

## 第 5 章

# 情感分析的应用

## 5.1 情感分析在股票预测中的应用



本节代码

股票价格波动的本质是对新信息的反应,在传统基于股市数值分析的基础上,研究新闻对股票市场的影响,有助于提高股票走势预测的准确率。本章节提出一种基于 CNN-BiLSTM 多特征融合的股票走势预测模型。该模型通过引入卷积神经网络和双向长短时记忆模型挖掘财经新闻中的新闻事件类型和新闻情感倾向,深度融合股市财务数据、新闻事件特征以及新闻情感特征,实现对股票走向的预测。

分别选取家用电器行业和通信行业的两支股票作为实验对象,验证所提模型对不同行业个股走势预测的可行性。实验结果表明,引入新闻事件和情感特征后,模型的预测准确率有提升,家用电器行业准确率提高了 11%,通信行业准确率提高了 22%。表明通过引入新闻事件类型和情感倾向能够提高股票走势预测的性能,此外,还评估了影响股票走势的因素,根据重要性对影响股票走势预测的特征进行了排序。

### 5.1.1 股票走势预测研究背景

股票走势预测是根据与股票相关的数据对股票价格的波动进行预判。现有研究表明,股票价格的走势并不遵循随机游走,在一定程度上可以被预测<sup>[1]</sup>。传统的量化投资属于早期的股票预测方法,通过结构化的、线性的历史交易数据对股票走势进行预测,主要采用线性回归、参数估计方法<sup>[2]</sup>。然而,股票价格的走势除了与历史交易数据相关,还容易受到非线性因素的影响,例如政策因素、投资心理以及突发事件等。如何统筹考虑线性的历史交易数据以及非线性的因素对股票价格的影响,进一步提高股票预测的准确率是近年来的研究热点。

随着信息技术的发展,媒体新闻信息作为反映与股票走势非线性因素相关的重要信号,被用于提高股票走势预测的准确率。大量研究表明,媒体新闻信息会对股票价格产生影响<sup>[3]</sup>。现有研究将媒体新闻信息中的情感极性(正面或者负面)作为反映市场状态的指标<sup>[4-6]</sup>。例如,Bollen 等人利用道琼斯指数分析 Twitter 的用户情感来进行股价预测<sup>[7]</sup>。现有研究证实了新闻媒体所蕴含的情绪信号对于股票市场价格有一定的预测能力,然而仍存在一些问题亟待解决。此外,新闻文本信息属于非结构化的数据,具备非线性、非平稳的特征,使得传统的量化投资分析方法并不适用。

股票走势预测是一个经典的基于时序数据的预测问题,近年来由于深度学习在文本处理方面的突出表现,越来越多学者尝试将深度学习用于解决基于时序数据的股票预测

问题。循环神经网络(Recurrent Neural Network, RNN)将时间概念融入网络结构中,天然地适用于处理时序数据。然而当输入的时序数据的序列太长时,RNN 会面临梯度消失问题。为了解决这一问题,长短时记忆模型(Long Short-Term Memory, LSTM)作为 RNN 的改进版本被提出。研究表明,LSTM 在处理基于时序数据的股票预测问题时性能优于 RNN 和传统的机器学习算法,例如随机森林(Random Forest)、支持向量机(Support Vector Machine, SVM)等<sup>[8-10]</sup>。大多现有基于深度学习的股票预测问题采用 LSTM 或者改进的 LSTM,例如 PSO-LSTM<sup>[9]</sup>、Stacked LSTM<sup>[10]</sup>、Time-Weighted LSTM<sup>[11]</sup>等,用于搭建股票预测模型。

现有基于 LSTM 的股票预测模型主要通过分析文本信息的情感极性特征,将媒体新闻的情感极性与历史交易数据作为股票预测算法的输入,需要考虑以下 3 方面的问题。

(1) 首先,财经新闻中的情感极性并不明显,大多是对客观事件的总结和报道,使得财经新闻的情感极性分析准确率并不高,影响股票走势的预测准确率<sup>[12]</sup>。

