp_lang = create_dataset(path, num_examples)
78     input_tensor, inp_lang_tokenizer = tokenize(inp_lang)
79     target_tensor, targ_lang_tokenizer = tokenize(targ_lang)
80     return input_tensor, target_tensor, inp_lang_tokenizer, targ_l
81     ang_tokenizer
82
83 # Try experimenting with the size of that dataset
84 num_examples = 30000
85 input_tensor, target_tensor, inp_lang, targ_lang =
86 load_dataset(data_path, num_examples)
87
88 # Calculate max_length of the target tensors
89 max_length_targ, max_length_inp = max_length(target_tensor), ma
90 x_length(input_tensor)
91

```

```

87 # Creating training and validation sets using an 80-20 split
88 input_tensor_train, input_tensor_val, target_tensor_train, target_tensor_val = train_test_split(input_tensor, target_tensor, test_size=0.2)
89
90 # Show length
91 len(input_tensor_train), len(target_tensor_train), len(input_tensor_val), len(target_tensor_val)
92
93 BUFFER_SIZE = len(input_tensor_train)
94 BATCH_SIZE = 64
95 steps_per_epoch = len(input_tensor_train)//BATCH_SIZE
96 embedding_dim = 256
97 units = 1024
98 vocab_inp_size = len(inp_lang.word_index)+1
99 vocab_tar_size = len(targ_lang.word_index)+1
100
101 dataset = tf.data.Dataset.from_tensor_slices((input_tensor_train, target_tensor_train)).shuffle(BUFFER_SIZE)
102 dataset = dataset.batch(BATCH_SIZE, drop_remainder=True)
103
104 "2. build model"
105 "2.1 encoder"
106 class Encoder(tf.keras.Model):
107     def __init__(self, vocab_size, embedding_dim, encoding_units, batch_size):
108         super(Encoder, self).__init__()
109         self.batch_size = batch_size
110         self.encoding_units = encoding_units
111         self.embedding = keras.layers.Embedding(vocab_size, embedding_dim)
112         self.gru = keras.layers.GRU(self.encoding_units, return_sequences=True, return_state=True, recurrent_initializer='glorot_uniform')
113
114     def call(self, x, hidden):
115         x = self.embedding(x)

```

```

116 output, state = self.gru(x, initial_state = hidden)
117 return output, state
118
119 def initialize_hidden_state(self):
120 return tf.zeros((self.batch_size, self.encoding_units))
121 encoder = Encoder(vocab_inp_size, embedding_dim, units, BATCH_SIZE)
122
123 "2.2 attention"
124 class BahdanauAttention(tf.keras.Model):
125     def __init__(self, units):
126         super(BahdanauAttention, self).__init__()
127         self.W1 = tf.keras.layers.Dense(units)
128         self.W2 = tf.keras.layers.Dense(units)
129         self.V = tf.keras.layers.Dense(1)
130
131     def call(self, query, values):
132         # hidden shape == (batch_size, hidden size)
133         # hidden_with_time_axis shape == (batch_size, 1, hidden size)
134         # we are doing this to perform addition to calculate the score
135         hidden_with_time_axis = tf.expand_dims(query, 1)
136
137         # score shape == (batch_size, max_length, 1)
138         # we get 1 at the last axis because we are applying score to self.V
139         # the shape of the tensor before applying self.V is (batch_size, max_length, units)
140         score = self.V(tf.nn.tanh(self.W1(values) + self.W2(hidden_with_time_axis)))
141
142         # attention_weights shape == (batch_size, max_length, 1)
143         attention_weights = tf.nn.softmax(score, axis=1)
144
145         # context_vector shape after sum == (batch_size, hidden_size)
146         context_vector = attention_weights * values
147         context_vector = tf.reduce_sum(context_vector, axis=1)

```

```

148
149     return context_vector, attention_weights
150 # attention_layer = BahdanauAttention(10)
151
152 "2.3 decoder"
153 class Decoder(tf.keras.Model):
154     def __init__(self, vocab_size, embedding_dim, decoding_units,
155                  batch_size):
156         super(Decoder, self).__init__()
157         self.batch_size = batch_size
158         self.decoding_units = decoding_units
159         self.embedding = keras.layers.Embedding(vocab_size, embedding_dim)
160         self.gru = keras.layers.GRU(self.decoding_units, return_sequences=True, return_state=True, recurrent_initializer='glorot_uniform')
161         self.fc = keras.layers.Dense(vocab_size)
162         # used for attention
163         self.attention = BahdanauAttention(self.decoding_units)
164
165     def call(self, x, hidden, encoding_output):
166         # enc_output shape == (batch_size, max_length, hidden_size)
167         context_vector, attention_weights = self.attention(hidden, encoding_output)
168
169         # x shape after passing through embedding == (batch_size, 1, embedding_dim)
170         x = self.embedding(x)
171
172         # x shape after concatenation == (batch_size, 1, embedding_dim + hidden_size)
173         x = tf.concat([tf.expand_dims(context_vector, 1), x], axis=-1)
174
175         # passing the concatenated vector to the GRU
176         output, state = self.gru(x)
177
178         # output shape == (batch_size * 1, hidden_size)

```



```

178 output = tf.reshape(output, (-1, output.shape[2]))
179
180 # output shape == (batch_size, vocab)
181 x = self.fc(output)
182 return x, state, attention_weights
183 decoder = Decoder(vocab_tar_size, embedding_dim, units, BATCH_SIZE)
184
185 "2.4 loss & optimizer"
186 optimizer = keras.optimizers.Adam()
187 # reduction='none': 损失函数如何聚合，此处不聚合
188 loss_object = keras.losses.SparseCategoricalCrossentropy(from_logits=True, reduction='none')
189
190 def loss_function(real, pred):
191     # padding 填充的是0
192     mask = tf.math.logical_not(tf.math.equal(real, 0))
193     loss_ = loss_object(real, pred)
194     mask = tf.cast(mask, dtype=loss_.dtype)
195     loss_ *= mask
196     return tf.reduce_mean(loss_)
197
198 checkpoint_dir = './10-1_checkpoints'
199 if not os.path.exists(checkpoint_dir):
200     os.mkdir(checkpoint_dir)
201 checkpoint_prefix = os.path.join(checkpoint_dir, "ckpt")
202 checkpoint = tf.train.Checkpoint(optimizer=optimizer, encoder=encoder, decoder=decoder)
203
204 "2.5 train"
205 @tf.function
206 def train_step(inp, targ, encoding_hidden):
207     loss = 0
208     with tf.GradientTape() as tape:
209         encoding_output, encoding_hidden = encoder(inp, encoding_hidden)

```

```

210
211     decoding_hidden = encoding_hidden
212     # eg: <start> i am here <end>
213     # 1. <start> -> i
214     # 2. i -> am
215     # 3. am -> here
216     # 4. here -> <end>
217
218     # (batch_size,1)
219     decoding_input = tf.expand_dims([targ_lang.word_index['<start>']] * BATCH_SIZE, 1)
220     # Teacher forcing - feeding the target as the next input
221     for t in range(1, targ.shape[1]):
222         # passing enc_output to the decoder
223         predictions, decoding_hidden, _ = decoder(decoding_input, decoding_hidden, encoding_output)
224         loss += loss_function(targ[:, t], predictions)
225         # using teacher forcing
226         decoding_input = tf.expand_dims(targ[:, t], 1)
227
228     batch_loss = (loss / int(targ.shape[0]))
229     variables = encoder.trainable_variables + decoder.trainable_variables
230     gradients = tape.gradient(loss, variables)
231     optimizer.apply_gradients(zip(gradients, variables))
232     return batch_loss
233
234     EPOCHS = 10
235     for epoch in range(EPOCHS):
236         start = time.time()
237         encoding_hidden = encoder.initialize_hidden_state()
238         total_loss = 0
239         for (batch, (inp, targ)) in enumerate(dataset.take(steps_per_epoch)):
240             batch_loss = train_step(inp, targ, encoding_hidden)
241             total_loss += batch_loss

```

