

第 5 章



Transformer 与人机畅聊

当读完本章时,应该能够:

- 理解机器是如何学习说话和学会说话的。
- 理解机器问答的技术路线。
- 理解机器问答语言模型及评价方法。
- 理解 Transformer 模型的原理与方法。
- 理解 Transformer 独特的注意力机制。
- 理解 Transformer 模型的定义与训练。
- 实战基于腾讯聊天数据集+Transformer 的聊天机器人设计。
- 实战基于 Web API 的聊天服务器设计。
- 实战 Android 版人机畅聊客户机设计。

5.1 项目动力



机器视觉与自然语言处理是人工智能的两大热点领域。对于机器视觉学习,读者较容易入门,但对于自然语言的处理,特别是对于机器翻译、机器问答、人机聊天等应用的理解与学习,难度较高。究其原因,读者往往很难独立去完成一个上述自然语言处理项目,一旦缺少了亲力亲为的实战体验,理论不能与实践紧密结合,不能相互印证,那么对理论的掌握往往就是空中楼阁。

本章案例将带领读者从零起步,以 Transformer 为建模基础,在腾讯发布的中文聊天数据集上,训练出一款会聊天的机器人程序。聊天机器人除了不知疲倦、有求必应、一对多并发服务的优点外,聊天机器人的未来将兼具多种风格,例如滔滔不绝的机器人、激情演讲的机器人、能做灵魂触碰的机器人、能心有灵犀的机器人、能聊出思想火花的机器人、

能与你一起头脑风暴的机器人、能自我学习提高的机器人……

总之,具备强大语言能力的机器人、具备思想能力的机器人,正沿着人类追求的航向,奋力前行。



5.2 机器问答技术路线

机器问答与机器聊天(本书特指人机聊天)都是自然语言处理领域的经典问题,二者并不完全相同。

机器问答一般常见于阅读理解的场景,给定参考语料,要求机器人根据语料回答问题,问题的答案往往是在语料之中做了显式陈述,或者根据语料可以推断的。

对于机器聊天而言,聊天的随机性决定了应答难度更大一些,不但需要及时切换话题,还需要通过简短的交谈揣摩对方的真实意图,否则会答非所问。

机器问答与机器聊天的界限有时也是模糊的,例如,电商平台部署的客服机器人,确实可以解决客户的若干问题,节省了大规模的人力资源。这些客服机器人基于其收集的大量客服会话数据训练而成,兼具机器问答与机器聊天的特点。

经常听到新入门学生询问如何开发一个机器问答或者机器聊天的项目,图 5.1 给出了项目实施路径的四部曲。

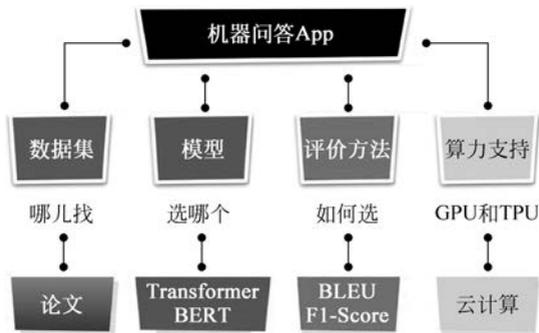


图 5.1 机器问答技术路线

第一步,解决数据集的问题。俗话说,巧妇难为无米之炊,强大的语料数据集对语言建模至关重要。基本原则是根据问题目标,选择适配的语料库。例如斯坦福大学发布的 CoQA、SQuAD 2.0 语料库,华盛顿大学牵头发布的 QuAC 语料库,百度发布的 DuConv 语料库,京东的 JDCC 客服语料库,腾讯的 NaturalConv 聊天语料库,清华大学的 KdConv 语料库等。这些经典的语料库都有论文解析,通过阅读论文可以得知语料库的规模、覆盖的领域、适合解决的问题。如果条件允许,可以考虑自建语料库。或者自建小规模语料库,用迁移学习的方法解决数据量不足的问题。

第二步,选择模型(算法)。工欲善其事,必先利其器。近年来,能够代表自然语言处理领域前沿热点的模型非 Transformer 和 BERT 莫属。研读前沿文献不难发现,活跃在各大语言建模挑战赛排行榜前列的模型,无不是以 Transformer 和 BERT 结构为核心的衍生品。图 5.2 给出了适合自然语言建模的经典技术,自底向上,反映了语言建模技术的

变迁与演进。

第三步,理解语言模型的评价方法。之所以强调这一点,是因为学生往往比较熟悉机器视觉领域的分类或回归问题的评价方法,这些方法在自然语言处理领域并不好用。以机器聊天或机器推理为例,可接受的正确答案往往不止一个。如何去评价这些语句长度不一、文字不一的句子,需要新方法。例如 BLEU 即是一种经典方法。

第四步,解决算力问题。普通的个人计算机,即使对于小规模数据集,也难以胜任模型训练。GPU 性能稍好的工作站,往往也难以满足中等建模问题对算力的需求。对于学生而言,需要寻找免费算力资源,才能不让自己的想法卡在最后一公里。幸运的是,百度、Google 等都提供了免费算力平台。

以本章的机器人聊天项目为例,在 3.7 节给出的计算机配置上,完成 38 个 epoch 的训练,需要 4 小时左右。在 Kaggle 平台提供的免费 TPUv3-8 上,只需要 12 分钟。速度提升 20 倍。

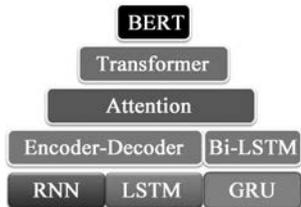


图 5.2 语言建模技术的演进

5.3 腾讯聊天数据集

2021 年,腾讯人工智能实验室发布了主题驱动的多轮中文聊天数据集,数据集名称为 NaturalConv,详情参见论文 *Naturalconv: A chinese dialogue dataset towards multi-turn topic-driven conversation* (WANG, X Y, LI C, ZHAO J, et al. 2021)。

NaturalConv 包含体育、娱乐、科技、游戏、教育和健康六个领域的 19.9K 场双人对话,合计 40 万条对话语句,每场对话平均包含 20.1 个句子,即对话双方进行了 10 次左右的表达。为了获取上述数据,腾讯支付了语料库整理人员 5 万美元的信息采集和加工整理费用。现在这个数据集可以在腾讯人工智能实验室官方网站免费下载。

NaturalConv 数据集包含的对话是以 6500 篇网络新闻为背景展开的。对话具备自然性,即对话双方围绕一个共同主题展开话语交流,属于自然情况下的聊天行为。

口语化是聊天的显著特点,这也是 NaturalConv 数据集设计时遵循的原则。表 5.1 给出了 NaturalConv 中的一场对话文本,A、B 两名学生关于当天的一场荷甲球赛展开的话题聊天。

表 5.1 A、B 两人对话样本抽样观察

轮次	A、B 两人的对话内容
A-1	嗨,你来得挺早啊。
B-1	是啊,你怎么来得这么晚?
A-2	昨晚我看了球赛,所以今早起晚了,也没吃饭。
B-2	现在这个点食堂应该有饭,你看什么球赛啊? 篮球吗?
A-3	不是,足球。
B-3	怪不得,足球时间长。
A-4	你知道吗? 每次都是普罗梅斯进球。
B-4	这个我刚才也看了新闻了,他好有实力啊。



续表

轮次	A、B 两人的对话内容
A-5	是啊,尤其是他那个帽子戏法,让我看得太惊心动魄了。
B-5	我一个同学在群里说了,每次聊天都离不开他,可见他的实力有多强大。
A-6	是啊,看来你那个同学和我是一样的想法。
B-6	我好不容易摆脱他的话题,你又来说出他的名字。
A-7	哈哈,你不懂我们对足球有多热爱。
B-7	我知道你热爱,我还记得你参加初中比赛还拿到冠军呢。你功不可没啊。
A-8	哈哈,还是你能记得我当时的辉煌。
B-8	没办法,咱俩从小一起长大的,彼此太了解了。
A-9	嗯,老师来了。
B-9	快打开课本,老师要检查。
A-10	嗯嗯,下课再聊。
B-10	嗯。

话题参照的新闻稿件如下。

新闻稿:北京时间今天凌晨,荷甲第4轮进行了两场补赛。阿贾克斯和埃因霍温均在家取得了胜利,两队7轮后同积17分,阿贾克斯以6个净胜球的优势领跑积分榜。0点30分,埃因霍温与格罗宁根的比赛开战,埃因霍温最终3:1在家获胜。2点45分,阿贾克斯主场与福图纳锡塔德之战开始。由于埃因霍温已经先获胜了,阿贾克斯必须获胜才能在积分榜上咬住对方。在整个上半场,阿贾克斯得势不得分,双方0:0互交白卷。在下半场中,阿贾克斯突然迎来了大爆发。在短短33分钟内,阿贾克斯疯狂打进5球,平均每6分钟就能取得1个进球。在第50分钟时,新援普罗梅斯为阿贾克斯打破僵局。塔迪奇左侧送出横传,普罗梅斯后点推射破门。第53分钟时亨特拉尔头球补射,内雷斯在门线前头球接力破门。第68分钟时,普罗梅斯近距离补射梅开二度。这名27岁的前场多面手,跑到场边来了一番尬舞。第77分钟时阿贾克斯收获第4球,客队后卫哈里斯在防传中时伸腿将球一捅,结果皮球恰好越过门将滚入网窝。在第83分钟时,普罗梅斯上演了帽子戏法,比分也最终被定格为5:0。在接到塔迪奇直传后,普罗梅斯禁区左侧反越位成功,他的单刀低射从门将裆下入网。普罗梅斯这次的庆祝动作是秀出三根手指,不过他手指从上到下抹过面部时的动作,很有点像是在擦鼻涕。

NaturalConv数据集与CMU DoG、Wizard of Wiki、DuConv、KdConv的对比如表5.2所示。NaturalConv涉及的主题更为广泛,双方对话的回合数更多,语料库规模更大。

表 5.2 与其他数据集对比

Dataset	Language	Document Type	Annotation Level	Topic	Avg. # turns	# uttrs
CMU DoG	English	Text	Sentence	Film	22.6	130k
Wizard of Wiki	English	Text	Sentence	Multiple	9.0	202k
DuConv	Chinese	Text&KG	Dialogue	Film	5.8	91k
KdConv	Chinese	Text&KG	Sentence	Film, Music, Travel	19.0	86k
NaturalConv	Chinese	Text	Dialogue	Sports, Ent, Tech, Games, Edu, Health	20.1	400k

NaturalConv 语料库中各主题的分布情况如表 5.3 所示。可以看到,体育类主题的文档数占比接近一半,健康类主题最少,只有 52 篇背景文档。

表 5.3 主题的分布情况

	Sports	Ent	Tech	Games	Edu	Health	Total
# document	3124	1331	1476	103	414	52	6500
# dialogues	9740	4403	4061	308	1265	142	19 919
# dialogues per document	3.1	3.3	2.8	3.0	3.1	2.7	3.0
# utterances	195 643	88 457	81 587	6180	25 376	2852	400 095
Avg. # utterances per dialogue	20.1	20.1	20.1	20.1	20.1	20.1	20.1
Avg. # tokens per utterance	12.0	12.4	12.3	12.1	12.6	12.5	12.2
Avg. # characters per utterance	17.8	18.1	18.6	17.8	18.1	18.3	18.1
Avg. # tokens per dialogue	241.1	248.2	247.5	242.9	248.3	251.1	244.8
Avg. # characters per dialogue	357.5	363.2	372.8	356.5	356.5	368.0	363.1

每篇文档相关的对话场次平均为 3 场。每场对话包含的语句平均为 20.1 个。平均每场对话包含的字数为 360 个左右。

可见,NaturalConv 语料库的主题分布并不均衡,这也可以解释,为什么后面训练的聊天机器人,偏爱于体育类的话题,或者对体育类的话题更“健谈”一些。语料库相当于语言模型的先天基因,决定了语言模型表现的倾向性。

5.4 Transformer 模型解析



Transformer 模型参见论文 *Attention is all you need* (VASWANI A, SHAZEER N, PARMAR N, et al. 2017)。Transformer 这个词的原意是“变形金刚”,作者用 Transformer 寓意其可伸缩性好,可以胜任的应用领域非常广泛。

