

# BERT

- [BERT](#)
  - [Intro:](#)
  - [Model architecture](#)
  - [Input/Output](#)
  - [Principle:](#)
    - [Embeddings](#)
    - [Pre-training](#)
    - [Fine-tuning](#)
  - [Comparison](#)
  - [代码调用](#)

## Intro:

- **BERT(Bidirectional Encoder Representation from Transformers)**是2018年10月由Google AI研究院发表的一篇名为《[BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding](#)》的文章提出的一种预训练模型，该模型在机器阅读理解顶级水平测试SQuAD1.1中表现出惊人的成绩: 全部两个衡量指标上全面超越人类，并且在11种不同NLP测试中创出SOTA表现，包括将GLUE基准推高至80.4% (绝对改进7.6%)，MNLI准确度达到86.7% (绝对改进5.6%)，成为NLP发展史上的里程碑式的模型成就。
- BERT的网络架构使用的是《Attention is all you need》中提出的多层Transformer结构。其最大的特点是抛弃了传统的RNN和CNN，通过Attention机制将任意位置的两个单词的距离转换成1，有效的解决了NLP中棘手的长期依赖问题。Transformer的结构在NLP领域中已经得到了广泛应用。

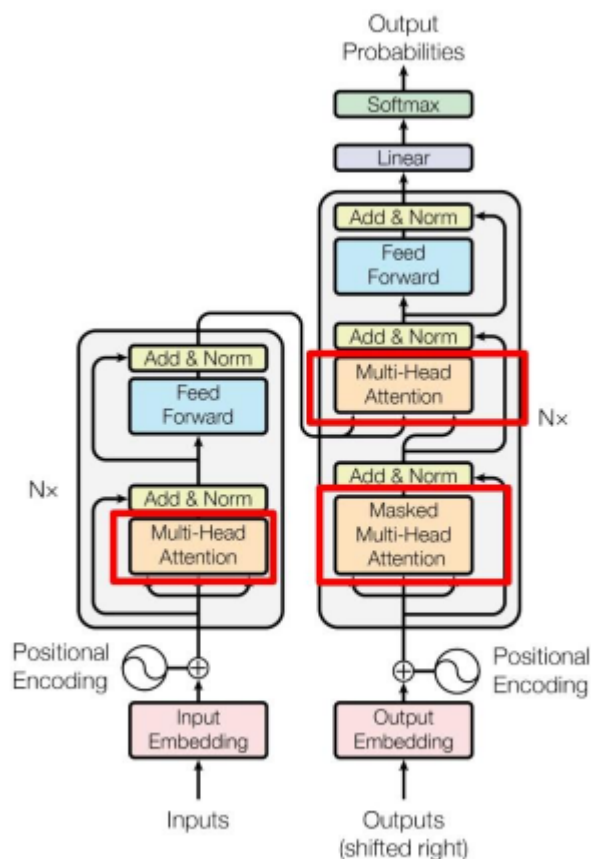


fig1. Transformer Architecture

## Model architecture

简单来讲，BERT就是由多层bidirectional Transformer encoder堆叠而成。Bidirectional指的是在MLM训练任务中利用上下文token做预测训练。

	L(# layers)	H(hidden size)	A(# heads)	Total Parameters
BERT(base)	12	768	12	110M
BERT(large)	24	1024	16	340M

