

机器翻译

背景介绍

机器翻译 (machine translation, MT) 是用计算机来实现不同语言之间翻译的技术。被翻译的语言通常称为源语言 (source language) ，翻译成的结果语言称为目标语言 (target language) 。机器翻译即实现从源语言到目标语言转换的过程，是自然语言处理的重要研究领域之一。

早期机器翻译系统多为基于规则的翻译系统，需要由语言学家编写两种语言之间的转换规则，再将这些规则录入计算机。该方法对语言学家的要求非常高，而且我们几乎无法总结一门语言会用到的所有规则，更何况两种甚至更多的语言。因此，传统机器翻译方法面临的主要挑战是无法得到一个完备的规则集合[1]。

为解决以上问题，统计机器翻译 (Statistical Machine Translation, SMT) 技术应运而生。在统计机器翻译技术中，转化规则是由机器自动从大规模的语料中学习得到的，而非我们人主动提供规则。因此，它克服了基于规则的翻译系统所面临的知识获取瓶颈的问题，但仍然存在许多挑战：1) 人为设计许多特征 (feature) ，但永远无法覆盖所有的语言现象；2) 难以利用全局的特征；3) 依赖于许多预处理环节，如词语对齐、分词或符号化 (tokenization) 、规则抽取、句法分析等，而每个环节的错误会逐步累积，对翻译的影响也越来越大。

近年来，深度学习技术的发展为解决上述挑战提供了新的思路。将深度学习应用于机器翻译任务的方法大致分为两类：1) 仍以统计机器翻译系统为框架，只是利用神经网络来改进其中的关键模块，如语言模型、调序模型等 (见图1的左半部分) ；2) 不再以统计机器翻译系统为框架，而是直接用神经网络将源语言映射到目标语言，即端到端的神经网络机器翻译 (End-to-End Neural Machine Translation, End-to-End NMT) (见图1的右半部分) ，简称为NMT模型。

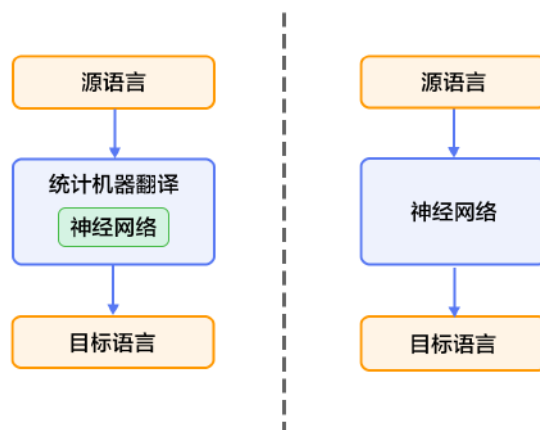


图1. 基于神经网络的机器翻译系统

本教程主要介绍NMT模型，以及如何用PaddlePaddle来训练一个NMT模型。

效果展示

以中英翻译（中文翻译到英文）的模型为例，当模型训练完毕时，如果输入如下已分词的中文句子：

- 1. 这些是希望的曙光和解脱的迹象。
- 1. 0 -5.36816 these are signs of hope and relief . <e>
- 2. 1 -6.23177 these are the light of hope and relief . <e>
- 3. 2 -7.7914 these are the light of hope and the relief of hope . <e>

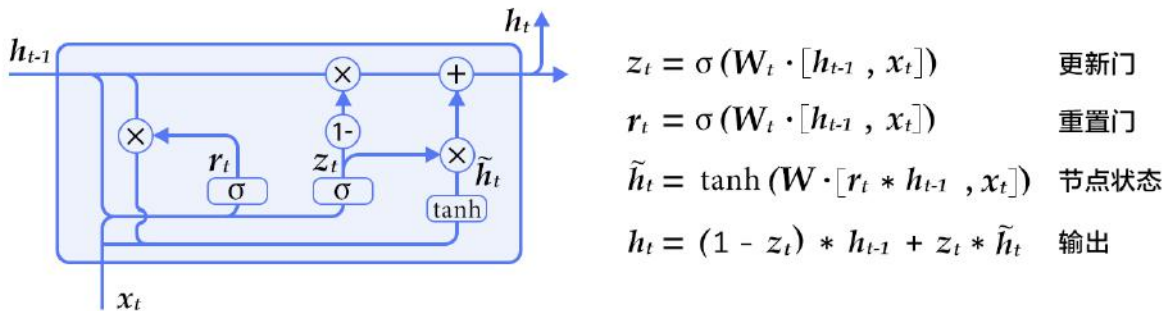


图2. GRU（门控循环单元）

一般来说，具有短距离依赖属性的序列，其重置门比较活跃；相反，具有长距离依赖属性的序列，其更新门比较活跃。另外，Chung等人[3]通过多组实验表明，GRU虽然参数更少，但是在多个任务上都和LSTM有相近的表现。

双向循环神经网络

我们已经在语义角色标注一章中介绍了一种双向循环神经网络，这里介绍Bengio团队在论文[2,4]中提出的另一种结构。该结构的目的是输入一个序列，得到其在每个时刻的特征表示，即输出的每个时刻都用定长向量表示到该时刻的上下文语义信息。

具体来说，该双向循环神经网络分别在时间维以顺序和逆序——即前向（forward）和后向（backward）——依次处理输入序列，并将每个时间步RNN的输出拼接成为最终的输出层。这样每个时间步的输出节点，都包含了输入序列中当前时刻完整的过去和未来的上下文信息。下图展示的是一个按时间步展开的双向循环神经网络。该网络包含一个前向和一个后向RNN，其中有六个权重矩阵：输入到前向隐层和后向隐层的权重矩阵（ W_1, W_3 ），隐层到隐层自己的权重矩阵（ W_2, W_5 ），前向隐层和后向隐层到输出层的权重矩阵（ W_4, W_6 ）。注意，该网络的前向隐层和后向隐层之间没有连接。