Transformer 模型结构如图 5.3 所示。左侧代表编码器,右侧代表解码器。编码器与解码器的结构非常相似。

编码器由 N 个结构重复的单元连接而成。每个单元包含两个残差块:第一个残差块以多头注意力模块为核心;第二个残差块以全连接网络为核心。

解码器的主体是由 N 个结构重复的单元连接而成,最后跟上一个全连接层和 Softmax 分类层。每个单元包含三个残差块:第一个残差块以带掩码的多头注意力模块为核心;第二个残差块是交叉多头注意力模块;第三个模块为全连接网络模块。

注意,Transformer 中用 Linear 表示的全连接网络,不包含激活函数,是一个线性结构。

编码器的输入层需要将序列的嵌入向量与序列的位置编码向量做叠加运算,然后并行输入到三个 Linear 网络,得到 Q 、 K 、 V 三个输入向量。

编码器输出层的输出将直接给到解码器各个单元的第二个残差块。

解码器除了接收来自编码器的输入,还有一个被称作 Output Embedding 的输入。在模型训练期间,Output Embedding 用样本的标签向量表示。在模型推理时,解码器是

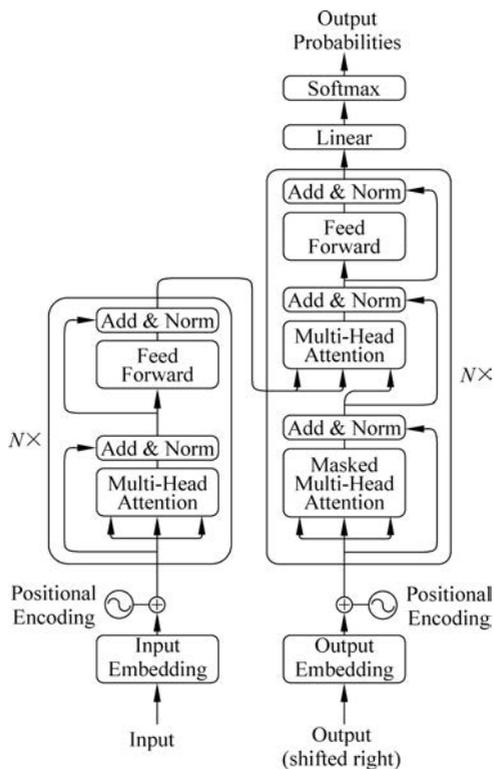


图 5.3 Transformer 模型结构(见彩插)

一个自回归结构,单步推理生成的输出,将回馈给 Output Embedding 作为解码器的输入,参入整个推理过程。

编码器与解码器之间的连接方式有很多,图 5.4 给出了 Transformer 的经典推荐方式。即编码器最后一个单元的输出,给到解码器各个单元的交叉多头注意力模块。

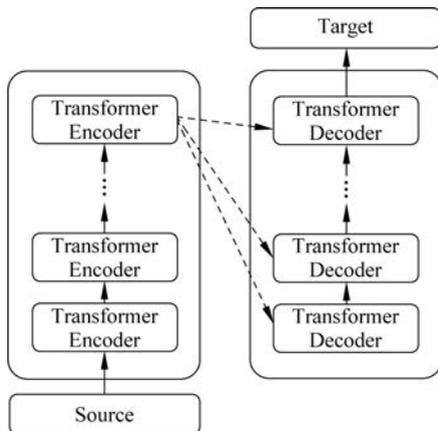


图 5.4 Transformer 的经典推荐方式

事实上,也可以让编码器各个单元的输出,只给到解码器的同层次单元,或者给到解

码器的所有单元。

Transformer 的核心运算体现在注意力机制上,图 5.5 给出了 Transformer 的单头注意力机制与多头注意力机制的计算逻辑。

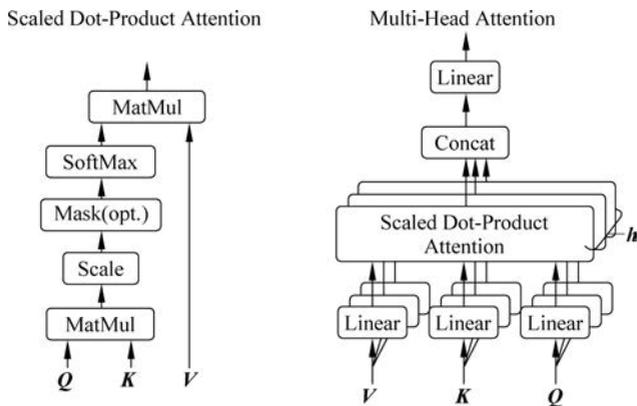


图 5.5 Transformer 的单头注意力机制与多头注意力机制的计算逻辑

Transformer 将编码的输入序列,通过全连接网络学习为 Q 、 K 、 V 三个嵌入向量。左图为基于点积的单头注意力逻辑。右图为 h 个多头注意力并行堆叠、合并计算的逻辑。

Q 、 K 之间完成注意力计算,形成注意力权重,用注意力权重乘以 V ,得到单个注意力模块的输出。

直观上看, Q 、 K 之间的注意力强度反映了序列中上下文之间的关系强弱,将这种强弱关系映射到向量 V 上,即可实现对输入序列的编码和特征提取。 Q 、 K 、 V 之间的计算逻辑与互动关系如图 5.6 所示。

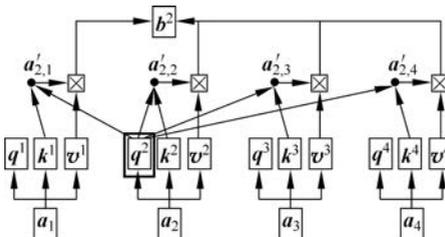


图 5.6 Q 、 K 、 V 之间的计算逻辑与互动关系

假设输入序列为 $a_1 a_2 a_3 a_4$,现在只考虑 a_2 对应的输出 b^2 ,通过注意力计算得到 b^2 的步骤解析如下。

- (1) a_2 通过三个独立的全连接网络学习,得到编码 q^2 、 k^2 、 v^2 , q^2 与 k^2 按照图 5.5 所示的点积注意力计算逻辑,经 Softmax 输出 q^2 与 k^2 之间的归一化关系权重向量 $a'_{2,2}$ 。
- (2) q^2 与 k^1 经注意力计算,Softmax 层输出 q^2 与 k^1 的归一化关系权重向量 $a'_{2,1}$ 。
- (3) 重复步骤(2),得到 $a'_{2,3}$ 和 $a'_{2,4}$ 。

至此, a_2 与序列 $a_1 a_2 a_3 a_4$ 中其他单词的关系(包括与自身的关系),已经通过 q^2 与 k^1 、 k^2 、 k^3 、 k^4 之间的注意力计算得到,表示为四个权重向量 $a'_{2,1}$ 、 $a'_{2,2}$ 、 $a'_{2,3}$ 和 $a'_{2,4}$ 。

(4) 得到 $\mathbf{b}^2 = \mathbf{a}'_{2,1} \times v^1 + \mathbf{a}'_{2,2} \times v^2 + \mathbf{a}'_{2,3} \times v^3 + \mathbf{a}'_{2,4} \times v^4$ 。现在, 可以认为向量 \mathbf{b}^2 是对向量 \mathbf{a}_2 施加上下文全局注意力后的新表示。

同样的方法可以得到 \mathbf{b}^1 、 \mathbf{b}^3 、 \mathbf{b}^4 。

从 $\mathbf{a}_1 \mathbf{a}_2 \mathbf{a}_3 \mathbf{a}_4$ 到 $\mathbf{b}^1 \mathbf{b}^2 \mathbf{b}^3 \mathbf{b}^4$, 依靠注意力机制, 完成了一次特征提取与变换。

下面通过一个例子演示注意力的计算过程。如图 5.7 所示(图片源自 Ketan Doshi 博客), 假设输入的序列为 You are welcome, 规定序列长度为 4, 所以后面需要补一个空位, 不妨用 PAD 表示。假定每个单词向量的编码长度为 3。

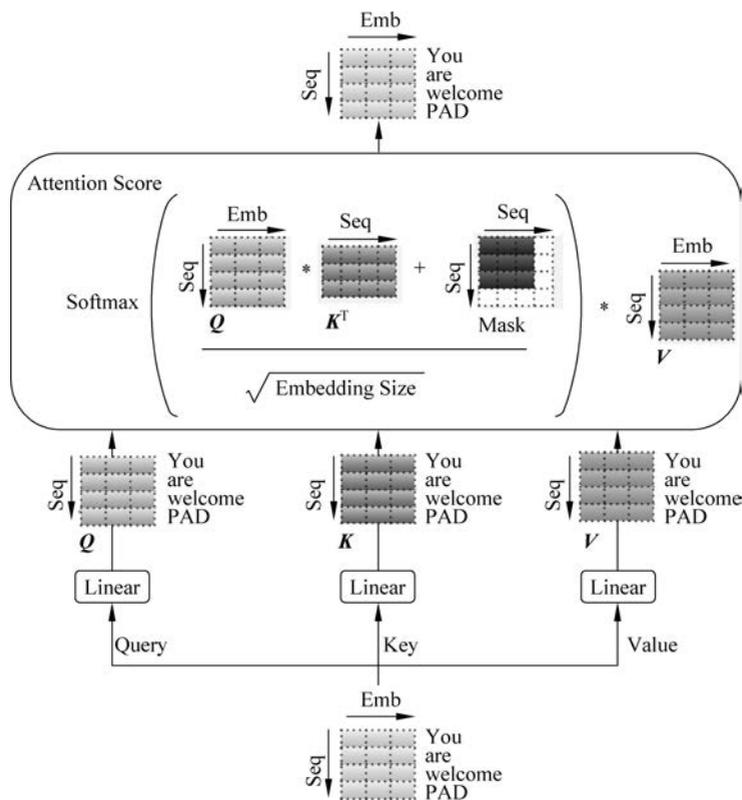


图 5.7 注意力计算举例

You are welcome PAD 构成了一个维度为 $(4, 3)$ 的特征矩阵, 分别送入三个 Linear 网络, 得到 \mathbf{Q} 、 \mathbf{K} 、 \mathbf{V} 三个特征编码矩阵, 维度均为 $(4, 3)$ 。

注意力的计算可以归纳为式(5.1)。

$$\text{Attention}(\mathbf{Q}, \mathbf{K}, \mathbf{V}) = \text{Softmax}\left(\frac{\mathbf{Q}\mathbf{K}^T}{\sqrt{d_k}}\right)\mathbf{V} \quad (5.1)$$

其中, d_k 表示 \mathbf{Q} 、 \mathbf{K} 序列中单词向量的长度。考虑到 \mathbf{Q} 与 \mathbf{K} 的点积运算, 有可能放大了特征的输出值, 所以分母除以 $\sqrt{d_k}$, 这也是 Scaled Dot-Product Attention 名称的由来。

Transformer 的相关参数设置如表 5.4 所示。

表 5.4 Transformer 的相关参数设置

参数名称	参数值
编码器单元数 N	6
解码器单元数 N	6
输入输出向量的长度 d_{model}	512
注意力头数 h	8
Q, K, V 的向量长度	$d_k = d_v = d_{\text{model}} / h = 64$

5.5 机器人项目初始化



用 PyCharm 在 TensorFlow_to_Android 项目下新建目录 MyRobot。在 MyRobot 目录下新建 models 和 transformer 两个子目录。

models 目录用于存放训练好的聊天机器人模型。transformer 目录用于存放数据集预处理程序和模型定义程序。

在 transformer 目录下新建子目录 dataset。将下载的腾讯自然语言聊天数据集 NaturalConv 解压后存放到 dataset 目录下。

数据集可在腾讯人工智能实验室官方网站免费下载。下载地址为 <https://ai.tencent.com/ailab/nlp/dialogue/#datasets>。

