

实验2 机器翻译

实验2 机器翻译

1.实验背景及简介

- 1.1实验背景
- 1.2实验简介

2.实验目的

3.实验步骤

- 3.1 数据上传与同步(OBS)
- 3.1 Seq2Seq
- 3.2 Transformer

4.实验原理

- 4.1 Seq2Seq
 - 4.1.1 编码器
 - 4.1.2 译码器
 - 4.1.3 训练模型
- 4.2 Transformer

5 实验结果

- 5.1 Seq2Seq
- 5.2 Transformer

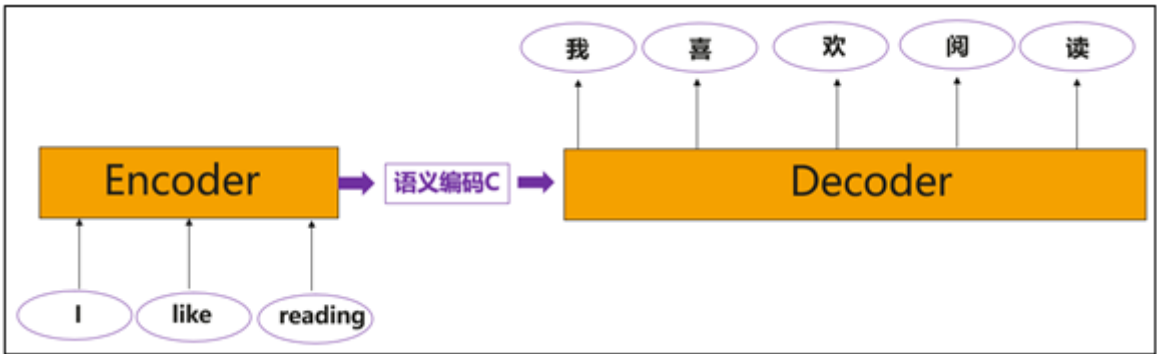
1.实验背景及简介

1.1实验背景

在实际生活中，机器翻译是人工智能技术比较广泛的一个应用。在机器翻译应用中，由于输入序列与输出序列的长度通常不一样。需要采用序列到序列（Sequence to Sequene, Seq2Seq）的映射框架来解决这一问题，它能将一个可变长序列映射到另一个可变长序列。本次实验将探索Seq2Seq基础模型在机器翻译中的应用，以及Attention注意力机制、Transformer模型对基础Seq2Seq模型的改进。

1.2实验简介

翻译任务在日常生活应用广泛，如手机中有各种翻译软件，可以满足人们交流、阅读的需求。本实验基于Seq2Seq编码器-解码器框架，结合GRU单元实现英文转中文的翻译任务，框架示意图如下：



GRU（门递归单元）是一种递归神经网络算法，就像LSTM（长短期存储器）一样。它是由Kyunghyun Cho、Bart van Merriënboer等在2014年的文章“使用RNN编码器-解码器学习短语表示用于统计机器翻译”中提出的。本文提出了一种新的神经网络模型RNN Encoder-Decoder，该模型由两个递归神经网络（RNN）组成，为了提高翻译任务的效果，我们还参考了“神经网络的序列到序列学习”和“联合学习对齐和翻译的神经机器翻译”。

2.实验目的

1. 掌握神经网络在NLP机器翻译领域的应用
2. 理解seq2seq解码器-编码器框架
3. 理解attention注意力机制，transformer架构

3.实验步骤

3.1 数据上传与同步(OBS)



- Seq2Seq

```
In [1]: import moxing as mox
mox.file.copy_parallel(src_url="s3://bucket-mzh/lab2/源代码-Seq2Seq/data/", dst_url='./data/')
path_to_file = "data/cmn.txt"

INFO:root:Using MoXing-v1.17.3-d858ff4a
INFO:root:Using OBS-Python-SDK-3.20.9.1
```

- Transformer

```
In [1]: import moxing as mox
mox.file.copy_parallel(src_url="s3://bucket-mzh/lab2/源代码-transformer/mt_transformer_mindspore_1.1/data/", dst_url='./data/')
mox.file.copy_parallel(src_url="s3://bucket-mzh/lab2/源代码-transformer/mt_transformer_mindspore_1.1/src/", dst_url='./src/')

INFO:root:Using MoXing-v1.17.3-d858ff4a
INFO:root:Using OBS-Python-SDK-3.20.9.1
```