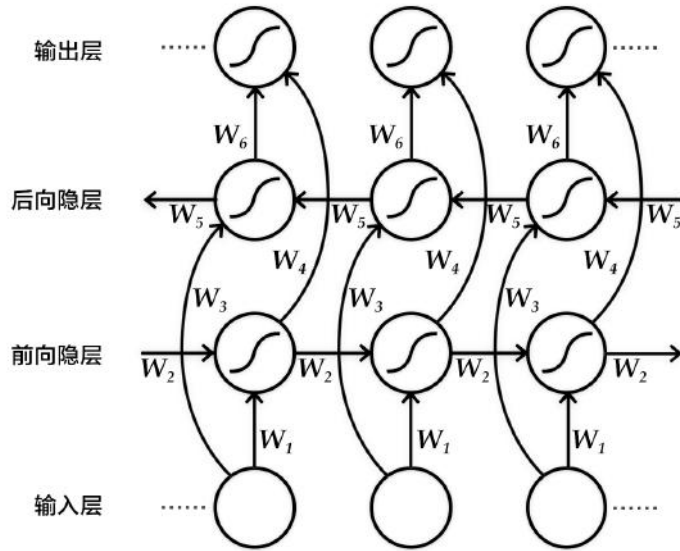


图3. 按时间步展开的双向循环神经网络

编码器-解码器框架

编码器-解码器 (Encoder-Decoder) [2] 框架用于解决由一个任意长度的源序列到另一个任意长度的目标序列的变换问题。即编码阶段将整个源序列编码成一个向量，解码阶段通过最大化预测序列概率，从中解码出整个目标序列。编码和解码的过程通常都使用RNN实现。

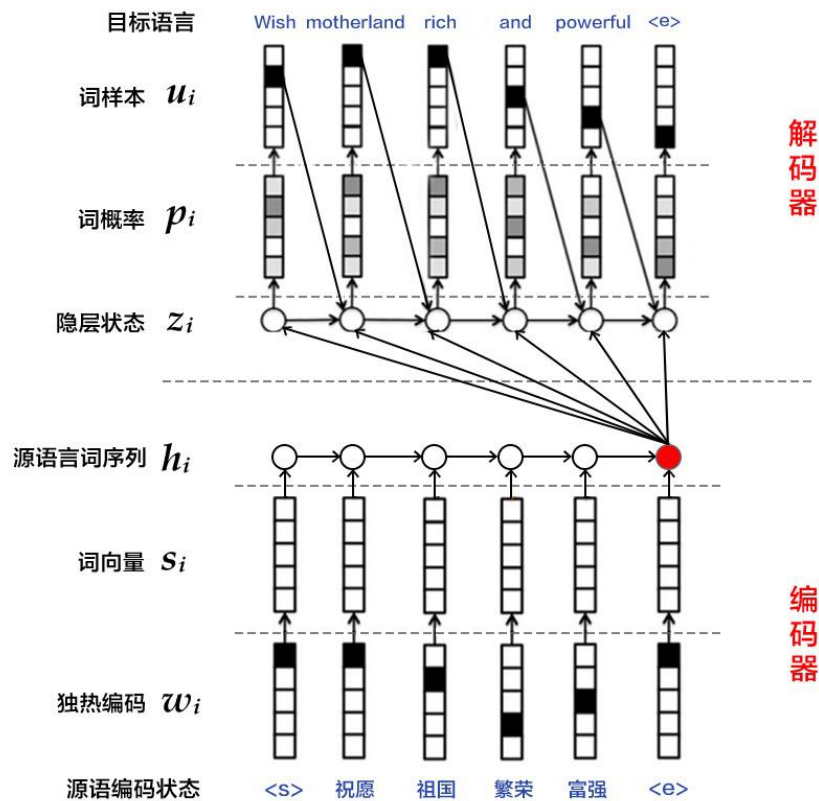


图4. 编码器-解码器框架

编码器

编码阶段分为三步：

1. one-hot vector表示：将源语言句子 $\mathbf{x} = \{x_1, x_2, \dots, x_T\}$ 的每个词 x_i 表示成一个列向量 $w_i \in \mathbb{R}^{|V|}, i = 1, 2, \dots, T$ 。这个向量 w_i 的维度与词汇表大小 $|V|$ 相同，并且只有一个维度上有值1（该位置对应该词在词汇表中的位置），其余全是0。
2. 映射到低维语义空间的词向量：one-hot vector表示存在两个问题，1) 生成的向量维度往往很大，容易造成维数灾难；2) 难以刻画词与词之间的关系（如语义相似性，也就是无法很好地表达语义）。因此，需再 one-hot vector映射到低维的语义空间，由一个固定维度的稠密向量（称为词向量）表示。记映射矩阵为 $C \in \mathbb{R}^{K \times |V|}$ ，用 $s_i = Cw_i$ 表示第 i 个词的词向量， K 为向量维度。
3. 用RNN编码源语言词序列：这一过程的计算公式为 $h_i = \mathcal{O}_\theta(h_{i-1}, s_i)$ ，其中 h_0 是一个全零的向量， \mathcal{O}_θ 是一个非线性激活函数，最后得到的 $\mathbf{h} = \{h_1, \dots, h_T\}$ 就是RNN依次读入源语言 T 个词的状态编码序列。整句话的向量表示可以采用 \mathbf{h} 在最后一个时间步 T 的状态编码，或使用时间维上的池化（pooling）结果。

