

# 2021年10月32日总结

## Deep Learning

### 1. 预训练模型:

- 往往用于NLP，几乎在所有的NLP任务中都得到了最好的结果。预训练模型就意味着把人类的语言知识，先学了一个东西，然后再代入到某个具体任务，就顺手了。预训练模型是迁移学习的一种应用，其中最关键的三个技术分别是Transformer、自监督学习和微调。

### 2. 无监督学习:

- 无监督学习广泛采用的方式是自动编码器，编码器将样本映射到隐层向量，解码器将这个隐层向量映射回样本空间，我们期待网络的输入和输出可以保持一致（理想情况，无损重构），同时隐层向量的维度远远小于输入样本的维度，以此达到了降维的目的，利用学习到的隐层向量（代替原始的输入样本）再进行聚类任务时将更加的简单高效。
- 如何学习隐层向量，称之为Representation Learning（表征/表示学习），但是这样简单的编码解码过程不能够学习到更多的语义特征。

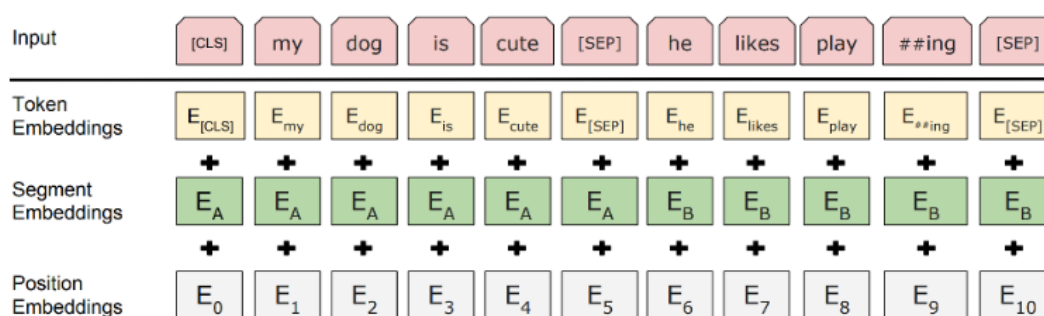
### 3. 自监督学习:

- 自监督学习主要是利用辅助任务（pretext）从大规模的无监督数据中挖掘自身的监督信息，通过这种构造的监督信息对网络进行训练，从而可以学习到对下游任务有价值的表征。（也就是说自监督学习的监督信息不是人工标注的，而是算法在大规模无监督数据中自动构造监督信息，来进行监督学习或训练）。
- 自监督学习的主要方法之一就是基于上下文(Context Based)

### 4. BERT

- 分为bert\_base和bert\_large，前者使用了12个encoder（与Transformer模型中的相同），后者使用了24个encoder。
- 输入部分：bert中的 input = token emb + segment emb + position emb，下图是个例子：

### BERT input representation



其中的CLS、SEP都是因为NSP任务（处理两个句子的关系，是个二分类任务），SEP用于隔开句子；CLS在预训练之后的输出向量用于做NSP二分类任务，但CLS向量不能代表整个句子的语义信息。

- 如何做预训练:

#### MLM-掩码语言模型

- 对句子进行一个mask，掩盖一些地方，比如“我爱吃饭-> 我爱mask饭”，然后预测mask，这样打破了文本，让模型进行文本重建。
- mask概率，bert是随机mask15%的单词，然后对mask的词，10%替换成其它，10%原封不动，80%替换成mask。

NSP:



1. python导入自己写的模块出错：在pycharm中，每个project都有一个sources root，这是寻找文件的默认路径，如果调用模块不再sources root 中就会报错，将模块的上层目录设定为sources root就可以了。
2. python的tqdm是一个快速，可扩展的进度条类。
3. 在使用conda创建新的虚拟环境时，要关闭网络代理，否则会报错。
4. tensorflow-gpu相当于是tensorflow的进阶版本，当没有配置好cuda和cudnn，或者没有gpu时，就会像tensorflow一样，只使用CPU运行。

## 问题

1. 这个是怎么运行的？

### Fine-Tune

We fine-tuned the model on 4\*P40 GPUs.

```
cd code2n1

lang=php #programming language
lr=5e-5
batch_size=64
beam_size=10
source_length=256
target_length=128
data_dir=./data/code2n1/CodeSearchNet
output_dir=model/$lang
train_file=$data_dir/$lang/train.jsonl
dev_file=$data_dir/$lang/valid.jsonl
eval_steps=1000 #400 for ruby, 600 for javascript, 1000 for others
train_steps=50000 #20000 for ruby, 30000 for javascript, 50000 for others
pretrained_model=microsoft/codebert-base #Roberta: roberta-base

python run.py --do_train --do_eval --model_type roberta --model_name_or_path $pretrained_model --train_filename $train_file --dev_file
```

- 2.