

一个应用实例，介绍在 NLP 应用中，从数据处理、到模型构建、再到模型训练、最后使用模型进行预测的完整过程。

## 9.6 应用实例：微博评论谣言分析——监测网络舆情

随着互联网的飞速发展，层出不穷的媒体、论坛、微博等平台已经成为了大众获取信息的主要渠道，这样一种新的信息传播形式，早已经深入了网友们的日常生活之中。习近平总书记指出“网民来自老百姓，老百姓上了网，民意也就上了网”。在新媒体时代，一方面，大量网民可以在网络上讨论和转播社会热点事件，网络成为表达民意的重要渠道，助推了舆情的产生；另一方面，网络的虚拟性和用户主体隐匿性意味着网络中不可避免的会存在大量与客观现实不符合的谣言。



应用实例

网络谣言的泛滥，可以对个人、社会和文化造成广泛的负面影响。网络谣言会传播虚假信息，误导人们相信不准确的事实，这可能导致人们做出错误的决策。网络谣言也常常伴随着对个人、组织或公司的恶意攻击，导致声誉受损。网络谣言还可能引发社交不安定，造成群体恐慌、骚乱或冲突。当人们无法分辨真假信息时，会导致社会的信任度下降，破坏了媒体和政府的权威，使人们对权威机构的信息产生怀疑。

因此，在新媒体时代，无论对于政府还是企业，网络舆情监测都应是一项非常重要的工作。党的十八大以来，我国高度重视网络舆情的治理，并做出一系列重大部署，在党的二十大报告中也明确指出要“加强全媒体传播体系建设，塑造主流舆论新格局。健全网络综合治理体系，推动形成良好网络生态。”

人工智能和深度学习技术提供了一种有效的方式来自动化、规模化地进行网络舆情检测，并提高其效率和准确性，在数据处理和信息抽取、情感分析、文本分类和主题建模、舆情预测、虚假信息和恶意评论的检测等方面发挥了重要的作用。可以帮助社交媒体平台、新闻机构和政府部门更有效地应对网络谣言的传播和危害，更好地维护信息的真实性和公共利益。

本案例采用循环神经网络对网络谣言进行识别，案例使用的数据集是从新浪微博不实信息举报平台抓取的中文谣言数据，数据集中共包含 1538 条谣言和 1849 条非谣言，存储于 `wb_rumor.txt` 文件中，文件的结构如图 9-26 所示。

```

0 【人间惨剧】不能忘却的历史：1984年12月3日晨，印度博帕尔市的联合碳化物有限公司一P-X项目发生氟化物泄漏，覆盖方圆60公里。2.5万人直接致死，55万人间接致死，另外
有20多万人永久残废的人间惨剧。现在，当地居民仍然要忍受癌症及胎儿畸形的折磨。@宁波发布
0 红薯+“墨水”=紫薯！！！无良奸商！！下图是一小块“来伊份”牌香脆紫薯浸泡一夜后脱落的“墨水”！食品啊，杀人！@逊别拉 @蔡康永 @李厚霖 @药大颠儿 @何昊 @乐嘉 @
刘大鸿 @源家姐姐 @Tracy---liu
0 【20日，洛阳一名约50岁裸女站在一辆法院的车上大喊“苍天啊”，引发市民围观与。据说该女子弟弟被法院人员杀害，上诉多年，但是多部门一直推诿扯皮。无奈之下，只好出此
下策。】-----点！“该有多大的绝望、有多大的冤屈才能让一个母亲不顾尊严、不顾羞耻而这样做！”@杂谈五味 @观点联播V
0 【超市手推车插入式摆放宛若同性爱姿势，专家建议取缔】教育专家王建立一直以敢说真话闻名，近日他指出，超市的手推车从后面一个插入的摆放方式，让他看到就脸红。
这样的摆放容易让人想入非非，孩子去一次超市就相当于看一次A片。王专家进一步指出，手推车还是更恶俗的群P，危害无穷，必须取缔！
1 【演出预告】青春版越剧《胭脂》，中国越剧最新一代男女合演，迈出外地巡演脚步，第一站：4月2日余姚市凤山街道永丰村；第二站：5月24日#绍兴大剧院#，走出杭城面向
全国，传递青春活力，敬请支持！一起期待下一站[心][给力]转播FC提供的仿#爱情公寓#胭脂预告，娱你一笑http://t.cn/zYV236d @新浪戏剧
1 一个人的成熟，在思想里。一个人的天真，在眼神里。成熟和天真，并非一对矛盾词。成熟是一种生活态度，天真是一种生活方式。每张面孔都会变老，但生活方式可以年轻一辈
子。上帝若是爱你，会让你做一生的小孩子。晚安！
1 不要想太多，你的想象力会 创造出个一开始并不存在的问题。
1 #肯定没人理# 小透明啥都不会就最后几分钟凑个热闹。转发的姑娘，接下去我每到一地方都会寄明信片给你，大概寄到明年光棍节为止吧。转发有效期是明天上午九点之前。
今年去过的地方有乌镇、杭州、武汉、沈阳、长白山、牡丹江、哈尔滨、北京、横店、温州、武义。我会私信你拿地址。
1 钱多了，友少了；车多了，路少了；娱多了，乐少了；说多了，假多了；真少了；城多了，村少了；虚多了，实少了；贫多了，廉少了；话多了，理少了；梦多了，觉少
了；雾多了，雪少了；草多了，树少了；恨多了，爱少了；黑多了，白少了；暗多了，明少了；衣多了，穿少了。
1 《好好先生作品集》在#那些黄金年代# 里，Warwick Goble（1862-1943）的风格可以算是绝对的“童叟无害”，我们几乎看不到他画沉重、阴郁的题材；即便是画了，也被他清