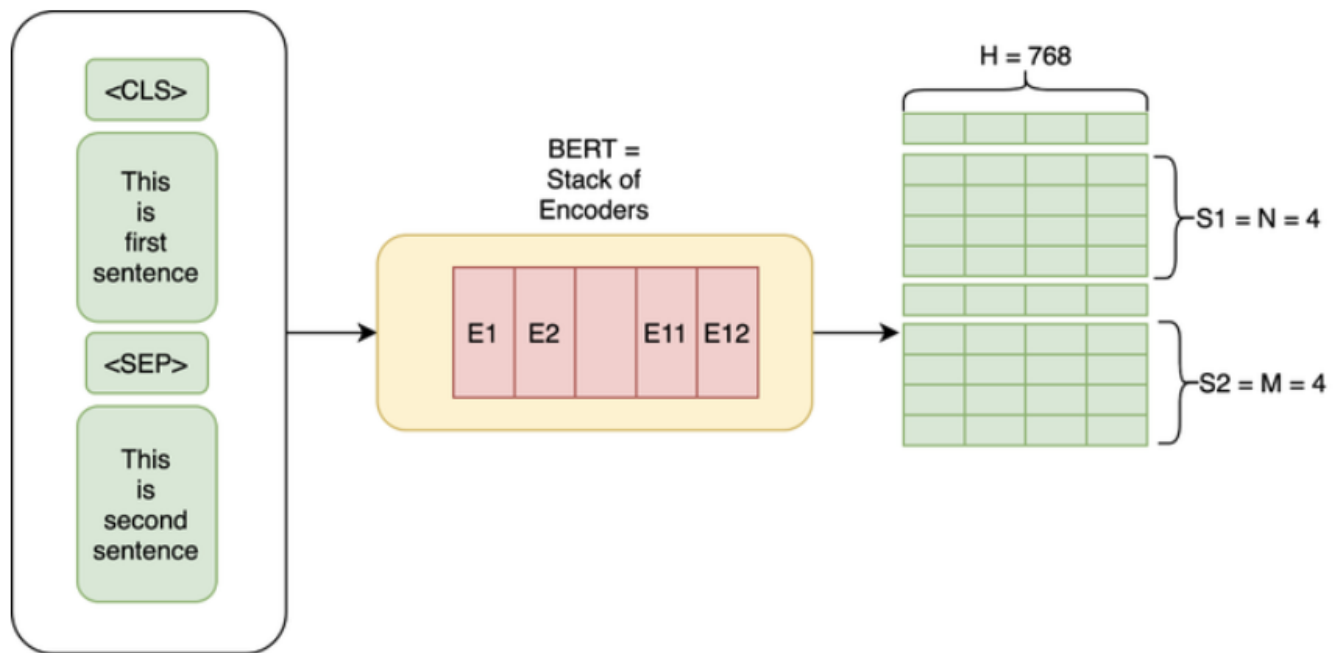


fig2.BERT simple architecture

## Input/Output

- One sentence or a pair of sentences(e.g.<Question, Answer>)
- 每个sequence的第一个token是[CLS]，用于分类，对于非分类模型，该符号可以省去。
- 每个sequence的最后一个token是[SEP]，表示分句符号，用于断开输入语料中的两个句子
- The final hidden vector of [CLS] token as  $C \in \mathbb{R}^H$ , the final hidden vector for the  $i^{th}$  input token is  $T_i \in \mathbb{R}^H$ . C可以作为整句话的语义表示，从而用于下游的分类任务等。[CLS]本身没有语义，因而经过12层self-attention相比其他词更能表示整句话的语义。

## Principle:

BERT在嵌入层以及预训练任务和模型微调方面有自己的独特创新，下面我们将分别介绍BERT如何对input做embedding以及BERT是怎么被预训练的，还有我们可能更加关心的如何运用模型。

## Embeddings

Input的embeddings是由词嵌入(token embeddings)、段落嵌入(segment embedding)以及位置嵌入(position embeddings)相加而得的。

- Token Embeddings: 词向量，通过建立字向量将每个word转化成一个一维向量，特别的英文单词会做更细粒度的切分(playing = play + ##ing)。

- 输入文本后，tokenizer会对文本进行tokenization处理，两个特殊的token会插入在文本开头[CLS]和文本结尾[SEP].
- Token Embeddings层会将每个word转化为768维向量
- Segment Embeddings:用于区别两种句子，bert能处理句子对语义相似任务，当两个句子被简单拼接后送入模型中，bert利用该embeddings区分两个句子
  - segment embeddings层有两种向量表示，前一个向量把0赋值给第一个句子的各个token，后一个向量把1赋给第二个句子的各个token
  - 文本分类任务只有一个句子，那么segment embeddings就全是0
- Position Embeddings:位置向量，与transformer中的三角函数固定编码不同，随机初始化后通过数据训练学习的，不仅可以标记位置还可以学习到这个位置有什么用