(2) 其次,如何将非结构化的文本信息特征与结构化的股票交易数据融合在一起。现有大多数研究将新闻文本的情感值与高维度的历史交易数据、公司财务数据直接拼接在一起,作为股票预测的输入<sup>[12-14]</sup>。这种方法很容易将情感信息淹没在高维度的结构化信息中。添加了情感维度的股票预测准确率甚至低于不添加情感维度的股票预测方法<sup>[15]</sup>。

(3) 最后,新闻文本的情感极性并不总是与股票走势的涨跌正相关。例如,新闻“中兴通讯高层大变动,少壮派受重用”的情感属性为正面,但并没有对中兴股票的上涨有积极影响。新闻文本情感极性可能与股票的涨跌负相关。与情感极性相比,新闻事件本身更能够代表媒体新闻对股票走势的影响。

为了解决上述 3 个问题,本节提出一种基于 CNN-BiLSTM 多特征融合的股票走势预测模型,通过融入新闻事件类型和情感极性提高股票走势预测的准确率。CNN 在新闻文本的事件特征提取上的性能更为突出<sup>[16]</sup>,而双向 Bi-LSTM 采用两个 LSTM 网络获取文本前向和后向的语境信息,更适合用于处理考虑上下文语境的情感极性判别,在情感分析上较单个 LSTM 能提升 3% 的性能<sup>[17]</sup>。因此,所提模型一方面采用从新闻报道中提取出的客观财经事件,例如中标事件、上市事件、停牌事件等。另一方面,采用 Bi-LSTM 对新闻报道的情感极性进行分析,计算新闻文本的情感分值。股票新闻特征包括新闻的事件类型和情感分值,与股票的财务数值特征一起作为 LSTM 网络的输入,利用历史的股票信息预测股票未来的涨跌情况。此外,所提模型通过对比实验研究影响股票走势预测模型的特征重要性。

### 5.1.2 相关研究工作

现有对股票走势的预测研究可以分为两类:一类是基于数值分析;另一类是融合数值和文本信息的股票预测模型。

早期量化投资分析属于传统的数值分析方法,主要基于历史交易数据、公司财务数据和宏观数据对股票走势进行预测<sup>[18]</sup>。Chen 等验证个股的历史股票数据用来预测个股的未来走势<sup>[19]</sup>。这些方法主要根据历史交易数据发现并描述数据随时间变化的规律。为

了挖掘大量数值信息的规律,传统的机器学习算法,例如 SVM、人工神经网络(Artificial Neural Networks, ANN)、朴素贝叶斯、随机森林等,被用于对大量的历史股票数据分析<sup>[20]</sup>。张玉川等利用 SVM 对个股涨跌进行预测,通过实证分析验证 SVM 对个股涨跌分类的有效性<sup>[21]</sup>。Karen 等对比 ANN 和 SVM 在股票走势预测上的性能,发现 ANN 在预测准确率上优于 SVM,前馈 ANN 由于能够同时预测股票走势的涨跌以及股票价格而被广泛采用<sup>[22]</sup>。然而,由于股票价格本质上是随机变量的噪声观测值,只采用历史交易数据进行分析具有一定局限,无法进一步提升预测效果。行为经济学理论指出,投资者面对复杂和不确定的决策问题时,很容易受到个人和社会环境情感状态的影响<sup>[23]</sup>。股票价格变化的根源是对新信息的反应,媒体新闻文章作为外生变量的信息,对短期价格预测有帮助<sup>[24]</sup>。

另一类融合股票数值信息和新闻信息的股票预测模型应运而生。Vanstone 等研究新闻及新闻的情感极性是否会对股票价格的预测起作用<sup>[24]</sup>。研究表明通过统计与股票相关的新闻文本数量以及 Twitter 条数,并将其作为股票价格预测的输入可以进一步提高预测准确率。Chen 等利用带有门循环单元的 RNN 模型研究新浪微博上的财经新闻的情感极性,并融合股价数值特征一起预测股票走势<sup>[25]</sup>。Manuel R. 等利用 CNN 和 RNN 研究融合新闻标题和技术指标的股票走势模型,证明新闻标题比新闻内容更有利于提高预测准确率<sup>[26]</sup>。岑咏华等考察新闻网站、股吧、博客等媒介信息所蕴含的情感信号对于股票市场的影响效应,发现投资者对于积极情感的反应更及时、更强烈<sup>[27]</sup>。