3.1 Seq2Seq

- 数据预处理

```
1 def prepare_data(data_path, vocab_save_path, max_seq_len):
2     with open(data_path, 'r', encoding='utf-8') as f:
3         data = f.read()
4
5     # 读取文本文件，按行分割，再将每行分割成语句对
6     data = data.split('\n')
7
8     # 截取前2000行数据进行训练
9     data = data[:2000]
```

```

10
11 # 分割每行中的中英文
12 en_data = [normalizeString(line.split('\t')[0]) for line in data]
13
14 ch_data = [line.split('\t')[1] for line in data]
15
16 # 利用集合，获得中英文词汇表
17 en_vocab = set(' '.join(en_data).split(' '))
18 id2en = [EOS] + [SOS] + list(en_vocab)
19 en2id = {c:i for i,c in enumerate(id2en)}
20 en_vocab_size = len(id2en)
21 # np.savetxt(os.path.join(vocab_save_path, 'en_vocab.txt'),
np.array(id2en), fmt='%s')
22
23 ch_vocab = set(' '.join(ch_data))
24 id2ch = [EOS] + [SOS] + list(ch_vocab)
25 ch2id = {c:i for i,c in enumerate(id2ch)}
26 ch_vocab_size = len(id2ch)
27 # np.savetxt(os.path.join(vocab_save_path, 'ch_vocab.txt'),
np.array(id2ch), fmt='%s')
28
29 # 将句子用词汇表id表示
30 en_num_data = np.array([[1] + [int(en2id[en]) for en in line.split('
')] + [0] for line in en_data])
31 ch_num_data = np.array([[1] + [int(ch2id[ch]) for ch in line] + [0] for
line in ch_data])
32
33 #将短句子扩充到统一的长度
34 for i in range(len(en_num_data)):
35     num = max_seq_len + 1 - len(en_num_data[i])
36     if(num >= 0):
37         en_num_data[i] += [0]*num
38     else:
39         en_num_data[i] = en_num_data[i][:max_seq_len] + [0]
40
41 for i in range(len(ch_num_data)):
42     num = max_seq_len + 1 - len(ch_num_data[i])
43     if(num >= 0):
44         ch_num_data[i] += [0]*num
45     else:
46         ch_num_data[i] = ch_num_data[i][:max_seq_len] + [0]
47
48
49 np.savetxt(os.path.join(vocab_save_path, 'en_vocab.txt'),
np.array(id2en), fmt='%s')
50
51 np.savetxt(os.path.join(vocab_save_path, 'ch_vocab.txt'),
np.array(id2ch), fmt='%s')
52
53 return en_num_data, ch_num_data, en_vocab_size, ch_vocab_size

```

- 设置超参数

```

1 from easydict import EasyDict as edict
2
3 # CONFIG
4 cfg = edict({

```

```

5     'en_vocab_size': 1154,
6     'ch_vocab_size': 1116,
7     'max_seq_length': 10,
8     'hidden_size': 1024,
9     'batch_size': 16,
10    'eval_batch_size': 1,
11    'learning_rate': 0.001,
12    'momentum': 0.9,
13    'num_epochs': 15,
14    'save_checkpoint_steps': 125,
15    'keep_checkpoint_max': 10,
16    'dataset_path': './preprocess',
17    'ckpt_save_path': './ckpt',
18    'checkpoint_path': './ckpt/gru-15_125.ckpt'
19 })

```

- GRU单元

```

1 class GRU(nn.Cell):
2     def __init__(self, config, is_training=True):
3         super(GRU, self).__init__()
4         if is_training:
5             self.batch_size = config.batch_size
6         else:
7             self.batch_size = config.eval_batch_size
8         self.hidden_size = config.hidden_size
9         self.weight_i, self.weight_h, self.bias_i, self.bias_h = \
10             gru_default_state(self.batch_size, self.hidden_size,
11 self.hidden_size)
12         self.rnn = P.DynamicGRUV2()
13         self.cast = P.Cast()
14
15     def construct(self, x, hidden):
16         x = self.cast(x, mstype.float16)
17         y1, h1, _, _, _ = self.rnn(x, self.weight_i, self.weight_h,
18 self.bias_i, self.bias_h, None, hidden)
19         return y1, h1