第3步也可以使用双向循环神经网络实现更复杂的句编码表示，具体可以用双向GRU实现。前向GRU按照词序列 (x_1, x_2, \dots, x_T) 的顺序依次编码源语言端词，并得到一系列隐层状态 (h_1, h_2, \dots, h_T) 。类似的，后向GRU按照 $(x_T, x_{T-1}, \dots, x_1)$ 的顺序依次编码源语言端词，得到 (h_1, h_2, \dots, h_T) 。最后对于词 x_i ，通过拼接两个GRU的结果得到它的隐层状态，即 $h_i = \left[\begin{matrix} \rightarrow h_i^T \\ \leftarrow h_i^T \end{matrix} \right]^T$ 。

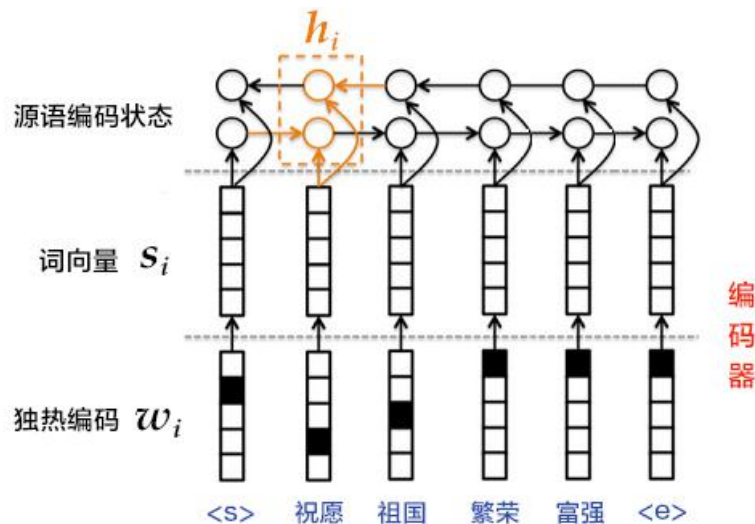


图5. 使用双向GRU的编码器

解码器

机器翻译任务的训练过程中，解码阶段的目标是最大化下一个正确的源语言词的概率。思路是：

1. 每一个时刻，根据源语言句子的编码信息（又叫上下文向量，context vector） c 、真实目标语言序列的第 i

个词 u_i 和 i 时刻RNN的隐层状态 z_i ，计算出下一个隐层状态 z_{i+1} 。计算公式如下：

$$z_{i+1} = \phi_{\theta}(c, u_i, z_i)$$

其中 ϕ_{θ} 是一个非线性激活函数； $c = qh$ 是源语言句子的上下文向量，在不使用注意力机制时，如果编码器的输出是源语言句子编码后的最后一个元素，则可以定义 $c = h_T$ ； u_i 是目标语言序列的第 i 个单词， u_0 是目标语言序列的开始标记 $\langle s \rangle$ ，表示解码开始； z_i 是 i 时刻解码RNN的隐层状态， z_0 是一个全零的向量。

2. 将 z_{i+1} 通过softmax归一化，得到目标语言序列的第 $i + 1$ 个单词的概率分布 p_{i+1} 。概率分布公式如下：

$$p(u_{i+1} | u_{<i+1>, \mathbf{x}}) = \text{softmax}(W_s z_{i+1} + b_z)$$

其中 $W_s z_{i+1} + b_z$ 是对每个可能的输出单词进行打分，再用softmax归一化就可以得到第 $i + 1$ 个词的概率 p_{i+1} 。

3. 根据 p_{i+1} 和 u_{i+1} 计算代价。

4. 重复步骤1~3，直到目标语言序列中的所有词处理完毕。

机器翻译任务的生成过程，通俗来讲就是根据预先训练的模型来翻译源语言句子。生成过程中的解码阶段和上述训练过程的有所差异，具体介绍请见[柱搜索算法](#)。

注意力机制

如果编码阶段的输出是一个固定维度的向量，会带来以下两个问题：1) 不论源语言序列的长度是5个词还是50个词，如果都用固定维度的向量去编码其中的语义和句法结构信息，对模型来说是一个非常高的要求，特别是对长句子序列而言；2) 直觉上，当人类翻译一句话时，会对与当前译文更相关的源语言片段上给予更多关注，且关注点会随着翻译的进行而改变。而固定维度的向量则相当于，任何时刻都对源语言所有信息给予了同等程度的关注，这是不合理的。因此，Bahdanau等人[4]引入注意力（attention）机制，可以对编码后的上下文片段进行解码，以此来解决长句子的特征学习问题。下面介绍在注意力机制下的解码器结构。

与简单的解码器不同，这里 z_i 的计算公式为：

$$z_{i+1} = \phi_{\theta}(c_i, u_i, z_i)$$

可见，源语言句子的编码向量表示为第 i 个词的上下文片段 c_i ，即针对每一个目标语言中的词 u_i ，都有一个特定的 c_i 与之对应。 c_i 的计算公式如下：

$$c_i = \sum_{j=1}^T a_{ij} h_j, a_i = [a_{i1}, a_{i2}, \dots, a_{iT}]$$

从公式中可以看出，注意力机制是通过在编码器中各时刻的RNN状态 h_j 进行加权平均实现的。权重 a_{ij} 表示目标语言中第 i 个词对源语言中第 j 个词的注意力大小， a_{ij} 的计算公式如下：

$$a_{ij} = \frac{\exp(e_{ij})}{\sum_{k=1}^T \exp(e_{ik})}$$

$$e_{ij} = \text{align}(z_i, h_j)$$

其中，*align*可以看作是一个对齐模型，用来衡量目标语言中第*i*个词和源语言中第*j*个词的匹配程度。具体而言，这个程度是通过解码RNN的第*i*个隐层状态 z_i 和源语言句子的第*j*个上下文片段 h_j 计算得到的。传统的对齐模型中，目标语言的每个词明确对应源语言的一个或多个词（hard alignment）；而在注意力模型中采用的是soft alignment，即任何两个目标语言和源语言词间均存在一定的关联，且这个关联强度是由模型计算得到的实数，因此可以融入整个NMT框架，并通过反向传播算法进行训练。