```

图9-26 训练集数据内容和格式

从图 9-26 可以看出，每一条数据包含两部分：第一部分为标签，值为 0 或 1，0 表示谣言，1 表示非谣言；第二部分为原文的文字内容。首先加载数据集，代码如下。

代码9-5 加载数据集

```

# 初始化两个空列表，用于存储标签和文本内容
labels = []
dataset = []

# 读取原始数据文件
file_path = "wb_rumor.txt"
with open(file_path, "r", encoding="utf-8") as file:
    for line in file:
        label, text = line.strip().split('\t')
        labels.append(int(label))
        dataset.append(text)

# 输出数据条目数
print(f"数据条目数为: {len(dataset)}")

# 输出示例数据
for i in range(10): # 输出前 10 条数据作为示例
    print(f"第{i}条数据, Label: {labels[i]}, Text: {dataset[i]}")

```

以上代码的运行结果如图 9-27 所示，显示了数据集数据总条目数以及前 10 条数据。

数据条目数为：3387

第0条数据, Label: 0, Text: 【人间惨剧】不能忘却的历史：1984年12月3日晨，印度博帕尔市的联合碳化物有限公司一P-X项目发生氟化物泄漏，覆盖方圆60公里。2.5万人直接致死，55万人间接致死，另外有20多万人永久残废的人间惨剧。现在，当地居民仍然要忍受癌症及胎儿畸形的折磨。@宁波发布

第1条数据, Label: 0, Text: 红薯+“墨水”=紫薯！！！无良奸商！！下图是一小块“来伊份”牌香脆紫薯浸泡一夜后脱落的“墨水”！食品啊，杀人！@逊别拉 @蔡康永 @李厚霖 @药大颠儿 @何昊 @乐嘉 @刘大鸿 @源家姐姐 @Tracy---liu

第2条数据, Label: 0, Text: 【20日，洛阳一名约50岁裸女站在一辆法院的车上大喊“苍天啊”，引发市民围观与。据说该女子弟弟被法院人员杀害，上诉多年，但是多部门一直推诿扯皮。无奈之下，只好出此下策。】-----点！“该有多大的绝望、有多大的冤屈才能让一个母亲不顾尊严、不顾羞耻而这样做！”@杂谈五味 @观点联播V

第3条数据, Label: 0, Text: 【超市手推车插入式摆放宛若同性爱姿势，专家建议取缔】教育专家王建立一直以敢说真话闻名，近日他指出，超市的手推车从后面一个插入的摆放方式，让他看到就脸红。这样的摆放容易让人想入非非，孩子去一次超市就相当于看一次A片。王专家进一步指出，手推车还是更恶俗的群P，危害无穷，必须取缔！

第4条数据, Label: 1, Text: 【演出预告】青春版越剧《胭脂》，中国越剧最新一代男女合演，迈出外地巡演脚步，第一站：4月2日余姚市凤山街道永丰村；第二站：5月24日#绍兴大剧院#，走出杭城面向全国，传递青春活力，敬请支持！一起期待下一站[心][给力]转播FC提供的仿#爱情公寓#胭脂预告，娱你一笑<http://t.cn/zYV236d> @新浪戏剧