```

- Encoder

```

1 class Encoder(nn.Cell):
2     def __init__(self, config, is_training=True):
3         super(Encoder, self).__init__()
4         self.vocab_size = config.en_vocab_size
5         self.hidden_size = config.hidden_size
6         if is_training:
7             self.batch_size = config.batch_size
8         else:
9             self.batch_size = config.eval_batch_size
10
11         self.trans = P.Transpose()
12         self.perm = (1, 0, 2)
13         self.embedding = nn.Embedding(self.vocab_size, self.hidden_size)
14         self.gru = GRU(config,
15 is_training=is_training).to_float(mstype.float16)

```

```

15         self.h = Tensor(np.zeros((self.batch_size,
self.hidden_size)).astype(np.float16))
16
17     def construct(self, encoder_input):
18         embeddings = self.embedding(encoder_input)
19         embeddings = self.trans(embeddings, self.perm)
20         output, hidden = self.gru(embeddings, self.h)
21         return output, hidden

```

- Decoder

```

1 class Decoder(nn.Cell):
2     def __init__(self, config, is_training=True, dropout=0.1):
3         super(Decoder, self).__init__()
4
5         self.vocab_size = config.ch_vocab_size
6         self.hidden_size = config.hidden_size
7         self.max_len = config.max_seq_length
8
9         self.trans = P.Transpose()
10        self.perm = (1, 0, 2)
11        self.embedding = nn.Embedding(self.vocab_size, self.hidden_size)
12        self.dropout = nn.Dropout(1-dropout)
13        self.attn = nn.Dense(self.hidden_size, self.max_len)
14        self.softmax = nn.Softmax(axis=2)
15        self.bmm = P.BatchMatMul()
16        self.concat = P.Concat(axis=2)
17        self.attn_combine = nn.Dense(self.hidden_size * 2,
self.hidden_size)
18
19        self.gru = GRU(config,
is_training=is_training).to_float(mstype.float16)
20        self.out = nn.Dense(self.hidden_size, self.vocab_size)
21        self.logsoftmax = nn.LogSoftmax(axis=2)
22        self.cast = P.Cast()
23    def construct(self, decoder_input, hidden, encoder_output):
24        embeddings = self.embedding(decoder_input)
25        embeddings = self.dropout(embeddings)
26        # calculate attn
27        attn_weights = self.softmax(self.attn(embeddings))
28        encoder_output = self.trans(encoder_output, self.perm)
29        attn_applied = self.bmm(attn_weights,
self.cast(encoder_output, mstype.float32))
30        output = self.concat((embeddings, attn_applied))
31        output = self.attn_combine(output)
32
33
34        embeddings = self.trans(embeddings, self.perm)
35        output, hidden = self.gru(embeddings, hidden)
36        output = self.cast(output, mstype.float32)
37        output = self.out(output)
38        output = self.logsoftmax(output)
39
40        return output, hidden, attn_weights

```

- Seq2Seq整体结构

```

1 class Seq2Seq(nn.Cell):
2     def __init__(self, config, is_train=True):
3         super(Seq2Seq, self).__init__()
4         self.max_len = config.max_seq_length
5         self.is_train = is_train
6
7         self.encoder = Encoder(config, is_train)
8         self.decoder = Decoder(config, is_train)
9         self.expanddims = P.ExpandDims()
10        self.squeeze = P.Squeeze(axis=0)
11        self.argmax = P.ArgMaxWithValue(axis=int(2), keep_dims=True)
12        self.concat = P.Concat(axis=1)
13        self.concat2 = P.Concat(axis=0)
14        self.select = P.Select()
15
16        def construct(self, src, dst):
17            encoder_output, hidden = self.encoder(src)
18            decoder_hidden = self.squeeze(encoder_output[self.max_len-
192:self.max_len-1:1, :, :])
19            if self.is_train:
20                outputs, _ = self.decoder(dst, decoder_hidden, encoder_output)
21            else:
22                decoder_input = dst[:,0:1:1]
23                decoder_outputs = ()
24                for i in range(0, self.max_len):
25                    decoder_output, decoder_hidden, _ =
26self.decoder(decoder_input,
27decoder_hidden, encoder_output)
28                    decoder_hidden = self.squeeze(decoder_hidden)
29                    decoder_output, _ = self.argmax(decoder_output)
30                    decoder_output = self.squeeze(decoder_output)
31                    decoder_outputs += (decoder_output,)
32                    decoder_input = decoder_output
33                outputs = self.concat(decoder_outputs)
34            return outputs