```

242     if batch % 100 == 0:
243         print('Epoch {} Batch {} Loss {:.4f}'.format(epoch + 1, batch,
            batch_loss.numpy()))
244     # saving (checkpoint) the model every 2 epochs
245     if (epoch + 1) % 2 == 0:
246         checkpoint.save(file_prefix = checkpoint_prefix)
247
248     print('Epoch {} Loss {:.4f}'.format(epoch + 1, total_loss / steps_per_epoch))
249     print('Time taken for 1 epoch {} sec\n'.format(time.time() - start))
250
251     "3. evaluation"
252     "3.1 give sentence, return translated results"
253     "3.2 visualize results (attention)"
254     def evaluate(sentence):
255         # attention_plot: (11,16)
256         attention_plot = np.zeros((max_length_targ, max_length_inp))
257         sentence = preprocess_sentence(sentence)
258         inputs = [inp_lang.word_index[i] for i in sentence.split(' ')]
259         inputs = keras.preprocessing.sequence.pad_sequences([inputs],
            maxlen=max_length_inp, padding='post')
260         inputs = tf.convert_to_tensor(inputs)
261
262         result = ''
263         hidden = tf.zeros((1, units))
264         encoding_out, encoding_hidden = encoder(inputs, hidden)
265         decoding_hidden = encoding_hidden
266         # eg: <start> -> A
267         # A -> B -> C -> D
268         # decoding_input.shape: (1,1)
269         decoding_input = tf.expand_dims([targ_lang.word_index['<start>']], 0)
270
271         for t in range(max_length_targ):
272             predictions, decoding_hidden, attention_weights = decoder(decoding_input, decoding_hidden, encoding_out)

```

