

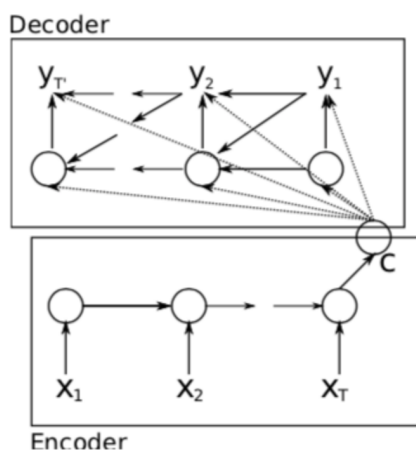
Table of Contents

- ▼ [1 序列到序列任务中的编码器-解码器架构 \(Seq2Seq with RNN Encoder-Decoder\)](#)
 - [1.1 RNN Encoder-Decoder网络架构](#)
 - [1.2 RNN Encoder-Decoder编码器Encoder原理](#)
 - [1.3 RNN Encoder-Decoder解码器Decoder原理](#)
 - [1.4 RNN Encoder-Decoder原理实现](#)
- ▼ [2 序列到序列任务中的注意力机制 \(Seq2Seq with Attention\)](#)
 - [2.1 Seq2Seq with Attention网络架构](#)
 - [2.2 Seq2Seq with Attention编码器Encoder原理](#)
 - [2.3 Seq2Seq with Attention解码器Decoder原理](#)
 - [2.4 Seq2Seq with Attention原理实现](#)
- ▼ [3 注意力机制Attention](#)
 - ▼ [3.1 柔性注意力机制 \(Soft Attention\)](#)
 - [3.1.1 注意力分布](#)
 - [3.1.2 加权平均](#)
 - [3.2 键值对注意力机制 \(Key-Value Pair Attention Mechanism\)](#)
 - [3.3 多头注意力机制 \(Multi-Head Attention Mechanism\)](#)
 - [3.4 自注意力模型 \(Self-Attention Model\)](#)

▼ 序列到序列任务中的注意力机制 (Seq2Seq with Attention)

▼ 1 序列到序列任务中的编码器-解码器架构 (Seq2Seq with RNN Encoder-Decoder)

▼ 1.1 RNN Encoder-Decoder网络架构



RNN Encoder-Decoder神经网络架构使用循环神经网络学习，将变长源序列 X 编码成定长向量表示 c ，并将学习的定长向量表示 c 解码成变长目标序列 Y 。模型的编码器和解码器被联合训练，以最大化给定源序列的目标序列的条件概率。

源文本序列: $X = (x_1, x_2, \dots, x_N)$ 其中, $x_i = (l_1, l_2, \dots, l_j, \dots, l_K)$, 其中 $l_j = I(i = j)$, $(j = 1, \dots, K)$ 。

目标文本序列: $Y = (y_1, y_2, \dots, y_M)$ 其中, $y_i = (l_1, l_2, \dots, l_j, \dots, l_K)$, 其中 $l_j = (i = j)$, $(j = 1, \dots, K)$

▼ 1.2 RNN Encoder-Decoder编码器Encoder原理

源文本单词的词嵌入表示： $e(\mathbf{x}_i) \in \mathbb{R}^{500}$

编码器的隐藏状态由1000个隐藏单元组成。

编码器隐藏状态初始化，在 $t = 0$ 时刻第 j 个隐藏单元

$$h_j^{\langle 0 \rangle} = 0$$

在 t 时刻第 j 个隐藏单元

$$h_j^{\langle t \rangle} = z_j h_j^{\langle t-1 \rangle} + (1 - z_j) \tilde{h}_j^{\langle t \rangle}$$

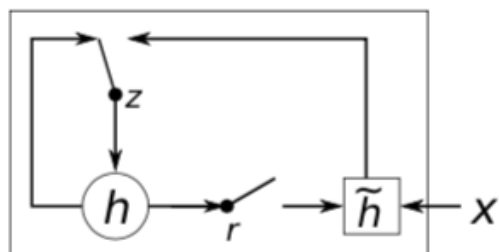
其中，

$$\tilde{h}_j^{\langle t \rangle} = \tanh\left(\left[\mathbf{W}_e(\mathbf{x}_t)\right]_j + \left[\mathbf{U}(\mathbf{r} \odot \mathbf{h}^{\langle t-1 \rangle})\right]_j\right)$$

$$z_j = \sigma\left(\left[\mathbf{W}_z e(\mathbf{x}_t)\right]_j + \left[\mathbf{U}_z \mathbf{h}^{\langle t-1 \rangle}\right]_j\right)$$

$$r_j = \sigma\left(\left[\mathbf{W}_r e(\mathbf{x}_t)\right]_j + \left[\mathbf{U}_r \mathbf{h}^{\langle t-1 \rangle}\right]_j\right)$$