第5条数据, Label: 1, Text: 一个人的成熟，在思想里。一个人的天真，在眼神里。成熟和天真，并非一对矛盾词。成熟是一种生活态度，天真是一种生活方式。每张面孔都会变老，但生活方式可以年轻一辈子。上帝若是爱你，会让你做一生的小孩子。晚安！

第6条数据, Label: 1, Text: 不要想太多。你的想象力会 创造出个一开始并不存在的问题。

第7条数据, Label: 1, Text: #肯定没人理# 小透明啥都不会就最后几分钟凑个热闹。转发的姑娘，接下去我每到一地方都会寄明信片给你，大概寄到明年光棍节为止吧。转发有效期是明天上午九点之前。今年去过的地方有乌镇、杭州、武汉、沈阳、长白山、牡丹江、哈尔滨、北京、横店、温州、武义。我会私信你拿地址。

第8条数据, Label: 1, Text: 钱多了，友少了；车多了，路少了；娱多了，乐少了；说多了，假多了；真少了；城多了，村少了；虚多了，实少了；贫多了，廉少了；话多了，理少了；梦多了，觉少了；雾多了，雪少了；草多了，树少了；恨多了，爱少了；黑多了，白少了；暗多了，明少了；衣多了，穿少了。

第9条数据, Label: 1, Text: 《好好先生作品集》在#那些黄金年代# 里，Warwick Goble（1862-1943）的风格可以算是绝对的“童叟无害”，我们几乎看不到他画沉重、阴郁的题材；即便是画了，也被他清新的“Fairy Style”给稀释了。该作品包收录了包括成名作《水宝宝》的9个绘本，文件内附电子书地址共179幅。戳<http://t.cn/8syw3qy>

图9-27 数据集信息

在图 9-27 中矩形框标注部分可以看到，文本中包含形如“@何炅 @乐嘉 @刘大鸿 @源家姐姐 @Tracy---liu”的转发标记、形如“#肯定没人理#”的话题标签和形如“http://t.cn/zYV23Gd”的 URL，这些对分析文本的是否为谣言没有帮助，我们可以利用正则表达式过滤掉这些内容，代码如下。



正则表达式

代码9-6 使用正则表达式过滤不相关内容

```
import re

# 定义一个函数来删除文本中的 URL
def filter_text(text):
    text = re.sub(r'@[^\w]+', '', text) # 去掉转发标记
    text = re.sub(r'#[^#]+', '', text) # 去掉话题标签
    text = re.sub(r'http[s]?://\S+', '', text) # 去掉 URL
    return text

# 对每个文本应用 filter_text 函数
dataset = [filter_text(text) for text in dataset]

# 输出示例数据
for i in range(10): # 输出前 10 条数据作为示例
    print(f"第{i}条数据, Label: {labels[i]}, Text: {dataset[i]}")
```

以上代码的输出结果如图 9-28 所示，可以看到，三种类型的不相关内容全部从文本中删除。

第0条数据, Label: 0, Text: 【人间惨剧】不能忘却的历史: 1984年12月3日晨, 印度博帕尔市的联合碳化物有限公司一P-X项目发生氰化物泄漏, 覆盖方圆60公里。2.5万人直接致死, 55万人间接致死, 另外有20多万人永久残废的人间惨剧。现在, 当地居民仍然要忍受癌症及胎儿畸形的折磨。

第1条数据, Label: 0, Text: 红薯+“墨水”=紫薯!!! 无良奸商!!! 下图是一小块“来伊份”牌香脆紫薯浸泡一夜后脱落的“墨水”! 食品啊, 杀人! ---liu

第2条数据, Label: 0, Text: 【20日, 洛阳一名约50岁裸女站在一辆法院的车上大喊“苍天啊”, 引发市民围观与。据说该女子弟弟被法院人员杀害, 上诉多年, 但是多部门一直推诿扯皮。无奈之下, 只好出此下策。】-----、点! “该有多大的绝望、有多大的冤屈才能让一个母亲不顾尊严、不顾羞耻而这样做!”

第3条数据, Label: 0, Text: 【超市手推车插入式摆放宛若同性爱姿势, 专家建议取缔】教育专家王建立一直以敢说真话闻名, 近日他指出, 超市的手推车从后面一个插入的摆放方式, 让他看到就脸红。这样的摆放容易让人想入非非, 孩子去一次超市就相当于看一次A片。王专家进一步指出, 手推车还是更恶俗的群P, 危害无穷, 必须取缔!