```

- 损失函数

```

1 class NLLLoss(_Loss):
2     '''
3     NLLLoss function
4     '''
5     def __init__(self, reduction='mean'):
6         super(NLLLoss, self).__init__(reduction)
7         self.one_hot = P.OneHot()
8         self.reduce_sum = P.ReduceSum()
9
10        def construct(self, logits, label):
11            label_one_hot = self.one_hot(label, F.shape(logits)[-1],
12F.scalar_to_array(1.0),
13F.scalar_to_array(0.0))
14            #print('NLLLoss label_one_hot:', label_one_hot, label_one_hot.shape)
15            #print('NLLLoss logits:', logits, logits.shape)
16            #print('xxx:', logits * label_one_hot)
17            loss = self.reduce_sum(-1.0 * logits * label_one_hot, (1,))
18            return self.get_loss(loss)

```

- 训练网络与优化器

```
1 network = Seq2Seq(cfg)
2 network = WithLossCell(network, cfg)
3 optimizer = nn.Adam(network.trainable_params(),
4   learning_rate=cfg.learning_rate, beta1=0.9, beta2=0.98)
5 model = Model(network, optimizer=optimizer)
```

- 构建模型

```
1 loss_cb = LossMonitor()
2 config_ck =
3   CheckpointConfig(save_checkpoint_steps=cfg.save_checkpoint_steps,
4     keep_checkpoint_max=cfg.keep_checkpoint_max)
5 ckpoint_cb = ModelCheckpoint(prefix="gru", directory=cfg.ckpt_save_path,
6   config=config_ck)
7 time_cb = TimeMonitor(data_size=ds_train.get_dataset_size())
8 callbacks = [time_cb, ckpoint_cb, loss_cb]
9
10 model.train(cfg.num_epochs, ds_train, callbacks=callbacks,
11   dataset_sink_mode=True)
```

3.2 Transformer

- 数据处理函数

```
1 def data_prepare(cfg, eval_idx):
2     tokenizer = tokenization.WhiteSpaceTokenizer(vocab_file=cfg.vocab_file)
3
4     writer_train = Filewriter(cfg.train_file_mindrecord, cfg.num_splits)
5     writer_eval = Filewriter(cfg.eval_file_mindrecord, cfg.num_splits)
6     data_schema = {"source_sos_ids": {"type": "int32", "shape": [-1]},
7       "source_sos_mask": {"type": "int32", "shape": [-1]},
8       "source_eos_ids": {"type": "int32", "shape": [-1]},
9       "source_eos_mask": {"type": "int32", "shape": [-1]},
10      "target_sos_ids": {"type": "int32", "shape": [-1]},
11      "target_sos_mask": {"type": "int32", "shape": [-1]},
12      "target_eos_ids": {"type": "int32", "shape": [-1]},
13      "target_eos_mask": {"type": "int32", "shape": [-1]}
14    }
15
16     writer_train.add_schema(data_schema, "transformer train")
17     writer_eval.add_schema(data_schema, "transformer eval")
18
19     index = 0
20     f_train = open(cfg.train_file_source, 'w', encoding='utf-8')
21     f_test = open(cfg.eval_file_source, 'w', encoding='utf-8')
22     f = open(cfg.input_file, "r", encoding='utf-8')
23     for s_line in f:
24         print("finish {}/{}".format(index, 23607), end='\r')
25
26         line = tokenization.convert_to_unicode(s_line)
27
28         source_line, target_line = line.strip().split("\t")
```

```

29     source_tokens = tokenizer.tokenize(source_line)
30     target_tokens = tokenizer.tokenize(target_line)
31
32     if len(source_tokens) >= (cfg.max_seq_length-1) or
len(target_tokens) >= (cfg.max_seq_length-1):
33         if cfg.clip_to_max_len:
34             source_tokens = source_tokens[:cfg.max_seq_length-1]
35             target_tokens = target_tokens[:cfg.max_seq_length-1]
36         else:
37             continue
38
39     index = index + 1
40     # print(source_tokens)
41     instance = create_training_instance(source_tokens, target_tokens,
cfg.max_seq_length)
42
43     if index in eval_idx:
44         f_test.write(s_line)
45         features = write_instance_to_file(writer_eval, instance,
tokenizer, cfg.max_seq_length)
46     else:
47         f_train.write(s_line)
48         features = write_instance_to_file(writer_train, instance,
tokenizer, cfg.max_seq_length)
49     f.close()
50     f_test.close()
51     f_train.close()
52     writer_train.commit()
53     writer_eval.commit()