$\sigma(\cdot)$ 为sigmoid函数， \odot 为向量元素乘法， $\mathbf{W}, \mathbf{W}_z, \mathbf{W}_r \in \mathbb{R}^{1000 \times 500}$ 和 $\mathbf{U}, \mathbf{U}_z, \mathbf{U}_r \in \mathbb{R}^{1000 \times 1000}$ 为权值矩阵。为了使方程齐整，省略了偏置项。



源文本最后第 N 时刻，编码器的隐藏状态计算完成，源文本的定长向量表示

$$\mathbf{c} = \tanh(\mathbf{V} \mathbf{h}^{\langle N \rangle})$$

其中， $\mathbf{V} \in \mathbb{R}^{1000 \times 1000}$ 为权值矩阵。

▼ 1.3 RNN Encoder-Decoder 解码器 Decoder 原理

解码器隐藏状态初始化，在 $t = 0$ 时刻

$$\mathbf{h}'^{\langle 0 \rangle} = \tanh(\mathbf{V}' \mathbf{c})$$

其中， $\mathbf{V}' \in \mathbb{R}^{1000 \times 1000}$ 为权值矩阵。

在 t 时刻第 j 个隐藏单元

$$h_j'^{\langle t \rangle} = z_j' h_j'^{\langle t-1 \rangle} + (1 - z_j') \tilde{h}_j'^{\langle t \rangle}$$

其中，

$$\begin{aligned}\tilde{h}_j^{(t)} &= \tanh\left(\left[\mathbf{W}'_e(\mathbf{y}_{t-1})\right]_j + r'_j\left[\mathbf{U}'_h\mathbf{h}^{(t-1)} + \mathbf{C}\mathbf{c}\right]\right) \\ z'_j &= \sigma\left(\left[\mathbf{W}'_z e(\mathbf{y}_{t-1})\right]_j + \left[\mathbf{U}'_z\mathbf{h}^{(t-1)}\right]_j + \left[\mathbf{C}_z\mathbf{c}\right]_j\right) \\ r'_j &= \sigma\left(\left[\mathbf{W}'_r e(\mathbf{y}_{t-1})\right]_j + \left[\mathbf{U}'_r\mathbf{h}^{(t-1)}\right]_j + \left[\mathbf{C}_r\mathbf{c}\right]_j\right)\end{aligned}$$

其中, $\mathbf{W}', \mathbf{W}'_z, \mathbf{W}'_r \in \mathbb{R}^{1000 \times 500}$ 和 $\mathbf{U}', \mathbf{U}'_z, \mathbf{U}'_r \in \mathbb{R}^{1000 \times 1000}$ 以及 $\mathbf{C}', \mathbf{C}'_z, \mathbf{C}'_r \in \mathbb{R}^{1000 \times 1000}$ 为权值矩阵。

目标文本单词的词嵌入表示: $e(\mathbf{y}_i) \in \mathbb{R}^{500}$, 且在 $t=0$ 时刻 $e(\mathbf{y}_0) = \mathbf{0}$ 。

在每个时刻 t , 解码器计算生成第 j 个单词的概率

$$p(y_{t,j} = 1 | \mathbf{y}_{t-1}, \dots, \mathbf{y}_1, X) = \frac{\exp(\mathbf{g}_j \mathbf{s}^{(t)})}{\sum_{j=1}^K \exp(\mathbf{g}_j \mathbf{s}^{(t)})}$$

其中, 最大输出单元 (maxout unit)

$$s_i^{(t)} = \max \left\{ s'_{2i-1}{}^{(t)}, s'_{2i}{}^{(t)} \right\}$$

且

$$\mathbf{s}'^{(t)} = \mathbf{O}_h \mathbf{h}^{(t)} + \mathbf{O}_y \mathbf{y}_{t-1} + \mathbf{O}_c \mathbf{c}$$

$\mathbf{O}_h, \mathbf{O}_c \in \mathbb{R}^{500 \times 1000}$ 和 $\mathbf{O}_y \in \mathbb{R}^{500 \times 500}$ 以及 $\mathbf{G} = [\mathbf{g}_1, \dots, \mathbf{g}_K] \in \mathbb{R}^{K \times 1000}$ 为权值矩阵。