第4条数据, Label: 1, Text: 【演出预告@】青春版越剧《胭脂》· 中国越剧最新一代男女合演, 迈出外地巡演脚步, 第一站: 4月2日余姚市凤山街道永丰村; 第二站: 5月24日, 走出杭城面向全国, 传递青春活力, 敬请支持! 一起期待下一站[心][给力]转播FC提供的仿胭脂预告, 娱乐一乐

第5条数据, Label: 1, Text: 一个人的成熟, 在思想里。一个人的天真, 在眼神里。成熟和天真, 并非一对矛盾词。成熟是一种生活态度, 天真是一种生活方式。每张面孔都会变老, 但生活方式可以年轻一辈子。上帝若是爱你, 会让你做一生的大孩子。晚安!

第6条数据, Label: 1, Text: 不要想太多, 你的想象力会 创造出个一开始并不存在的问题。

第7条数据, Label: 1, Text: 小透明啥都不会就最后几分钟凑个热闹。转发的姑娘, 接下去我每到一个地方都会寄明信片给你, 大概寄到明年光棍节为止吧。转发有效期是明天上午九点之前。今年去过的地方有乌镇、杭州、武汉、沈阳、长白山、牡丹江、哈尔滨、北京、横店、温州、武义。我会私信你拿地址。

第8条数据, Label: 1, Text: 钱多了, 友少了; 车多了, 路少了; 娱多了, 乐少了; 说多了, 做少了; 假多了, 真少了; 城多了, 村少了; 虚多了, 实少了; 贫多了, 廉少了; 话多了, 理少了; 梦多了, 觉少了; 雾多了, 雾少了; 草多了, 树少了; 爱多了, 爱少了; 黑多了, 白少了; 暗多了, 明少了; 衣多了, 穿少了。

第9条数据, Label: 1, Text: 《好好先生作品集》在 里, Warwick Goble (1862-1943) 的风格可以算是绝对的“重男轻女”, 我们几乎看不到他画沉重、阴郁的题材; 即便是画了, 也被他清新的“Fairy Style”给稀释了。该作品包收录了包括成名作《水宝宝》的9个绘本, 文件内附电子书地址共179幅。戳→

图9-28 过滤文本中不相关内容

接下来，我们将数据集划分为训练集和测试集，比例为 7:3，代码如下。

代码9-7 拆分数据集

```
# 定义数据集列表
train_dataset = []
train_labels = []
test_dataset = []
test_labels = []

# 数据集的总长度
total_samples = len(labels)
```

```
# 计算各部分数据集的长度
train_size = int(total_samples * 0.7)
test_size = total_samples - train_size

# 将数据划分为训练集、验证集和测试集
train_dataset = dataset[:train_size]
train_labels = labels[:train_size]

test_dataset = dataset[train_size:]
test_labels = labels[train_size:]

print(f"训练集数据量为: {len(train_dataset)}")
print(f"测试集数据量为: {len(test_dataset)}")
```

以上代码的输出结果如图 9-29 所示，训练集包含 2370 条数据，测试集包含 1017 条数据。

训练集数据量为: 2370  
测试集数据量为: 1017

图9-29 拆分数据集

然后，我们建立词汇表，代码如下。

代码9-8 建立词汇表

```
import jieba
from torchtext.vocab import build_vocab_from_iterator

# 对迭代器中的文本进行分词处理
tokenized_iterator = (jieba.cut(sample) for sample in train_dataset)

# 使用 `build_vocab_from_iterator` 函数构建词汇表
vocab = build_vocab_from_iterator(tokenized_iterator, specials=["<unk>"])
vocab.set_default_index(vocab["<unk>"]) # 设置默认索引，如果找不到单词，则会选择默认索引

# 输出词汇表长度
print("词汇表长度:", len(vocab))

# 输出示例词汇对应编码
print(vocab(['问题', '发生']))
```

以上代码中，`tokenized_iterator = (jieba.cut(sample) for sample in train_dataset)`使用了一个生成器表达式，对 `train_dataset` 中的每一个样本进行 jieba 分词。生成器表达式将会创建一个迭代器，该迭代器在每次迭代时都会返回一个分词后的样本。