```

273 # storing the attention weights to plot later on
274 # attention_weights: (1,16,1) -> 16维向量
275 attention_weights = tf.reshape(attention_weights, (-1, ))
276 attention_plot[t] = attention_weights.numpy()
277 # prediction.shape: (batch_size, vocab_size) (1, 4935)
278 predicted_id = tf.argmax(predictions[0]).numpy()
279 result += targ_lang.index_word[predicted_id] + ' '
280 if targ_lang.index_word[predicted_id] == '<end>':
281     return result, sentence, attention_plot
282 # the predicted ID is fed back into the model
283 decoding_input = tf.expand_dims([predicted_id], 0)
284 return result, sentence, attention_plot
285
286 # function for plotting the attention weights
287 def plot_attention(attention, sentence, predicted_sentence):
288     fig = plt.figure(figsize=(10,10))
289     ax = fig.add_subplot(1, 1, 1)
290     ax.matshow(attention, cmap='viridis')
291     fontdict = {'fontsize': 14}
292     ax.set_xticklabels([''] + sentence, fontdict=fontdict, rotation
n=90)
293     ax.set_yticklabels([''] + predicted_sentence, fontdict=fontdict,
t)
294     plt.show()
295
296 def translate(sentence):
297     result, sentence, attention_plot = evaluate(sentence)
298     print('Input: %s' % (sentence))
299     print('Predicted translation: {}'.format(result))
300     attention_plot = attention_plot[:len(result.split(' ')),
:len(sentence.split(' '))]
301     plot_attention(attention_plot, sentence.split(' '), result.sp
it(' '))
302
303 checkpoint.restore(tf.train.latest_checkpoint(checkpoint_dir))
304 translate(u'hace mucho frio aqui.')

```

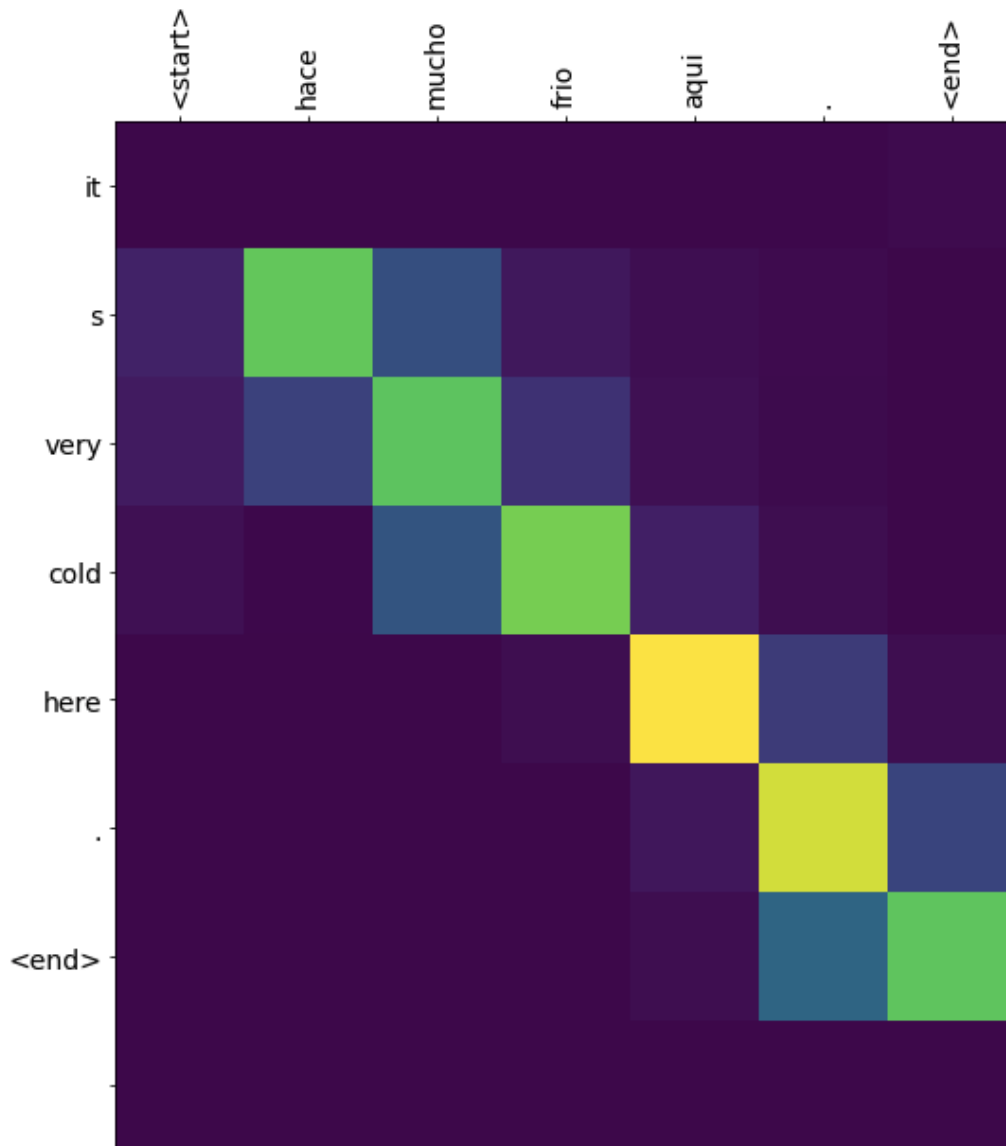
```
translate(u'hace mucho frio aqui.')
```

Input: <start> hace mucho frio aqui . <end>

Predicted translation: it s very cold here . <end>