```

- 定义数据加载函数

```

1  def load_dataset(batch_size=1, data_file=None):
2      """
3      Load mindrecord dataset
4      """
5      ds = de.MindDataset(data_file,
6                          columns_list=["source_eos_ids", "source_eos_mask",
7                                       "target_sos_ids", "target_sos_mask",
8                                       "target_eos_ids", "target_eos_mask"],
9                          shuffle=False)
10     type_cast_op = deC.TypeCast(mstype.int32)
11     ds = ds.map(input_columns="source_eos_ids", operations=type_cast_op)
12     ds = ds.map(input_columns="source_eos_mask", operations=type_cast_op)
13     ds = ds.map(input_columns="target_sos_ids", operations=type_cast_op)
14     ds = ds.map(input_columns="target_sos_mask", operations=type_cast_op)
15     ds = ds.map(input_columns="target_eos_ids", operations=type_cast_op)
16     ds = ds.map(input_columns="target_eos_mask", operations=type_cast_op)
17     # apply batch operations
18     ds = ds.batch(batch_size, drop_remainder=True)
19     ds.channel_name = 'transformer'
20     return ds

```

- 训练函数

```

1  def train(cfg):

```



```

2      """
3      Transformer training.
4      """
5
6      train_dataset = load_dataset(cfg.batch_size, data_file=cfg.data_path)
7
8      netwithloss = TransformerNetworkWithLoss(transformer_net_cfg, True)
9
10     if cfg.checkpoint_path:
11         parameter_dict = load_checkpoint(cfg.checkpoint_path)
12         load_param_into_net(netwithloss, parameter_dict)
13
14     lr =
Tensor(create_dynamic_lr(schedule="constant*rsqrt_hidden*linear_warmup*rsqr
t_decay",
15     training_steps=train_dataset.get_dataset_size()*cfg.epoch_size,
16     learning_rate=cfg.lr_schedule.learning_rate,
17     warmup_steps=cfg.lr_schedule.warmup_steps,
18     hidden_size=transformer_net_cfg.hidden_size,
19     start_decay_step=cfg.lr_schedule.start_decay_step,
20                                     min_lr=cfg.lr_schedule.min_lr),
mstype.float32)
21     optimizer = Adam(netwithloss.trainable_params(), lr)
22
23     callbacks = [TimeMonitor(train_dataset.get_dataset_size()),
LossCallBack()]
24     if cfg.enable_save_ckpt:
25         ckpt_config =
CheckpointConfig(save_checkpoint_steps=cfg.save_checkpoint_steps,
26
keep_checkpoint_max=cfg.save_checkpoint_num)
27         ckpoint_cb = ModelCheckpoint(prefix=cfg.save_checkpoint_name,
directory=cfg.save_checkpoint_path, config=ckpt_config)
28         callbacks.append(ckpt_cb)
29
30     if cfg.enable_lossscale:
31         scale_manager =
DynamicLossScaleManager(init_loss_scale=cfg.init_loss_scale_value,
32
scale_factor=cfg.scale_factor,
33
scale_window=cfg.scale_window)
34         update_cell = scale_manager.get_update_cell()
35         netwithgrads =
TransformerTrainOneStepWithLossScaleCell(netwithloss,
optimizer=optimizer, scale_update_cell=update_cell)
36     else:
37         netwithgrads = TransformerTrainOneStepCell(netwithloss,
optimizer=optimizer)
38
39     netwithgrads.set_train(True)
40     model = Model(netwithgrads)

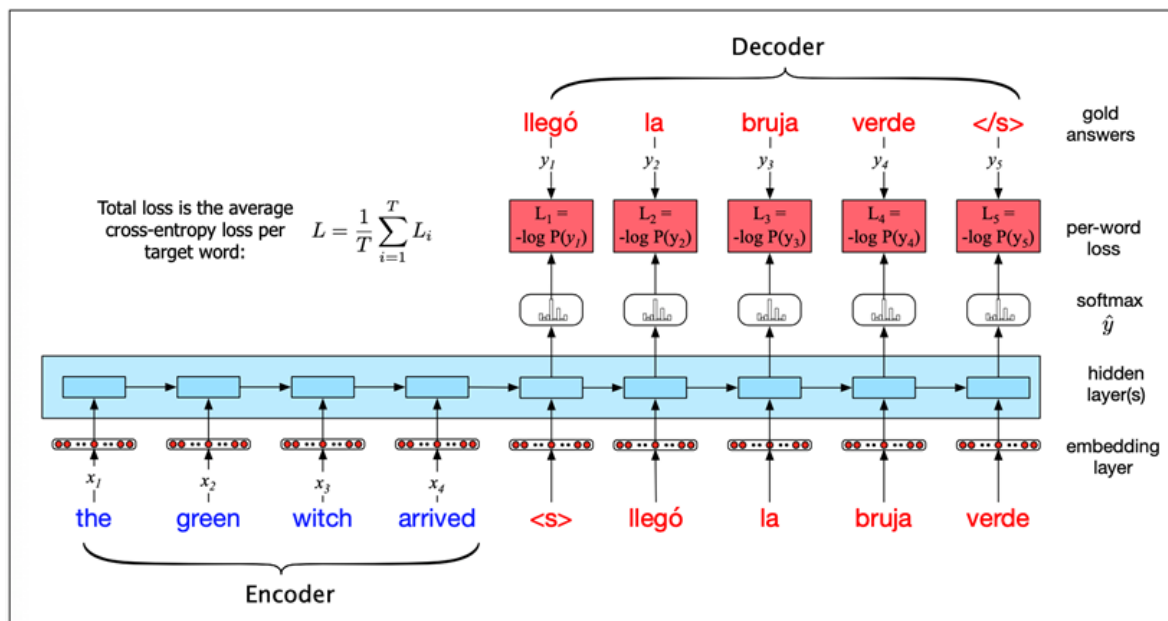
```

```
41 model.train(cfg.epoch_size, train_dataset, callbacks=callbacks,
dataset_sink_mode=cfg.enable_data_sink)
```