通过使用 `torchtext` 的 `build_vocab_from_iterator` 函数从训练数据集构建词汇表。`TorchText` 是 `PyTorch` 的一个 NLP 库，提供了一些主要的 NLP 构建块，如预训练词向量、文本处理工具以及常用的 NLP 数据集。`build_vocab_from_iterator` 函数从分词后的样本中构建词汇表，它接收一个迭代

器作为输入，并将迭代器中的所有唯一的词语收集到词汇表中。参数 `specials=["<unk>"]` 指定了一些特殊的词语，这些词语将会被添加到词汇表中。在这个例子中，只添加了一个特殊词语 "<unk>"，用于表示词汇表中不存在的词语。`vocab.set_default_index` 设置了词汇表的默认索引。如果在将文本转换为词汇索引时遇到词汇表中不存在的词语，那么就会返回默认索引。

代码 9-7 的运行结果如图 9-30 所示，显示了词汇表的长度以及“问题”和“发生”两个词的索引值。

词汇表长度：22815  
[219, 161]

图9-30 词汇表信息

接下来，我们将数据集中的文本转换为词汇表中对应的索引形式，代码如下。

代码9-9 转换文本

```
train_dataset_indices = []
test_dataset_indices = []

# 转换索引
for data in train_dataset:
    train_dataset_indices.append([vocab[token] for token in
jieba.cut(data)])
for data in test_dataset:
    test_dataset_indices.append([vocab[token] for token in jieba.cut(data)])

# 输出样例
print(train_dataset[0])
print(train_dataset_indices[0])
```

以上代码的输出结果如图 9-31 所示，将训练集和测试集中的文本数据转换为其在词汇表中的索引形式。

【人间惨剧】不能忘却的历史：1984年12月3日晨，印度博帕尔市的联合碳化物有限公司一P-X项目发生氟化物泄漏，覆盖方圆60公里。2.5万人直接致死，55万人间接致死，另外有20多万人永久残疾的人间惨剧。现在，当地居民仍然要忍受癌症及胎儿畸形的折磨。

[12, 1385, 2188, 13, 153, 8267, 2, 554, 11, 6734, 53, 99, 29, 54, 4521, 1, 738, 7578, 1575, 2, 1793, 9203, 1596, 71, 3261, 62, 2737, 1172, 161, 8841, 886  
9, 1, 3163, 8575, 579, 1688, 3, 4997, 1519, 365, 1627, 1, 2342, 1519, 3878, 1627, 1, 2127, 18, 177, 1904, 17, 6053, 4599, 2, 1385, 2188, 3, 150, 1, 561,  
1573, 4063, 50, 2521, 901, 303, 2659, 6199, 2, 4465, 3]

图9-31 文本转换为索引值

由于循环神经网络接受的输入文本必须是相同长度，因此我们需要统一数据集中文本的长度，代码如下。

代码9-10 统一文本长度

```
import torch
from torch.nn.utils.rnn import pad_sequence

# 将列表元素转换为指定长度的张量
def convert_list(ls, max_length):
    processed_data = []
```

```

for item in ls:
    if len(item) < max_length:
        # 如果长度不足指定长度, 进行填充
        padded_item = torch.tensor(item + [0] * (max_length - len(item)))
    else:
        # 如果长度超过指定长度, 进行截断
        padded_item = torch.tensor(item[:max_length])
    processed_data.append(padded_item)
return processed_data

# 使用 pad_sequence 函数将列表转换为张量
padded_train_dataset = pad_sequence(convert_list(train_dataset_indices,
100), batch_first=True)
padded_test_dataset = pad_sequence(convert_list(test_dataset_indices, 100),
batch_first=True)
print(padded_train_dataset[0])

```

以上代码的输出结果如图 9-32 所示, 将数据集中的每一个文本序列都处理为长度为 100 的张量。如果原始序列的长度少于 100, 就通过添加 0 来进行填充; 如果原始序列的长度超过 100, 就通过截断来进行处理。

```

tensor([ 12, 1385, 2188, 13, 153, 8267, 2, 554, 11, 6734, 53, 99,
        29, 54, 4521, 1, 738, 7578, 1575, 2, 1793, 9203, 1596, 71,
        3261, 62, 2737, 1172, 161, 8841, 8869, 1, 3163, 8575, 579, 1688,
        3, 4997, 1519, 365, 1627, 1, 2342, 1519, 3878, 1627, 1, 2127,
        18, 177, 1904, 17, 6053, 4599, 2, 1385, 2188, 3, 150, 1,
        561, 1573, 4063, 50, 2521, 901, 303, 2659, 6199, 2, 4465, 3,
        0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0,
        0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0,
        0, 0, 0, 0])