E:\Anaconda3\lib\site-packages\ipykernel\_launcher.py:51: UserWarning: FixedFormatter should only be used together with FixedLocator

E:\Anaconda3\lib\site-packages\ipykernel\_launcher.py:52: UserWarning: FixedFormatter should only be used together with FixedLocator



## 2.Transformer

# 机器翻译

## 模型思想-Attention

- ◆ 基于attention的seq2seq
  - ◆ 去除定长编码瓶颈，信息无损从Encoder传到Decoder
- ◆ 但是
  - ◆ 采用GRU，计算依然有瓶颈，并行度不高
  - ◆ 只有Encoder和Decoder之间有attention

# 机器翻译

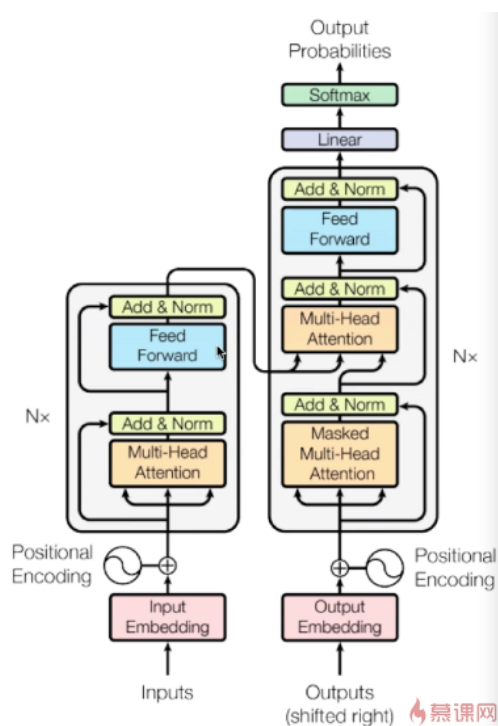
## 模型思想-Attention

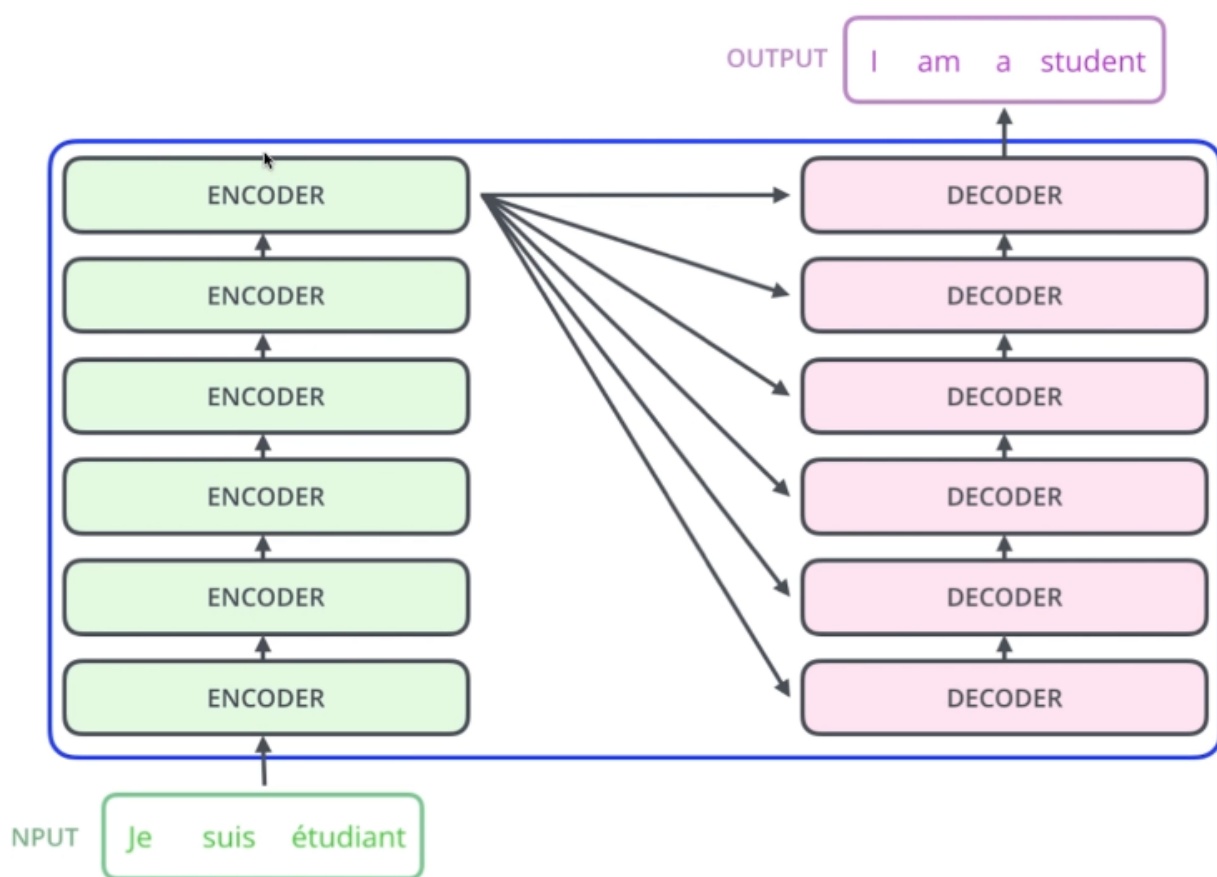
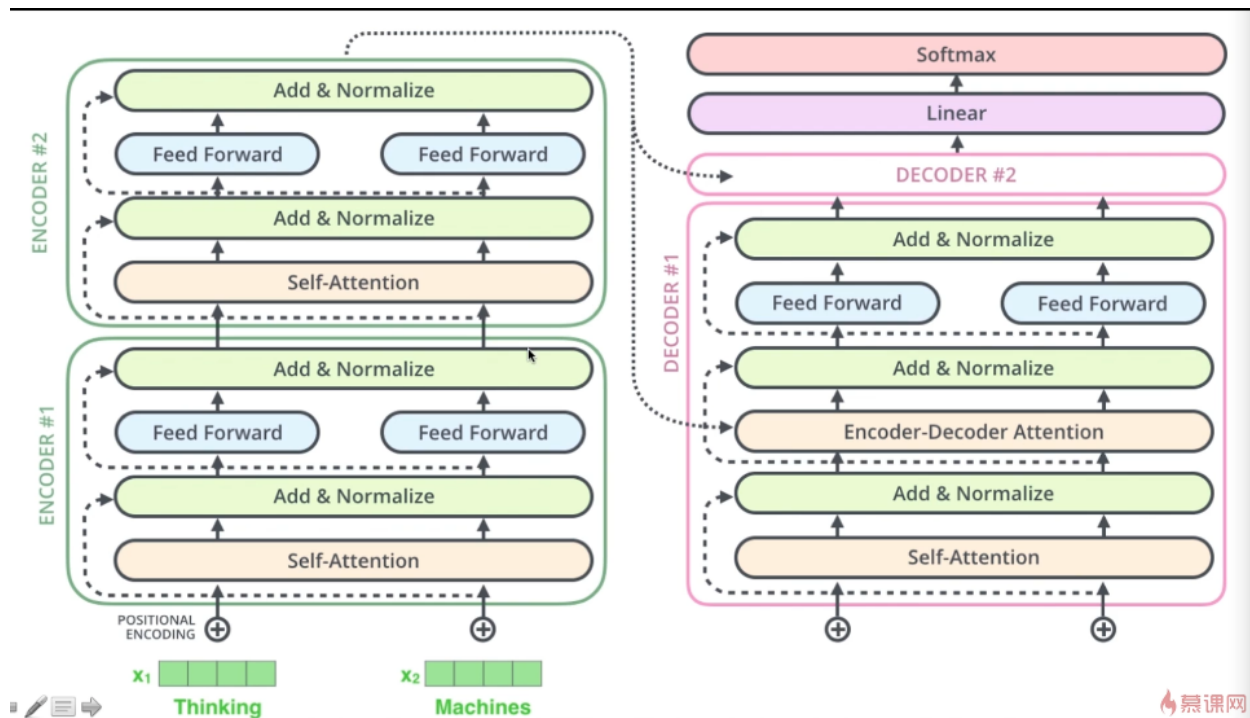
- ◆ 能否去掉RNN？
- ◆ 能否给输入和输出分别加上self attention？
- **GRU中顺序处理序列，而Transformer中不知道句子的顺序，因此需要添加Positional Encoding**

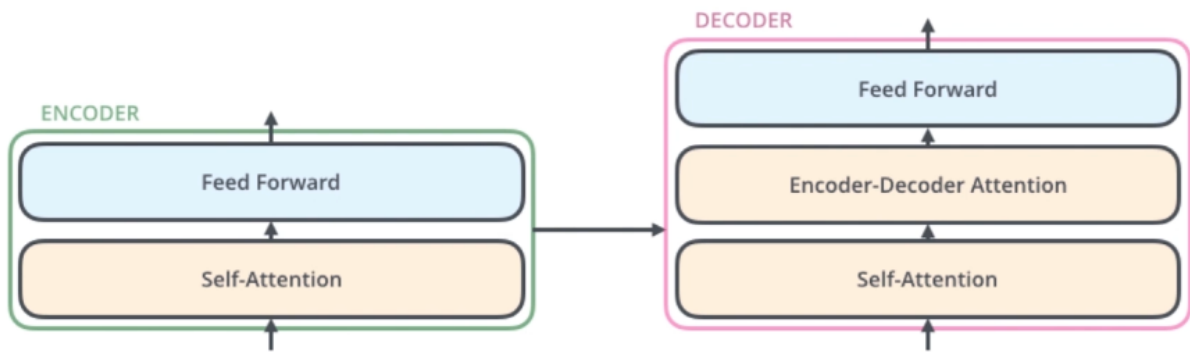
## Transformer模型

### 模型结构-Transformer

- ◆ Encoder-Decoder
- ◆ 多层Encoder-Decoder
- ◆ 位置编码
- ◆ 多头注意力
  - ◆ 缩放点积注意力
- ◆ Add & norm







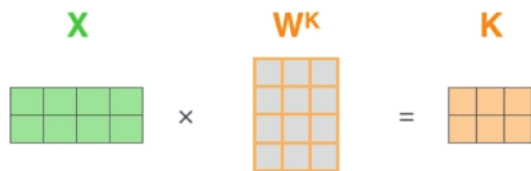
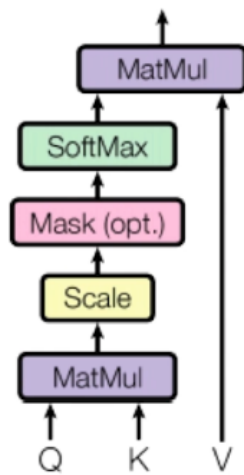
## Transformer模型

### 模型结构-Attention

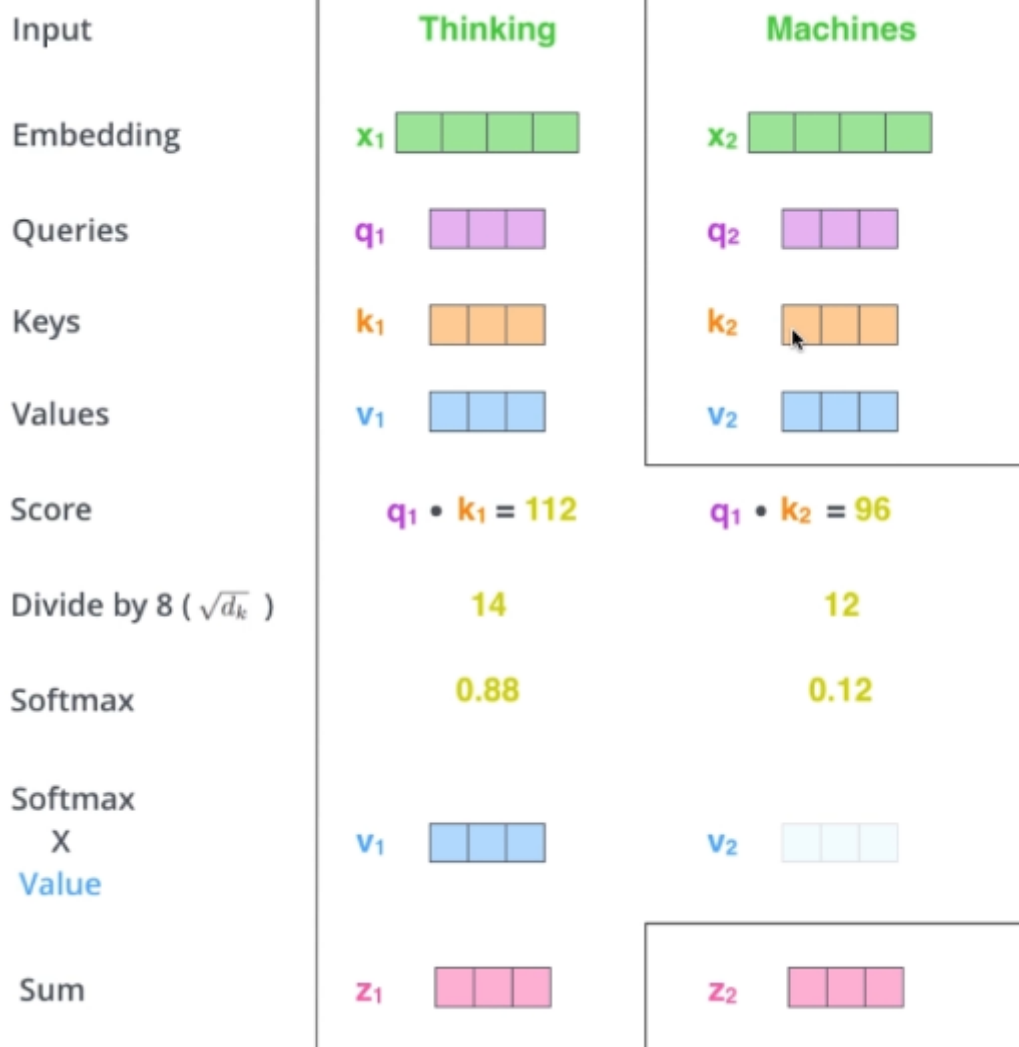
◆ 缩放点积attention

$$\text{Attention}(Q, K, V) = \text{softmax}\left(\frac{QK^T}{\sqrt{d_k}}\right)V$$

### Scaled Dot-Product Attention

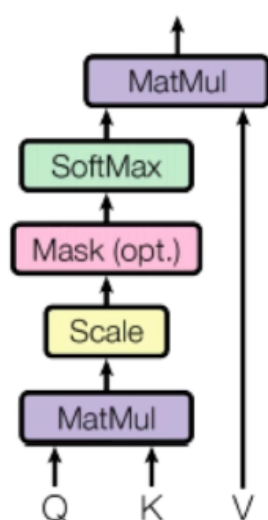




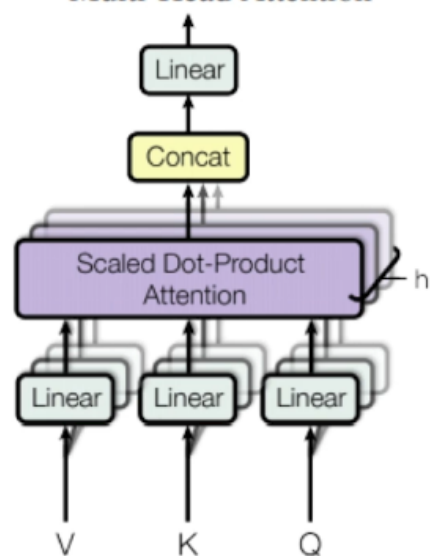


## Transformer模型

Scaled Dot-Product Attention



Multi-Head Attention



# Transformer模型

## 模型结构-Transformer

### ◆ 位置编码

$$PE_{(pos, 2i)} = \sin(pos/10000^{2i/d_{\text{model}}})$$

$$PE_{(pos, 2i+1)} = \cos(pos/10000^{2i/d_{\text{model}}})$$

### ◆ Pos = 1, d = 128

### ◆ [sin(1/10000^(0/128)), cos(1/10000^(1/128)), sin(1/10000^(2/128)), cos(1/10000^(3/128)), ...]



# Transformer模型

## 模型结构-Decoder

- ◆ Train的时候并行化
- ◆ Inference的时候仍然要序列式完成
- ◆ Self attention时前词不能见后词
  - ◆ Mask来实现

