

# 基于 RNN Attention 的 生成式聊天模型

邢超 (Chao Xing)

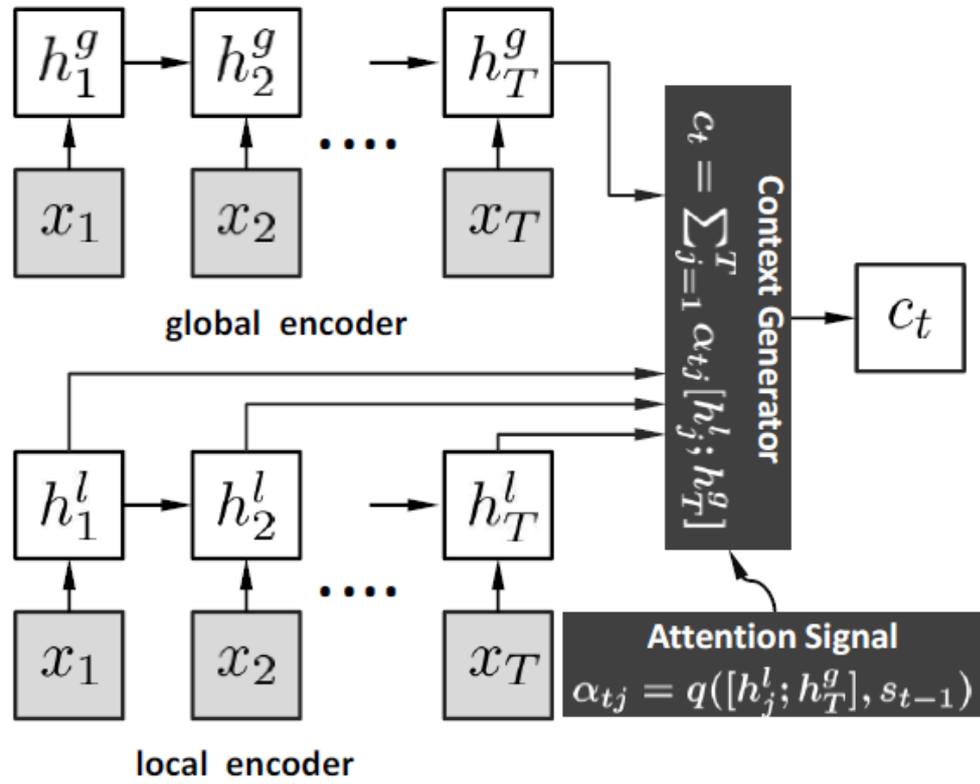
白子薇 (Ziwei Bai)

2017/1/11

简介 .....	3
数据准备 .....	4
使用说明 .....	5
数据准备脚本 .....	5
训练脚本 .....	5
测试脚本 .....	5
实验结果展示 .....	6

## 简介

本项目基于<sup>[1]</sup>进行序列对序列网络进行复现。使用的模型如图 1 所示。



图一 序列对序列模型图

## 数据处理

### 处理规则

数据的收集与初步处理请参考《聊天数据收集与清理建议》

为了方便起见,我们将对话对中的 Q 或 post 统一称为 post, A 或 response 统一称为 response。

为了使数据更加适合本模型,我们只保留符合以下某一条件的数据,并按条件不同分组:

- post 和 response 的词数量都大小 3, 小于等于 5
- post 和 response 的词数量都大小 5, 小于等于 10
- post 的词数量大于 10, 小于等于 15; response 的词的数量大于 5, 小于等于 10
- post 的词数量大于 15, 小于等于 20; response 的词的数量大于 5, 小于等于 10

### 数据处理程序

#### 程序

```
bucket.py
```

#### 运行示例

```
python buckey.py ../data/post ../data/response
```

## 使用说明

### 数据准备脚本

- 程序名称
  - run\_prepare.sh
- 功能
  - 对数据进行分 bucket
  - 每组随机采样 500 句作为测试数据
  - 对训练数据进行补 MASK，根据之前对数据的分组，将每组数据补到该组数据词数量最大值。
  - 将通过词表示的补 MASK 后数据转化为通过词的 index 表示的数据
- 备注
  - 建议训练时，sample 100w 句进行训练

### 训练脚本

- 程序名称
  - run\_train.sh
- 功能
  - 训练模型并保存
- 使用示例
  - sh run\_train.sh 0
  - 参数 '0' 代表实验运行指定的 GPU

### 测试脚本

- 程序名称
  - run\_test.sh
- 功能
  - 根据训练得到的模型进行测试，测试模型在不同组数据下的表现，并保存结果
- 使用示例
  - sh run\_test.sh 0
  - 参数 '0' 代表实验运行指定的 GPU

## 核心程序说明

本次提交的 python 核心脚本分别为: `simple_train_process.py`, `simple_test_process.py`, `toolbox`。

其中:

`rnn_local.py` 为训练脚本, 参数为:

-post:接受训练数据的 post 文件, 文件内容是 post 数据中每个词的 index

-response: 接受训练数据的 response 文件, 文件内容是 response 数据中每个词的 index

-vec:接受一个词向量矩阵,词向量需要进行 mean variance normalization

-dict: 接受一个表示词和 index 对应关系的字典, 文件类型是.pkl 文件, 通过 cPickle 库读取。

-output: 输出模型存放的位置以及名称。

-lr: 学习速率

-batch:每个 batch 的大小。本文中所述的 RNN Attention 模型, 采用 mini-batch 方法更新参数。

-epoch 模型训练的轮数。

-encoder: 编码器的大小

-decoder: 解码器的大小

使用示例:

```
python rnn_local.py -post ../data/data_all/post_filter -res ../data/data_all/response_filter
  -dict ../data/data_all/recut_word_to_index.pkl -vec ../data/data_all/word_matrix.npy -e
ncoder 512 -decoder 512 -epoch 101 -batch 500 -show 1 -output ../model-new/mode
l -lr 1e-5 >> run.log 2>> err.log
```

`rnn_local.predict.py` 为训练脚本, 参数为:

-post:接受训练数据的 post 文件, 文件内容是分好词的 post 数据

-res: 接受训练数据的 response 文件, 文件内容是分好词的 response 数据

-post-len:post 补齐 mask 后的长度

-res-len: response 补 end 和 mask 后的长度

-dict: 接受一个表示词和 index 对应关系的字典, 文件类型是.pkl 文件, 通过 cPickle 库读取。

-result: 输出结果存放位置及名称

-encoder: 编码器的大小

-decoder: 解码器的大小

-model: 预测所使用的模型

使用示例:

```
python ../src/rnn_local.predict.py -post /work5/baizw/Chatting_Model_RNN/data/data_n
ew/P2R_v1_sample/P2R_post_0 -res /work5/baizw/Chatting_Model_RNN/data/data_new/
P2R_v1_sample/P2R_response_0 -post-len 6 -res-len 6 -vec ../data/word_matrix.npy -di
ct ../data/recut_word_to_index.pkl ../data/recut_index_to_word.pkl -model /work4/chao
s/RNN-ATT/model_1-10/$model -encoder 512 -decoder 512 -result ../result2/0_$model
```

## 实验设置

### 词向量准备

我们选择谷歌的开源工具 `word2vec` 对语料进行训练。

- 语料分布
  - 汇联 1.7G 语料，对话数据
- 分词
  1. 使用 `jieba` 分词工具，通过结巴分词中的默认词表对数据进行分词，允许使用 HMM 模型,得到一个词表 `vocab1`
  2. 将数据按字分开，得到另一个词表 `vocab2`
  3. 将 `vocab1` 中词频前 3W 的词，和 `vocab2` 中所有词合在一起，构成一个新的词表 `vocab`
  4. 使用 `jieba` 分词工具，根据 `vocab` 再次对数据进行分词，不允许使用 HMM 模型，我们训练词向量及后续模型的训练和测试都依赖于这一套分词结果。（该词表已提交）

- `word2vec` 训练参数设置如下表所示

Window	Iteration	Cbow	Size	Sample	negative	hs	binary
8	20	0 (skip-gram)	200	1e-4	5	0	0

### 参数设置

encoder:512

decoder:512

## 实验结果展示(model-290)

ID	Post	生成的 response	标准 response
1	但我悄悄告诉你	我再也不相信爱情了	那都不重要
2	安全感是自己给的	这句话说的真好	没安全感是吧
3	说你爱我	我爱你你永远爱	我我我我说不出口
4	这套戏要看有惊喜	好厉害惊喜你都知道	早前看了宣传片还不错
5	所以说不要惹女人 女人是很可怕滴	我靠这句话太有水平了	话说有几个女淫能做到
6	他在保护自己的朋友	看来他不够爱你啊	看来他不够爱你啊
7	中午啦吃饭吃饭吃饭吃饭	吃饭的孩纸们伤不起啊	这是早饭还是昨天的晚饭
8	你太聪明了这都被你发现了不简单	我优点多着呢这只是冰山一角	还行还行排的上前三
9	据说喜欢在句末发这个的孩子有一颗玲珑婉转的心	我也是我也是	我有一颗玲珑心有谁懂得珍惜
10	你一生中看过的小说中最希望哪部能改编为电影	爱神的黑白羽翼盗墓笔记	太多了眼下希望是百年孤独
11	这是一个从上厕所到坐电车都有准备的女人	这桀骜的女人真危险	还真方便佩服这位大姐太有才啦
12	无论我此时是多么的彷徨迷茫最终我都要过上自己想要的生活	我也要这样的生活	感觉讲到心里去了你们呢
13	我就喜欢看哈林还没转身前纠结像吃了大便的那种表情	我也喜欢我也喜欢	怎么办又想看快男又想看好声音
14	罗斯训练中的这个球有让你更怀念他的折叠暴扣吗	我也想训练求方法	弹跳明显弱了一成啊难道保守实力
15	一起来填空每天摸着大肚肚我最想跟宝宝说的话是	我想我们应该会成为好朋友的	快点跟我变身滚出来

- [1] Shang L, Lu Z, Li H. Neural responding machine for short-text conversation[J]. arXiv preprint arXiv:1503.02364, 2015.