```

图9-32 统一文本序列长度

处理好文本后, 我们创建数据加载器, 代码如下。

#### 代码9-11 创建数据加载器

```

# 导入需要的库
from torch.utils.data import Dataset, DataLoader

# 自定义的文本数据集类, 继承自 PyTorch 的 Dataset 类
class TextDataset(Dataset):
    def __init__(self, texts, labels):
        self.texts = texts # 文本数据
        self.labels = labels # 相应的标签数据

    def __getitem__(self, index):
        text = self.texts[index] # 获取文本数据
        label = self.labels[index] # 获取标签数据
        return text, label # 返回文本和标签

    def __len__(self):

```



```

        return len(self.texts) # 返回数据集的长度

# 创建自定义的 Dataset 对象，分别为训练集和测试集
train_dataset = TextDataset(padded_train_dataset,
torch.tensor(train_labels)) # 使用文本数据和标签数据初始化训练集对象
test_dataset = TextDataset(padded_test_dataset, torch.tensor(test_labels))
# 使用文本数据和标签数据初始化测试集对象

# 设置批大小
batch_size = 128

# 创建 DataLoader 对象，用于批量加载数据
train_data_loader = DataLoader(train_dataset, batch_size=batch_size,
shuffle=True) # 用于训练集的数据加载器，shuffle=True 表示打乱数据
test_data_loader = DataLoader(test_dataset, batch_size=batch_size,
shuffle=False) # 用于测试集的数据加载器，shuffle=False 表示不打乱数据

```

以上代码用来处理文本数据并创建用于训练和测试的数据加载器（DataLoader）。代码中首先导入了 PyTorch 中用于处理数据的 Dataset 和 DataLoader 类。TextDataset 是一个自定义的数据集类，继承自 PyTorch 的 Dataset 类，用于存储文本数据和标签数据。其中，\_\_init\_\_(self, texts, labels) 是类的构造函数，接受文本数据（texts）和相应的标签数据（labels）作为参数，并将它们存储在类的实例变量中；\_\_getitem\_\_(self, index) 方法用于获取数据集中指定索引 index 的数据，它返回文本数据和相应的标签数据；len\_\_(self) 方法返回数据集的总长度，即数据样本的数量。接下来，两个 TextDataset 对象被创建，分别用于训练集和测试集。它们被初始化时传入了文本数据（padded\_train\_dataset 和 padded\_test\_dataset）和相应的标签数据（train\_labels 和 test\_labels）。batch\_size = 128 定义了批大小，用于指定每次从数据集中加载多少样本，该值将用于创建数据加载器。最后，使用 DataLoader 类创建了两个数据加载器对象 train\_data\_loader 和 test\_data\_loader，分别用于训练集和测试集的数据加载。

接下来构建 LSTM 网络模型，代码如下。

#### 代码9-12 构建 LSTM 网络模型

```

# 本段代码定义 LSTM 模型，与后面定义 GRU 模型的代码只可执行其一
import torch.nn as nn
import torch.optim as optim

# 定义 LSTM 模型
class LSTMClassifier(nn.Module):
    def __init__(self, input_size, hidden_size, num_classes):
        super(LSTMClassifier, self).__init__()
        self.embedding = nn.Embedding(input_size, hidden_size)
        self.lstm = nn.LSTM(hidden_size, hidden_size, batch_first=True,
bidirectional=True)

```

```

self.dropout = nn.Dropout(0.5)
self.fc = nn.Linear(hidden_size, num_classes)

def forward(self, x):
    x = self.embedding(x)
    _, (h, _) = self.lstm(x)
    h = self.dropout(h[-1])
    out = self.fc(h)
    return out

# 定义模型参数
input_size = len(vocab)
hidden_size = 128
num_classes = 3

device = torch.device("cuda" if torch.cuda.is_available() else "cpu")

# 创建模型实例
model = LSTMClassifier(input_size, hidden_size, num_classes)
model.to(device)
model.to(device)