在 MyRobot 目录下新建主程序 main.py, 负责模型的训练、保存、评估和测试。

初始化后的项目结构如图 5.8 所示。

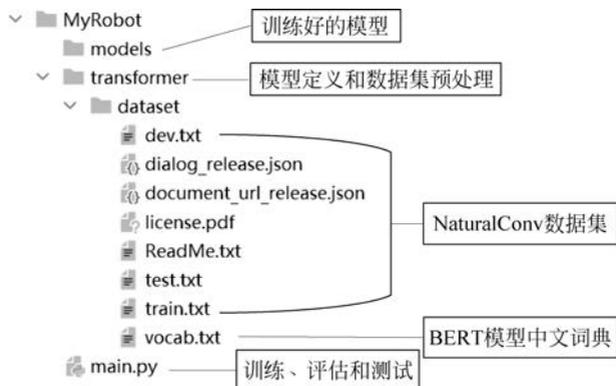


图 5.8 MyRobot 项目初始结构

dataset 目录下的 vocab.txt 是 BERT 中文模型的词典文件。这个文件不属于腾讯自然语言聊天数据集, 需要从 BERT 官方网站下载中文模型, 从中抽取中文词典文件。

BERT 中文预训练模型可从官方网站下载。下载地址为 <https://github.com/google-research/bert>。

本章案例采用 BERT 的分词模型对中文语句分词。



5.6 数据集预处理与划分

模型训练之前,需要首先准备好数据,让数据能够直接“喂入”模型进行训练,为了提高模型“喂入”的效率,往往还要设计数据加载模式。

在 transformer 目录下新建 dataset.py 程序,如程序源码 P5.1 所示。

程序源码 P5.1 dataset.py 数据集预处理与划分

```

1 import codecs
2 import json
3 import re
4 import tensorflow as tf
5 from tensorflow.keras.preprocessing.sequence import pad_sequences
6 # 析取数据集(训练集、验证集、测试集),将所有的"问"与"答"分开
7 def extract_conversations(hparams, data_list, dialog_list):
8     inputs, outputs = [], [] # 问答列表
9     for dialog in dialog_list:
10         if dialog['dialog_id'] in data_list:
11             if len(dialog['content']) % 2 == 0:
12                 i = 0
13                 for line in dialog['content']:
14                     if (i % 2 == 0):
15                         inputs.append(line) # "问"列表
16                     else:
17                         outputs.append(line) # "答"列表
18                     i += 1
19                 # 限定样本总数
20                 # if len(inputs) >= hparams.total_samples:
21                     # return inputs, outputs
22     return inputs, outputs
23 # 分词,过滤掉超过长度的句子,短句补齐
24 def tokenize_and_filter(hparams, inputs, outputs, tokenizer):
25     tokenized_inputs, tokenized_outputs = [], []
26     for (sentence1, sentence2) in zip(inputs, outputs):
27         sentence1 = tokenizer.tokenize(sentence1) # 分词
28         sentence1 = tokenizer.convert_tokens_to_ids(sentence1) # ids
29         sentence2 = tokenizer.tokenize(sentence2)
30         sentence2 = tokenizer.convert_tokens_to_ids(sentence2)
31         sentence1 = hparams.start_token + sentence1 + hparams.end_token
32         sentence2 = hparams.start_token + sentence2 + hparams.end_token
33         if len(sentence1) <= hparams.max_length and len(sentence2) <= hparams.max_
34             length:
35             tokenized_inputs.append(sentence1)
36             tokenized_outputs.append(sentence2)
37     # 补齐
38     tokenized_inputs = pad_sequences(tokenized_inputs, \
39                                     maxlen = hparams.max_length, padding = 'post')

```

```
40     tokenized_outputs = pad_sequences(tokenized_outputs, \
41                                     maxlen=hparams.max_length, padding='post')
42     return tokenized_inputs, tokenized_outputs
43 # 读文件
44 def get_data(datafile):
45     with open(f'{datafile}', 'r') as f:
46         data_list = f.readlines()
47         for i in range(len(data_list)):
48             data_list[i] = re.sub(r'\n', '', data_list[i])
49     return data_list
50 # 返回训练集和验证集
51 def get_dataset(hparams, tokenizer, dialog_file, train_file, valid_file):
52     dialog_list = json.loads(codecs.open(f"{dialog_file}", "r", "utf-8").read())
53     print(dialog_list[0])
54     train_list = get_data(f'{train_file}')
55     train_questions, train_answers = extract_conversations(hparams, train_list, dialog_list)
56     train_questions, train_answers = tokenize_and_filter(hparams, \
57                                                         list(train_questions), list(train_answers), tokenizer)
58     # 构建训练集
59     train_dataset = tf.data.Dataset.from_tensor_slices((
60         {
61             'inputs': train_questions,
62             # 解码器使用正确的标签作为输入
63             'dec_inputs': train_answers[:, :-1] # 去掉最后一个元素或 END_TOKEN
64         },
65         {
66             'outputs': train_answers[:, 1:] # 去掉 START_TOKEN
67         },
68     ))
69     train_dataset = train_dataset.cache()
70     train_dataset = train_dataset.shuffle(len(train_questions))
71     train_dataset = train_dataset.batch(hparams.batchSize)
72     train_dataset = train_dataset.prefetch(tf.data.AUTOTUNE)
73     # 构建验证集
74     valid_list = get_data(f'{valid_file}')
75     valid_questions, valid_answers = extract_conversations(\
76                                     hparams, valid_list, dialog_list)
77     valid_questions, valid_answers = tokenize_and_filter(hparams, \
78                                                         list(valid_questions), list(valid_answers), tokenizer)
79     valid_dataset = tf.data.Dataset.from_tensor_slices((
80         {
81             'inputs': valid_questions,
82             # 解码器使用正确的标签作为输入
83             'dec_inputs': valid_answers[:, :-1] # 去掉最后一个元素或 END_TOKEN
84         },
85         {
86             'outputs': valid_answers[:, 1:] # 去掉 START_TOKEN
87         },
88     ))
```

```

89     valid_dataset = valid_dataset.cache()
90     valid_dataset = valid_dataset.shuffle(len(valid_questions))
91     valid_dataset = valid_dataset.batch(hparams.batchSize)
92     valid_dataset = valid_dataset.prefetch(tf.data.AUTOTUNE)
93     return train_dataset, valid_dataset

```

程序源码解析参见本节视频教程。



5.7 定义 Transformer 输入层编码

模型定义的完整流程如图 5.9 所示。本节首先完成输入层的编码定义。后续各节分步完成注意力机制、编码器、解码器的模块定义,最后合成为模型的整体定义。

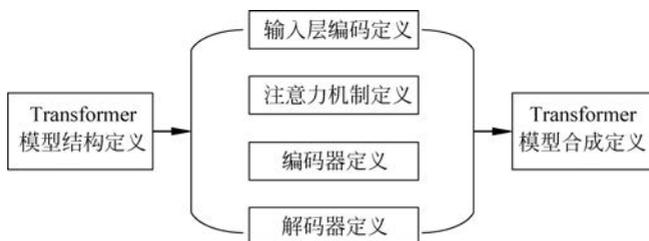


图 5.9 Transformer 模型定义的完整流程

在 transformer 目录下新建模型定义程序 model.py。关于编码器输入层的定义逻辑如程序源码 P5.2 所示。

程序源码 P5.2 model.py 之输入层定义

```

1 import matplotlib.pyplot as plt
2 import tensorflow as tf
3 from tensorflow.keras import Input, Model
4 from tensorflow.keras.layers import Dense, Lambda, Embedding, Dropout, \
5     add, LayerNormalization
6 from tensorflow.keras.utils import plot_model
7 # 定义掩码矩阵
8 def create_padding_mask(x):
9     # 找出序列中的 padding, 设置掩码值为 1
10    mask = tf.cast(tf.math.equal(x, 0), tf.float32)
11    # (batch_size, 1, 1, sequence length)
12    return mask[:, tf.newaxis, tf.newaxis, :]
13 # 测试语句
14 print(create_padding_mask(tf.constant([[1, 2, 0, 3, 0], [0, 0, 0, 4, 5]])))
15 # 解码器的前向掩码
16 def create_look_ahead_mask(x):
17    seq_len = tf.shape(x)[1]
18    look_ahead_mask = 1 - tf.linalg.band_part(tf.ones((seq_len, seq_len)), -1, 0)
19    padding_mask = create_padding_mask(x)
20    return tf.maximum(look_ahead_mask, padding_mask)

```

```
21 # 测试
22 print(create_look_ahead_mask(tf.constant([[1, 2, 0, 4, 5]])))
23 # 位置编码类
24 class PositionalEncoding(tf.keras.layers.Layer):
25     def __init__(self, position, d_model):
26         super(PositionalEncoding, self).__init__()
27         self.pos_encoding = self.positional_encoding(position, d_model)
28     def get_config(self):
29         config = super(PositionalEncoding, self).get_config()
30         config.update({
31             'position': self.position,
32             'd_model': self.d_model,
33         })
34         return config
35     def get_angles(self, position, i, d_model):
36         angles = 1 / tf.pow(10000, (2 * (i // 2)) / tf.cast(d_model, tf.float32))
37         return position * angles
38     def positional_encoding(self, position, d_model):
39         angle_rads = self.get_angles( \
40             position = tf.range(position, dtype = tf.float32)[:, tf.newaxis], \
41             i = tf.range(d_model, dtype = tf.float32)[tf.newaxis, :], \
42             d_model = d_model)
43         # 奇数位置用正弦函数
44         sines = tf.math.sin(angle_rads[:, 0::2])
45         # 偶数位置用余弦函数
46         cosines = tf.math.cos(angle_rads[:, 1::2])
47         pos_encoding = tf.concat([sines, cosines], axis = -1)
48         pos_encoding = pos_encoding[tf.newaxis, ...]
49         return tf.cast(pos_encoding, tf.float32)
50     def call(self, inputs):
51         return inputs + self.pos_encoding[:, :tf.shape(inputs)[1], :]
52 # 测试
53 sample_pos_encoding = PositionalEncoding(50, 512)
54 plt.pcolormesh(sample_pos_encoding.pos_encoding.numpy()[0], cmap = 'RdBu')
55 plt.xlabel('Depth')
56 plt.xlim((0, 512))
57 plt.ylabel('Position')
58 plt.colorbar()
59 plt.show()
```

假定输入层向量的维度为(50,512),即序列长度为50(也即有50个单词),每个单词的嵌入向量长度为512,则整个输入层向量的位置编码及其几何分布如图5.10所示。对于中文而言,分词以单个中文字符为单位,所以位置编码是针对单个中文字符而言的。对于英文,分词以英文单词为单位。函数曲线的变化趋势体现了位置编码函数对不同位置编码的差异性。

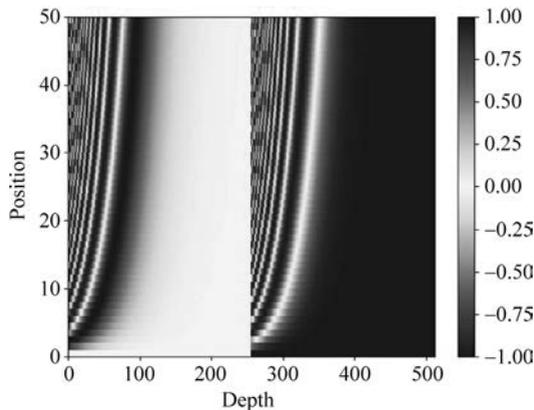


图 5.10 位置编码及其几何分布



5.8 定义 Transformer 注意力机制

先定义单头注意力函数,实现单头注意力机制的计算,然后合成多头注意力计算模块。编码逻辑如程序源码 P5.3 所示。

程序源码 P5.3 model.py 之注意力机制

```

1  # 计算注意力
2  def scaled_dot_product_attention(query, key, value, mask):
3      matmul_qk = tf.matmul(query, key, transpose_b=True)
4      # 计算 qk
5      depth = tf.cast(tf.shape(key)[-1], tf.float32)
6      logits = matmul_qk / tf.math.sqrt(depth)
7      # 添加掩码以将填充标记归零
8      if mask is not None:
9          logits += (mask * -1e9)
10     # 在最后一个轴上实施 softmax
11     attention_weights = tf.nn.softmax(logits, axis=-1)
12     output = tf.matmul(attention_weights, value)
13     return output
14 # 定义多头注意力类,继承了 Layer 类
15 class MultiHeadAttention(tf.keras.layers.Layer):
16     def __init__(self, d_model, num_heads, name="multi_head_attention"):
17         super(MultiHeadAttention, self).__init__(name=name)
18         self.num_heads = num_heads
19         self.d_model = d_model
20         assert d_model % self.num_heads == 0
21         self.depth = d_model // self.num_heads
22         self.query_dense = Dense(units=d_model)
23         self.key_dense = Dense(units=d_model)
24         self.value_dense = Dense(units=d_model)
25         self.dense = Dense(units=d_model)