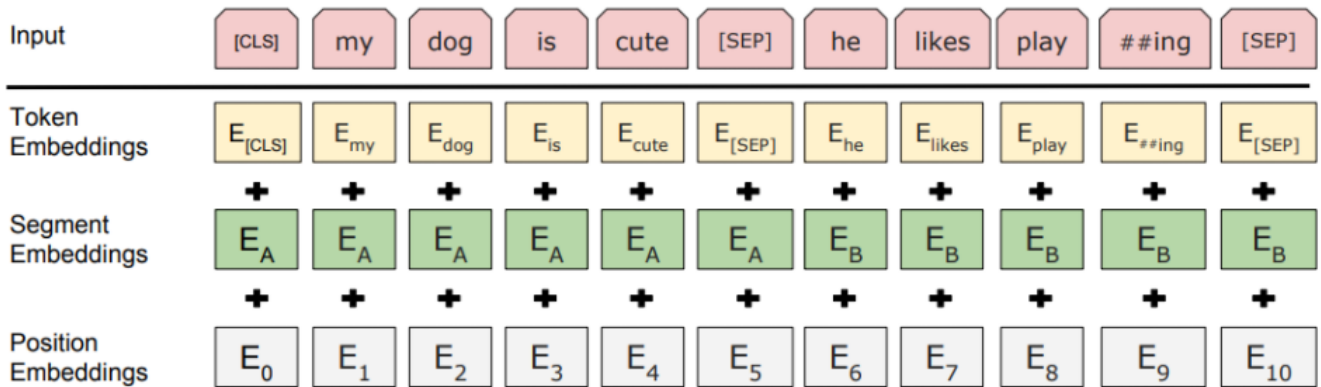


fig3.  $Input\ Embeddings = Token\ Embeddings + Segment\ Embeddings + Position\ Embeddings$

## Pre-training

摒弃传统利用left-to-right language model和right-to-left language model 来预训练bert，而是用两个特殊的unsupervised任务来初始化参数：MLM，NSP。预训练采用BooksCorpus(800M words) and English Wikipedia(2500M words)。