```

以上代码的运行结果如图 9-33 所示，定义了一个 LSTM 分类器模型，其中 `nn.Embedding` 是一个嵌入层，负责将输入的单词索引转换为词向量。这层的输入大小是词汇表的大小（`len(vocab)`），输出大小是词向量的大小（即隐藏层大小）。`nn.LSTM` 是 LSTM 层，它会处理嵌入层的输出，并将其传递到下一层。这层的输入大小是嵌入向量的大小，输出大小也是隐藏层大小。此外，我们设置了 `batch_first=True` 以便 LSTM 可以正确处理我们的输入。随后是一个 Dropout 层，防止过拟合。最后是一个全连接层，它会将 LSTM 层的输出转换为我们最终需要的类别预测。这个层的输入大小是 LSTM 的输出大小，输出大小是类别的数量。

```

LSTMClassifier(
  (embedding): Embedding(22815, 128)
  (lstm): LSTM(128, 128, batch_first=True, bidirectional=True)
  (dropout): Dropout(p=0.5, inplace=False)
  (fc): Linear(in_features=128, out_features=3, bias=True)
)

```

图9-33 LSTM 分类器模型结构

处理好以上所有的工作，可以开始进行训练，代码如下。

#### 代码9-13 模型训练

```

# 定义损失函数和优化器
learning_rate = 0.001
criterion = nn.CrossEntropyLoss()
optimizer = optim.Adam(model.parameters(), lr=learning_rate)
num_epochs = 15

# 在训练集和验证集上计算准确率的函数

```



```
def calculate_accuracy(outputs, labels):
    _, predicted = torch.max(outputs, 1)
    total = labels.size(0)
    correct = (predicted == labels).sum().item()
    accuracy = correct / total
    return accuracy

# 训练模型
for epoch in range(num_epochs):
    train_loss = 0.0
    train_total = 0
    train_correct = 0

    # 训练集上的准确率
    model.train()
    for batch_text_indices, batch_labels in train_data_loader:
        batch_text_indices = batch_text_indices.to(device)
        batch_labels = batch_labels.to(device)

        # 将数据传入模型进行前向传播
        outputs = model(batch_text_indices)

        # 计算损失
        loss = criterion(outputs, batch_labels)
        train_loss += loss.item() * batch_labels.size(0)

        # 计算准确率
        accuracy = calculate_accuracy(outputs, batch_labels)
        train_correct += int(accuracy * batch_labels.size(0))
        train_total += batch_labels.size(0)

        # 反向传播和优化
        optimizer.zero_grad()
        loss.backward()
        optimizer.step()

    train_accuracy = train_correct / train_total
    train_loss = train_loss / train_total

    # 打印训练过程中的损失值和准确率
    print(f"Epoch [{epoch+1}/{num_epochs}], Train Loss: {train_loss:.4f},
    Train Accuracy: {train_accuracy:.4f}")
```

以上代码的输出结果如图 9-34 所示，从训练结果可知，模型在训练集上的准确率最高达到 0.9996。

```

Epoch [1/15], Train Loss: 0.9472, Train Accuracy: 0.5173
Epoch [2/15], Train Loss: 0.7143, Train Accuracy: 0.5827
Epoch [3/15], Train Loss: 0.6456, Train Accuracy: 0.6329
Epoch [4/15], Train Loss: 0.5546, Train Accuracy: 0.7173
Epoch [5/15], Train Loss: 0.4474, Train Accuracy: 0.8017
Epoch [6/15], Train Loss: 0.3322, Train Accuracy: 0.8595
Epoch [7/15], Train Loss: 0.2381, Train Accuracy: 0.9114
Epoch [8/15], Train Loss: 0.1817, Train Accuracy: 0.9342
Epoch [9/15], Train Loss: 0.1112, Train Accuracy: 0.9667
Epoch [10/15], Train Loss: 0.0785, Train Accuracy: 0.9772
Epoch [11/15], Train Loss: 0.0475, Train Accuracy: 0.9882
Epoch [12/15], Train Loss: 0.0288, Train Accuracy: 0.9924
Epoch [13/15], Train Loss: 0.0184, Train Accuracy: 0.9962
Epoch [14/15], Train Loss: 0.0111, Train Accuracy: 0.9987
Epoch [15/15], Train Loss: 0.0072, Train Accuracy: 0.9996