```

```

26 def get_config(self):
27     config = super(MultiHeadAttention, self).get_config()
28     config.update({
29         'num_heads':self.num_heads,
30         'd_model':self.d_model,
31     })
32     return config
33 def split_heads(self, inputs, batch_size):
34     inputs = Lambda(lambda inputs:tf.reshape(inputs, \
35         shape = (batch_size, -1, self.num_heads, self.depth)))(inputs)
36     return Lambda(lambda inputs: tf.transpose(inputs, perm = [0, 2, 1, 3]))(inputs)
37 def call(self, inputs):
38     query, key, value, mask = inputs['query'], inputs['key'], inputs[
39         'value'], inputs['mask']
40     batch_size = tf.shape(query)[0]
41     # 线性层变换
42     query = self.query_dense(query)
43     key = self.key_dense(key)
44     value = self.value_dense(value)
45     # 分头
46     query = self.split_heads(query, batch_size)
47     key = self.split_heads(key, batch_size)
48     value = self.split_heads(value, batch_size)
49     # 定义缩放的点积注意力
50     scaled_attention = scaled_dot_product_attention(query, key, value, mask)
51     scaled_attention = Lambda(lambda scaled_attention: tf.transpose(
52         scaled_attention, perm = [0, 2, 1, 3]))(scaled_attention)
53     # 堆叠注意力头
54     concat_attention = Lambda(lambda scaled_attention: tf.reshape( \
55         scaled_attention, (batch_size, -1, self.d_model)))(scaled_attention)
56     # 多头注意力最后一层
57     outputs = self.dense(concat_attention)
58     return outputs

```

5.9 定义 Transformer 编码器



先完成编码器一个单元的定义,然后由多个单元合成整个编码器。编码逻辑如程序源码 P5.4 所示。

程序源码 P5.4 model.py 之编码器定义

```

1 # 编码器中的一层,即编码器的一个单元定义
2 def encoder_layer(units, d_model, num_heads, dropout, name = "encoder_layer"):
3     inputs = tf.keras.Input(shape = (None, d_model), name = "inputs")
4     padding_mask = tf.keras.Input(shape = (1, 1, None), name = "padding_mask")
5     attention = MultiHeadAttention( \
6         d_model, num_heads, name = "attention")({ \

```

```
7         'query': inputs,
8         'key': inputs,
9         'value': inputs,
10        'mask': padding_mask
11    })
12    attention = Dropout(rate = dropout)(attention)
13    add_attention = add([inputs, attention])
14    attention = LayerNormalization(epsilon = 1e - 6)(add_attention)
15    outputs = Dense(units = units, activation = 'relu')(attention)
16    outputs = Dense(units = d_model)(outputs)
17    outputs = Dropout(rate = dropout)(outputs)
18    add_attention = add([attention, outputs])
19    outputs = LayerNormalization(epsilon = 1e - 6)(add_attention)
20    return Model(inputs = [inputs, padding_mask], outputs = outputs, name = name)
21 # 测试
22 sample_encoder_layer = encoder_layer(
23     units = 512,
24     d_model = 128,
25     num_heads = 4,
26     dropout = 0.3,
27     name = "sample_encoder_layer")
28 plot_model(sample_encoder_layer, to_file = 'encoder_layer.png', show_shapes = True)
29 # 定义编码器,由多个单元合成编码器
30 def encoder(vocab_size,
31            num_layers,
32            units,
33            d_model,
34            num_heads,
35            dropout,
36            name = "encoder"):
37     inputs = Input(shape = (None, ), name = "inputs")
38     padding_mask = Input(shape = (1, 1, None), name = "padding_mask")
39     embeddings = Embedding(vocab_size, d_model)(inputs)
40     embeddings * = Lambda(lambda d_model: tf.math.sqrt(tf.cast(d_model, tf.float32)))(d_model)
41     embeddings = PositionalEncoding(vocab_size, d_model)(embeddings)
42     outputs = Dropout(rate = dropout)(embeddings)
43     for i in range(num_layers):
44         outputs = encoder_layer(
45             units = units,
46             d_model = d_model,
47             num_heads = num_heads,
48             dropout = dropout,
49             name = "encoder_layer_{}".format(i),
50             )([outputs, padding_mask])
51     return Model(inputs = [inputs, padding_mask], outputs = outputs, name = name)
52 # 编码器测试
53 sample_encoder = encoder(
54     vocab_size = 21128,
55     num_layers = 2,
```

```
56     units = 512,  
57     d_model = 128,  
58     num_heads = 4,  
59     dropout = 0.3,  
60     name = "sample_encoder")  
61 plot_model(sample_encoder, to_file = 'encoder.png', show_shapes = True)
```

测试编码器程序,观察生成的编码器逻辑结构,与 Transformer 论文解析的结构做对照,在实践中灵活配置编码器的相关参数,可得到适配问题需求的编码器。

5.10 定义 Transformer 解码器



解码器的设计思路与编码器类似。先完成解码器一个单元的定义,然后由多个单元合成整个解码器。编码逻辑如程序源码 P5.5 所示。

程序源码 P5.5 model.py 之解码器定义

```
1 # 定义解码器中的一层,一个解码单元  
2 def decoder_layer(units, d_model, num_heads, dropout, name = "decoder_layer"):  
3     inputs = Input(shape = (None, d_model), name = "inputs")  
4     enc_outputs = Input(shape = (None, d_model), name = "encoder_outputs")  
5     look_ahead_mask = Input(shape = (1, None, None), name = "look_ahead_mask")  
6     padding_mask = Input(shape = (1, 1, None), name = 'padding_mask')  
7     attention1 = MultiHeadAttention(  
8         d_model, num_heads, name = "attention_1")(inputs = {  
9         'query': inputs,  
10        'key': inputs,  
11        'value': inputs,  
12        'mask': look_ahead_mask  
13    })  
14    add_attention = tf.keras.layers.add([attention1, inputs])  
15    attention1 = tf.keras.layers.LayerNormalization(epsilon = 1e - 6)(add_attention)  
16    attention2 = MultiHeadAttention(  
17        d_model, num_heads, name = "attention_2")(inputs = {  
18        'query': attention1,  
19        'key': enc_outputs,  
20        'value': enc_outputs,  
21        'mask': padding_mask  
22    })  
23    attention2 = Dropout(rate = dropout)(attention2)  
24    add_attention = add([attention2, attention1])  
25    attention2 = LayerNormalization(epsilon = 1e - 6)(add_attention)  
26    outputs = Dense(units = units, activation = 'relu')(attention2)  
27    outputs = Dense(units = d_model)(outputs)  
28    outputs = Dropout(rate = dropout)(outputs)  
29    add_attention = add([outputs, attention2])  
30    outputs = LayerNormalization(epsilon = 1e - 6)(add_attention)
```

```
31     return Model(
32         inputs = [inputs, enc_outputs, look_ahead_mask, padding_mask],
33         outputs = outputs,
34         name = name)
35 # 测试
36 sample_decoder_layer = decoder_layer(
37     units = 512,
38     d_model = 128,
39     num_heads = 4,
40     dropout = 0.3,
41     name = "sample_decoder_layer")
42 plot_model(sample_decoder_layer, to_file = 'decoder_layer.png', show_shapes = True)
43 # 定义解码器,合成多个解码单元
44 def decoder(vocab_size,
45            num_layers,
46            units,
47            d_model,
48            num_heads,
49            dropout,
50            name = 'decoder'):
51     inputs = Input(shape = (None, ), name = 'inputs')
52     enc_outputs = Input(shape = (None, d_model), name = 'encoder_outputs')
53     look_ahead_mask = Input(shape = (1, None, None), name = 'look_ahead_mask')
54     padding_mask = Input(shape = (1, 1, None), name = 'padding_mask')
55     embeddings = Embedding(vocab_size, d_model)(inputs)
56     embeddings *= Lambda(lambda d_model: tf.math.sqrt( \
57         tf.cast(d_model, tf.float32)))(d_model)
58     embeddings = PositionalEncoding(vocab_size, d_model)(embeddings)
59     outputs = Dropout(rate = dropout)(embeddings)
60     for i in range(num_layers):
61         outputs = decoder_layer(
62             units = units,
63             d_model = d_model,
64             num_heads = num_heads,
65             dropout = dropout,
66             name = 'decoder_layer_{}'.format(i),
67         )(inputs = [outputs, enc_outputs, look_ahead_mask, padding_mask])
68     return Model(
69         inputs = [inputs, enc_outputs, look_ahead_mask, padding_mask],
70         outputs = outputs,
71         name = name)
72 # 解码器测试
73 sample_decoder = decoder(
74     vocab_size = 21128,
75     num_layers = 2,
76     units = 512,
77     d_model = 128,
78     num_heads = 4,
79     dropout = 0.3,
```

```
80     name = "sample_decoder")
81     plot_model(sample_decoder, to_file = 'decoder.png', show_shapes = True)
```

测试解码器程序,观察生成的解码器逻辑结构,与 Transformer 论文解析的结构做对照,在实践中灵活配置解码器的相关参数,可得到适配问题需求的解码器。

5.11 Transformer 模型合成



在前面分步完成的各个模块的基础上,定义 Transformer 的完整模型。编程逻辑如程序源码 P5.6 所示。

程序源码 P5.6 model.py 之 Transformer 定义

```
1  # 定义 Transformer 模型
2  def transformer(vocab_size,
3                 num_layers,
4                 units,
5                 d_model,
6                 num_heads,
7                 dropout,
8                 name = "transformer"):
9      inputs = Input(shape = (None, ), name = "inputs")
10     dec_inputs = Input(shape = (None, ), name = "dec_inputs")
11     enc_padding_mask = Lambda(
12         create_padding_mask, output_shape = (1, 1, None),
13         name = 'enc_padding_mask')(inputs)
14     # 解码器第一个注意力块的前向掩码
15     look_ahead_mask = Lambda(
16         create_look_ahead_mask,
17         output_shape = (1, None, None),
18         name = 'look_ahead_mask')(dec_inputs)
19     # 对编码器输出到解码器第二个注意力块的内容掩码
20     dec_padding_mask = Lambda(
21         create_padding_mask, output_shape = (1, 1, None),
22         name = 'dec_padding_mask')(inputs)
23     enc_outputs = encoder(
24         vocab_size = vocab_size,
25         num_layers = num_layers,
26         units = units,
27         d_model = d_model,
28         num_heads = num_heads,
29         dropout = dropout,
30         )(inputs = [inputs, enc_padding_mask])
31     dec_outputs = decoder(
32         vocab_size = vocab_size,
33         num_layers = num_layers,
34         units = units,
```

```

35         d_model = d_model,
36         num_heads = num_heads,
37         dropout = dropout,
38     )(inputs = [dec_inputs, enc_outputs, look_ahead_mask, dec_padding_mask])
39     outputs = Dense(units = vocab_size, name = "outputs")(dec_outputs)
40     return Model(inputs = [inputs, dec_inputs], outputs = outputs, name = name)
41 # 测试
42 sample_transformer = transformer(
43     vocab_size = 21128,
44     num_layers = 4,
45     units = 512,
46     d_model = 128,
47     num_heads = 4,
48     dropout = 0.3,
49     name = "sample_transformer")
50 plot_model(sample_transformer, to_file = 'transformer.png', show_shapes = True)

```

查看生成的 Transformer 结构图,理解 Transformer 的逻辑运算过程。实践中可灵活调整相关参数配置,以与目标问题相适配。

Transformer 模型定义期间,为了测试各模块程序的逻辑,在每一个模块后面都编写了测试语句。模型程序 model.py 可以迭代运行,观察测试语句的输出结果,加强对模型的理解。其中的 plot_model 函数可以绘制模型结构图,但是需要安装图形支持包 graphviz,根据程序运行时的相关提示,完成环境配置。

Transformer 各模块的逻辑解析及其测试,参见视频教程。



5.12 模型结构与参数配置

从本节到 5.16 节,分步完成聊天机器人模型的训练与评估。打开 5.5 节创建的主程序 main.py,首先完成库的导入和 Transformer 聊天机器人的参数配置和结构定义。编程逻辑如程序源码 P5.7 所示。