4.实验原理

4.1 Seq2Seq

在自然语言处理的很多应用中，输入和输出都可以是不定长序列。以机器翻译为例，输入可以是一段不定长的英语文本序列，输出可以是一段不定长的德语文本序列，如下面图中所示。



当输入和输出都是不定长序列时，我们可以使用编码器—解码器（encoder-decoder）或者 seq2seq 模型。序列到序列模型，简称 seq2seq 模型。这两个模型本质上都用到了两个循环神经网络，分别叫做**编码器和解码器**。编码器用来分析输入序列，解码器用来生成输出序列。两个循环神经网络是共同训练的。

4.1.1 编码器

编码器的作用是把一个不定长的输入序列变换成一个定长的背景变量 c ，并在该背景变量中编码输入序列信息。常用的编码器是循环神经网络。假设输入序列是 x_1, \dots, x_T ，例如 x_i 是输入句子中的第 i 个词。在时间步 t ，循环神经网络将输入 x_t 的特征向量 x_t 和上个时间步的隐藏状态 h_{t-1} 变换为当前时间步的隐藏状态 h_t 。用函数 f 表达循环神经网络隐藏层的变换：

$$h_t = f(x_t, h_{t-1})$$

编码器通过自定义函数 q 将各个时间步的隐藏状态变换为背景变量：

$$c = q(h_1, \dots, h_T)$$

4.1.2 译码器

编码器输出的背景变量 c 编码了整个输入序列 x_1, \dots, x_T 的信息。给定训练样本中的输出序列 $y_1, y_2, \dots, y_{T'}$ ，对每个时间步 t' （符号与输入序列或编码器的时间步 t 有区别），解码器输出 $y_{t'}$ 的条件概率将基于之前的输出序列 $y_1, \dots, y_{t'-1}$ 和背景变量 c ，即：

$$P(y_{t'} | y_1, \dots, y_{t'-1}, c)$$

可以使用另一个循环神经网络作为解码器。在输出序列的时间步 t' ，解码器将上一时间步的输出 $y_{t'-1}$ 以及背景变量 c 作为输入，并将它们与上一时间步的隐藏状态 $s_{t'-1}$ 变换为当前时间步的隐藏状态 $s_{t'}$ 。因此，我们可以用函数 g 表达解码器隐藏层的变换：

$$s_{t'} = g(y_{t'-1}, c, s_{t'-1})$$

4.1.3 训练模型

根据最大似然估计，我们可以最大化输出序列基于输入序列的条件概率：

$$\begin{aligned} P(y_1, \dots, y_{t'-1} \mid x_1, \dots, x_T) &= \prod_{t'=1}^{T'} P(y_{t'} \mid y_1, \dots, y_{t'-1}, x_1, \dots, x_T) \\ &= \prod_{t'=1}^{T'} P(y_{t'} \mid y_1, \dots, y_{t'-1}, c) \end{aligned}$$