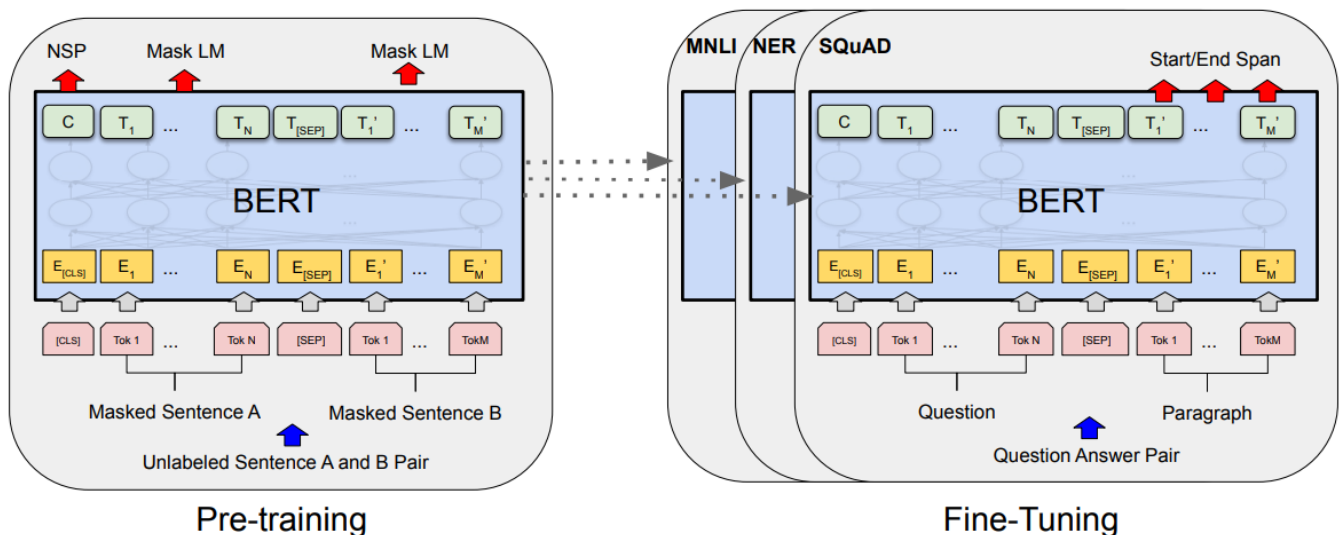


fig4. Two stages of BERT

- Masked LM:利用以一定的机制掩盖input中的target tokens的方式训练模型利用上下文推理预测被掩盖词的能力(类似完型填空)
  - sequence的15%会被选择为target tokens，对每个被选择的位置上的token，通过一系列实验测试得到最为理想的mask strategy:
    - 80%的概率mask target token，即tokenize为[MASK]

- 10%的概率不mask该target token, 即不对input不做任何修改
  - 10%的概率替换该target token为其他随机token
- 训练过程就是通过用上述策略修改input后, 将final hidden vector中对应被mask的token  $T_i$  扔到基于词汇表的softmax函数里进行预测, 训练采用cross entropy loss
- 为什么要用这样的策略:如果句子中的target tokens被100% masked的话, fine-tuning过程中模型就会碰到没有见过的单词, 而加入随机token和保持不变的意义是为了缓解上述问题
- 该训练任务的缺点是有时候会随机mask一些本具有连续意义的词, 使得模型不容易学得词的语义信息, 这在google后续发布的BERT-WWM中得到解决

- 80% of the time: Replace the word with the [MASK] token, e.g., my dog is hairy → my dog is [MASK]
- 10% of the time: Replace the word with a random word, e.g., my dog is hairy → my dog is apple
- 10% of the time: Keep the word unchanged, e.g., my dog is hairy → my dog is hairy. The purpose of this is to bias the representation towards the actual observed word.

*fig5.Masking strategy*

- Next Sentence Prediction: 许多重要的下游任务, 如问答(QA)和自然语言推理(NLI), 都是基于对两个句子之间关系的理解, 而语言建模并没有直接捕捉到这些关系。为了训练一个理解句子关系的模型, 我们预先训练了一个可以从任何单语语料库中生成的二值化下一个句子预测任务。判断句子B是否是句子A的下文。如果是的话输出'IsNext', 否则输出'NotNext'
  - 训练数据的生成方式是从平行语料中随机抽取的连续两句话, 其中50%保留抽取的两句话, 它们符合IsNext关系, 另外50%的第二句话是随机从预料中提取的, 它们的关系是NotNext的。
  - 利用之前提到的输出层[CLS], 用于做分类任务
  - 在此后的研究 (论文《Crosslingual language model pretraining》等) 中发现, NSP任务可能并不是必要的, 消除NSP损失在下游任务的性能上能够与原始BERT持平或略有提高。针对这一点, 后续的RoBERTa、ALBERT、spanBERT都移去了NSP任务。

```

Input = [CLS] the man went to [MASK] store [SEP]
        he bought a gallon [MASK] milk [SEP]

Label = IsNext

Input = [CLS] the man [MASK] to the store [SEP]
        penguin [MASK] are flight ##less birds [SEP]

Label = NotNext
  