程序源码 P5.7 main.py 之模型结构与参数配置

```

1 from tensorflow.keras.callbacks import ModelCheckpoint, EarlyStopping
2 from tensorflow.keras.losses import SparseCategoricalCrossentropy
3 from tensorflow import multiply, minimum
4 from tensorflow.keras.optimizers.schedules import LearningRateSchedule
5 from tensorflow.keras.metrics import sparse_categorical_accuracy
6 from tensorflow.python.ops.math_ops import rsqrt
7 from tensorflow.keras.optimizers import Adam
8 from bert.tokenization.bert_tokenization import FullTokenizer
9 import numpy as np
10 # 用 BLEU 方法评估模型
11 from nltk.translate.bleu_score import sentence_bleu
12 from transformer.model import *

```

```
13 from transformer.dataset import *
14 if __name__ == "__main__":
15     dialog_list = json.loads(\
16         codecs.open("transformer/dataset/dialog_release.json", \
17             "r", "utf-8").read())
18     print(dialog_list[0])           # 第一条数据
19     # 以下参数可根据需要调整,为了演示,可以将相关参数调低一些
20     # 最大句子长度
21     MAX_LENGTH = 40
22     # 最大样本数量
23     MAX_SAMPLES = 120000           # 可根据需要调节
24     BATCH_SIZE = 64                # 批处理大小
25     # Transformer 参数定义
26     NUM_LAYERS = 2                 # 编码器解码器 block 重复数,论文中是 6
27     D_MODEL = 128                  # 编码器解码器宽度,论文中是 512
28     NUM_HEADS = 4                  # 注意力头数,论文中是 8
29     UNITS = 512                    # 全连接网络宽度,论文中输入输出为 512
30     DROPOUT = 0.1                  # 与论文一致
31     VOCAB_SIZE = 21128              # BERT 词典长度
32     START_TOKEN = [VOCAB_SIZE]     # 序列起始标志
33     END_TOKEN = [VOCAB_SIZE + 1]   # 序列结束标志
34     VOCAB_SIZE = VOCAB_SIZE + 2    # 加上开始与结束标志后的词典长度
35     EPOCHS = 50                    # 训练代数
36     bert_vocab_file = 'transformer/dataset/vocab.txt'
37     tokenizer = FullTokenizer(bert_vocab_file)
38     # 聊天模型参数配置与结构定义
39     model = transformer(
40         vocab_size = VOCAB_SIZE,
41         num_layers = NUM_LAYERS,
42         units = UNITS,
43         d_model = D_MODEL,
44         num_heads = NUM_HEADS,
45         dropout = DROPOUT)
46     model.summary()
```

为了满足教学演示需要,程序源码 P5.7 中将 Transformer 编码器与解码器的单元数均缩减为 2,其他参数也有相应缩减,模型可训练参数总量为 9 060 746 个。

模型采用了 BERT 分词方法,故需要安装 BERT 模型库。安装命令为:

```
pip install bert-for-tf2
```

模型采用了两种评价方法:一是计算模型的回归损失;二是计算 BLEU 得分。需要安装 BLEU 函数库。安装命令为:

```
pip install nltk
```



5.13 学习率动态调整

为了优化模型训练过程,加快模型收敛速度,指定了学习率动态调整策略,编码逻辑如程序源码 P5.8 所示。

程序源码 P5.8 main.py 之学习率动态调整

```
1 # 学习率动态调整
2 class CustomSchedule(LearningRateSchedule):
3     def __init__(self, d_model, warmup_steps = 4000):
4         super(CustomSchedule, self).__init__()
5         self.d_model = tf.constant(d_model, dtype = tf.float32)
6         self.warmup_steps = warmup_steps
7     def get_config(self):
8         return {"d_model": self.d_model, "warmup_steps": self.warmup_steps}
9     def __call__(self, step):
10        arg1 = rsqrt(step)
11        arg2 = step * (self.warmup_steps ** -1.5)
12        return multiply(rsqrt(self.d_model), minimum(arg1, arg2))
13 # 测试
14 sample_learning_rate = CustomSchedule(d_model = 256)
15 plt.plot(sample_learning_rate(tf.range(200000, dtype = tf.float32)))
16 plt.ylabel("Learning Rate")
17 plt.xlabel("Train Step")
18 plt.show()
19 learning_rate = CustomSchedule(D_MODEL) # 学习率
```

学习率变化曲线如图 5.11 所示,训练初期学习率采取上升策略以加快训练速度,训练中期保持学习率为一个稳定值,训练后期对学习率采取衰减策略,以期寻找最优解。

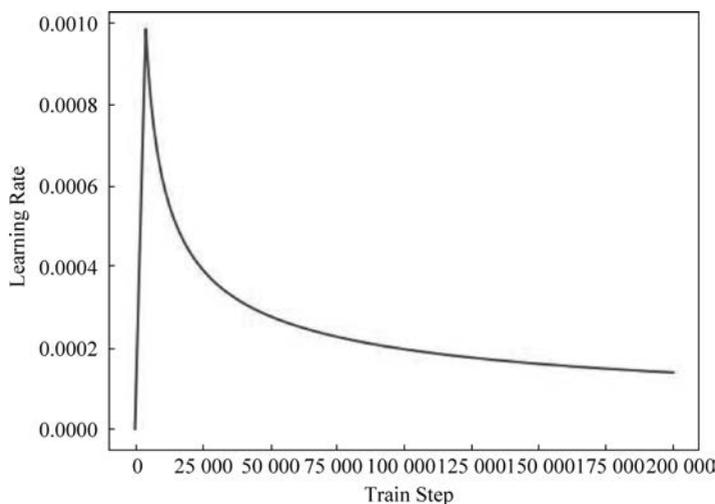


图 5.11 学习率变化曲线

当然,图 5.11 显示本案例跳过了学习率恒定的阶段,在达到最高值后直接开始衰减。

5.14 模型训练过程



模型训练之前,指定模型采用的优化算法为 Adam,定义分类交叉熵损失函数,并定义准确率评价标准,完成模型编译。训练过程中,保存可能取得的最优模型的权重,用提前终止回调函数控制模型训练进程。编码逻辑如程序源码 P5.9 所示。

程序源码 P5.9 main.py 之模型训练过程

```
1 # 定义损失函数
2 def loss_function(y_true, y_pred):
3     y_true = tf.reshape(y_true, shape = (-1, MAX_LENGTH - 1))
4     loss = SparseCategoricalCrossentropy(
5         from_logits = True, reduction = 'none')(y_true, y_pred)
6     mask = tf.cast(tf.not_equal(y_true, 0), tf.float32)
7     loss = tf.multiply(loss, mask)
8     return tf.reduce_mean(loss)
9 # 自定义准确率函数
10 def accuracy(y_true, y_pred):
11     # 调整标签的维度为:(batch_size, MAX_LENGTH - 1)
12     y_true = tf.reshape(y_true, shape = (-1, MAX_LENGTH - 1))
13     return sparse_categorical_accuracy(y_true, y_pred)
14 # 优化算法
15 optimizer = Adam(learning_rate, beta_1 = 0.9, beta_2 = 0.98, epsilon = 1e-9)
16 model.compile(optimizer = optimizer, loss = loss_function, metrics = [accuracy])
17 # 定义回调函数:保存最优模型
18 checkpoint = ModelCheckpoint("robot_weights.h5",
19                               monitor = "val_loss",
20                               mode = "min",
21                               save_best_only = True,
22                               save_weights_only = True,
23                               verbose = 1)
24 # 定义回调函数:提前终止训练
25 earlystop = EarlyStopping(monitor = 'val_loss',
26                            min_delta = 0,
27                            patience = 10,
28                            verbose = 1,
29                            restore_best_weights = True)
30 # 将回调函数组织为回调列表
31 callbacks = [earlystop, checkpoint]
32 dialog_file = 'transformer/dataset/dialog_release.json'
33 train_file = 'transformer/dataset/train.txt'
34 valid_file = 'transformer/dataset/dev.txt'
35 class Hparams():
36     def __init__(self, start_token, end_token, batchSize, total_samples, max_length):
```

```

37         self.start_token = start_token
38         self.end_token = end_token
39         self.batchSize = batchSize
40         self.total_samples = total_samples
41         self.max_length = max_length
42     hparams = Hparams
43     hparams.start_token = START_TOKEN
44     hparams.end_token = END_TOKEN
45     hparams.total_samples = MAX_SAMPLES
46     hparams.batchSize = BATCH_SIZE
47     hparams.max_length = MAX_LENGTH
48     # 加载并划分数据集
49     train_dataset, valid_dataset = get_dataset(hparams, tokenizer, dialog_file,
50     train_file, valid_file)
51     # 模型训练
52     history = model.fit(train_dataset, epochs = EPOCHS, validation_data = valid_dataset,
53     callbacks = callbacks)

```

执行程序源码 P5.9, 开始模型训练。训练结束后, 在当前目录下会保存最佳模型的权重文件 robot_weights.h5。

如果计算机配置不够, 可以先不要考虑模型可用性, 适当降低模型参数配置, 先运行并通过项目逻辑。

本项目在 Kaggle 服务器上的训练结果可以参见链接 <https://www.kaggle.com/code/upsunny/naturalconv-chatbot/notebook>。

如果训练模型的计算机配置过低, 无法完成模型训练时, 可以先将本书素材包中的 robot_weights_l2.h5 模型复制到 MyRobot 的 models 目录下。按照视频教程演示的方法, 调用预训练模型完成后续测试任务。



5.15 损失函数与准确率曲线

绘制模型损失函数曲线与准确率曲线, 有助于观察模型的过拟合情况, 判断模型的泛化能力。编码逻辑如程序源码 P5.10 所示。

程序源码 P5.10 main.py 之损失函数与准确率曲线

```

1     # 损失函数曲线
2     plt.figure(figsize = (12, 6))
3     x = range(1, len(history.history['loss']) + 1)
4     plt.plot(x, history.history['loss'])
5     plt.plot(x, history.history['val_loss'])
6     plt.xticks(x)
7     plt.ylabel('Loss')
8     plt.xlabel('Epoch')
9     plt.legend(['train', 'test'])

```

```
10 plt.title('Loss over training epochs')
11 plt.savefig('loss.png')
12 plt.show()
13 # 准确率曲线
14 plt.figure(figsize = (12, 6))
15 plt.plot(x, history.history['accuracy'])
16 plt.plot(x, history.history['val_accuracy'])
17 plt.ylabel('Accuracy')
18 plt.xlabel('Epoch')
19 plt.xticks(x)
20 plt.legend(['train', 'test'])
21 plt.title('Accuracy over training epochs')
22 plt.savefig('acc.png')
23 plt.show()
```

损失函数曲线如图 5.12 所示。模型在第 14 代之前,损失保持了较快的下降速度。从第 20 代开始,模型优化的幅度不明显,逐渐呈现过拟合趋势。

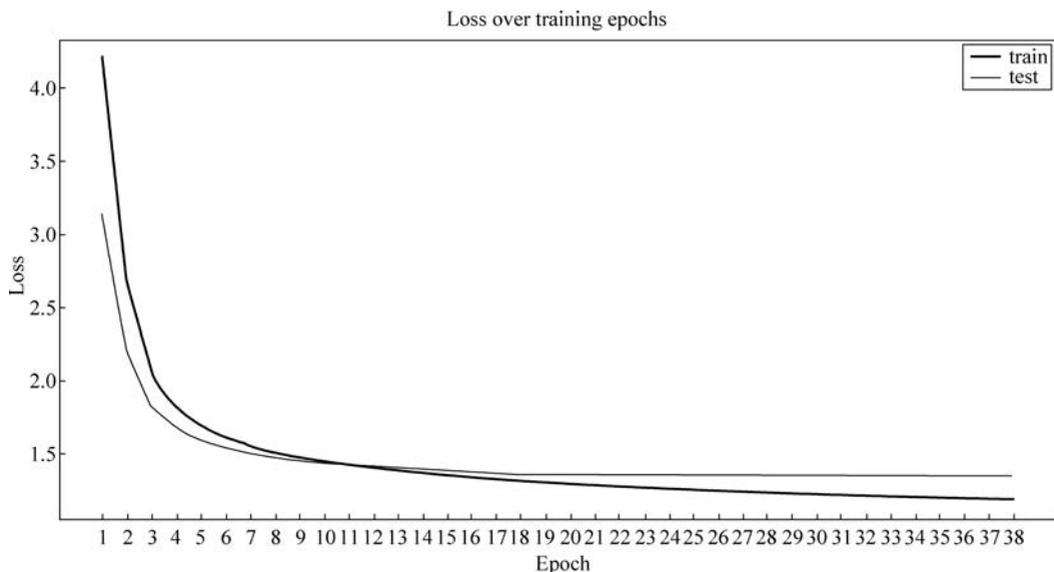


图 5.12 损失函数曲线

准确率曲线如图 5.13 所示,与损失函数曲线基本保持了一致的判断。在第 14 代之前,模型准确率保持较快的增长,第 20 代之后,训练集上的准确率仍保持缓慢增长,验证集上的准确率则几乎保持不变,模型逐渐呈现过拟合趋势。