并得到该输出序列的损失：

$$-\log P(y_1, \dots, y_{t'-1} \mid x_1, \dots, x_T) = -\sum_{t'=1}^{T'} \log P(y_{t'} \mid y_1, \dots, y_{t'-1}, c)$$

假设解码器的输出是一段文本序列。设输出文本词典 Y （包含特殊符号“”）的大小为 $|Y|$ ，输出序列的最大长度为 T' 。所有可能的输出序列一共有 $O(|Y|^{T'})$ 种。这些输出序列中所有特殊符号“”后面的子序列将被舍弃。

4.2 Transformer

Transformer网络如下图所示，其中左边为编码网络，右边为解码网络。

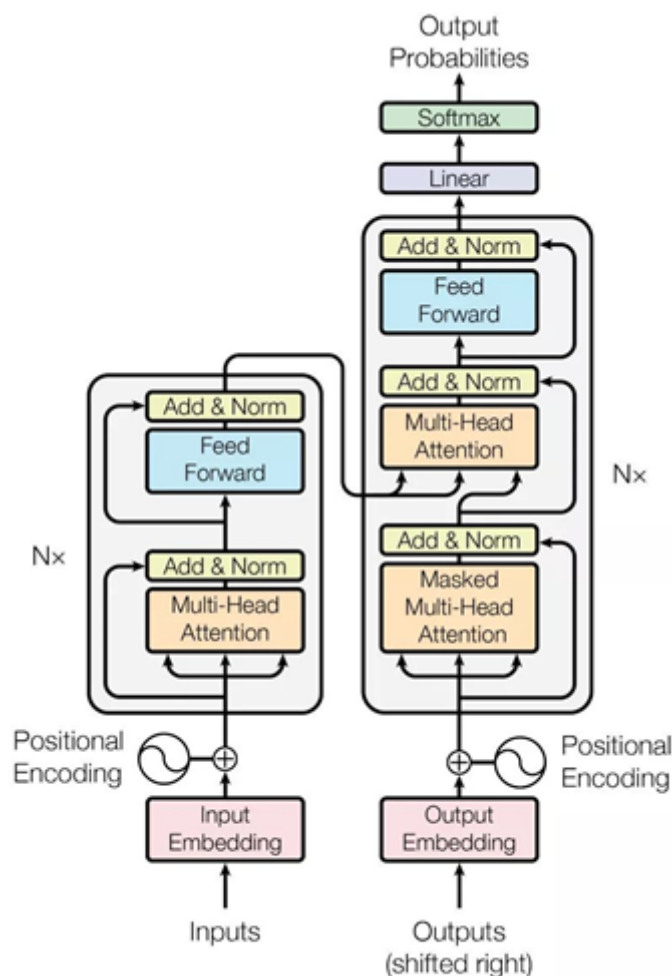


Figure 1: The Transformer - model architecture.

每一个编码器在结构上都是一样的，但它们的权重参数是不同的。每一个编码器里面，可以分为2层（Self-Attention 层、前馈神经网络）。输入编码器的文本数据，首先会经过一个 Self Attention 层，这个层处理一个词的时候，不仅会使用这个词本身的信息，也会使用句子中其他词的信息。接下来，Self Attention 层的输出会经过前馈神经网络。同理，解码器也具有这两层，但是这两层中间还插入了一个 Encoder-Decoder Attention 层。

在self-attention中，每个单词有3个不同的向量，它们分别是Query向量（Q），Key向量（K）和Value向量（V）。它们是通过3个不同的权值矩阵由嵌入向量乘以三个不同的权值矩阵，self-attention的公式如下：