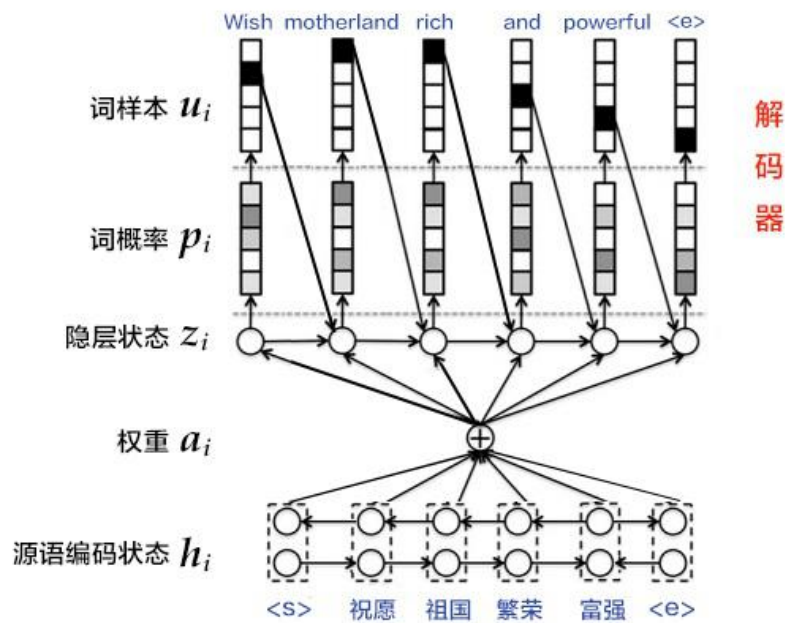


图6. 基于注意力机制的解码器

柱搜索算法

柱搜索（beam search）是一种启发式图搜索算法，用于在图或树中搜索有限集合中的最优扩展节点，通常用在解空间非常大的系统（如机器翻译、语音识别）中，原因是内存无法装下图或树中所有展开的解。如在机器翻译任务中希望翻译“<s>你好<e>”，就算目标语言字典中只有3个词（<s>，<e>，hello），也可能生成无限句话（hello循环出现的次数不定），为了找到其中较好的翻译结果，我们可采用柱搜索算法。

柱搜索算法使用广度优先策略建立搜索树，在树的每一层，按照启发代价（heuristic cost）（本教程中，为生成词的log概率之和）对节点进行排序，然后仅留下预先确定的个数（文献中通常称为beam width、beam size、柱宽度等）的节点。只有这些节点会在下一层继续扩展，其他节点就被剪掉了，也就是说保留了质量较高的节点，剪枝了质量较差的节点。因此，搜索所占用的空间和时间大幅减少，但缺点是无法保证一定获得最优解。

使用柱搜索算法的解码阶段，目标是最大化生成序列的概率。思路是：

1. 每一个时刻，根据源语言句子的编码信息 \mathbf{c} 、生成的第 i 个目标语言序列单词 \mathbf{u}_i 和 i 时刻RNN的隐层状态 \mathbf{z}_i ，计算出下一个隐层状态 \mathbf{z}_{i+1} 。
2. 将 \mathbf{z}_{i+1} 通过 `softmax` 归一化，得到目标语言序列的第 $i + 1$ 个单词的概率分布 \mathbf{p}_{i+1} 。
3. 根据 \mathbf{p}_{i+1} 采样出单词 \mathbf{u}_{i+1} 。
4. 重复步骤1~3，直到获得句子结束标记 `<e>` 或超过句子的最大生成长度为止。

注意： \mathbf{z}_{i+1} 和 \mathbf{p}_{i+1} 的计算公式同[解码器](#)中的一样。且由于生成时的每一步都是通过贪心法实现的，因此并不能保证得到全局最优解。

数据准备

下载与解压缩

本教程使用WMT-14数据集中的[bitexts\(after selection\)](#)作为训练集，[dev+test data](#)作为测试集和生成集。

在Linux下，只需简单地运行以下命令：

```
1. cd data
2. ./wmt14_data.sh
```

文件夹名	法英平行语料文件	文件数	文件大小
train	ccb2_pc30.src, ccb2_pc30.trg, etc	12	3.55G
test	ntst1213.src, ntst1213.trg	2	1636k
gen	ntst14.src, ntst14.trg	2	864k

- `XXX.src` 是源法语文件，`XXX.trg` 是目标英语文件，文件中的每行存放一个句子
- `XXX.src` 和 `XXX.trg` 的行数一致，且两者任意第 i 行的句子之间都有着——对应的关系。

用户自定义数据集（可选）

如果您想使用自己的数据集，只需按照如下方式组织，并将它们放在 `data` 目录下：

```
1. user_dataset
2. |— train
3. |   |— train_file1.src
4. |   |— train_file1.trg
5. |   |— ...
6. |— test
```

```
7. | | | test_file1.src
8. | | | test_file1.trg
9. | | | ...
10. | | gen
11. | | | gen_file1.src
12. | | | gen_file1.trg
13. | | | ...
```

- 一级目录 `user_dataset` : 用户自定义的数据集名字。
- 二级目录 `train`、`test` 和 `gen` : 必须使用这三个文件夹名字。
- 三级目录 : 存放源语言到目标语言的平行语料库文件, 后缀名必须使用 `.src` 和 `.trg`。

数据预处理

我们的预处理流程包括两步 :

- 将每个源语言到目标语言的平行语料库文件合并为一个文件 :
- 合并每个 `XXX.src` 和 `XXX.trg` 文件为 `XXX`。
- `XXX` 中的第 i 行内容为 `XXX.src` 中的第 i 行和 `XXX.trg` 中的第 i 行连接, 用 `\t` 分隔。
- 创建训练数据的“源字典”和“目标字典”。每个字典都有 **DICTSIZE** 个单词, 包括: 语料中词频最高的 (`DICTSIZE - 3`) 个单词, 和 3 个特殊符号 `<s>` (序列的开始)、`<e>` (序列的结束) 和 `<unk>` (未登录词)。