这些现有的融合股票数值信息和新闻信息的股票预测方法主要考虑提取新闻的情感极性作为股票数值信息的补充。然而,新闻是对客观事实的描述,情感极性大多比较隐晦,使得情感特征对于预测准确率的提升并不明显。考虑到这一局限,现有研究<sup>[24-26]</sup>大多采用情感表达较为明显的 Twitter、微博文本作为新闻信息的来源。Zhao 等提出使用隐含狄利克雷分布主题模型提取出微博文本的关键词,再基于关键词分析微博文本的情感特征作为股票预测的输入<sup>[28]</sup>。区别于现有工作,考虑到新闻事件比新闻情感更能代表新闻媒体信息对股票走势的影响,本节采用融合新闻事件和情感特征的多特征融合方法,对股票的数值特征进行补充,进一步提升股票预测的准确率。

### 5.1.3 基于新闻事件和情感特征的股票预测模型

股票走势预测是一个二分类问题,当涨跌幅高于阈值时,判为上涨样本;反之则为下跌样本。设采样间隔为  $a$  天,根据过去  $t-a$  至  $t-1$  天的数据,预测第  $t$  天的股票涨跌。

所提模型的架构如图 5-1 所示。基本思路是:

(1) 分别爬取股票相关的财务数据和新闻标题数据,对财务数据和新闻标题数据进行预处理,得到股票数据库和新闻数据库。

(2) 划分个股新闻事件,对新闻进行事件类型标注,利用 CNN 训练标注数据得到新闻事件分类器。

(3) 对新闻文本进行情感标注,利用 Bi-LSTM 建立新闻情感分类器。

(4) 将(2)和(3)训练得到的股票新闻特征和股票数值特征作为 LSTM 网络的输入,

比较不同特征的融合方式,并预测股票走势的涨跌。

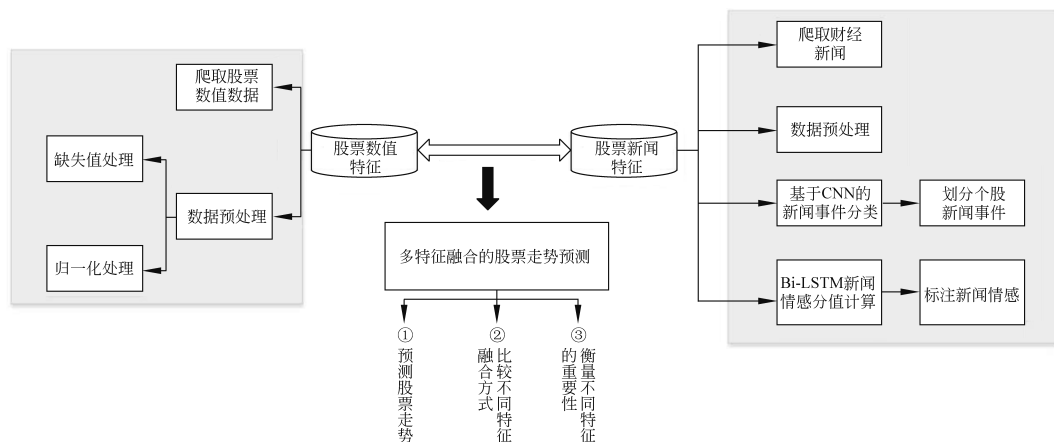


图 5-1 融合新闻事件和情感特征的股票走势预测流程图

### 1. 股票财务数值特征提取

相关研究表明公司的市盈率、市净率、主力资金净流入等指标跟个股的走势相关<sup>[19,29-30]</sup>。因此,股票财务数值特征选取公司的财务数据(如市盈率、市净率),公司的资金流向数据(如主力资金净流入、换手率)和股票数据(如开盘价、收盘价)。考虑到大盘指数和板块指数对股票走势的影响,也选取了个股的大盘指数和板块指数作为股票财务数值指标。