```

*fig6.Example input and corresponding output in NSP task*

Fine-tuning

微调基于预训练好的模型，对具体的下游任务，进行参数的微调。该步骤体现了预训练的强大作用，它使得预训练大模型得以在相关的各种下游任务中高效地取得理想表现。

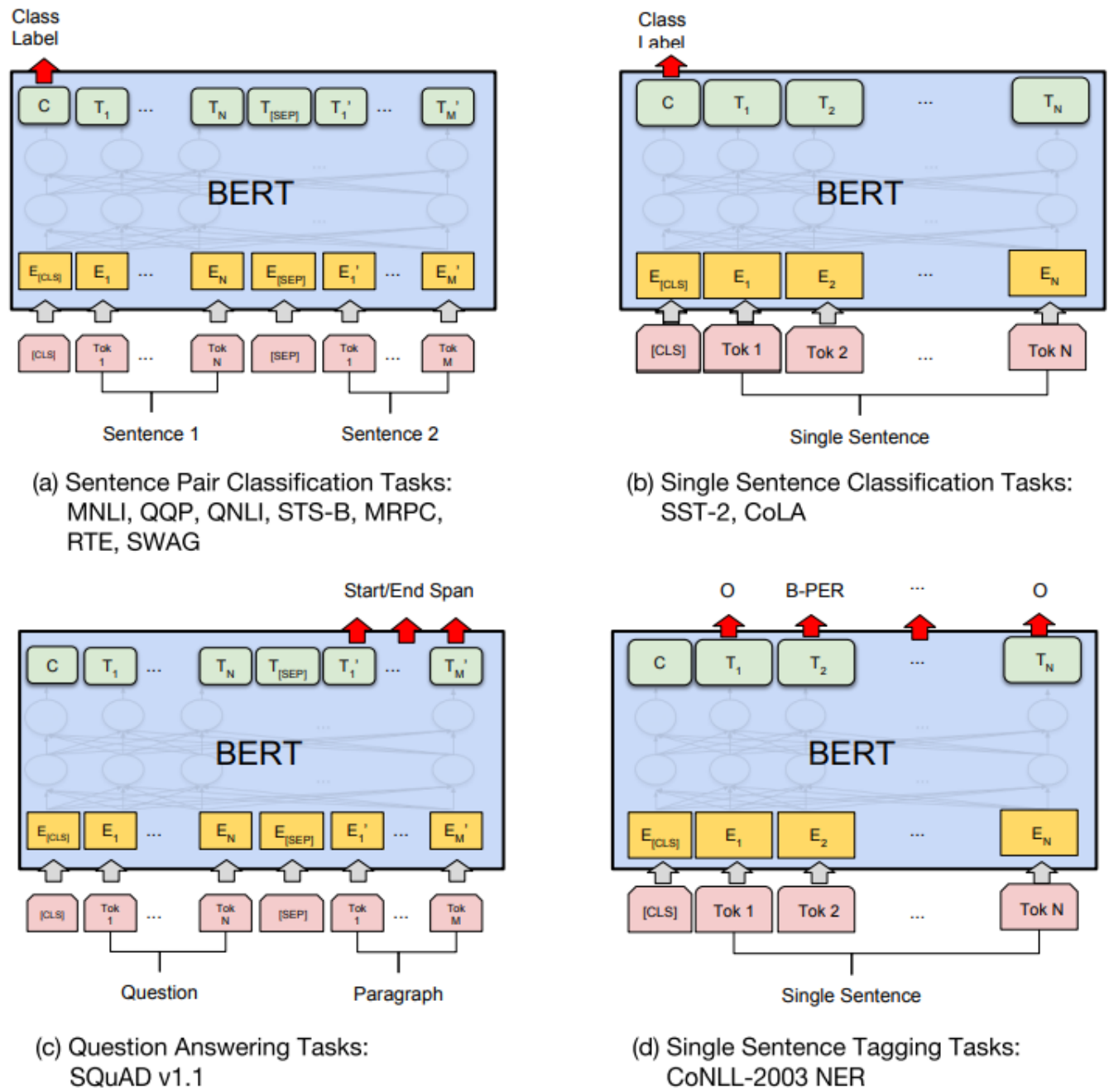


fig7.Fine tuning tasks

模型出处在上方的NLP的四个领域的11个经典benchmark实验中进行测试，均取得显著提升，达到SOTA水平

# Comparison

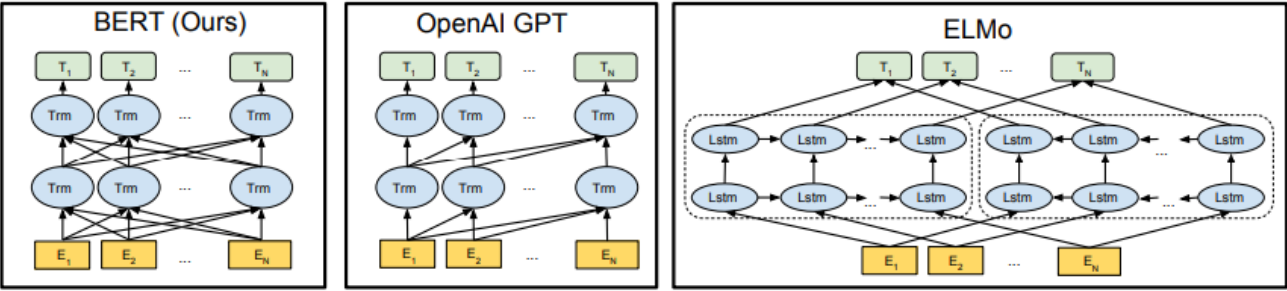


Figure 3: Differences in pre-training model architectures. BERT uses a bidirectional Transformer. OpenAI GPT uses a left-to-right Transformer. ELMo uses the concatenation of independently trained left-to-right and right-to-left LSTMs to generate features for downstream tasks. Among the three, only BERT representations are jointly conditioned on both left and right context in all layers. In addition to the architecture differences, BERT and OpenAI GPT are fine-tuning approaches, while ELMo is a feature-based approach.

分析比较BERT，OpenAI GPT 和 ELMo

	BERT	OpenAI GPT	ELMo