```
1 import matplotlib as mpl
2 import matplotlib.pyplot as plt
3 %matplotlib inline
4 import numpy as np
5 import sklearn
6 import pandas as pd
7 import os
8 import sys
9 import time
10 import tensorflow as tf
11 from tensorflow import keras
```

```

12 print(tf.__version__)
13 print(sys.version_info)
14 for module in mpl, np, pd, sklearn, tf, keras:
15     print(module.__name__, module.__version__)
16
17 # 1. loads data
18 # 2. preprocesses data -> dataset
19 # 3. tools
20 # 3.1 generates position embedding
21 # 3.2 create mask. (a. padding, b. decoder)
22 # 3.3 scaled_dot_product_attention
23 # 4. builds model
24 # 4.1 MultiheadAttention
25 # 4.2 EncoderLayer
26 # 4.3 DecoderLayer
27 # 4.4 EncoderModel
28 # 4.5 DecoderModel
29 # 4.6 Transformer
30 # 5. optimizer & loss
31 # 6. train step -> train
32 # 7. Evaluate and Visualize
33
34 "1. loads data"
35 import tensorflow_datasets as tfds
36 examples, info = tfds.load('ted_hrlr_translate/pt_to_en', with_info=True, as_supervised = True)
37 train_examples, val_examples = examples['train'], examples['validation']
38
39 "2. preprocesses data -> dataset"
40 en_tokenizer = tfds.deprecated.text.SubwordTextEncoder.build_from_corpus((en.numpy() for pt, en in train_examples), target_vocab_size = 2 ** 13)
41 pt_tokenizer = tfds.deprecated.text.SubwordTextEncoder.build_from_corpus((pt.numpy() for pt, en in train_examples), target_vocab_size = 2 ** 13)
42 buffer_size = 20000

```

```

43 batch_size = 64
44 max_length = 40
45
46 def encode_to_subword(pt_sentence, en_sentence):
47     # pt_tokenizer.vocab_size: 用以指代<start>
48     # pt_tokenizer.vocab_size + 1: 用以指代<end>
49     pt_sequence = [pt_tokenizer.vocab_size] \
50     + pt_tokenizer.encode(pt_sentence.numpy()) \
51     + [pt_tokenizer.vocab_size + 1]
52     en_sequence = [en_tokenizer.vocab_size] \
53     + en_tokenizer.encode(en_sentence.numpy()) \
54     + [en_tokenizer.vocab_size + 1]
55     return pt_sequence, en_sequence
56
57 def filter_by_max_length(pt, en):
58     return tf.logical_and(tf.size(pt) <= max_length, tf.size(en) <
59     = max_length)
60
61 # dataset的map无法直接调用python函数，需要转为py_function
62 def tf_encode_to_subword(pt_sentence, en_sentence):
63     return tf.py_function(encode_to_subword, [pt_sentence, en_sentence], [tf.int64, tf.int64])
64
65 train_dataset = train_examples.map(tf_encode_to_subword)
66 train_dataset = train_dataset.filter(filter_by_max_length)
67 # [-1], [-1]: 数据由两个维度，每个维度都扩展到当前维度下最高的值
68 train_dataset =
69 train_dataset.shuffle(buffer_size).padded_batch(batch_size, padded
70 _shapes=([-1], [-1]))
71
72 valid_dataset = val_examples.map(tf_encode_to_subword)
73 valid_dataset = valid_dataset.filter(filter_by_max_length).padd
74 ed_batch(batch_size, padded_shapes=([-1], [-1]))

```

```

1 "3. tools"
2 "3.1 generates position embedding"
3 # PE(pos, 2i) = sin(pos / 10000^(2i/d_model))

```

```

4 # PE(pos, 2i+1) = cos(pos / 10000^(2i/d_model))
5
6 # pos.shape: [sentence_length, 1] 词在句子中的位置
7 # i.shape : [1, d_model] 词在embed中的位置
8 # result.shape: [sentence_length, d_model]
9 def get_angles(pos, i, d_model):
10     angle_rates = 1 / np.power(10000,(2 * (i // 2)) / np.float32(d_model))
11     return pos * angle_rates
12
13 def get_position_embedding(sentence_length, d_model):
14     angle_rads = get_angles(np.arange(sentence_length)[: , np.newaxis], np.arange(d_model)[np.newaxis, :], d_model)
15     # sines.shape: [sentence_length, d_model / 2]
16     # cosines.shape: [sentence_length, d_model / 2]
17     sines = np.sin(angle_rads[:, 0::2])
18     cosines = np.cos(angle_rads[:, 1::2])
19     # position_embedding.shape: [sentence_length, d_model]
20     position_embedding = np.concatenate([sines, cosines], axis = -1)
21     # position_embedding.shape: [1, sentence_length, d_model]
22     position_embedding = position_embedding[np.newaxis, ...]
23     return tf.cast(position_embedding, dtype=tf.float32)
24
25 "3.2 create mask. (a. padding, b. decoder)"
26 # 1. padding mask, 2. look ahead
27
28 # batch_data.shape: [batch_size, seq_len]
29 def create_padding_mask(batch_data):
30     padding_mask = tf.cast(tf.math.equal(batch_data, 0), tf.float32)
31     # [batch_size, 1, 1, seq_len]
32     return padding_mask[:, tf.newaxis, tf.newaxis, :]
33
34 # attention_weights.shape: [3,3]
35 # [[1, 0, 0],

```

```

36 # [4, 5, 0], 第二个单词分别与第一、二、三个单词的attention
37 # [7, 8, 9]]
38 # 由于mask时padding为1, 其余为0; 所以需要创建一个上三角矩阵
39 def create_look_ahead_mask(size):
40     mask = 1 - tf.linalg.band_part(tf.ones((size, size)), -1, 0)
41     return mask # (seq_len, seq_len)
42
43 # create_look_ahead_mask(3)
44 # <tf.Tensor: shape=(3, 3), dtype=float32, numpy=
45 # array([[0., 1., 1.],
46 # [0., 0., 1.],
47 # [0., 0., 0.]], dtype=float32)>
48
49 "3.3 scaled_dot_product_attention"
50 def scaled_dot_product_attention(q, k, v, mask):
51     """
52     Args:
53     - q: shape == (... , seq_len_q, depth)
54     - k: shape == (... , seq_len_k, depth)
55     - v: shape == (... , seq_len_v, depth_v)
56     - seq_len_k == seq_len_v
57     - mask: shape == (... , seq_len_q, seq_len_k)
58     Returns:
59     - output: weighted sum
60     - attention_weights: weights of attention
61     """
62
63     # matmul_qk.shape: (... , seq_len_q, seq_len_k)
64     # transpose_b = True: 第二个矩阵做转置
65     matmul_qk = tf.matmul(q, k, transpose_b = True)
66     dk = tf.cast(tf.shape(k)[-1], tf.float32)
67     scaled_attention_logits = matmul_qk / tf.math.sqrt(dk)
68     if mask is not None:
69         # 使得在softmax后值趋近于0 (一个数加上一个极小值, softmax后趋向于0)

```