1.4 RNN Encoder-Decoder原理实现

```
In [1]: 1 import numpy as np
        2 import tensorflow as tf
        3
        4 tf.__version__
```

executed in 6.54s, finished 14:57:18 2019-09-10

```
Out[1]: '1.12.0'
```

In [18]:

```
1  tf.reset_default_graph()
2
3  char_arr = [c for c in 'SEPabcdefghijklmnopqrstuvwxyz']
4  num_dic = {n: i for i, n in enumerate(char_arr)}
5
6  seq_data = [['man', 'women'], ['black', 'white'], ['king', 'queen'], [
7              'girl', 'boy'], ['up', 'down'], ['high', 'low']]
8
9  n_step = 5
10 n_hidden = 128
11 n_class = len(num_dic)
12
13
14 def make_batch(seq_data):
15     input_batch, output_batch, target_batch = [], [], []
16
17     for seq in seq_data:
18         for i in range(2):
19             seq[i] = seq[i] + 'P' * (n_step - len(seq[i]))
20
21             input = [num_dic[n] for n in seq[0]]
22             output = [num_dic[n] for n in ('S' + seq[1])]
23             target = [num_dic[n] for n in (seq[1] + 'E')]
24
25             input_batch.append(np.eye(n_class)[input])
26             output_batch.append(np.eye(n_class)[output])
27
28             target_batch.append(target)
29
30     return input_batch, output_batch, target_batch
31
32
33 enc_input = tf.placeholder(tf.float32, [None, None, n_class])
34 dec_input = tf.placeholder(tf.float32, [None, None, n_class])
35 targets = tf.placeholder(tf.int64, [None, None])
36
37 with tf.variable_scope('encoder'):
38     enc_cell = tf.nn.rnn_cell.BasicRNNCell(n_hidden)
39     enc_cell = tf.nn.rnn_cell.DropoutWrapper(enc_cell, output_keep_prob=0.5)
40     _, enc_states = tf.nn.dynamic_rnn(enc_cell, enc_input, dtype=tf.float32)
41
42 with tf.variable_scope('decoder'):
43     dec_cell = tf.nn.rnn_cell.BasicRNNCell(n_hidden)
44     dec_cell = tf.nn.rnn_cell.DropoutWrapper(dec_cell, output_keep_prob=0.5)
45     outputs, _ = tf.nn.dynamic_rnn(
46         dec_cell, dec_input, initial_state=enc_states, dtype=tf.float32)
47
48 model = tf.layers.dense(outputs, n_class, activation=None)
49
50 cost = tf.reduce_mean(tf.nn.sparse_softmax_cross_entropy_with_logits(
51     logits=model, labels=targets))
52 optimizer = tf.train.AdamOptimizer(0.001).minimize(cost)
53
54 sess = tf.Session()
55 sess.run(tf.global_variables_initializer())
56 input_batch, output_batch, target_batch = make_batch(seq_data)
57
58 for epoch in range(5000):
59     _, loss = sess.run([optimizer, cost], feed_dict={
60         enc_input: input_batch, dec_input: output_batch, targets: ta
61     })
62     if (epoch + 1) % 1000 == 0:
63         print('Epoch:', '%04d' % (epoch + 1), 'cost = ', '{:.6f}'.format(loss))
64
65 def translate(word):
66     seq_data = [word, 'P' * len(word)]
67
68     input_batch, output_batch, _ = make_batch([seq_data])
69     prediction = tf.argmax(model, 2)
70
71     result = sess.run(prediction, feed_dict={
```

```

72         enc_input: input_batch, dec_input: output_batch})
73
74     decoded = [char_arr[i] for i in result[0]]
75     end = decoded.index('E')
76     translated = ''.join(decoded[:end])
77
78     return translated.replace('P', ' ')
79
80
81 print('test')
82 print('man ->', translate('man'))
83 print('mans ->', translate('mans'))
84 print('king ->', translate('king'))
85 print('black ->', translate('black'))
86 print('up ->', translate('up'))