模型设置了提前结束的条件,如果连续 10 代的损失函数不下降,则提前终止模型训练。这就是为什么设置了 50 代的训练,却在第 38 代终止训练的原因。

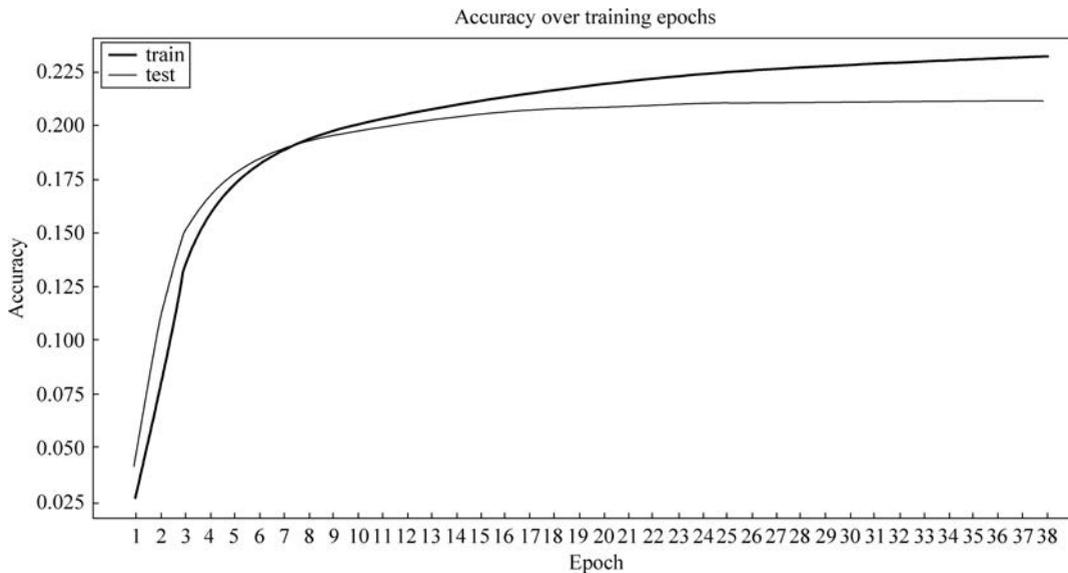


图 5.13 准确率曲线



5.16 聊天模型评估与测试

用 5.15 节训练好的模型做随机对话测试,并用 BLEU 评分观察预测结果。编码逻辑如程序源码 P5.11 所示。

程序源码 P5.11 main.py 之聊天模型评估与测试

```

1  # 加载训练好的模型
2  model.load_weights('models/robot_weights_12.h5')
3  # 用模型做聊天推理,A,B两人聊天,输入A的句子,得到B的回应
4  def evaluate(sentence):
5      sentence = tokenizer.tokenize(sentence)
6      sentence = START_TOKEN + tokenizer.convert_tokens_to_ids(sentence) + END_TOKEN
7      sentence = tf.expand_dims(sentence, axis = 0)
8      output = tf.expand_dims(START_TOKEN, 0)
9      for i in range(MAX_LENGTH):
10         predictions = model(inputs = [sentence, output], training = False)
11         # 选择最后一个输出
12         predictions = predictions[:, -1:, :]
13         predicted_id = tf.cast(tf.argmax(predictions, axis = -1), tf.int32)
14         # 如果是 END_TOKEN 则结束预测
15         if tf.equal(predicted_id, END_TOKEN[0]):
16             break
17         # 把已经得到的预测值串联起来,作为解码器的新输入
18         output = tf.concat([output, predicted_id], axis = -1)
19     return tf.squeeze(output, axis = 0)
20 # 模拟聊天间的问答,输入问话,输出回答

```

```
21 def predict(question):
22     prediction = evaluate(question)           # 调用模型推理
23     predicted_answer = tokenizer.convert_ids_to_tokens(
24         np.array([i for i in prediction if i < VOCAB_SIZE - 2]))
25     print(f'问话者: {question}')
26     print(f'答话者: {" ".join(predicted_answer)}')
27     return predicted_answer
28 # 几组随机测试
29 output1 = predict('嗨,你好呀.')           # 训练集中的样本
30 print("")
31 output2 = predict('昨晚的比赛你看了吗?') # 随机问话 1
32 print("")
33 output3 = predict('你最喜欢的人是谁?')   # 随机问话 2
34 print("")
35 output4 = predict('真热,下点儿雨就好了') # 随机问话 3
36 print("")
37 output5 = predict('这个老师讲课怎么样?') # 随机问话 4
38 print("")
39 output6 = predict('今天收获大吗?')       # 随机问话 5
40 print("")
41 # 多轮对话测试,自问自答
42 sentence = '你最近听说过«中国女排»这部电影吗?'
43 for _ in range(5):
44     sentence = " ".join(predict(sentence))
45     print("")
46 reference = '是呀,我觉得这部«中国女排»应该能有很高的收视率.'
47 pred_sentence = predict(reference)
48 reference = tokenizer.tokenize(reference)
49 # 1-gram BLEU 计算
50 BLEU_1 = sentence_bleu([reference], pred_sentence, weights = (1, 0, 0, 0))
51 print(f"\n BLEU-1 评分: {BLEU_1}")
52 # 2-gram BLEU 计算
53 BLEU_2 = sentence_bleu([reference], pred_sentence, weights = (0.5, 0.5, 0, 0))
54 print(f"\n BLEU-2 评分: {BLEU_2}")
55 # 3-gram BLEU 计算
56 BLEU_3 = sentence_bleu([reference], pred_sentence, weights = (0.33, 0.33, 0.33, 0))
57 print(f"\n BLEU-3 评分:: {BLEU_3}")
58 # 4-gram BLEU 计算
59 BLEU_4 = sentence_bleu([reference], pred_sentence, weights = (0.25, 0.25, 0.25, 0.25))
60 print(f"\n BLEU-4 评分:: {BLEU_4}")
61 # 5-gram BLEU 计算
62 BLEU_5 = sentence_bleu([reference], pred_sentence, weights = (0.2, 0.2, 0.2, 0.2, 0.2))
63 print(f"\n BLEU-5 评分:: {BLEU_5}")
```

程序源码 P5.11 测试结果如表 5.5 所示。假定问话者由人扮演,答话者由机器扮演。自问自答部分由机器同时扮演两个角色。

表 5.5 聊天模型推理抽样观察

简单问答：(一问一答模式)

问话者：嗨，你好呀。

答话者：你好。

问话者：昨晚的比赛你看了吗？

答话者：看了，昨天我看了，昨天晚上的比赛，我看了。

问话者：你最喜欢的人是谁？

答话者：我喜欢的就是那个梅西，他的球技也很好。

问话者：真热，下点儿雨就好了。

答话者：是啊，我也是刚刚看到这个新闻，你看了吗？

问话者：这个老师讲课怎么样？

答话者：我们学校的教师，学校也是很重要的。

问话者：今天收获大吗？

答话者：今天休班。

自问自答：(连续推理测试)

问话者：你最近听说过《中国女排》这部电影吗？

答话者：没有哎，我最近没怎么关注电影。

问话者：没有哎，我最近没怎么关注电影。

答话者：你看了吗？

问话者：你看了吗？

答话者：看了，这部电影的预告片很好看。

问话者：看了，这部电影的预告片很好看。

答话者：是啊，这部电影的主演是谁啊？

问话者：是啊，这部电影的主演是谁啊？

答话者：这部电影是《中国机长》，叫《我的祖国》。

注意：表中的对话解析参见视频教程。

模型还对下面这组问答给出了 BLEU 评分。

问话者：是呀，我觉得这部《中国女排》应该能拿下很高的收视率。

答话者：是呀，这次的世界杯的表现也是非常不错的。

BLEU-1 评分：0.22517932221294598

BLEU-2 评分：0.1332178835716084

BLEU-3 评分：0.09227103858589292

BLEU-4 评分：1.8955151000606497e-78

BLEU-5 评分：1.8662507507148366e-124

事实上，受限于答话句子的长度，BLEU-4 与 BLEU-5 评分没有实际意义。BLEU-1 的得分值表明模型具备一定的可用性与参考性。

直观看，机器的回答是有些跑题，甚至答非所问。但是似乎前后又有一定联系。因为前者说收视率很高，是一个正面评价。后者给出的是对世界杯的正面评价。或许，机器人并不知道如何理解和接续问话者的表达，只是根据自己建模得到的经验做了一个力所能及的回答。

至于为什么这个回答与世界杯有关,而不是与电影有关,一是问话者的语言中包含“中国女排”,这可以理解为体育相关的话题;二是在 5.3 节已经指明,给定的数据集偏重体育语料,会使得训练的模型偏爱体育表达。

事实上,对人类之间的交流而言,这完全不是问题,因为其中的“这部”“收视率”等字眼表明谈论的《中国女排》是一部电影。

5.17 聊天模型部署到服务器



将训练好的 Transformer 聊天模型部署到 Web 服务器上,可以实现一对多的聊天服务。在第 1 章已经搭建了一个基于 Flask 的 Web API 服务框架。在此基础上,迭代追加支持机器人聊天的 Web API 设计。

打开 Server 目录下的 app.py 程序,追加聊天机器人的服务逻辑,如程序源码 P5.12 所示。

程序源码 P5.12 app.py 之聊天模型部署到服务器

```
1 # Transformer 参数定义
2 # 最大句子长度
3 MAX_LENGTH = 40
4 NUM_LAYERS = 2 # 编码器解码器 block 重复数,论文中是 6
5 D_MODEL = 128 # 编码器解码器宽度,论文中是 512
6 NUM_HEADS = 4 # 注意力头数,论文中是 8
7 UNITS = 512 # 全连接网络宽度,论文中输入输出为 512
8 DROPOUT = 0.1 # 与论文一致
9 VOCAB_SIZE = 21128 # 词典长度
10 START_TOKEN = [VOCAB_SIZE] # 序列起始标志
11 END_TOKEN = [VOCAB_SIZE + 1] # 序列结束标志
12 VOCAB_SIZE = VOCAB_SIZE + 2 # 加上开始与结束标志后的词典长度
13 # 分词器
14 bert_vocab_file = '../MyRobot/transformer/dataset/vocab.txt'
15 tokenizer = FullTokenizer(bert_vocab_file)
16 # 模型
17 robot_model = transformer(
18     vocab_size=VOCAB_SIZE,
19     num_layers=NUM_LAYERS,
20     units=UNITS,
21     d_model=D_MODEL,
22     num_heads=NUM_HEADS,
23     dropout=DROPOUT)
24 # 加载权重文件.注意,上面的模型参数必须与权重文件对应的结构一致
25 robot_model.load_weights('../MyRobot/models/robot_weights_12.h5')
26 # 用模型做聊天推理,A,B两人聊天,输入A的句子,得到B的回应
27 def robot_evaluate(sentence):
28     sentence = tokenizer.tokenize(sentence)
29     sentence = START_TOKEN + tokenizer.convert_tokens_to_ids(sentence) + END_TOKEN
30     sentence = tf.expand_dims(sentence, axis=0)
```

```

31     output = tf.expand_dims(START_TOKEN, 0)
32     for i in range(MAX_LENGTH):
33         predictions = robot_model(inputs = [sentence, output], training = False)
34         # 选择最后一个输出
35         predictions = predictions[:, -1:, :]
36         predicted_id = tf.cast(tf.argmax(predictions, axis = -1), tf.int32)
37         # 如果是 END_TOKEN 则结束预测
38         if tf.equal(predicted_id, END_TOKEN[0]):
39             break
40         # 把已经得到的预测值串联起来,作为解码器的新输入
41         output = tf.concat([output, predicted_id], axis = -1)
42     return tf.squeeze(output, axis = 0)
43 # 输入问话,输出回答
44 def robot_predict(question):
45     prediction = robot_evaluate(question)
46     predicted_answer = tokenizer.convert_ids_to_tokens(
47         np.array([i for i in prediction if i < VOCAB_SIZE - 2]))
48     return "".join(predicted_answer)
49 # 机器人聊天 API
50 @app.route('/robot', methods = ['post'])
51 def robot():
52     message = request.get_json(force = True)
53     question = message['question']
54     answer = robot_predict(question)
55     response = {
56         'answer': answer
57     }
58     print(response)
59     return jsonify(response), 200