对财务数值数据首先进行数据缺失值处理:如果某一天的数据缺少某个指标(停牌等原因),则将该天的数据去掉。此外,考虑到不同财务数值指标的性质不同,数量级也不同,直接使用指标的原始数值可能导致数值较高的指标在训练中有更为主导的作用,削弱数值较低的指标的影响,因此对财务数值数据进行 z-score<sup>[31]</sup> 标准化处理,保证不同指标数据之间的可比性。

设  $D_j = (d_{1j}, d_{2j}, \dots, d_{ij}, \dots, d_{nj})$  表示第  $j$  个财务指标在  $T_n$  天内的数值组成的向量,  $d_{ij}$  表示第  $j$  个财务指标在第  $i$  天的数值,将  $D_j$  中的每一个数值进行 z-score 标准化处理为

$$d_{ij} = \frac{d_{ij} - \mu_j}{\sigma_j} \tag{5-1}$$

其中,  $\mu_j$  和  $\sigma_j$  表示财务指标  $D_j$  中所有数值的均值和标准差。表 5-1 描述了由  $p$  个股票指标在  $T_n$  天内构成的财务特征矩阵,  $D_1$  到  $D_p$  表示  $p$  个财务特征。

表 5-1 股票财务特征矩阵

	$D_1$	$D_2$	...	$D_j$	...	$D_p$
$T_1$	$d_{11}$	$d_{12}$	...	$d_{1j}$	...	$d_{1p}$
...	...	...	...	...	...	...

续表

	$D_1$	$D_2$	...	$D_j$	...	$D_p$
$T_i$	$d_{i1}$	$d_{i2}$	...	$d_{ij}$	...	$d_{ip}$
...	...	...	...	...	...	...
$T_n$	$d_{n1}$	$d_{n2}$	...	$d_{nj}$	...	$d_{np}$

## 2. 基于 CNN 模型的新闻事件分类及特征提取

新闻事件特征提取从新闻标题中提取出客观的金融事件。首先对新闻标题进行数据预处理：利用结巴分词器进行分词和去停用词，同时引入自定义的停用词和金融词典提高分词准确率。自定义金融词典包括常见金融词汇、A 股上市公司代码及简称、A 股上市公司实际控制人及高管姓名等。

根据客观事件划分金融新闻事件。参考国泰安 (CSMAR) 经济数据库中新闻词条的关键字段作为事件划分的分类依据，一共划分了 82 个新闻事件类型。表 5-2 列举了部分新闻事件。

表 5-2 部分新闻事件类型

事件类别	事件名称
交易类	停牌 复牌 资金流入 资金流出 大宗交易 股价倒挂 创新高
股权类	挂牌 借壳 举牌 收购并购 资产重组 资产冻结 股权转让
投融资类	投资 投建 中标 发行债券 发行股票 可转债 募资 质押 分红
公司事务类	注册资本变更 快速发展 战略合作 拓展业务 高管减持或离职
外部事件类	登上龙虎榜 交易所处罚 评级利好 评级下调 政策利好

基于 CNN 搭建新闻事件分类器。新闻事件分类器的训练过程与基于 CNN 的文本分类过程类似，详细过程可参考笔者前期研究工作<sup>[32]</sup>，其主要思路是：根据划分的 82 个新闻事件对股票新闻进行标注，训练 CNN 模型，模型包括输入层、卷积层、池化层和全连接层操作，每一层的输出是下一层的输入。首先采用 Word2vec 训练新闻标题，将得到的词向量矩阵作为卷积层的输入，卷积层利用滤波器对新闻标题的词向量矩阵进行卷积操作，产生特征图。池化层对特征图进行采样，抽取每个特征图中最重要的特征传入全连接层。最后，全连接层通过 softmax 函数获得新闻标题最终分类结果，输出新闻标题的事件类型。