```

70 scaled_attention_logits += (mask * -1e9)
71 # attention_weights.shape: (... , seq_len_q, seq_len_k)
72 attention_weights = tf.nn.softmax(scaled_attention_logits, axis
73 s = -1)
74 # output.shape: (... , seq_len_q, depth_v)
74 output = tf.matmul(attention_weights, v)
75 return output, attention_weights
76
77 def print_scaled_dot_product_attention(q, k, v):
78     temp_out, temp_att = scaled_dot_product_attention(q, k, v, Non
79 e)
80     print("Attention weights are:")
81     print(temp_att)
82     print("Output is:")
83     print(temp_out)

```

```

1 "4. builds model"
2 "4.1 MultiheadAttention"
3 class MultiHeadAttention(keras.layers.Layer):
4     """
5     理论上:
6     x -> Wq0 -> q0
7     x -> Wk0 -> k0
8     x -> Wv0 -> v0
9
10    实战中:
11    q -> Wq0 -> q0
12    k -> Wk0 -> k0
13    v -> Wv0 -> v0
14
15    实战中技巧:
16    q -> Wq -> Q -> split -> q0, q1, q2...
17    """
18    def __init__(self, d_model, num_heads):
19        super(MultiHeadAttention, self).__init__()
20        self.num_heads = num_heads

```

```

21 self.d_model = d_model
22 assert self.d_model % self.num_heads == 0
23
24 self.depth = self.d_model // self.num_heads
25
26 self.WQ = keras.layers.Dense(self.d_model)
27 self.WK = keras.layers.Dense(self.d_model)
28 self.WV = keras.layers.Dense(self.d_model)
29
30 self.dense = keras.layers.Dense(self.d_model)
31
32 def split_heads(self, x, batch_size):
33     # x.shape: (batch_size, seq_len, d_model)
34     # d_model = num_heads * depth
35     # x -> (batch_size, num_heads, seq_len, depth)
36
37     x = tf.reshape(x,
38                     (batch_size, -1, self.num_heads, self.depth))
39     return tf.transpose(x, perm=[0, 2, 1, 3])
40
41 def call(self, q, k, v, mask):
42     batch_size = tf.shape(q)[0]
43
44     q = self.WQ(q) # q.shape: (batch_size, seq_len_q, d_model)
45     k = self.WK(k) # k.shape: (batch_size, seq_len_k, d_model)
46     v = self.WV(v) # v.shape: (batch_size, seq_len_v, d_model)
47
48     # q.shape: (batch_size, num_heads, seq_len_q, depth)
49     q = self.split_heads(q, batch_size)
50     # k.shape: (batch_size, num_heads, seq_len_k, depth)
51     k = self.split_heads(k, batch_size)
52     # v.shape: (batch_size, num_heads, seq_len_v, depth)
53     v = self.split_heads(v, batch_size)
54
55     # scaled_attention_outputs.shape: (batch_size, num_heads, seq_
    len_q, depth)

```



```

56 # attention_weights.shape: (batch_size, num_heads, seq_len_q,
    seq_len_k)
57 scaled_attention_outputs, attention_weights = \
58 scaled_dot_product_attention(q, k, v, mask)
59
60 # scaled_attention_outputs.shape: (batch_size, seq_len_q, num_
    heads, depth)
61 scaled_attention_outputs = tf.transpose(
62 scaled_attention_outputs, perm = [0, 2, 1, 3])
63 # concat_attention.shape: (batch_size, seq_len_q, d_model)
64 concat_attention = tf.reshape(scaled_attention_outputs,
65 (batch_size, -1, self.d_model))
66
67 # output.shape : (batch_size, seq_len_q, d_model)
68 output = self.dense(concat_attention)
69 return output, attention_weights
70
71 def feed_forward_network(d_model, dff):
72 # dff: dim of feed forward network.
73 return keras.Sequential([
74 keras.layers.Dense(dff, activation='relu'),
75 keras.layers.Dense(d_model)
76 ])
77
78 "4.2 EncoderLayer"
79 class EncoderLayer(keras.layers.Layer):
80     """
81     x -> self attention -> add & normalize & dropout
82     -> feed_forward -> add & normalize & dropout
83     """
84     def __init__(self, d_model, num_heads, dff, rate=0.1):
85         super(EncoderLayer, self).__init__()
86         self.mha = MultiHeadAttention(d_model, num_heads)
87         self.ffn = feed_forward_network(d_model, dff)
88
89         self.layer_norm1 = keras.layers.LayerNormalization(epsilon = 1
    e-6)

```

```

90 self.layer_norm2 = keras.layers.LayerNormalization(epsilon = 1
e-6)
91
92 self.dropout1 = keras.layers.Dropout(rate)
93 self.dropout2 = keras.layers.Dropout(rate)
94
95 def call(self, x, training, encoder_padding_mask):
96 # x.shape : (batch_size, seq_len, dim=d_model)
97 # attn_output.shape: (batch_size, seq_len, d_model)
98 # out1.shape : (batch_size, seq_len, d_model)
99 attn_output, _ = self.mha(x, x, x, encoder_padding_mask)
100 attn_output = self.dropout1(attn_output, training=training)
101 out1 = self.layer_norm1(x + attn_output)
102
103 # ffn_output.shape: (batch_size, seq_len, d_model)
104 # out2.shape : (batch_size, seq_len, d_model)
105 ffn_output = self.ffn(out1)
106 ffn_output = self.dropout2(ffn_output, training=training)
107 out2 = self.layer_norm2(out1 + ffn_output)
108 return out2
109
110 "4.3 DecoderLayer"
111 class DecoderLayer(keras.layers.Layer):
112     """
113     x -> self attention -> add & normalize & dropout -> out1
114     out1 , encoding_outputs -> attention -> add & normalize & dropout -> out2
115     out2 -> ffn -> add & normalize & dropout -> out3
116     """
117     def __init__(self, d_model, num_heads, dff, rate = 0.1):
118         super(DecoderLayer, self).__init__()
119
120         self.mha1 = MultiHeadAttention(d_model, num_heads)
121         self.mha2 = MultiHeadAttention(d_model, num_heads)
122
123         self.ffn = feed_forward_network(d_model, dff)
124

```

```

125 self.layer_norm1 = keras.layers.LayerNormalization(epsilon = 1e-6)
126 self.layer_norm2 = keras.layers.LayerNormalization(epsilon = 1e-6)
127 self.layer_norm3 = keras.layers.LayerNormalization(epsilon = 1e-6)
128
129 self.dropout1 = keras.layers.Dropout(rate)
130 self.dropout2 = keras.layers.Dropout(rate)
131 self.dropout3 = keras.layers.Dropout(rate)
132
133 def call(self, x, encoding_outputs, training,
134 decoder_mask, encoder_decoder_padding_mask):
135     # decoder_mask: 由look_ahead_mask和decoder_padding_mask合并而来
136     # x.shape: (batch_size, target_seq_len, d_model)
137     # encoding_outputs.shape: (batch_size, input_seq_len, d_model)
138
139     # attn1, out1.shape : (batch_size, target_seq_len, d_model)
140     attn1, attn_weights1 = self.mha1(x, x, x, decoder_mask)
141     attn1 = self.dropout1(attn1, training = training)
142     out1 = self.layer_norm1(attn1 + x)
143
144     # attn2, out2.shape : (batch_size, target_seq_len, d_model)
145     attn2, attn_weights2 = self.mha2(out1, encoding_outputs, encoding_outputs, encoder_decoder_padding_mask)
146     attn2 = self.dropout2(attn2, training = training)
147     out2 = self.layer_norm2(attn2 + out1)
148
149     # ffn_output, out3.shape: (batch_size, target_seq_len, d_model)
150     ffn_output = self.ffn(out2)
151     ffn_output = self.dropout3(ffn_output, training=training)
152     out3 = self.layer_norm3(ffn_output + out2)
153     return out3, attn_weights1, attn_weights2
154
155 "4.4 EncoderModel"
156 class EncoderModel(keras.layers.Layer):

```