```

待完成 Android 客户机后,再与服务器做联合测试。



5.18 Android 项目初始化

新建 Android 项目,模板选择 Empty Activity,项目参数设定如图 5.14 所示。项目

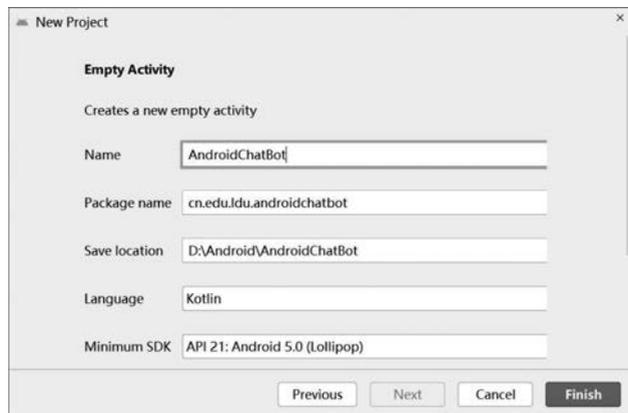


图 5.14 项目参数设定

名称为 AndroidChatBot,包名称为 cn.edu.ldu.androidchatbot,编程语言为 Kotlin,SDK 最小版本号设置为 API 21: Android 5.0(Lollipop),单击 Finish,完成项目初始化。

选择项目根目录,右击,借助快捷菜单命令 New→Package 分别新建 ui、utils、network、data 四个子目录。各子目录的功能及其包含的程序如表 5.6 所示。

表 5.6 项目各子目录的功能及其包含的程序

子目录名	程序名	功能
ui	MainActivity	主控界面逻辑控制
	MessagingAdapter	主控界面数据适配器
utils	Constant	定义全局性常量对象
	Time	定义时间戳对象
network	ApiService	定义网络访问服务接口
data	Answer	匹配服务器应答消息结构的实体类
	Message	消息实体类

依照表 5.6 的提示,依次完成各个程序模块的创建。MainActivity 在项目初始化时已自动生成,将其移动到 ui 目录下即可。项目初始结构如图 5.15 所示。

其他一些简单的初始化工作包括:

(1) 定义实体类 Answer。

```
data class Answer(val answer:String)
```

(2) 定义实体类 Message。

```
data class Message(val message:String, val id:String, val time:String)
```

(3) 在清单文件中声明 Internet 访问权限。

```
<uses-permission android:name="android.permission.INTERNET" />
```

(4) 支持 HTTP 通信。

考虑到本案例采用 HTTP 通信,还需要在 application 节点中添加支持 HTTP 的属性。

```
android:usesCleartextTraffic="true"
```

(5) 添加模块依赖。

在模块配置文件开头的 plugins{} 节点后面添加 Kotlin 扩展语句。

```
apply plugin: 'kotlin-android-extensions'
```

在末尾的 dependencies{} 节点中追加以下依赖库:

```
// RecyclerView
implementation("androidx.recyclerview:recyclerview:1.2.1")
```



图 5.15 项目初始结构

```
// For control over item selection of both touch and mouse driven selection
implementation("androidx.recyclerview:recyclerview-selection:1.1.0")
// Coroutines
implementation 'org.jetbrains.kotlin:kotlin-coroutines-core:1.5.0'
implementation 'org.jetbrains.kotlin:kotlin-coroutines-android:1.5.0'
// Retrofit2
implementation "com.squareup.retrofit2:retrofit:2.9.0"
implementation "com.squareup.retrofit2:converter-scalars:2.9.0"
implementation 'com.squareup.retrofit2:converter-gson:2.9.0'
// Glide
implementation 'com.github.bumptech.glide:glide:4.11.0'
annotationProcessor 'com.github.bumptech.glide:compiler:4.11.0'
```

注意,修改 build.gradle 文件之后,需要同步更新项目配置。系统设计过程中,可根据需要,随时为模块引入相关依赖库。

(6) 定义工具性对象类。

Constant 定义两个常量标识符,用于区别消息发送者和接收者。

```
package cn.edu.ldu.androidchatbot.utils
object Constant {
    const val SEND_ID = "SEND_ID"
    const val RECEIVE_ID = "RECEIVE_ID"
}
```

Time 定义时间戳,表示消息收发时间。

```
package cn.edu.ldu.androidchatbot.utils
import java.sql.Date
import java.sql.Timestamp
import java.text.SimpleDateFormat
object Time {
    fun timeStamp(): String {
        val timeStamp = Timestamp(System.currentTimeMillis())
        val sdf = SimpleDateFormat("HH:mm")
        val time = sdf.format(Date(timeStamp.time))
        return time.toString()
    }
}
```

(7) 定义网络通信模块。

```
package cn.edu.ldu.androidchatbot.network
import okhttp3.RequestBody
import okhttp3.ResponseBody
import retrofit2.Response
import retrofit2.http.Body
```

```
import retrofit2.http.POST
interface ApiService {
    @POST("/robot")
    suspend fun getAnswer(@Body body: RequestBody): Response < ResponseBody >
}
```

修改程序名称为“我的聊天机器人”。接下来,转入界面设计和主控逻辑设计。

5.19 Android 聊天界面设计



聊天界面是程序的主控界面,核心控件是构建聊天列表的 RecyclerView。主控界面由布局文件 activity_main.xml 定义。RecyclerView 中的单行布局由 message_item.xml 定义。

为了美化界面风格,定义两个形状控件,用于渲染聊天消息的背景。send_box.xml 定义的形状用于衬托和突出用户发送的消息。round_box.xml 定义的形状作为发送按钮的背景轮廓。receive_box.xml 定义的形状用于衬托和突出用户收到的消息。

相关资源列表如图 5.16 所示。

选择 drawable 目录,右击,在弹出的快捷菜单中执行 New→Drawable Recourse File 命令,设置文件名称为 receive_box,根元素为 shape,如图 5.17 所示,单击 OK 按钮,创建形状资源文件。

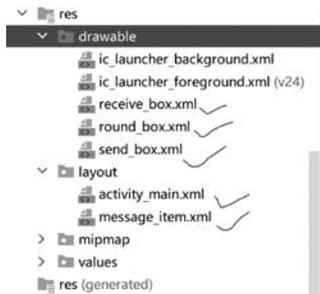


图 5.16 与界面布局相关的资源列表

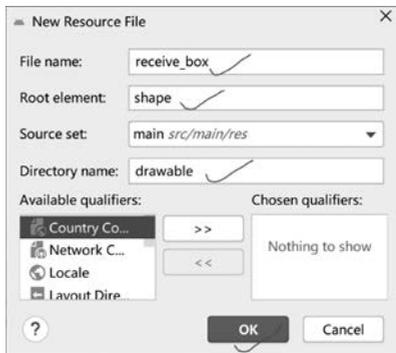


图 5.17 创建 receive_box.xml 文件

替换 receive_box.xml 文件内容。

```
<?xml version="1.0" encoding="utf-8"?>
<shape xmlns:android="http://schemas.android.com/apk/res/android">
    <corners android:topLeftRadius="0dp"
        android:topRightRadius="20dp"
        android:bottomLeftRadius="20dp"
        android:bottomRightRadius="20dp"/>
</shape>
```

这是一个左上角为直角、其他三个角为圆角的矩形框。
用类似的方法定义 send_box.xml。

```
<?xml version = "1.0" encoding = "utf - 8"?>
< shape xmlns:android = "http://schemas.android.com/apk/res/android">
    < corners android:topLeftRadius = "20dp"
        android:topRightRadius = "20dp"
        android:bottomLeftRadius = "20dp"
        android:bottomRightRadius = "0dp" />
</shape >
```

这是一个右下角为直角、其他三个角为圆角的矩形框。

再定义 round_box.xml,这是一个四个角都是圆角的矩形框,用于修饰主界面下方文本框和按钮所在的区域。

```
<?xml version = "1.0" encoding = "utf - 8"?>
< shape xmlns:android = "http://schemas.android.com/apk/res/android">
    < corners android:radius = "20dp"></corners >
</shape >
```

用程序源码 P5.13 所示的脚本,完成主控界面布局。

程序源码 P5.13 activity_main.xml 主控界面布局脚本

```
1 <?xml version = "1.0" encoding = "utf - 8"?>
2 < RelativeLayout xmlns:android = "http://schemas.android.com/apk/res/android"
3     xmlns:tools = "http://schemas.android.com/tools"
4     android:layout_width = "match_parent"
5     android:layout_height = "match_parent"
6     tools:context = ". ui. MainActivity">
7     < LinearLayout
8         android:id = "@ + id/ll_layout_bar"
9         android:layout_width = "match_parent"
10        android:layout_height = "wrap_content"
11        android:layout_alignParentBottom = "true"
12        android:background = "# E4E4E4"
13        android:orientation = "horizontal">
14        < EditText
15            android:id = "@ + id/input_box"
16            android:inputType = "textShortMessage"
17            android:layout_width = "match_parent"
18            android:layout_height = "wrap_content"
19            android:layout_margin = "10dp"
20            android:layout_weight = ". 5"
21            android:background = "@drawable/round_box"
22            android:backgroundTint = "@android:color/white"
23            android:hint = "输入消息..."
24            android:padding = "10dp"
```

```
25         android:singleLine = "true" />
26     <Button
27         android:id = "@ + id/btn_send"
28         android:layout_width = "match_parent"
29         android:layout_height = "match_parent"
30         android:layout_margin = "10dp"
31         android:layout_weight = "1"
32         android:background = "@drawable/round_box"
33         android:backgroundTint = "#26A69A"
34         android:text = "发送"
35         android:textColor = "@android:color/white" />
36 </LinearLayout >
37 <androidx.recyclerview.widget.RecyclerView
38     android:id = "@ + id/chat_view"
39     android:layout_width = "match_parent"
40     android:layout_height = "match_parent"
41     android:layout_above = "@ id/ll_layout_bar"
42     android:layout_below = "@ + id/dark_divider"
43     tools:itemCount = "20"
44     tools:listitem = "@layout/message_item" />
45 <View
46     android:layout_width = "match_parent"
47     android:layout_height = "10dp"
48     android:background = "#42A5F5"
49     android:id = "@ + id/dark_divider" />
50 </RelativeLayout >
```

还有最后一项工作,定义 RecyclerView 中每一行的结构。

选择 layout 目录,右击,在弹出的快捷菜单中执行 New→Layout Resource File 命令,如图 5.18 所示,设定文件名称为 message_item,单击 OK 按钮。