通过 CNN 新闻分类器对特定股票每天的新闻进行统计。每条新闻输入到新闻分类器后将输出一个事件类型，再统计每天各个事件出现的频率，得到表 5-3 所示的新闻事件特征矩阵。 $S_1$  到  $S_q$  表示  $q$  个新闻事件，这里  $q=82$ ， $s_{ij}$  表示在日期  $T_i$  出现事件特征  $S_j$  的频次。

表 5-3 新闻事件特征矩阵

	$S_1$	$S_2$	...	$S_j$	...	$S_q$
$T_1$	$s_{11}$	$s_{12}$	...	$s_{1j}$	...	$s_{1q}$
...	...	...	...	...	...	...
$T_i$	$s_{i1}$	$s_{i2}$	...	$s_{ij}$	...	$s_{iq}$
...	...	...	...	...	...	...
$T_n$	$s_{n1}$	$s_{n2}$	...	$s_{nj}$	...	$s_{nq}$

### 3. 新闻文本情感特征提取

新闻事件是对客观事件的描述,而新闻文本的情感表征新闻事件的情感极性,例如消极或者积极。新闻文本的情感极性判断需要考虑上下文的语境信息,Bi-LSTM 采用 2 个 LSTM 网络进行训练,一个训练序列从文本前面开始,一个训练序列从文本后面开始,这两个训练序列连接到同一个输出层。Bi-LSTM 能够整合每个点的过去和未来信息,相比于单个 LSTM 在文本情感极性判断上更有优势。图 5-2 为利用 Bi-LSTM 计算新闻标题情感分值的过程。

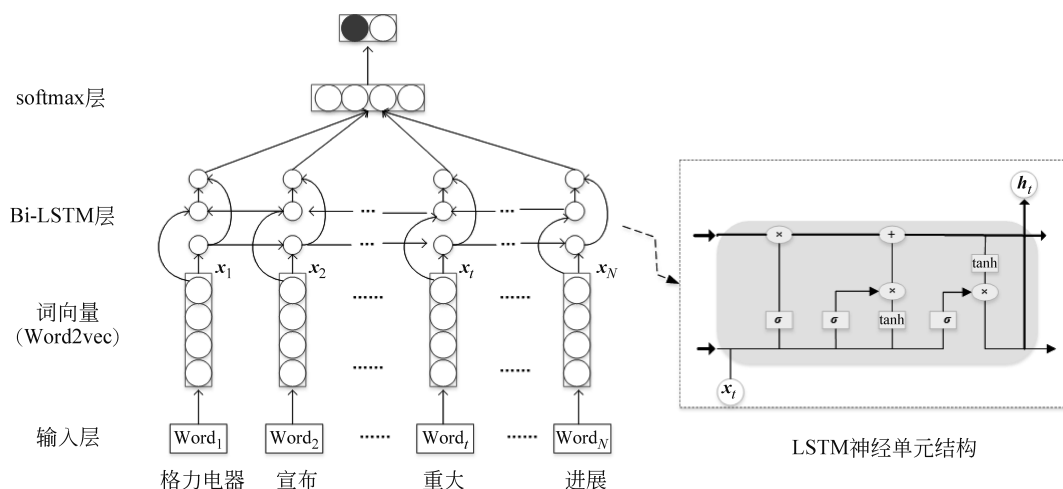


图 5-2 基于 Bi-LSTM 的新闻情感分析模型

对数据预处理后的每一条新闻标题  $d$ , 最长裁剪  $N$  个词语, 并设置 LSTM 的处理步长为  $N$ 。对于长度不足  $N$  的新闻标题, 进行左置零补齐。针对每一个采样时刻 ( $t \leq N$ ), 将词语  $w_t$  通过 Word2vec 训练得到的词向量  $x_t$  输入到一个包含  $L$  个神经元的 LSTM 神经网络层。该神经网络层输出一个维度为  $L$  的隐含状态向量  $h_t$ 。每一个神经元设置三种门限结构, 即遗忘门、输入门和输出门, 基于过去的隐含状态向量  $h_{t-1}$  和当前输入  $x_t$ , 决定需要遗忘哪些信息、输入哪些新的信息以及对哪些新的记忆信息编码输出得到  $h_t$ 。具体地, 在时刻  $t$ , LSTM 层的计算如公式(5-2)~公式(5-6)所示<sup>[33]</sup>。