```

157 def __init__(self, num_layers, input_vocab_size, max_length, d_model, num_heads, dff, rate=0.1):
158     super(EncoderModel, self).__init__()
159     self.d_model = d_model
160     self.num_layers = num_layers
161     self.max_length = max_length
162
163     self.embedding =
keras.layers.Embedding(input_vocab_size, self.d_model)
164     # position_embedding.shape: (1, max_length, d_model)
165     self.position_embedding = get_position_embedding(max_length, self.d_model)
166
167     self.dropout = keras.layers.Dropout(rate)
168     self.encoder_layers = [
169     EncoderLayer(d_model, num_heads, dff, rate)
170     for _ in range(self.num_layers)]
171
172     def call(self, x, training, encoder_padding_mask):
173     # x.shape: (batch_size, input_seq_len)
174     input_seq_len = tf.shape(x)[1]
175     tf.debugging.assert_less_equal(input_seq_len,
self.max_length, "input_seq_len should be less or equal to self.max
_length")
176
177     # x.shape: (batch_size, input_seq_len, d_model)
178     x = self.embedding(x)
179     # 扩大x的值, 使得与position_embedding相加后, x起的作用大一些
180     x *= tf.math.sqrt(tf.cast(self.d_model, tf.float32))
181     x += self.position_embedding[:, :input_seq_len, :]
182     x = self.dropout(x, training = training)
183
184     for i in range(self.num_layers):
185     x = self.encoder_layers[i](x, training, encoder_padding_mask)
186     # x.shape: (batch_size, input_seq_len, d_model)
187     return x
188

```

```

189 "4.5 DecoderModel"
190 class DecoderModel(keras.layers.Layer):
191     def __init__(self, num_layers, target_vocab_size, max_length, d_model, num_heads, dff, rate=0.1):
192         super(DecoderModel, self).__init__()
193         self.num_layers = num_layers
194         self.max_length = max_length
195         self.d_model = d_model
196         self.embedding = keras.layers.Embedding(target_vocab_size, d_model)
197         self.position_embedding = get_position_embedding(max_length, d_model)
198
199         self.dropout = keras.layers.Dropout(rate)
200         self.decoder_layers = [DecoderLayer(d_model, num_heads, dff, rate) for _ in range(self.num_layers)]
201
202     def call(self, x, encoding_outputs, training, decoder_mask, encoder_decoder_padding_mask):
203         # x.shape: (batch_size, output_seq_len)
204         output_seq_len = tf.shape(x)[1]
205         tf.debugging.assert_less_equal(output_seq_len, self.max_length, "output_seq_len should be less or equal to self.max_length")
206
207         attention_weights = {}
208         # x.shape: (batch_size, output_seq_len, d_model)
209         x = self.embedding(x)
210         x *= tf.math.sqrt(tf.cast(self.d_model, tf.float32))
211         x += self.position_embedding[:, :output_seq_len, :]
212
213         x = self.dropout(x, training = training)
214
215         for i in range(self.num_layers):
216             x, attn1, attn2 = self.decoder_layers[i](x, encoding_outputs, training, decoder_mask, encoder_decoder_padding_mask)
217             attention_weights['decoder_layer{}_attn1'.format(i+1)] = attn1
218             attention_weights['decoder_layer{}_attn2'.format(i+1)] = attn2

```

```

219 # x.shape: (batch_size, output_seq_len, d_model)
220 return x, attention_weights
221
222 "4.6 Transformer"
223 class Transformer(keras.Model):
224     def __init__(self, num_layers, input_vocab_size, target_vocab_size, max_length, d_model, num_heads, dff, rate=0.1):
225         super(Transformer, self).__init__()
226
227         self.encoder_model = EncoderModel(num_layers,
            input_vocab_size, max_length, d_model, num_heads, dff, rate)
228
229         self.decoder_model = DecoderModel(num_layers, target_vocab_size, max_length, d_model, num_heads, dff, rate)
230
231         self.final_layer = keras.layers.Dense(target_vocab_size)
232
233     def call(self, inp, tar, training, encoder_padding_mask, decoder_mask, encoder_decoder_padding_mask):
234         # encoding_outputs.shape: (batch_size, input_seq_len, d_model)
235         encoding_outputs = self.encoder_model(inp, training, encoder_padding_mask)
236
237         # decoding_outputs.shape: (batch_size, output_seq_len, d_model)
238         decoding_outputs, attention_weights = self.decoder_model(tar, encoding_outputs, training, decoder_mask, encoder_decoder_padding_mask)
239
240         # predictions.shape: (batch_size, output_seq_len, target_vocab_size)
241         predictions = self.final_layer(decoding_outputs)
242
243         return predictions, attention_weights
244

```

```

1 "5. optimizer & loss"
2 num_layers = 4
3 d_model = 128

```

```

4 dff = 512
5 num_heads = 8
6
7 input_vocab_size = pt_tokenizer.vocab_size + 2
8 target_vocab_size = en_tokenizer.vocab_size + 2
9
10 dropout_rate = 0.1
11
12 transformer = Transformer(num_layers,input_vocab_size,target_vocab_size,max_length,d_model, num_heads, dff, dropout_rate)
13
14
15 learning_rate = 0.001
16 optimizer = keras.optimizers.Adam(learning_rate,beta_1 = 0.9,beta_2 = 0.98,epsilon = 1e-9)
17
18 loss_object = keras.losses.SparseCategoricalCrossentropy(from_logits = True, reduction = 'none')
19
20 def loss_function(real, pred):
21     mask = tf.math.logical_not(tf.math.equal(real, 0))
22     loss_ = loss_object(real, pred)
23
24     mask = tf.cast(mask, dtype=loss_.dtype)
25     loss_ *= mask
26
27     return tf.reduce_mean(loss_)
28
29 def create_masks(inp, tar):
30     """
31     Encoder:
32     - encoder_padding_mask (self attention of EncoderLayer)
33     Decoder:
34     - look_ahead_mask (self attention of DecoderLayer)
35     - encoder_decoder_padding_mask (encoder-decoder attention of DecoderLayer)
36     - decoder_padding_mask (self attention of DecoderLayer)

```

```

37     """
38     encoder_padding_mask = create_padding_mask(inp)
39     encoder_decoder_padding_mask = create_padding_mask(inp)
40
41     look_ahead_mask = create_look_ahead_mask(tf.shape(tar)[1])
42     decoder_padding_mask = create_padding_mask(tar)
43     decoder_mask =
44         tf.maximum(decoder_padding_mask, look_ahead_mask)
45     return encoder_padding_mask, decoder_mask, encoder_decoder_padi
46         ding_mask

```

```

1  "6. train step -> train"
2  train_loss = keras.metrics.Mean(name = 'train_loss')
3  train_accuracy = keras.metrics.SparseCategoricalAccuracy(name =
4      'train_accuracy')
5
6  @tf.function
7  def train_step(inp, tar):
8      tar_inp = tar[:, :-1]
9      tar_real = tar[:, 1:]
10
11     encoder_padding_mask, decoder_mask, encoder_decoder_padding_ma
12     sk = create_masks(inp, tar_inp)
13
14     with tf.GradientTape() as tape:
15         predictions, _ = transformer(inp, tar_inp, True, encoder_paddin
16             g_mask, decoder_mask, encoder_decoder_padding_mask)
17         loss = loss_function(tar_real, predictions)
18
19         gradients = tape.gradient(loss, transformer.trainable_variable
20             s)
21         optimizer.apply_gradients(zip(gradients, transformer.trainable
22             _variables))
23         train_loss(loss)
24         train_accuracy(tar_real, predictions)
25

```