预处理可以使用 `preprocess.py` :

```
1. python preprocess.py -i INPUT [-d DICTSIZE] [-m]
```

- `-i INPUT` : 输入的原始数据集路径。
- `-d DICTSIZE` : 指定的字典单词数, 如果没有设置, 字典会包含输入数据集中的所有单词。
- `-m --mergeDict` : 合并“源字典”和“目标字典”, 即这两个字典的内容完全一样。

本教程的具体命令如下 :

```
1. python preprocess.py -i data/wmt14 -d 30000
```

请耐心等待几分钟的时间, 您会在屏幕上看到 :

```
1. concat parallel corpora for dataset
2. build source dictionary for train data
3. build target dictionary for train data
4. dictionary size is 30000
```

预处理好的数据集存放在 `data/pre-wmt14` 目录下 :


```

1.  pre-wmt14
2.  |— train
3.  |   |— train
4.  |— test
5.  |   |— test
6.  |— gen
7.  |   |— gen
8.  |— train.list
9.  |— test.list
10. |— gen.list
11. |— src.dict
12. |— trg.dict

```

- `train`、`test` 和 `gen`：分别包含了法英平行语料库的训练、测试和生成数据。其每个文件的每一行以“\t”分为两列，第一列是法语序列，第二列是对应的英语序列。
- `train.list`、`test.list` 和 `gen.list`：分别记录了 `train`、`test` 和 `gen` 文件夹中的文件路径。
- `src.dict` 和 `trg.dict`：源（法语）和目标（英语）字典。每个字典都含有30000个单词，包括29997个最高频单词和3个特殊符号。

提供数据给PaddlePaddle

我们通过 `dataproducer.py` 将数据提供给PaddlePaddle。具体步骤如下：

1. 首先，引入PaddlePaddle的PyDataProvider2包，并定义三个特殊符号。

```

1.  from paddle.trainer.PyDataProvider2 import *
2.  UNK_IDX = 2      #未登录词
3.  START = "<s>"    #序列的开始
4.  END = "<e>"     #序列的结束

```

2. 其次，使用初始化函数 `hook`，分别定义了训练模式和生成模式下的数据输入格式（`input_types`）。

- 训练模式：有三个输入序列，其中“源语言序列”和“目标语言序列”作为输入数据，“目标语言的下一个词序列”作为标签数据。
- 生成模式：有两个输入序列，其中“源语言序列”作为输入数据，“源语言序列编号”作为输入数据的编号（该输入非必须，可以省略）。

`hook` 函数中的 `src_dict_path` 是源语言字典路径，`trg_dict_path` 是目标语言字典路

径，`is_generating`（训练或生成模式）是从模型配置中传入的对象。`hook` 函数的具体调用方式请见[训练模型配置说明](#)。

```

1.  def hook(settings, src_dict_path, trg_dict_path, is_generating, file_list,
2.         **kwargs):
3.      # job_mode = 1: 训练模式; 0: 生成模式
4.      settings.job_mode = not is_generating
5.
6.      def fun(dict_path): # 根据字典路径加载字典

```

```

7.         out_dict = dict()
8.         with open(dict_path, "r") as fin:
9.             out_dict = {
10.                 line.strip(): line_count
11.                 for line_count, line in enumerate(fin)
12.             }
13.         return out_dict
14.
15.     settings.src_dict = fun(src_dict_path)
16.     settings.trg_dict = fun(trg_dict_path)
17.
18.     if settings.job_mode:                                #训练模式
19.         settings.input_types = {
20.             'source_language_word':                    #源语言序列
21.                 integer_value_sequence(len(settings.src_dict)),
22.             'target_language_word':                    #目标语言序列
23.                 integer_value_sequence(len(settings.trg_dict)),
24.             'target_language_next_word':                #目标语言的下一个词序列
25.                 integer_value_sequence(len(settings.trg_dict))
26.         }
27.     else:                                                #生成模式
28.         settings.input_types = {
29.             'source_language_word':                    #源语言序列
30.                 integer_value_sequence(len(settings.src_dict)),
31.             'sent_id':                                #源语言序列编号
32.                 integer_value_sequence(len(open(file_list[0], "r").readlines()))
33.         }

```

3. 最后，使用 `process` 函数打开文本文件 `file_name`，读取每一行，将行中的数据转换成与 `input_types` 一致的格式，再用 `yield` 关键字返回给PaddlePaddle进程。具体来说，

- 在源语言序列的每句话前面补上开始符号 `<s>`、末尾补上结束符号 `<e>`，得到 `"source_language_word"`；
- 在目标语言序列的每句话前面补上 `<s>`，得到 `"target_language_word"`；
- 在目标语言序列的每句话末尾补上 `<e>`，作为目标语言的下一个词序列（`"target_language_next_word"`）。

```

1.     def _get_ids(s, dictionary): # 获得源语言序列中的每个单词在字典中的位置
2.         words = s.strip().split()
3.         return [dictionary[START]] + \
4.                 [dictionary.get(w, UNK_IDX) for w in words] + \
5.                 [dictionary[END]]
6.
7.     @provider(init_hook=hook, pool_size=50000)
8.     def process(settings, file_name):
9.         with open(file_name, 'r') as f:
10.            for line_count, line in enumerate(f):
11.                line_split = line.strip().split('\t')
12.                if settings.job_mode and len(line_split) != 2:
13.                    continue
14.                src_seq = line_split[0]