$$f_t = \sigma(W_f \cdot [h_{t-1}, x_t] + b_f) \quad (5-2)$$

$$i_t = \sigma(W_i \cdot [h_{t-1}, x_t] + b_i) \quad (5-3)$$

$$C_t = f_t \Theta C_{t-1} + i_t \Theta \tanh(W_c \cdot [h_{t-1}, x_t] + b_c) \quad (5-4)$$

$$o_t = \sigma(W_o \cdot [h_{t-1}, x_t] + b_o) \quad (5-5)$$

$$h_t = o_t \Theta \tanh(C_t) \quad (5-6)$$

首先,遗忘门基于 Sigmoid 函数  $\sigma$  判断过去的记忆对于当前的记忆状态,有多大意义值得保留,通过公式(5-2)计算生成系数  $f_t$ 。接着,输入门判断当前的词输入向量  $x_t$  在多大意义上值得保留,根据公式(5-3)生成系数  $i_t$ 。然后,神经元根据公式(5-4)更新生成当前时刻的状态  $C_t$ 。最后,输出门基于  $C_t$  判断新的记忆在多大意义上值得输出,输出的隐含状态  $h_t$  由公式(5-6)表示。公式(5-2)~(5-6)中,  $W_*$  和  $b_*$  分别表示权重矩阵和偏置向量,  $\Theta$  为乘操作。

经过上述计算,得到新闻标题  $d$  在  $t=N$  时刻上的隐含状态编码  $h_N$ ,基于  $h_N$  通过 softmax 函数得到  $d$  在不同情感类别  $\{1, -1\}$ , 即  $\{\text{正向}, \text{负向}\}$  上的概率分布向量为

$$y_d = \text{Softmax}(W_o \cdot h_N + b_o) \quad (5-7)$$

基于  $y_d$ , 计算新闻标题  $d$  的情感倾向 (Sentiment Orientation) 为

$$so_d = (1, -1) \cdot y_d \quad (5-8)$$

其中,  $so \in [-1, 1]$ , 当  $so_d > 0$ , 情感倾向为积极(正向), 否则情感倾向为消极(负向)。

设在日期  $T_i$  一共有  $m_i$  条新闻, 每条新闻的情感极性用  $so_j$  来表示,  $j=1, 2, \dots, m_i$ ,  $so_j \in [-1, 1]$ 。则在  $T_i$  当天的新闻情感总分值计算为

$$S_i = \frac{1}{m_i} \sum_{j=1}^{m_i} so_j \quad (5-9)$$

将新闻情感向量  $S = [S_1, S_2, \dots, S_i, \dots, S_n]^T$  作为一系列特征添加到表 5-3 所示的新闻事件特征矩阵中, 得到股票新闻特征矩阵。

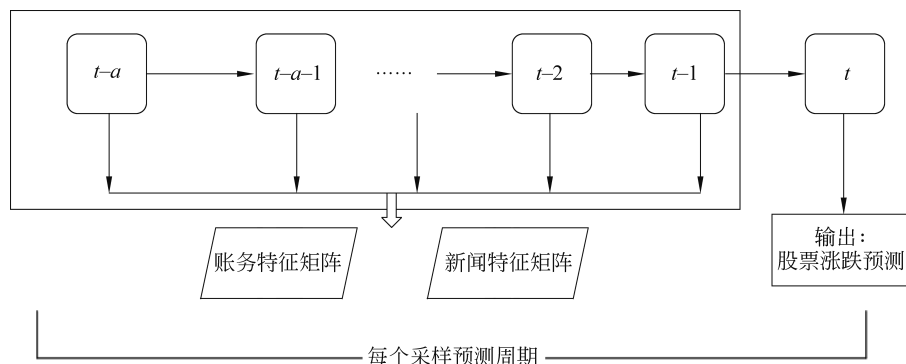


图 5-3 股票预测模型采样周期示意图

#### 4. 股票走势预测

股票预测将财务特征矩阵、股票新闻特征矩阵按照日期合并, 作为股票走势预测模型的输入。股票走势预测作为前向时序建模预测, 不用考虑后向序列的影响, 因此采用 LSTM 而非 Bi-LSTM 作为股票涨跌预测。