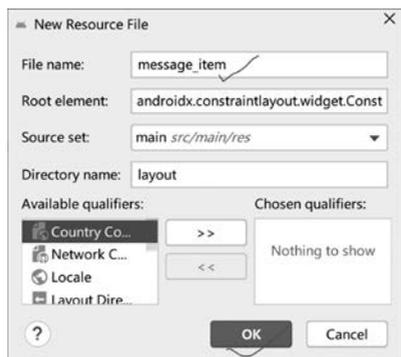


图 5.18 创建单行布局文件

用程序源码 P5.14 所示的脚本程序,定义单行布局。

程序源码 P5.14 message_item.xml 单行布局文件

```
1 <?xml version = "1.0" encoding = "utf - 8"?>
2 <RelativeLayout xmlns:android = "http://schemas.android.com/apk/res/android"
```

```

3     android:layout_width = "match_parent"
4     android:layout_height = "wrap_content">
5     <TextView
6         android:id = "@ + id/tv_message"
7         android:layout_width = "200dp"
8         android:layout_height = "wrap_content"
9         android:layout_margin = "4dp"
10        android:background = "@drawable/send_box"
11        android:backgroundTint = "# 26A69A"
12        android:padding = "14dp"
13        android:text = "发出的消息"
14        android:textColor = "@android:color/white"
15        android:textSize = "18sp"
16        android:layout_alignParentEnd = "true" />
17    <TextView
18        android:visibility = "visible"
19        android:id = "@ + id/tv_bot_message"
20        android:layout_width = "200dp"
21        android:layout_height = "wrap_content"
22        android:layout_margin = "4dp"
23        android:background = "@drawable/receive_box"
24        android:backgroundTint = "# FF7043"
25        android:padding = "14dp"
26        android:text = "收到的消息"
27        android:textColor = "@android:color/white"
28        android:textSize = "18sp"
29        android:layout_alignParentStart = "true" />
30 </RelativeLayout >

```

运行主程序,观察主控界面效果。



5.20 Android 聊天逻辑设计

主控逻辑包括两部分:一部分放在 ActivityMain 中;另一部分放在 MessagingAdapter 中。在 ui 包目录下新建消息适配器程序 MessagingAdapter,编码逻辑如程序源码 P5.15 所示。

程序源码 P5.15 MessagingAdapter 消息适配器

```

1 package cn.edu.ldu.androidchatbot.ui
2 import android.view.LayoutInflater
3 import android.view.View
4 import android.view.ViewGroup
5 import androidx.recyclerview.widget.RecyclerView
6 import cn.edu.ldu.androidchatbot.R
7 import cn.edu.ldu.androidchatbot.data.Message
8 import cn.edu.ldu.androidchatbot.utils.Constant.RECEIVE_ID

```

```
9 import cn.edu.ldu.androidchatbot.utils.Constant.SEND_ID
10 import kotlinx.android.synthetic.main.message_item.view.*
11 class MessagingAdapter: RecyclerView.Adapter<MessagingAdapter.MessageViewHolder>() {
12     var messageList = mutableListOf<Message>()
13     inner class MessageViewHolder(itemView: View): RecyclerView.ViewHolder(itemView) {
14         init {
15             itemView.setOnClickListener {
16                 messageList.removeAt(adapterPosition)
17                 notifyItemRemoved(adapterPosition)
18             }
19         }
20     }
21     override fun onCreateViewHolder(parent: ViewGroup,
22                                     viewType: Int): MessageViewHolder {
23         return MessageViewHolder(
24             LayoutInflater.from(parent.context).inflate(
25                 R.layout.message_item, parent, false))
26     }
27     override fun onBindViewHolder(holder: MessageViewHolder, position: Int) {
28         val currentMessage = messageList[position]
29         when (currentMessage.id) {
30             SEND_ID -> {
31                 holder.itemView.tv_message.apply {
32                     text = currentMessage.message
33                     visibility = View.VISIBLE
34                 }
35                 holder.itemView.tv_bot_message.visibility = View.GONE
36             }
37             RECEIVE_ID -> {
38                 holder.itemView.tv_bot_message.apply {
39                     text = currentMessage.message
40                     visibility = View.VISIBLE
41                 }
42                 holder.itemView.tv_message.visibility = View.GONE
43             }
44         }
45     }
46     override fun getItemCount(): Int {
47         return messageList.size
48     }
49     fun insertMessage(message: Message) {
50         this.messageList.add(message)
51         notifyItemInserted(messageList.size)
52         notifyDataSetChanged()
53     }
54 }
```

ActivityMain 编码逻辑如程序源码 P5.16 所示,负责与用户的交互,包括收发消息及界面滚动。收发消息均通过后台协程实现。

程序源码 P5.16 ActivityMain 主控逻辑

```

1 package cn.edu.ldu.androidchatbot.ui
2 import androidx.appcompat.app.AppCompatActivity
3 import android.os.Bundle
4 import androidx.recyclerview.widget.LinearLayoutManager
5 import cn.edu.ldu.androidchatbot.R
6 import cn.edu.ldu.androidchatbot.data.Answer
7 import cn.edu.ldu.androidchatbot.data.Message
8 import cn.edu.ldu.androidchatbot.network.ApiService
9 import cn.edu.ldu.androidchatbot.utils.Constant.RECEIVE_ID
10 import cn.edu.ldu.androidchatbot.utils.Constant.SEND_ID
11 import cn.edu.ldu.androidchatbot.utils.Time
12 import com.google.gson.Gson
13 import kotlinx.android.synthetic.main.activity_main.*
14 import kotlinx.coroutines.Dispatchers
15 import kotlinx.coroutines.GlobalScope
16 import kotlinx.coroutines.launch
17 import kotlinx.coroutines.withContext
18 import okhttp3.MediaType
19 import okhttp3.RequestBody
20 import org.json.JSONObject
21 import retrofit2.Retrofit
22 import retrofit2.converter.gson.GsonConverterFactory
23 class MainActivity : AppCompatActivity() {
24     // 腾讯服务器教学演示地址
25     private val BASE_URL = "http://120.53.107.28"
26     // private val BASE_URL = "http://192.168.0.103:5000" // 本地服务器地址
27     private val retrofit = Retrofit.Builder() // 初始化 Retrofit 框架
28         .addConverterFactory(GsonConverterFactory.create())
29         .baseUrl(BASE_URL)
30         .build()
31         .create(ApiService::class.java)
32     private lateinit var adapter: MessagingAdapter
33     var messageList = mutableListOf<Message>() // 消息列表
34     override fun onCreate(savedInstanceState: Bundle?) {
35         super.onCreate(savedInstanceState)
36         setContentView(R.layout.activity_main)
37         recyclerView() // 列表视图
38         clickEvents() // 单击事件
39         customMessage("你好,很高兴见到你!") // 欢迎语
40     }
41     private fun clickEvents() {
42         btn_send.setOnClickListener { // 发送消息单击事件
43             sendMessage()
44         }
45         input_box.setOnClickListener { // 输入消息事件
46             GlobalScope.launch {
47                 withContext(Dispatchers.Main){

```

```
48         chat_view.scrollToPosition(adapter.itemCount - 1)
49     }
50 }
51 }
52 }
53 override fun onStart() {
54     super.onStart()
55     GlobalScope.launch {
56         withContext(Dispatchers.Main) {
57             chat_view.scrollToPosition(adapter.itemCount - 1)
58         }
59     }
60 }
61 private fun recyclerView() { // 视图绑定到消息适配器
62     adapter = MessagingAdapter()
63     chat_view.adapter = adapter
64     chat_view.layoutManager = LinearLayoutManager(applicationContext)
65 }
66 private fun sendMessage() { // 发送消息
67     val message = input_box.text.toString()
68     val timeStamp = Time.timeStamp()
69     if (message.isNotEmpty()) {
70         messageList.add(Message(message, SEND_ID, timeStamp))
71         input_box.setText("")
72         adapter.insertMessage(Message(message, SEND_ID, timeStamp))
73         chat_view.scrollToPosition(adapter.itemCount - 1)
74         botResponse(message)
75     }
76 }
77 private fun botResponse(message: String) { // 接收消息
78     val timeStamp = Time.timeStamp()
79     GlobalScope.launch(Dispatchers.Main) {
80         val request = JSONObject()
81         request.put("question", message)
82         val body: RequestBody =
83             RequestBody.create(
84                 MediaType.parse("application/json"),
85                 request.toString()
86             )
87         val response = retrofit.getAnswer(body)
88         if (response.isSuccessful) {
89             val json: String = response.body()!!.string()
90             var gson = Gson()
91             var reply = gson.fromJson(
92                 json,
93                 Answer::class.java
94             )
95             messageList.add(Message(reply.answer, RECEIVE_ID, timeStamp))
96             adapter.insertMessage(Message(reply.answer, RECEIVE_ID, timeStamp))
97             chat_view.scrollToPosition(adapter.itemCount - 1)
98         } else { }
99     }
```

```

100     }
101     private fun customMessage(message: String) {           // 自定义欢迎消息
102         GlobalScope.launch {
103             val timeStamp = Time.timeStamp()
104             withContext(Dispatchers.Main){
105                 messageList.add(Message(message, RECEIVE_ID, timeStamp))
106                 adapter.insertMessage(Message(message, RECEIVE_ID, timeStamp))
107                 chat_view.scrollToPosition(adapter.itemCount - 1)
108             }
109         }
110     }
111 }

```

程序解析参见本节视频教程。



5.21 客户机与服务器联合测试

现在可以开始项目联合测试了。先运行 Web 服务器,再运行客户机。可以用模拟器测试,也可以用真机测试。既可以本地测试,也可以远程测试。

如果做本地测试,首先运行 Server 目录下的服务器程序 app.py,观察控制台上输出的服务器地址,将程序源码 P5.16 中第 26 行语句中的地址修改为 Web API 运行的地址,注释掉第 25 行中的远程服务器地址,然后做客户机与服务器的本地联合测试。

如果做远程测试,可以采用书中的腾讯服务器地址,直接运行 Android 客户机程序,在模拟器与真机上与远程服务器做人机畅聊测试。图 5.19 给出了一组随机对话,显示了



(a) 与本地服务器测试



(b) 与远程服务器测试

图 5.19 真机与本地服务器和远程服务器通信的测试效果

用真机分别与本地服务器和远程服务器通信的测试效果。

至此,本章实现的聊天机器人已经初露锋芒。无论其聊天水平怎么样,对于一个仅有900万训练参数的小模型,依靠一个超小规模语料库,实现的语言对答还是令人惊讶的。今天的一小步,孕育着未来的一大步。

事实上,图5.19展示的人机聊天,是故意增加了难度的。首先,语料库中关于健康和航天类的语料是偏少的。其次,对于第一组聊天,涉及奥密克戎这个新生词汇,超过了模型的“知识范畴”;对于第二组聊天,航天的话题无疑也超过了模型的“知识范畴”。

可以试试体育类的话题,机器人的回答则会生动得多,有趣得多。做完这个项目,逗着自己缔造的机器人玩玩,成就感与幸福感如期而至。原来,是如此容易让机器说话的呀。更多测试细节,参见本节视频教程。

5.22 小结



本章以人机畅聊的境界追求为动力,遵循机器问答的技术设计路线,完成了Transformer聊天模型+Web API+Android聊天客户机的项目设计。本章项目对于提升读者在自然语言智能领域的理论水平和实践能力,具备非常好的教学示范效果。

沿着语料库分析,Transformer模型解析,聊天模型建模、训练、评估,Web服务器模型部署以及Android客户机设计这些环环相扣的步骤,读者可以体验到学习过程中的循序渐进,体验到量变到质变,体验到从混沌到顿悟、从顿悟到彻悟的华丽蜕变。

机器是如何学习说话的?机器是如何学会说话的?机器说话的本质是什么?现阶段局限是什么?所有这些不再是秘密。

5.23 习题

1. 简要描述机器问答与机器聊天的不同。
2. 简要描述实施机器问答项目的方法与步骤。
3. 腾讯自然语言聊天数据集 NaturalConv 有哪些特点?
4. Transformer 的创新点有哪些?
5. 简要解析 Transformer 单头注意力的计算逻辑。
6. Transformer 模型适合哪些场景的应用?
7. 结合项目实战,谈谈为什么本章实现的聊天模型偏爱体育方面的表达。
8. 根据 Transformer 多头注意力机制的编程设计,绘制其计算逻辑流程图。
9. 绘制 Transformer 单层编码器的工作流程图。
10. 绘制 Transformer 单层解码器的工作流程图。
11. 模型训练期间,采用动态学习率调整策略有何优势?
12. Transformer 聊天机器人中的损失函数是如何定义的?

13. 描述 Transformer 聊天机器人在 Web 服务器上的部署流程。
14. 简述 Android 聊天机器人的客户机界面布局设计。
15. 绘图说明 Android 聊天机器人的客户机运行逻辑。
16. 为什么 Android 客户机的收发消息流程需要定义到协程中？
17. 结合本章项目实战,谈谈机器是如何学习说话的、机器是如何学会说话的、机器说话的本质是什么。