```

executed in 1m 17.2s, finished 16:15:30 2019-09-10

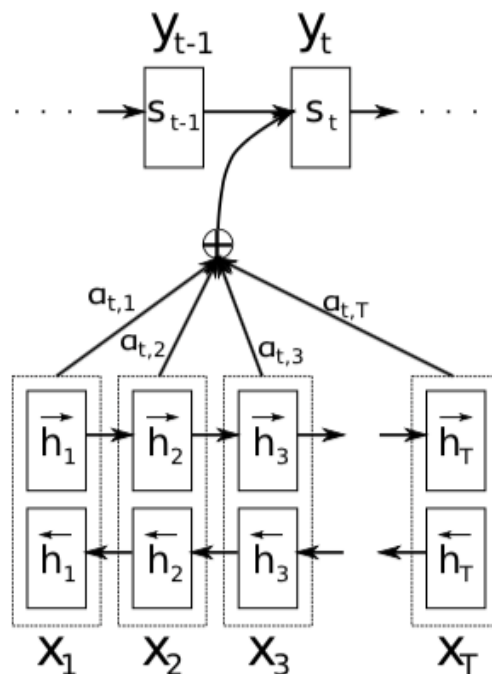
```

Epoch: 1000 cost = 0.000840
Epoch: 2000 cost = 0.000314
Epoch: 3000 cost = 0.000181
Epoch: 4000 cost = 0.000044
Epoch: 5000 cost = 0.000027
test
man -> women
mans -> women
king -> queen
black -> white
up -> down

```

2 序列到序列任务中的注意力机制 (Seq2Seq with Attention)

2.1 Seq2Seq with Attention网络架构



seq2seq with Attention神经网络架构中，编码器采用双向循环神经网络学习将输入序列 \mathbf{x} 编码成每个时刻的上下文向量（注意力分布） c_i ，解码器学习将上下文向量 c_i 解码为输出序列 \mathbf{y} 。

源文本序列： $\mathbf{x} = (x_1, \dots, x_{T_x})$ ，其中 $x_i \in \mathbb{R}^{K_x}$ 为one-of-K编码， K_x 为源语言词表长度， T_x 为源语料长度。

目标文本序列： $\mathbf{y} = (y_1, \dots, y_{T_y})$ ，其中 $y_i \in \mathbb{R}^{K_y}$ 为one-of-K编码， K_y 为目标语言词表长度， T_y 为目标语料长度。

2.2 Seq2Seq with Attention编码器Encoder原理

编码器Encoder采用双向循环神经网络，前向状态计算

$$\vec{h}_i = \begin{cases} (1 - z_i) \circ \vec{h}_{i-1} + z_i \circ \vec{h}_{i-1} & , \text{if } i > 0 \\ 0 & , \text{if } i = 0 \end{cases}$$

其中,

$$\vec{h}_i = \tanh\left(\vec{W}Ex_i + \vec{U}\left[\vec{r}_i \circ \vec{h}_{i-1}\right]\right) \vec{z}_i = \sigma\left(\vec{W}_zEx_i + \vec{U}_z\vec{h}_{i-1}\right) \vec{r}_i = \sigma\left(\vec{W}_rEx_i + \vec{U}_r\vec{h}_{i-1}\right)$$

$E \in \mathbb{R}^{m \times K_x}$ 为词嵌入矩阵, m 为词嵌入维度。 $\vec{W}, \vec{W}_z, \vec{W}_r \in \mathbb{R}^{n \times m}$ 和 $\vec{U}, \vec{U}_z, \vec{U}_r \in \mathbb{R}^{n \times n}$ 为权值矩阵, n 为隐藏单元数。 $\sigma(\cdot)$ 通常为sigmoid函数。

后向状态 $\left(\overleftarrow{h}_1, \dots, \overleftarrow{h}_{T_x}\right)$ 计算相同。与权值矩阵不同, 我们在前向和后向RNN网络中共享词嵌入矩阵 E 。

将前向和后向状态关联起来得到注释 $(h_1, h_2, \dots, h_{T_x})$,

其中,

$$h_i = \begin{bmatrix} \vec{h}_i \\ \overleftarrow{h}_i \end{bmatrix}$$

2.3 Seq2Seq with Attention解码器Decoder原理

解码器Decoder隐层转态 s_i 由解码器注释 h_i 计算的注意力分布 c_i 得到

$$s_i = (1 - z_i) \circ s_{i-1} + z_i \circ \tilde{s}_i$$

其中,

$$\tilde{s}_i = \tanh\left(WEy_{i-1} + U\left[r_i \circ s_{i-1}\right] + Cc_i\right) z_i = \sigma\left(W_zEy_{i-1} + U_zs_{i-1} + C_r c_i\right) r_i = \sigma\left(W_rEy_{i-1} + U_rs_{i-1} + C_r c_i\right)$$

$E \in \mathbb{R}^{m \times K_y}$ 为目标语言的词嵌入矩阵, m 为词嵌入维度。 $W, W_z, W_r \in \mathbb{R}^{n \times m}$ 和 $U, U_z, U_r \in \mathbb{R}^{n \times n}$ 以及 $C, C_z, C_r \in \mathbb{R}^{n \times 2n}$