```

图9-34 LSTM 分类器模型训练结果

训练好模型后，我们可以在测试集上测试预测的准确率，代码如下。

代码9-14 预测准确率

```

# 在测试集上进行预测
model.eval()
test_total = 0
test_correct = 0
with torch.no_grad():
    for batch_text_indices, batch_labels in test_data_loader:
        batch_text_indices = batch_text_indices.to(device)
        batch_labels = batch_labels.to(device)

        # 将数据传入模型进行前向传播
        outputs = model(batch_text_indices)

        # 计算准确率
        accuracy = calculate_accuracy(outputs, batch_labels)
        test_correct += int(accuracy * batch_labels.size(0))
        test_total += batch_labels.size(0)

test_accuracy = test_correct / test_total

# 打印测试集的准确率
print(f'Test - Accuracy: {test_accuracy:.4f}')
```

以上代码的输出结果如图 9-35 所示，在测试集上预测的准确率为 73.91%。

```
Test - Accuracy: 0.7837
```

图9-35 LSTM 模型在测试集上的准确率

我们可以尝试通过修改模型、调整超参数、数据增广等方式提升模型预测的准确率。

## 9.7 本章小结

本章主要介绍了循环神经网络（RNN）以及其在自然语言处理（NLP）中的应用。

首先，我们探讨了序列数据的特点以及循环神经网络的基本概念，讲解了循环神经网络在处理序列数据时的工作原理，包括前向传播、反向传播、权重共享以及梯度消失和梯度爆炸问题。接着，我们讨论了为解决循环神经网络的一些问题而提出的改进方案，如长短期记忆网络（LSTM）和门控循环单元（GRU）。然后，我们深入到自然语言处理的主题，了解了分词、文本向量化、词嵌入等基础概念，并讨论了循环神经网络和大模型在自然语言处理中的应用。最后，以微博评论为例，具体展示了如何使用循环神经网络来处理实际的自然语言处理问题。

通过本章的学习，我们了解到了循环神经网络在序列数据处理，尤其是在自然语言处理中的强大能力。

## 9.8 课后习题

### 一、单项选择题

1. 在循环神经网络中，信息的传递和持久化是通过什么机制实现的？（ ）
  - A. 循环连接
  - B. 反馈连接
  - C. 正向传播
  - D. 反向传播
2. 循环神经网络最适合处理哪种类型的数据？（ ）
  - A. 序列数据
  - B. 图像数据
  - C. 音频数据
  - D. 自然语言数据
3. 循环神经网络中的神经元模型是如何建立和进行前向传播和反向传播的？（ ）
  - A. 通过引入循环连接
  - B. 通过设置卷积核的宽度
  - C. 通过使用带有自反馈的神经元构造网络

D. 通过增加一个循环回路来维护隐藏状态向量

4. 什么是循环神经网络中的隐藏状态向量？（ ）

A. 一个神经元的状态向量

B. 一个时间步的网络状态向量

C. 一个序列中所有时间步的网络状态向量

D. 一个循环神经网络的状态向量

5. 循环神经网络与全连接神经网络和卷积神经网络有什么不同？（ ）

A. 循环神经网络具有记忆能力，可以捕捉数据点之间的顺序关系

B. 循环神经网络只能处理序列数据，而全连接神经网络和卷积神经网络可以处理各种类型的数据

C. 循环神经网络的隐藏层中的神经元具有循环连接，而卷积神经网络的神经元之间没有循环连接

D. 循环神经网络的结构比全连接神经网络和卷积神经网络更复杂，需要更多的参数进行训练

## 二、填空题

1. \_\_\_\_\_是一种强大而灵活的神经网络架构，专门设计用于处理序列数据。

2. 循环神经网络通过在网络内部引入\_\_\_\_\_，使得信息可以在网络中传递并持久化。

3. 循环神经网络在\_\_\_\_\_、\_\_\_\_\_、\_\_\_\_\_等任务上取得了卓越的成就。

4. 循环神经网络可以捕捉到数据点之间的\_\_\_\_\_关系和\_\_\_\_\_相关性，从而更好地理解  
和预测序列数据的特征和趋势。

5. \_\_\_\_\_是一种适用于处理序列数据的深度学习模型，可以捕捉数据点之间的长期依赖关系。

## 三、简答题

1. 什么是循环神经网络（RNN）？它为什么在处理序列数据方面表现优秀？

2. 简述循环神经网络的基本结构，并说明其如何实现信息的传递和记忆。

3. 什么是长短期记忆网络（LSTM）？它如何解决循环神经网络中的梯度消失和梯度爆炸问题？

4. 什么是自然语言处理？列举几个自然语言处理的典型应用场景。