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

5 实验结果

5.1 Seq2Seq

- 训练

```
In [15]: loss_cb = LossMonitor()
         config_ck = CheckpointConfig(save_checkpoint_steps=config.save_checkpoint_steps, keep_checkpoint_max=config.keep_checkpoint_max)
         ckpoint_cb = ModelCheckpoint(prefix="gru", directory=config.ckpt_save_path, config=config_ck)
         time_cb = TimeMonitor(data_size=ds_train.get_dataset_size())
         callbacks = [time_cb, ckpoint_cb, loss_cb]

         model.train(config.num_epochs, ds_train, callbacks=callbacks, dataset_sink_mode=True)

epoch: 1 step: 125, loss is 2.901339
epoch time: 77857.532 ms, per step time: 622.860 ms
epoch: 2 step: 125, loss is 2.4917648
epoch time: 11361.552 ms, per step time: 90.892 ms
epoch: 3 step: 125, loss is 2.2108736
epoch time: 11393.523 ms, per step time: 91.148 ms
epoch: 4 step: 125, loss is 1.6616195
epoch time: 11329.437 ms, per step time: 90.635 ms
epoch: 5 step: 125, loss is 1.5475513
epoch time: 11329.803 ms, per step time: 90.638 ms
epoch: 6 step: 125, loss is 1.2502322
epoch time: 11373.320 ms, per step time: 90.987 ms
epoch: 7 step: 125, loss is 0.97378296
epoch time: 11316.333 ms, per step time: 90.531 ms
epoch: 8 step: 125, loss is 0.3811617
epoch time: 11423.199 ms, per step time: 91.386 ms
epoch: 9 step: 125, loss is 0.2629779
epoch time: 11470.042 ms, per step time: 91.760 ms
epoch: 10 step: 125, loss is 0.27675834
epoch time: 11495.968 ms, per step time: 91.968 ms
epoch: 11 step: 125, loss is 0.2847989
epoch time: 11563.872 ms, per step time: 92.511 ms
epoch: 12 step: 125, loss is 0.10026628
epoch time: 11553.943 ms, per step time: 92.432 ms
epoch: 13 step: 125, loss is 0.20696855
epoch time: 11676.776 ms, per step time: 93.414 ms
epoch: 14 step: 125, loss is 0.10150813
epoch time: 11522.073 ms, per step time: 92.177 ms
epoch: 15 step: 125, loss is 0.15692167
epoch time: 11475.951 ms, per step time: 91.808 ms
```

- 测试

```
In [23]: translate('i need a bottle of juice')
         translate('Tom is a good teacher')
```

```
English ['i', 'need', 'a', 'bottle', 'of', 'juice']
中文 我需要一张邮票。
English ['tom', 'is', 'a', 'good', 'teacher']
中文 汤姆是个真的。
```

5.2 Transformer

- 训练

10. 启动训练

```
In [11]: train(train_cfg)
time: 436852, epoch: 15, step: 8832, outputs are [2.9468496]
time: 436886, epoch: 15, step: 8833, outputs are [2.8108442]
time: 436920, epoch: 15, step: 8834, outputs are [2.9115794]
time: 436954, epoch: 15, step: 8835, outputs are [2.9215643]
time: 436988, epoch: 15, step: 8836, outputs are [2.813668]
time: 437022, epoch: 15, step: 8837, outputs are [2.9565923]
time: 437057, epoch: 15, step: 8838, outputs are [2.885069]
time: 437091, epoch: 15, step: 8839, outputs are [2.8448014]
time: 437125, epoch: 15, step: 8840, outputs are [3.0417776]
time: 437159, epoch: 15, step: 8841, outputs are [3.0115216]
time: 437193, epoch: 15, step: 8842, outputs are [3.072069]
time: 437227, epoch: 15, step: 8843, outputs are [3.1246238]
time: 437261, epoch: 15, step: 8844, outputs are [2.9308126]
time: 437295, epoch: 15, step: 8845, outputs are [3.0508318]
time: 437329, epoch: 15, step: 8846, outputs are [3.0099447]
time: 437363, epoch: 15, step: 8847, outputs are [3.260268]
time: 437397, epoch: 15, step: 8848, outputs are [3.2783072]
time: 437431, epoch: 15, step: 8849, outputs are [3.5525923]
time: 437466, epoch: 15, step: 8850, outputs are [3.7860718]
epoch time: 21620.916 ms, per step time: 36.985 ms
```

• 测试

```
[17]: evaluate(eval_cfg)
```

```
source: Hop in .           跳 进 来 。
result: 在大部份的时候受到了一个普通的东西。
source: I ' m up .         我已经起来了。
result: 我起床了。
source: Why me ?           为什么是我？
result: 为什么我为什么为什么为什么我要？
source: Be fair .          公平点。
result: 公平等公平运公平。
source: Be kind .          友善点。
result: 好运点点点点穿点点点好运。
source: Be nice .          和气点。
result: 友善点。
```