架构	bidirectional transformer	left-to-right transformer	left-to-right and right-to-left LSTM
方法	fine-tuning	fine-tuning	feature-based

BERT与近年大热的GPT的区别：

- BERT利用transformer encoder侧网络，利用上下文信息编码token，GPT利用decoder侧网络，left-to-right的架构使得它适用于文本生成。
- GPT仅在BooksCorpus语料库上进行训练，而BERT在BooksCorpus和Wikipedia上进行训练
- GPT仅在fine-tuning时引入[CLS]和[SEP]，而BERT在预训练阶段学习[SEP],[CLS]以及A/B embeddings
- GPT was trained for 1M steps with a batch size of 32000 words, BERT was trained fro 1M steps with a batch size of 128000 words.
- GPT在所有fine-tuning实验中都统一采用 $5e^{-5}$ 的学习率，而BERT需要根据情况选取最优的

## Conclusion

BERT基于Transformer架构，采用多头自注意力机制进行信息挖掘，并利用大量无标签语言数据进行预训练，通过自学习捕捉语言模式。预训练的的目的是让模型学习到通用的语言表示，这些表示可以用于多种自然语言处理任务，如情感分析、命名实体识别、问答等。在预训练完成后，BERT模型需要进行微调以适配特定的任务，该过程将预训练模型输入特定任务的数据，并通过一套较小的网络结构进行微调。利用预训练好的BERT(如加载Hugging Face仓库)，我们可以很方便的在我们所需要的任务上运用强大的BERT。



## 代码调用（以文本分类为例）

Hugging Face为我们提供了方便的BERT接口，我们在下面的BERT运用示例中仅需要调用的是，Transformers 库内的BertModel和BertTokenizer模块。尚未安装对应第三方库的，运行 `pip install transformers` 即可。