将财务特征矩阵、新闻特征矩阵直接拼接在一起作为 LSTM 模型的输入可能会导致梯度消失问题<sup>[12]</sup>,因此本节采用 2 个 LSTM 神经元,将财务特征矩阵和新闻特征矩阵分别输入到两个 LSTM 中,再将结果进行向量合并,输入到全连接神经网络中,最后输出涨跌结果。当采样间隔为  $a$  天,将根据过去  $t-a$  至  $t-1$  天的数据,预测第  $t$  天的股票涨跌,如图 5-3 所示。

## 5.1.4 实验分析

### 1. 实验数据集

为了评估所提模型对不同行业和板块的适用性,实验分别选取家用电器行业和通信行业的两支股票:格力电器(000651.SZ)和中兴通讯(000063.SZ)作为实验对象。家用电器行业属于周期性行业,受到国内或国际经济周期性波动的影响较大。与此相反,通信行业属于防御性行业,受到经济周期性衰退和繁荣的影响较小,在股价表现上较为稳定。

实验一共选取了 12 个股票指标构成财务特征矩阵,分别是:开盘价、最低价、最高价、收盘价、主力资金净流入量、换手率、涨跌、涨跌幅、市盈率、市净率、深证 A 指以及板块指数。

财务数值数据主要来源于 Wind 金融数据库。Wind 是国内以金融证券数据为核心的大型金融工程和财经数据仓库,以科学的核查手段和先进的管理方法确保数据的准确性在 99.95% 以上。在 Wind 数据库中抓取 2010—02—02 到 2020—02—28 十年间,格力电器共 18 766 条新闻数据,29 352 条财务数据;中兴通讯共 18 796 条新闻数据,29 364 条财务数据。新闻事件分类器的标注数据来源为国泰安经济金融数据库中的新闻,一共对 27 800 条新闻数据进行事件标注,用于事件分类器模型的训练。

### 2. 新闻事件和情感分类器实验结果

#### (1) 参数设置和模型输入。

新闻事件分类器和情感分析模型的性能受到词向量维度、卷积窗口大小、迭代次数、过滤器数量等因素影响。在实验中通过十折交叉验证评估模型表现,选定最合适的参数组合。在本节数据集上表现最好的参数组合如表 5-4 所示。表中 Null 表示不需要这个参数设置。

表 5-4 模型参数设置

参 数	CNN	Bi-LSTM
词向量维度	300	300
卷积核个数	96	Null
卷积核大小	3,4,5	Null
Dropout	0.5	0.5



续表

参 数	CNN	Bi-LSTM
Batch_size	128	128
迭代次数	10	20
标题截取长度	Null	15
单层 LSTM 神经元个数	Null	[256,256]

## (2) 评价指标。

计算精确率(Precision)、召回率(Recall)以及  $F1$  值( $F1$ -measure)作为实验评价指标。第  $i$  类新闻事件分类的精确率、召回率和值分别用  $Pre_i$ 、 $R_i$  和  $F1_i$  表示,计算公式为

$$Pre_i = \frac{TP_i}{(TP_i + FP_i)} \quad (5-10)$$

$$R_i = \frac{TP_i}{(TP_i + FN_i)} \quad (5-11)$$

$$F1_i = \frac{2 \times Pre_i \times R_i}{Pre_i + R_i} \quad (5-12)$$

其中,  $TP_i$  为被正确分类到第  $i$  类的样本数量,  $FP_i$  为被错误分类到第  $i$  类的样本数量,  $FN_i$  为原本属于第  $i$  类,但是被分到其他类别的样本数量。