```

```

15.         src_ids = _get_ids(src_seq, settings.src_dict)
16.
17.         if settings.job_mode:
18.             trg_seq = line_split[1]
19.             trg_words = trg_seq.split()
20.             trg_ids = [settings.trg_dict.get(w, UNK_IDX) for w in trg_words]
21.
22.             # 如果任意一个序列长度超过80个单词，在训练模式下会移除这条样本，以防止RNN过深。
23.             if len(src_ids) > 80 or len(trg_ids) > 80:
24.                 continue
25.             trg_ids_next = trg_ids + [settings.trg_dict[END]]
26.             trg_ids = [settings.trg_dict[START]] + trg_ids
27.             yield {
28.                 'source_language_word': src_ids,
29.                 'target_language_word': trg_ids,
30.                 'target_language_next_word': trg_ids_next
31.             }
32.         else:
33.             yield {'source_language_word': src_ids, 'sent_id': [line_count]}

```

注意：由于本示例中的训练数据有3.55G，对于内存较小的机器，不能一次性加载进内存，所以推荐使用 `pool_size` 变量来设置内存中暂存的数据条数。

模型配置说明

数据定义

1. 首先，定义数据集路径和源/目标语言字典路径，并用 `is_generating` 变量定义当前配置是训练模式（默认）还是生成模式。该变量接受从命令行传入的参数，使用方法见[应用命令与结果](#)。

```

1. import os
2. from paddle.trainer_config_helpers import *
3.
4. data_dir = "./data/pre-wmt14" # 数据集路径
5. src_lang_dict = os.path.join(data_dir, 'src.dict') # 源语言字典路径
6. trg_lang_dict = os.path.join(data_dir, 'trg.dict') # 目标语言字典路径
7. is_generating = get_config_arg("is_generating", bool, False) # 配置模式

```

2. 其次，通过 `define_py_data_sources2` 函数从 `dataproducer.py` 中读取数据，并用 `args` 变量传入源/目标语言的字典路径以及配置模式。

```

1. if not is_generating:
2.     train_list = os.path.join(data_dir, 'train.list')
3.     test_list = os.path.join(data_dir, 'test.list')
4. else:
5.     train_list = None
6.     test_list = os.path.join(data_dir, 'gen.list')

```

```

7.
8.  define_py_data_sources2(
9.     train_list,
10.    test_list,
11.    module="dataproducer",
12.    obj="process",
13.    args={
14.        "src_dict_path": src_lang_dict, # 源语言字典路径
15.        "trg_dict_path": trg_lang_dict, # 目标语言字典路径
16.        "is_generating": is_generating # 配置模式
17.    })

```

算法配置

```

1.  settings(
2.     learning_method = AdamOptimizer(),
3.     batch_size = 50,
4.     learning_rate = 5e-4)

```

本教程使用默认的SGD随机梯度下降算法和Adam学习方法，并指定学习率为5e-4。注意：生成模式下的 `batch_size = 50`，表示同时生成50条序列。

模型结构

1. 首先，定义了一些全局变量。

```

1.  source_dict_dim = len(open(src_lang_dict, "r").readlines()) # 源语言字典维度
2.  target_dict_dim = len(open(trg_lang_dict, "r").readlines()) # 目标语言字典维度
3.  word_vector_dim = 512 # dimension of word vector # 词向量维度
4.  encoder_size = 512 # 编码器中的GRU隐层大小
5.  decoder_size = 512 # 解码器中的GRU隐层大小
6.
7.  if is_generating:
8.     beam_size=3 # # 柱搜索算法中的宽度
9.     max_length=250 # 生成句子的最大长度
10.    gen_trans_file = get_config_arg("gen_trans_file", str, None) # 生成后的文件

```

2. 其次，实现编码器框架。分为三步：

2.1 传入已经在 `dataproducer.py` 转换成one-hot vector表示的源语言序列 \mathbf{w} 。

```

1.  src_word_id = data_layer(name='source_language_word', size=source_dict_dim)

```

2.2 将上述编码映射到低维语言空间的词向量 \mathbf{s} 。

```

1.  src_embedding = embedding_layer(
2.     input=src_word_id,

```

```

3.     size=word_vector_dim,
4.     param_attr=ParamAttr(name='_source_language_embedding'))

```

2.3 用双向GRU编码源语言序列，拼接两个GRU的编码结果得到**h**。

```

1.     src_forward = simple_gru(input=src_embedding, size=encoder_size)
2.     src_backward = simple_gru(
3.         input=src_embedding, size=encoder_size, reverse=True)
4.     encoded_vector = concat_layer(input=[src_forward, src_backward])

```

3. 接着，定义基于注意力机制的解码器框架。分为三步：

3.1 对源语言序列编码后的结果（见2.3），过一个前馈神经网络（Feed Forward Neural Network），得到其映射。

```

1.     with mixed_layer(size=decoder_size) as encoded_proj:
2.         encoded_proj += full_matrix_projection(input=encoded_vector)

```

3.2 构造解码器RNN的初始状态。由于解码器需要预测时序目标序列，但在0时刻并没有初始值，所以我们希望对其进行初始化。这里采用的是将源语言序列逆序编码后的最后一个状态进行非线性映射，作为该初始值，即 $\mathbf{c}_0 = \mathbf{h}_T$ 。

```

1.     backward_first = first_seq(input=src_backward)
2.     with mixed_layer(
3.         size=decoder_size,
4.         act=TanhActivation(), ) as decoder_boot:
5.         decoder_boot += full_matrix_projection(input=backward_first)

```

3.3 定义解码阶段每一个时间步的RNN行为，即根据当前时刻的源语言上下文向量 \mathbf{c}_i 、解码器隐层状态 \mathbf{z}_i 和目标语言中第*i*个词 \mathbf{u}_i ，来预测第*i* + 1个词的概率 p_{i+1} 。