对BERT相关模型，或者感兴趣如何配置自定义的BERT框架以及自己预训练BERT，可以参考Huggface的[Bert模块配置详情](#)。

下面的代码块展示了，BertModel和BertTokenizer模块中一些方法的运用

```
1 from transformers import BertModel, BertTokenizer
2 BERT_PATH = 'bert-base-cased' #区别大小写的Bert_case, 可以指定预训练的Bert模型配置文件路径
3 tokenizer = BertTokenizer.from_pretrained(BERT_PATH)
4 print(tokenizer.tokenize('My dog is cute, he likes playing.'))
5 #['My', 'dog', 'is', 'cute', ',', 'he', 'likes', 'playing', '.']
6 bert = BertModel.from_pretrained(BERT_PATH)
7
8 example_text = 'I will watch Memento tonight'
9 bert_input = tokenizer(example_text, padding='max_length', #padding用于将每个sequence填充到
    指定的最大长度
10                        max_length = 10,                    #max_length用于指定每个sequence的最
    大长度
11                        truncation=True,                    #truncation为true代表每个序列中超过最大长
    度的标记将被截断
12                        return_tensors="pt") #return_tensors设置返回的张量类型, 'pt'-
    >pytorch, 'ft'->tensorflow
13 # ----- bert_input -----
14 print(tokenizer.tokenize(example_text))
15 #['I', 'will', 'watch', 'Me', '##mento', 'tonight']
16
17 print(bert_input['input_ids']) #每个token的id, 可以用decode来解码id得到token
18 #tensor([[ 101,   146,  1209,  2824,  2508, 26173,  3568,   102,     0,     0]])
19
20 print(tokenizer.encode(example_text)) #预训练词汇表中的索引编号
21 #[101, 146, 1209, 2824, 2508, 26173, 3568, 102]
22
23 print(tokenizer.decode(bert_input.input_ids[0]))
24 #[CLS] I will watch Memento tonight [SEP] [PAD] [PAD]
25
26 print(bert_input['token_type_ids']) #用于标识每个token哪个sequence, 也就是position
    embeddings
27 #tensor([[0, 0, 0, 0, 0, 0, 0, 0, 0, 0]])
28
29 print(bert_input['attention_mask']) #如果为0则为[PAD], 否则代表真实word或[CLS][SEP],
    用于标记哪些位置是实际输入哪些是填充输入, 防止[PAD]参与attention机制
30
```

下面我们以BBC新闻文本分类任务为例，展示BERT模型的一个简单用法，我们利用BERT模型学习表征，并在BERT输出处接一个简单的线性全连接层来做分类任务，最后能达到不错的分类效果。

```
1 #数据集划分
2 data_train, data_valid= np.split(df_train.sample(frac=1, random_state=42),
    [int(0.8*len(df_train))], axis=0)
3
```

```

4  #构建数据集类，以便后续采用DataLoader
5  tokenizer = BertTokenizer.from_pretrained('bert-base-cased')
6  labels = {'business':0,
7           'entertainment':1,
8           'sport':2,
9           'tech':3,
10          'politics':4
11          }
12
13  class Dataset(torch.utils.data.Dataset):
14      def __init__(self, df):
15          self.labels = [labels[label] for label in df['category']]
16          self.texts = [tokenizer(text,
17                               padding='max_length',
18                               max_length = 512,
19                               truncation=True,
20                               return_tensors="pt")
21                        for text in df['text']]
22
23      def __len__(self):
24          return len(self.labels)
25
26      def get_batch_labels(self, idx):
27          # Fetch a batch of labels
28          return np.array(self.labels[idx])
29
30      def get_batch_texts(self, idx):
31          # Fetch a batch of inputs
32          return self.texts[idx]
33
34      def __getitem__(self, idx):
35          batch_texts = self.get_batch_texts(idx)
36          batch_y = self.get_batch_labels(idx)
37          return batch_texts, batch_y
38
39  #构建基于Bert的分类器，根据前面介绍，我们简单的在bert的final hidden layer之后加一个线性全连接网络，利用hidden layer的[CLS]作为线性网络的输入即可用于分类任务
40  class BertClassifier(nn.Module):
41      def __init__(self, dropout=0.5):
42          super(BertClassifier, self).__init__()
43          self.bert = BertModel.from_pretrained('bert-base-cased')
44          self.dropout = nn.Dropout(dropout)
45          self.linear = nn.Linear(768, 5) #将768维语意表征向量映射到5维的类别概率向量上
46          self.relu = nn.ReLU()
47
48      def forward(self, input_id, mask):
49          _, pooled_output = self.bert(input_ids= input_id,
50                                     attention_mask=mask, return_dict=False)
51          dropout_output = self.dropout(pooled_output)
52          linear_output = self.linear(dropout_output)
53          final_layer = self.relu(linear_output)
54          return final_layer
55
56  #训练
57  def train(model, train_data, val_data, learning_rate, epochs):
58      # 通过Dataset类获取训练和验证集