## (3) 实验结果。

表 5-5 列举了新闻事件分类的实验结果,为了验证基于 CNN 的新闻事件分类器的性能,将其与 SVM 算法以及基于最大熵(Maximum Entropy, Maxent)算法的新闻事件分类性能做比较。SVM 在传统的机器学习算法中表现优异,而 Maxent 是线性对数模型,具有较好的分类能力。实验将新闻事件数据集 90% 的数据共 25 020 条新闻作为训练集,10% 的数据共 2780 条新闻作为测试集。

表 5-5 新闻事件分类精确率对比

	SVM	Maxent	CNN
训练集	90.78%	72.02%	93%
测试集	85.23%	69.42%	87.7%

基于 CNN 新闻事件分类器在训练集中精确率达 93%,测试集中精确率达 87.7%,优于基于 SVM 算法和基于 Maxent 算法的新闻事件分类精确率。为了分析基于 CNN 的新闻事件分类器对不同新闻事件的分类效果。表 5-6 列举了各类新闻事件的准确率、召回率和  $F1$  值,由于篇幅关系仅列出部分的事件类型。高召回率和高  $F1$  值表示新闻事件预测的覆盖率和准确率都比较好。

从表 5-6 可以看出高召回率和高  $F1$  值的事件类型大多属于公司公告类,公司公告的内容在不同公司之间的差异性不大,公告具有一定的模板性,分类效果较好。分类性能较差的新闻事件大部分是对价格走势的预测,走势好或者走势不好在不同行业、不同公司之

间新闻的内容差别比较大,因此识别起来准确率较低。

表 5-6 新闻事件分类的性能统计表

分类性能较好的新闻事件类型示例				分类性能较差的新闻事件类型示例			
新闻事件	精确率	召回率	F1 值	新闻事件	精确率	召回率	F1 值
登上龙虎榜	1.00	1.00	1.00	业绩下降	0.64	0.58	0.61
停牌	0.98	1.00	0.99	政策利好	0.81	0.65	0.72
工商变更	1.00	1.00	1.00	资本变更	1	0.22	0.36
中标	1.00	1.00	1.00	聘请高管	0.50	0.40	0.44
可转债	0.97	0.97	0.97	业绩增长	0.68	0.73	0.71
质押	1.00	1.00	1.00	预计下滑	0.67	0.61	0.64
交易所问询	0.94	1.00	0.97	利差消息	0.42	0.47	0.44
退市	1.00	1.00	1.00	利好消息	0.46	0.65	0.54

表 5-7 为新闻情感分类的实验结果。基于 Bi-LSTM 的新闻情感分类在训练集上的精确率为 99%,在测试集上精确率为 91%;基于 SVM 的新闻情感分类性能次之,在训练集和测试集上的精确率分别为 86.64%和 81.13%;基于 Maxnet 的新闻情感分类效果最差。利用训练得到的新闻情感分类模型分别对格力电器和中兴通讯的新闻数据集进行情感分类、计算新闻情感分值,生成得到新闻的情感向量矩阵。

表 5-7 新闻情感分类精确率对比

	SVM	Maxent	Bi-LSTM
训练集	86.6%	82.8%	99.0%
测试集	81.1%	76.1%	91.0%

### 3. 股票走势预测结果

为了验证引入新闻信息后对股票预测模型性能的提高,实验将对对比测试以下算法的性能。

- (1) 未引入新闻特征的 LSTM: 仅利用股票财务特征进行股票走势预测;
- (2) 引入新闻事件的 LSTM: 利用股票财务特征和如表 5-3 所示的新闻事件特征进行股票走势预测;
- (3) 新闻事件/情感融合的 LSTM: 利用股票财务特征、新闻事件特征以及情感特征进行股票走势预测。

(4) 新闻事件/情感融合的 GBDT: 基于梯度提升决策树(Gradient Boosting Decision Tree, GBDT)模型<sup>[34]</sup>进行股票预测,输入与上述模型(3)相同,对比 LSTM 和