为权值矩阵, n 为隐藏单元数。隐层初始状态 $s_0 = \tanh\left(W_s \overleftarrow{h}_1\right)$, 其中 $W_s \in \mathbb{R}^{n \times n}$ 。

每个时刻的上下文向量（注意力分布） c_i 的计算

$$c_i = \sum_{j=1}^{T_x} a_{ij} h_j$$

其中,

$$a_{ij} = \frac{\exp(e_{ij})}{\sum_{k=1}^{T_x} \exp(e_{ik})} e_{ij} = v_a^T \tanh\left(W_a s_{i-1} + U_a h_j\right)$$

h_j 为源文本序列第 j 个注释。 $v_a \in \mathbb{R}^n$, $W_a \in \mathbb{R}^{n' \times n}$, $U_a \in \mathbb{R}^{n' \times 2n}$ 为权值矩阵。

使用解码器状态 s_{i-1} , 上下文 c_i 和上时刻生成单词 y_{i-1} 定义目标单词 y_i 的概率

$$p(y_i | s_i, y_{i-1}, c_i) \propto \exp(y_i^\top W_o t_i)$$

其中,

$$t_i = \left[\max \left\{ \tilde{t}_{i, 2j-1}, \tilde{t}_{i, wj} \right\} \right]_{j=1, \dots, l}^\top$$

$\tilde{t}_{i,k}$ 是向量 \tilde{t}_i 的第 k 个元素,

$$\tilde{t}_i = U_o s_{i-1} + V_o E y_{i-1} + C_o c_i$$

$W_o \in \mathbb{R}^{K_y \times l}$, $U_o \in \mathbb{R}^{2l \times n}$, $C_o \in \mathbb{R}^{2l \times 2n}$ 是权值矩阵。

▼ 2.4 Seq2Seq with Attention原理实现

```
In [4]: 1 %matplotlib inline
2
3 import numpy as np
4 import tensorflow as tf
5 import matplotlib.pyplot as plt
```

executed in 343ms, finished 15:20:23 2019-09-10

```
In [5]: 1 tf.reset_default_graph()
```

executed in 6ms, finished 15:20:39 2019-09-10

```
In [6]: 1 sentences = ['ich mochte ein bier P', 'S i want a beer', 'i want a beer E']
2
3 word_list = " ".join(sentences).split()
4 word_list = list(set(word_list))
5 word_dict = {w: i for i, w in enumerate(word_list)}
6 number_dict = {i: w for i, w in enumerate(word_list)}
7 n_class = len(word_dict)
```

executed in 12ms, finished 15:20:48 2019-09-10

```
In [7]: 1 n_step = 5
2 n_hidden = 128
```

executed in 9ms, finished 15:20:59 2019-09-10

```
In [8]: 1 def make_batch(sentences):
2     input_batch = [np.eye(n_class)[word_dict[n] for n in sentences[0].split()]]
3     output_batch = [np.eye(n_class)[word_dict[n] for n in sentences[1].split()]]
4     target_batch = [word_dict[n] for n in sentences[2].split()]
5
6     return input_batch, output_batch, target_batch
```

executed in 8ms, finished 15:21:10 2019-09-10

```
In [9]: 1 enc_inputs = tf.placeholder(tf.float32, [None, None, n_class])
2 dec_inputs = tf.placeholder(tf.float32, [None, None, n_class])
3 targets = tf.placeholder(tf.int64, [1, n_step])
```

executed in 11ms, finished 15:21:19 2019-09-10

```
In [10]: 1 attn = tf.Variable(tf.random_normal([n_hidden, n_hidden]))
2 out = tf.Variable(tf.random_normal([n_hidden * 2, n_class]))
```

executed in 26ms, finished 15:21:29 2019-09-10

```
In [11]: 1 def get_att_score(dec_output, enc_output):
2         score = tf.squeeze(tf.matmul(enc_output, attn), 0)
3         dec_output = tf.squeeze(dec_output, [0, 1])
4
5         return tf.tensordot(dec_output, score, 1)
6
7 def get_att_weight(dec_output, enc_outputs):
8     attn_scores = []
9     enc_outputs = tf.transpose(enc_outputs, [1, 0, 2])
10    for i in range(n_step):
11        attn_scores.append(get_att_score(dec_output, enc_outputs[i]))
12
13    return tf.reshape(tf.nn.softmax(attn_scores), [1, 1, -1])
```