- decoder_mem记录了前一个时间步的隐层状态 \mathbf{z}_i ，其初始状态是decoder_boot。
- context通过调用simple_attention函数，实现公式 $\mathbf{c}_i = \sum_j \mathbf{1}^T \mathbf{a}_{ij} \mathbf{h}_j$ 。其中，enc_vec是 \mathbf{h}_j ，enc_proj是 \mathbf{h}_j 的映射（见3.1），权重 \mathbf{a}_{ij} 的计算已经封装在simple_attention函数中。
- decoder_inputs融合了 \mathbf{c}_i 和当前目标词current_word（即 \mathbf{u}_i ）的表示。
- gru_step通过调用gru_step_layer函数，在decoder_inputs和decoder_mem上做了激活操作，即实现公式 $\mathbf{z}_{i+1} = \phi_{\theta}(\mathbf{c}_i, \mathbf{u}_i, \mathbf{z}_i)$ 。
- 最后，使用softmax归一化计算单词的概率，将out结果返回，即实现公式 $p(\mathbf{u}_i | \mathbf{u}_{<i}, \mathbf{x}) = \text{softmax}(\mathbf{W}_s \mathbf{z}_i + \mathbf{b}_z)$ 。

```

1.     def gru_decoder_with_attention(enc_vec, enc_proj, current_word):
2.         decoder_mem = memory(
3.             name='gru_decoder', size=decoder_size, boot_layer=decoder_boot)
4.
5.         context = simple_attention(
6.             encoded_sequence=enc_vec,
7.             encoded_proj=enc_proj,

```

```

8.         decoder_state=decoder_mem, )
9.
10.        with mixed_layer(size=decoder_size * 3) as decoder_inputs:
11.            decoder_inputs += full_matrix_projection(input=context)
12.            decoder_inputs += full_matrix_projection(input=current_word)
13.
14.        gru_step = gru_step_layer(
15.            name='gru_decoder',
16.            input=decoder_inputs,
17.            output_mem=decoder_mem,
18.            size=decoder_size)
19.
20.        with mixed_layer(
21.            size=target_dict_dim, bias_attr=True,
22.            act=SoftmaxActivation()) as out:
23.            out += full_matrix_projection(input=gru_step)
24.        return out

```

4. 训练模式与生成模式下的解码器调用区别。

4.1 定义解码器框架名字，和 `gru_decoder_with_attention` 函数的前两个输入。注意：这两个输入使用 `StaticInput`，具体说明可见[StaticInput文档](#)。

```

1.     decoder_group_name = "decoder_group"
2.     group_input1 = StaticInput(input=encoded_vector, is_seq=True)
3.     group_input2 = StaticInput(input=encoded_proj, is_seq=True)
4.     group_inputs = [group_input1, group_input2]

```

4.2 训练模式下的解码器调用：

- 首先，将目标语言序列的词向量 `trg_embedding`，直接作为训练模式下的 `current_word` 传给 `gru_decoder_with_attention` 函数。
- 其次，使用 `recurrent_group` 函数循环调用 `gru_decoder_with_attention` 函数。
- 接着，使用目标语言的下一个词序列作为标签层 `lbl`，即预测目标词。
- 最后，用多类交叉熵损失函数 `classification_cost` 来计算损失值。

```

1.     if not is_generating:
2.         trg_embedding = embedding_layer(
3.             input=data_layer(
4.                 name='target_language_word', size=target_dict_dim),
5.             size=word_vector_dim,
6.             param_attr=ParamAttr(name='_target_language_embedding'))
7.         group_inputs.append(trg_embedding)
8.
9.         decoder = recurrent_group(
10.            name=decoder_group_name,
11.            step=gru_decoder_with_attention,
12.            input=group_inputs)
13.
14.         lbl = data_layer(name='target_language_next_word', size=target_dict_dim)

```

```
15.     cost = classification_cost(input=decoder, label=lbl)
16.     outputs(cost)
```

4.3 生成模式下的解码器调用：

- 首先，在序列生成任务中，由于解码阶段的RNN总是引用上一时刻生成出的词的词向量，作为当前时刻的输入，因此，使用 `GeneratedInput` 来自动完成这一过程。具体说明可见[GeneratedInput文档](#)。
- 其次，使用 `beam_search` 函数循环调用 `gru_decoder_with_attention` 函数，生成出序列id。
- 最后，使用 `seqtext_printer_evaluator` 函数，根据目标字典 `trg_lang_dict`，打印出完整的句子保存在 `gen_trans_file` 中。

```
1.     else:
2.         trg_embedding = GeneratedInput(
3.             size=target_dict_dim,
4.             embedding_name='_target_language_embedding',
5.             embedding_size=word_vector_dim)
6.         group_inputs.append(trg_embedding)
7.
8.         beam_gen = beam_search(
9.             name=decoder_group_name,
10.            step=gru_decoder_with_attention,
11.            input=group_inputs,
12.            bos_id=0,
13.            eos_id=1,
14.            beam_size=beam_size,
15.            max_length=max_length)
16.
17.         seqtext_printer_evaluator(
18.             input=beam_gen,
19.             id_input=data_layer(
20.                 name="sent_id", size=1),
21.             dict_file=trg_lang_dict,
22.             result_file=gen_trans_file)
23.         outputs(beam_gen)
```

注意：我们提供的配置在Bahdanau的论文[4]上做了一些简化，可参考[issue #1133](#)。

训练模型

可以通过以下命令来训练模型：

```
1.     ./train.sh
```

其中 `train.sh` 的内容为：

```
1.     paddle train \
2.     --config='seqToseq_net.py' \
```

```
3.  --save_dir='model' \  
4.  --use_gpu=false \  
5.  --num_passes=16 \  
6.  --show_parameter_stats_period=100 \  
7.  --trainer_count=4 \  
8.  --log_period=10 \  
9.  --dot_period=5 \  
10. 2>&1 | tee 'train.log'
```