```



```

58 train, val = Dataset(train_data), Dataset(val_data)
59 # DataLoader根据batch_size获取数据，训练时选择打乱样本
60 train_dataloader = torch.utils.data.DataLoader(train, batch_size=2, shuffle=True)
61 val_dataloader = torch.utils.data.DataLoader(val, batch_size=2)
62 # 判断是否使用GPU
63 use_cuda = torch.cuda.is_available()
64 device = torch.device("cuda" if use_cuda else "cpu")
65 # 定义损失函数和优化器
66 criterion = nn.CrossEntropyLoss()
67 optimizer = Adam(model.parameters(), lr=learning_rate)
68
69 if use_cuda:
70     model = model.cuda()
71     criterion = criterion.cuda()
72 # 开始进入训练循环
73 for epoch_num in range(epochs):
74     # 定义两个变量，用于存储训练集的准确率和损失
75     total_acc_train = 0
76     total_loss_train = 0
77     model.train() #设置模型为训练模式
78     # 进度条函数tqdm
79     for train_input, train_label in tqdm(train_dataloader):
80         train_label = train_label.type(torch.LongTensor).to(device)
81         mask = train_input['attention_mask'].to(device)
82         input_id = train_input['input_ids'].squeeze(1).to(device)
83
84         optimizer.zero_grad()#清空残留梯度
85     # 通过模型得到输出
86     output = model(input_id, mask)
87     # 计算损失
88     batch_loss = criterion(output, train_label)
89     total_loss_train += batch_loss.item()
90     # 计算精度
91     acc = (output.argmax(dim=1) == train_label).sum().item()
92     total_acc_train += acc
93     # 模型更新
94     batch_loss.backward()
95     optimizer.step()
96     # ----- 验证模型 -----
97     # 定义两个变量，用于存储验证集的准确率和损失
98     total_acc_val = 0
99     total_loss_val = 0
100    model.eval() #设置模型为验证模式
101    # 不需要计算梯度
102    with torch.no_grad():
103        # 循环获取数据集，并用训练好的模型进行验证
104        for val_input, val_label in val_dataloader:
105            # 如果有GPU，则使用GPU，接下来的操作同训练
106            val_label = val_label.type(torch.LongTensor).to(device)
107            mask = val_input['attention_mask'].to(device)
108            input_id = val_input['input_ids'].squeeze(1).to(device)
109
110            output = model(input_id, mask)
111
112            batch_loss = criterion(output, val_label)
113            total_loss_val += batch_loss.item()

```

```

114
115         acc = (output.argmax(dim=1) == val_label).sum().item()
116         total_acc_val += acc
117
118     print(
119         f'''Epochs: {epoch_num + 1}
120         | Train Loss: {total_loss_train / len(train_data): .3f}
121         | Train Accuracy: {total_acc_train / len(train_data): .3f}
122         | Val Loss: {total_loss_val / len(val_data): .3f}
123         | Val Accuracy: {total_acc_val / len(val_data): .3f}'''

```

```

1 EPOCHS = 5
2 model = BertClassifier()
3 LR = 1e-6
4 train(model, data_train, data_valid, LR, EPOCHS)

```

100%|██████████| 596/596 [02:38<00:00, 3.76it/s]

Epochs: 1

```

| Train Loss: 0.795
| Train Accuracy: 0.265
| Val Loss: 0.751
| Val Accuracy: 0.389

```

100%|██████████| 596/596 [02:44<00:00, 3.62it/s]

Epochs: 2

```

| Train Loss: 0.644
| Train Accuracy: 0.574
| Val Loss: 0.494
| Val Accuracy: 0.758

```

100%|██████████| 596/596 [02:52<00:00, 3.46it/s]

Epochs: 3

```

| Train Loss: 0.311
| Train Accuracy: 0.904
| Val Loss: 0.194
| Val Accuracy: 0.966

```

100%|██████████| 596/596 [02:51<00:00, 3.48it/s]

Epochs: 4

```

| Train Loss: 0.144
| Train Accuracy: 0.977
| Val Loss: 0.116
| Val Accuracy: 0.980

```

100%|██████████| 596/596 [02:50<00:00, 3.50it/s]

Epochs: 5

```

| Train Loss: 0.084
| Train Accuracy: 0.985
| Val Loss: 0.079
| Val Accuracy: 0.977

```