```

21 epochs = 2 # 0
22 for epoch in range(epochs):
23     start = time.time()
24     train_loss.reset_states()
25     train_accuracy.reset_states()
26
27     for (batch, (inp, tar)) in enumerate(train_dataset):
28         train_step(inp, tar)
29         if batch % 100 == 0:
30             print('Epoch {} Batch {} Loss {:.4f} Accuracy {:.4f}'.format(e
epoch + 1, batch, train_loss.result(), train_accuracy.result()))
31
32             print('Epoch {} Loss {:.4f} Accuracy {:.4f}'.format(epoch + 1,
train_loss.result(), train_accuracy.result()))
33             print('Time take for 1 epoch: {} secs\n'.format(time.time() -
start))

```

```

1 "7. Evaluate and Visualize"
2 """
3 eg: A B C D -> E F G H.
4 Train: A B C D, E F G -> F G H
5 Eval: A B C D -> E
6     A B C D, E -> F
7     A B C D, E F -> G
8     A B C D, E F G -> H
9 """
10 def evaluate(inp_sentence):
11     input_id_sentence = [pt_tokenizer.vocab_size] + pt_tokenizer.e
ncode(inp_sentence) + [pt_tokenizer.vocab_size + 1]
12     # encoder_input.shape: (1, input_sentence_length)
13     encoder_input = tf.expand_dims(input_id_sentence, 0)
14
15     # decoder_input.shape: (1, 1)
16     decoder_input = tf.expand_dims([en_tokenizer.vocab_size], 0)
17
18     for i in range(max_length):

```

```

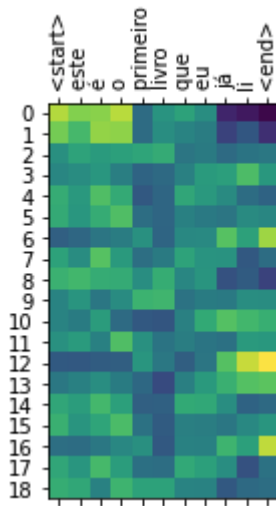
19 encoder_padding_mask, decoder_mask, encoder_decoder_padding_ma
sk = create_masks(encoder_input, decoder_input)
20 # predictions.shape: (batch_size, output_target_len, target_vo
cab_size)
21 predictions, attention_weights = transformer(encoder_input, dec
oder_input, False, encoder_padding_mask, decoder_mask, encoder_decoder
_padding_mask)
22 # predictions.shape: (batch_size, target_vocab_size)
23 predictions = predictions[:, -1, :]
24
25 predicted_id = tf.cast(tf.argmax(predictions, axis = -1), tf.in
t32)
26
27 if tf.equal(predicted_id, en_tokenizer.vocab_size + 1):
28     return tf.squeeze(decoder_input, axis = 0), attention_weights
29
30 decoder_input = tf.concat([decoder_input, [predicted_id]], axis
= -1)
31 return tf.squeeze(decoder_input, axis = 0), attention_weights
32
33 def plot_encoder_decoder_attention(attention, input_sentence, re
sult, layer_name):
34     fig = plt.figure(figsize = (16, 8))
35
36     input_id_sentence = pt_tokenizer.encode(input_sentence)
37
38     # attention.shape: (num_heads, tar_len, input_len)
39     attention = tf.squeeze(attention[layer_name], axis = 0)
40
41     for head in range(attention.shape[0]):
42         ax = fig.add_subplot(2, 4, head + 1)
43
44         # :-1: <end>的id没有加入attention矩阵
45         ax.matshow(attention[head][: -1, :])
46
47         fontdict = {'fontsize': 10}
48
49         ax.set_xticks(range(len(input_id_sentence) + 2))
50         ax.set_yticks(range(len(result)))

```

```

51
52 ax.set_ylim(len(result) - 1.5, -0.5)
53
54 ax.set_xticklabels(['<start>'] + [pt_tokenizer.decode([i]) for
i in input_id_sentence] + ['<end>'], fontdict = fontdict, rotation
= 90)
55 ax.set_yticklabels([en_tokenizer.decode([i]) for i in result if
i < en_tokenizer.vocab_size], fontdict = fontdict)
56 ax.set_xlabel('Head {}'.format(head + 1))
57 plt.tight_layout()
58 plt.show()
59
60 def translate(input_sentence, layer_name = ''):
61     result, attention_weights = evaluate(input_sentence)
62
63     predicted_sentence = en_tokenizer.decode([i for i in result if
i < en_tokenizer.vocab_size])
64
65     print("Input: {}".format(input_sentence))
66     print("Predicted translation: {}".format(predicted_sentence))
67
68     if layer_name:
69         plot_encoder_decoder_attention(attention_weights, input_senten
ce, result, layer_name)
70
71 translate('este é o primeiro livro que eu já li', layer_name =
'decoder_layer4_att2')
72 # Input: este é o primeiro livro que eu já li
73 # Predicted translation: it 's a lot of the world 's going to b
e a lot of the world .

```

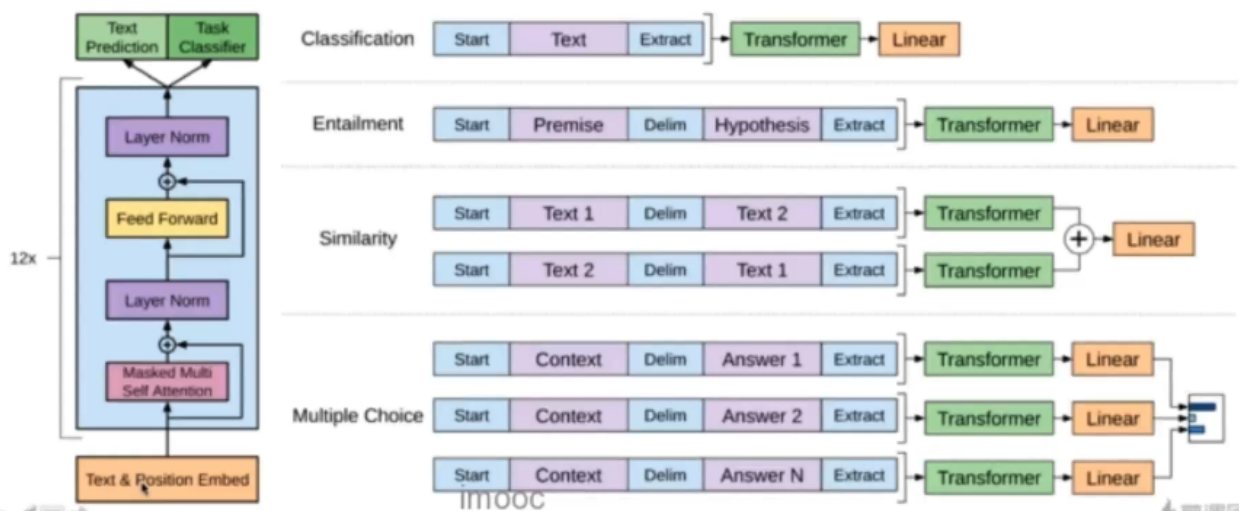


## Transformer模型

### ◆ GPT模型

- ◆ Generative Pre-training
- ◆ Transformer的Decoder部分
- ◆ 预训练+fine-tune到具体任务

## Generative Pre-Training模型



# Transformer模型

## ◆ Bert模型

- ◆ Bidirectional Encoder Representations from Transformers
- ◆ 预训练+fine-tune到特殊任务
- ◆ 双向网络
- ◆ 随机mask
- ◆ 预测下一个句子

## Bert模型