- config: 设置神经网络的配置文件。
- save_dir: 设置保存模型的输出路径。
- use_gpu: 是否使用GPU训练，这里使用CPU。
- num_passes: 设置passes的数量。PaddlePaddle中的一个pass表示对数据集中所有样本的一次完整训练。
- show_parameter_stats_period: 这里每隔100个batch显示一次参数统计信息。
- trainer_count: 设置CPU线程数或者GPU设备数。
- log_period: 这里每隔10个batch打印一次日志。
- dot_period: 这里每个5个batch打印一个点". "。

训练的损失函数每隔10个batch打印一次，您将会看到如下消息：

```
1.  I0719 19:16:45.952062 15563 TrainerInternal.cpp:160] Batch=10 samples=500 AvgCost=198  
.475 CurrentCost=198.475 Eval: classification_error_evaluator=0.737155 CurrentEval:  
classification_error_evaluator=0.737155  
2.  I0719 19:17:56.707319 15563 TrainerInternal.cpp:160] Batch=20 samples=1000  
AvgCost=157.479 CurrentCost=116.483 Eval: classification_error_evaluator=0.698392 Cu  
rrentEval: classification_error_evaluator=0.659065  
3.  .....
```

- AvgCost：从第0个batch到当前batch的平均损失值。
- CurrentCost：当前batch的损失值。
- classification_error_evaluator(Eval)：从第0个评估到当前评估中，每个单词的预测错误率。
- classification_error_evaluator(CurrentEval)：当前评估中，每个单词的预测错误率。

当classification_error_evaluator的值低于0.35时，模型就训练成功了。

应用模型

下载预训练的模型

由于NMT模型的训练非常耗时，我们在50个物理节点（每节点含有2颗6核CPU）的集群中，花了5天时间训练了16个pass，其中每个pass耗时7个小时。因此，我们提供了一个预先训练好的模型（pass-00012）供大家直接下载使用。该模型大小为205MB，在所有16个模型中有最高的BLEU评估值26.92。下载并解压模型的命令如下：


```
1. cd pretrained
2. ./wmt14_model.sh
```

应用命令与结果

可以通过以下命令来进行法英翻译：

```
1. ./gen.sh
```

其中 `gen.sh` 的内容为：

```
1. paddle train \
2. --job=test \
3. --config='seqToseq_net.py' \
4. --save_dir='pretrained/wmt14_model' \
5. --use_gpu=true \
6. --num_passes=13 \
7. --test_pass=12 \
8. --trainer_count=1 \
9. --config_args=is_generating=1,gen_trans_file="gen_result" \
10. 2>&1 | tee 'translation/gen.log'
```

与训练命令不同的参数如下：

- job：设置任务的模式为测试。
- save_dir：设置存放预训练模型的路径。
- num_passes和test_pass：加载第 i 轮[`test_pass, num_passes - 1`]轮的模型参数，这里只加载 `data/wmt14_model/pass-00012`。
- config_args：将命令行中的自定义参数传递给模型配置。`is_generating=1` 表示当前为生成模式，`gen_trans_file="gen_result"` 表示生成结果的存储文件。

翻译结果请见[效果展示](#)。

BLEU评估

BLEU(Bilingual Evaluation understudy)是一种广泛使用的机器翻译自动评测指标，由IBM的watson研究中心于2002年提出[5]，基本出发点是：机器译文越接近专业翻译人员的翻译结果，翻译系统的性能越好。其中，机器译文与人工参考译文之间的接近程度，采用句子精确度（precision）的计算方法，即比较两者的n元词组相匹配的个数，匹配的个数越多，BLEU得分越好。

[Moses](#) 是一个统计学的开源机器翻译系统，我们使用其中的 [multi-bleu.perl](#) 来做BLEU评估。下载脚本的命令如下：

```
1. ./moses_bleu.sh
```

BLEU评估可以使用 `eval_bleu` 脚本如下，其中FILE为需要评估的文件名，BEAMSIZE为柱宽度，默认使用 `data/wmt14/gen/ntst14.trg` 作为标准的翻译结果。

```
1. ./eval_bleu.sh FILE BEAMSIZE
```

本教程的具体命令如下：

```
1. ./eval_bleu.sh gen_result 3
```

您会在屏幕上看到：

```
1. BLEU = 26.92
```

总结

端到端的神经网络机器翻译是近几年兴起的一种全新的机器翻译方法。本章中，我们介绍了NMT中典型的“编码器-解码器”框架和“注意力”机制。由于NMT是一个典型的Seq2Seq (Sequence to Sequence , 序列到序列) 学习问题，因此，Seq2Seq中的query改写 (query rewriting)、摘要、单轮对话等问题都可以用本教程的模型来解决。

参考文献

1. Koehn P. [Statistical machine translation](#) [M]. Cambridge University Press, 2009.
2. Cho K, Van Merriënboer B, Gulcehre C, et al. [Learning phrase representations using RNN encoder-decoder for statistical machine translation](#) [C]//Proceedings of the 2014 Conference on Empirical Methods in Natural Language Processing (EMNLP), 2014: 1724-1734.
3. Chung J, Gulcehre C, Cho K H, et al. [Empirical evaluation of gated recurrent neural networks on sequence modeling](#) [J]. arXiv preprint arXiv:1412.3555, 2014.
4. Bahdanau D, Cho K, Bengio Y. [Neural machine translation by jointly learning to align and translate](#) [C]//Proceedings of ICLR 2015, 2015.
5. Papineni K, Roukos S, Ward T, et al. [BLEU: a method for automatic evaluation of machine translation](#) [C]//Proceedings of the 40th annual meeting on association for computational linguistics. Association for Computational Linguistics, 2002: 311-318.



本教程由PaddlePaddle创作，采用[知识共享 署名-非商业性使用-相同方式共享 4.0 国际 许可协议](#)进行许可。