executed in 12ms, finished 15:21:41 2019-09-10

```
In [12]: 1 model = []
2         # Attention = []
```

executed in 8ms, finished 15:21:51 2019-09-10

```
In [13]: 1 with tf.variable_scope('encode'):
2         enc_cell = tf.nn.rnn_cell.BasicRNNCell(n_hidden)
3         enc_cell = tf.nn.rnn_cell.DropoutWrapper(enc_cell, output_keep_prob=0.5)
4         enc_outputs, enc_hidden = tf.nn.dynamic_rnn(enc_cell, enc_inputs, dtype=tf.float32)
```

executed in 144ms, finished 15:22:03 2019-09-10

```
In [14]: 1 with tf.variable_scope('decode'):
2         dec_cell = tf.nn.rnn_cell.BasicRNNCell(n_hidden)
3         dec_cell = tf.nn.rnn_cell.DropoutWrapper(dec_cell, output_keep_prob=0.5)
4
5         inputs = tf.transpose(dec_inputs, [1, 0, 2])
6         hidden = enc_hidden
7         for i in range(n_step):
8             dec_output, hidden = tf.nn.dynamic_rnn(dec_cell,
9                                                     tf.expand_dims(inputs[i], 1),
10                                                     initial_state=hidden,
11                                                     dtype=tf.float32,
12                                                     time_major=True)
13             attn_weights = get_att_weight(dec_output, enc_outputs)
14             # Attention.append(tf.squeeze(attn_weights))
15
16             context = tf.matmul(attn_weights, enc_outputs)
17             dec_output = tf.squeeze(dec_output, 0)
18             context = tf.squeeze(context, 1)
19
20             model.append(tf.matmul(tf.concat((dec_output, context), 1), out))
```

executed in 664ms, finished 15:22:18 2019-09-10

```
In [15]: 1 # trained_attn = tf.stack([Attention[0], Attention[1], Attention[2], Attention[3],
2         model = tf.transpose(model, [1, 0, 2])
3         prediction = tf.argmax(model, 2)
4         cost = tf.reduce_mean(tf.nn.sparse_softmax_cross_entropy_with_logits(logits=model,
5         optimizer = tf.train.AdamOptimizer(0.001).minimize(cost)
```

executed in 1.34s, finished 15:22:33 2019-09-10


```
In [16]: 1 with tf.Session() as sess:
2         init = tf.global_variables_initializer()
3         sess.run(init)
4         for epoch in range(2000):
5             input_batch, output_batch, target_batch = make_batch(sentences)
6             _, loss, attention = sess.run([optimizer, cost, trained_attn],
7                                           feed_dict={enc_inputs: input_batch,
8                                                     dec_inputs: output_batch,
9                                                     targets: target_batch})
10
11             if (epoch + 1) % 400 == 0:
12                 print('Epoch:', '%04d' % (epoch + 1), 'cost=', '{:.6f}'.format(loss))
13
14             predict_batch = [np.eye(n_class)[word_dict[n] for n in 'P P P P P'.split()]]
15             result = sess.run(prediction, feed_dict={enc_inputs: input_batch,
16                                                     dec_inputs: predict_batch})
17             print(sentences[0].split(), '->', [number_dict[n] for n in result[0]])
```

executed in 20.8s, finished 15:23:05 2019-09-10

```
Epoch: 0400 cost= 0.000000
Epoch: 0800 cost= 0.000008
Epoch: 1200 cost= 0.000002
Epoch: 1600 cost= 0.002508
Epoch: 2000 cost= 0.000000
['ich', 'mochte', 'ein', 'bier', 'P'] -> ['i', 'want', 'a', 'beer', 'E']
```

3 注意力机制Attention

3.1 柔性注意力机制 (Soft Attention)

输入信息 $X = [x_1, \dots, x_N]$

注意力机制的计算:

1. 在输入信息上计算注意力分布;
2. 根据注意力分布计算输入信息的加权平均。

3.1.1 注意力分布

给定一个和任务相关的查询向量 \mathbf{q} , 用注意力变量 $z \in [1, N]$ 表示被选择信息的索引位置, 即 $z = i$ 表示选择了第 i 个输入信息。其中, 查询向量 \mathbf{q} 可以是动态生成的, 也可以是可学习的参数。

软性注意力的注意力分布

在给定输入信息 X 和查询变量 \mathbf{q} 下, 选择第 i 个输入信息的概率

$$\begin{aligned} \alpha_i &= p(z = i | X, \mathbf{q}) \\ &= \text{softmax}(s(\mathbf{x}_i, \mathbf{q})) \\ &= \frac{\exp(s(\mathbf{x}_i, \mathbf{q}))}{\sum_{j=1}^N \exp(s(\mathbf{x}_j, \mathbf{q}))} \end{aligned}$$

其中, α_i 称为注意力分布, $s(\mathbf{x}_i, \mathbf{q})$ 称为注意力打分函数。

注意力打分函数

- 加性模型 $s(\mathbf{x}_i, \mathbf{q}) = \mathbf{v}^\top \tanh(W\mathbf{x}_i + U\mathbf{q})$
- 点积模型 $s(\mathbf{x}_i, \mathbf{q}) = \mathbf{x}_i^\top \mathbf{q}$

- 缩放点积模型 $s(\mathbf{x}_i, \mathbf{q}) = \frac{\mathbf{x}_i^\top \mathbf{q}}{\sqrt{d}}$
- 双线性模型 $s(\mathbf{x}_i, \mathbf{q}) = \mathbf{x}_i^\top \mathbf{W} \mathbf{q}$

其中， \mathbf{W} , \mathbf{U} , \mathbf{v} 为可学习的网络参数， d 为输入信息的维度。

加性模型和点积模型的复杂度近似，但点积模型可利用矩阵乘积，计算效率跟高。当输入信息的维度 d 比较高，点积模型值方差较大，导致softmax函数的梯度较小，缩放点积模型可以解决。双线性模型是泛化的点积模型。若假设

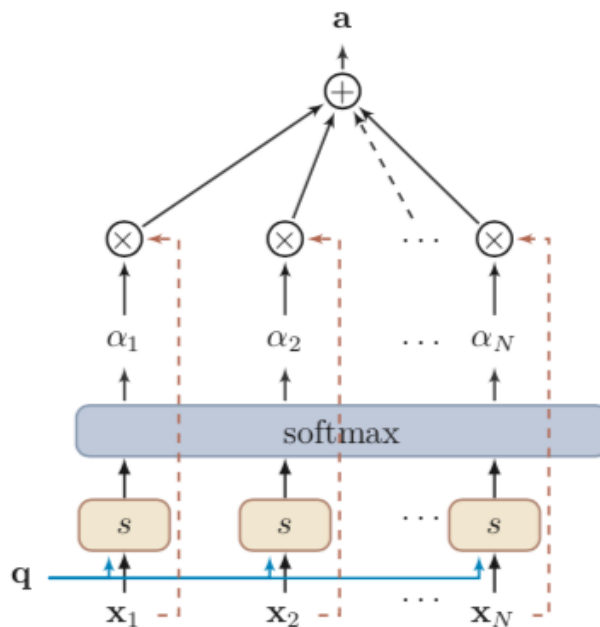
$\mathbf{W} = \mathbf{U}^\top \mathbf{V}$ ，则 $s(\mathbf{x}_i, \mathbf{q}) = \mathbf{x}_i^\top \mathbf{U}^\top \mathbf{V} \mathbf{q} = (\mathbf{U} \mathbf{x}_i)^\top (\mathbf{V} \mathbf{q})$ ，即分别对 \mathbf{x}_i 和 \mathbf{q} 进行线性变换后进行点积。相比点积模型，双线性模型在计算相似度是引入了非对称性。

注意力分布 α_i 可解释为在给定相关查询 \mathbf{q} 时，第 i 个信息受关注的程度。

3.1.2 加权平均

注意力函数

$$\begin{aligned} att(X, \mathbf{q}) &= \sum_{i=1}^N \alpha_i \mathbf{x}_i \\ &= \mathbb{E}_{z \sim p(z|X, \mathbf{q})}[\mathbf{x}] \end{aligned}$$



3.2 键值对注意力机制 (Key-Value Pair Attention Mechanism)

输入信息 $(K, V) = [(\mathbf{k}_1, \mathbf{v}_1), \dots, (\mathbf{k}_N, \mathbf{v}_N)]$ ，其中键用来计算注意力分布 α_i ，值用来计算聚合信息。

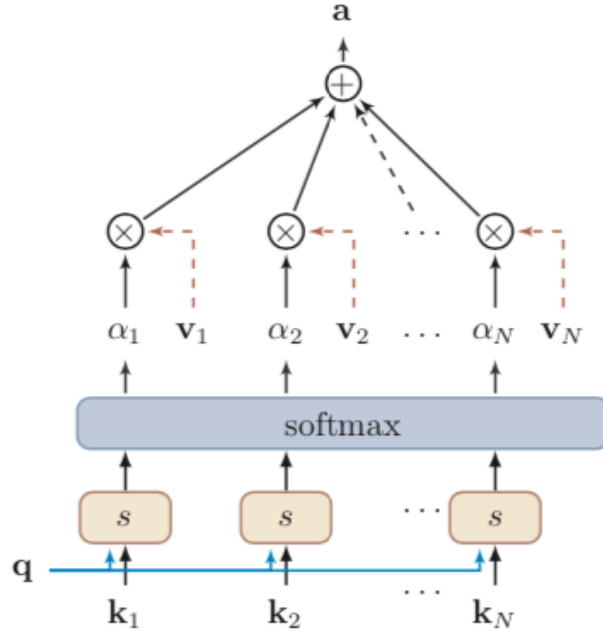
给定任务相关查询向量 \mathbf{q} 时，注意力分布

$$\alpha_i = \frac{\exp(s(\mathbf{k}_i, \mathbf{q}))}{\sum_{j=1}^N \exp(s(\mathbf{k}_j, \mathbf{q}))}$$

注意力函数

$$\begin{aligned} att((K, V), \mathbf{q}) &= \sum_{i=1}^N \alpha_i \mathbf{v}_i \\ &= \sum_{i=1}^N \frac{\exp(s(\mathbf{k}_i, \mathbf{q}))}{\sum_{j=1}^N \exp(s(\mathbf{k}_j, \mathbf{q}))} \mathbf{v}_i \end{aligned}$$

其中， $s(\mathbf{k}_i, \mathbf{q})$ 为打分函数。当 $K = V$ 时，键值对注意力机制等价于柔性注意力机制。



▼ 3.3 多头注意力机制 (Multi-Head Attention Mechanism)

多个查询 $Q = [\mathbf{q}_1, \dots, \mathbf{q}_M]$ 平行的计算从输入信息中选取多个信息。每个注意力关注输入信息的不同部分。

$$att((K, V), Q) = att((K, V), \mathbf{q}_1) \oplus \dots \oplus att((K, V), \mathbf{q}_M)$$

其中， \oplus 为向量拼接。

▼ 3.4 自注意力模型 (Self-Attention Model)

输入序列 $X = [\mathbf{x}_1, \dots, \mathbf{x}_N] \in \mathbb{R}^{d_1 \times N}$

输出序列 $H = [\mathbf{h}_1, \dots, \mathbf{h}_N] \in \mathbb{R}^{d_2 \times N}$

通过线性变换得到向量序列：

$$Q = W_Q X \in \mathbb{R}^{d_3 \times N} \quad K = W_K X \in \mathbb{R}^{d_3 \times N} \quad V = W_V X \in \mathbb{R}^{d_2 \times N}$$

其中， Q 为查询向量序列， K 为键向量序列， V 为值向量序列， W_Q, W_K, W_V 分别为可学习参数矩阵。

预测输出向量

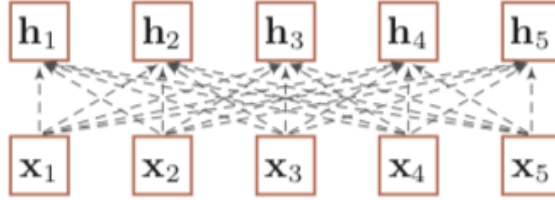
$$\begin{aligned}
\hat{\mathbf{h}}_i &= att(K, V, \mathbf{q}_i) \\
&= \sum_{j=1}^N \alpha_{ij} \mathbf{v}_j \\
&= \sum_{j=1}^N softmax(s(\mathbf{k}_j, \mathbf{q}_i)) \mathbf{v}_j
\end{aligned}$$

其中， $i, j \in [1, N]$ 为输出和输入向量序列的位置，连接权重 α_{ij} 由注意力机制动态生成。

若使用缩放点积模型作为打分函数，则输出向量序列

$$\begin{aligned}
H_{d_2 \times N} &= softmax\left(\frac{K^T Q}{\sqrt{d_3}}\right) V_{d_2 \times N} \\
&= softmax\left(\frac{K^T Q}{\sqrt{d_3}}\right) W_V X
\end{aligned}$$

其中，softmax为按列归一化的函数。



自注意力模型计算的权重 α_{ij} 只依赖 \mathbf{q}_i 和 \mathbf{k}_j 的相关性，而忽略了输入信息的位置信息。因此自注意力模型一般需要加入位置编码信息来进行修正。