

doi:10.3772/j.issn.1000-0135.2015.011.008

网络股评对股市走势的影响:基于文本情感分析的方法¹⁾

王洪伟¹ 张 对¹ 郑丽娟² 陆 颀³

(1. 同济大学经济与管理学院, 上海 200092; 2. 聊城大学商学院, 聊城 252000;

3. 上海通用汽车销售有限公司, 上海 201206)

摘要 为研究网络股评是否影响股民决策以及股市走势的问题, 本文收集新浪股吧中 30 只沪深 300 指数的股票的网络股评, 计算情感指数, 建立 ARMA-GARCHX 模型与 ARMAX-GARCH 模型, 从而分析股评中情感因素与股市走势之间的关系。本研究对证券投资机构 and 个人的投资决策具有参考价值。研究显示, 通过情感分析技术, 从网络股评信息中提取投资者情感是可行的, 且具有数据样本量大、时效性强、更接近投资者的真实情感等优点; 同时发现, 将投资者情感作为两个模型的输入项来预测股票价格的波动是可行的, 并且 ARMAX-GARCH 模型效果更好, 但模型中情感多为当期情感; ARMA-GARCHX 模型模拟效果不如 ARMAX-GARCH 模型, 但情感多为前期情感, 有更好的预测能力。

关键词 网络股评 情感指数 股市走势 ARMA-GARCHX 模型 ARMAX-GARCH 模型

The Effect of Online Comments on Stock Trends by Sentiment Analysis

Wang Hongwei¹, Zhang Dui¹, Zheng Lijuan² and Lu Ting³

(1. School of Economics and Management, Tongji University, Shanghai 200092;

2. School of Business, University of Liaocheng, Liaocheng 252000;

3. Saic General Motors Sales Co., LTD, Shanghai 201206)

Abstract This paper is to examine whether online stock comments could affect the stock trend by constructing ARMA-GARCHX model and ARMAX-GARCH model. The relationship between the sentiment of investors toward certain stocks and the stock trend is expected to be found for enabling better investment decisions. Our model is featured with three advantages: bigger sample size, better timeliness, and closer to true emotion of investors. The investors' sentiment is set as an external variable, which is then combined into the two proposed models. The empirical study is conducted with online comments on the 30 selected stocks of CSI 300 Index. The result shows that ARMAX-GARCH model outperforms ARMA-GARCHX model in term of prediction accuracy. In most situations, however, ARMAX-GARCH model prefer current sentiments to previous ones. By contrast, ARMA-GARCHX model uses more past sentiments, thus being more applicable to stock forecast.

Keywords online stock comments, sentiment index, stock trend, ARMA-GARCHX model, ARMAX-GARCH model

收稿日期: 2014 年 12 月 1 日

作者简介: 王洪伟, 男, 1973 年, 博士, 教授, 博导, 主要研究方向: 商务智能与情感计算, E-mail: hwwang@tongji.edu.cn。张对, 女, 1990 年, 硕士研究生, 主要研究方向: 金融工程。郑丽娟, 女, 1983 年, 博士, 主要研究方向: 商务智能与文本挖掘。陆颀, 男, 1980 年, 博士研究生, 主要研究方向: 商务智能。

1) 基金项目: 国家自然科学基金项目(71371144), 上海市哲学社会科学规划课题一般项目(2013BGL004)。

1 引言

股市是经济的晴雨表,揭示股市运行规律,对投资决策具有重要意义。影响股市走势的因素包括:宏观经济因素、行业发展因素、公司内部因素、股民情感因素^[1]。前三者短期内改变较少,难以反映股市实时波动的情况。而股民情感因素,因为受消息、股市走势、股民心理等因素影响,具有易变性,并在一定程度上反映影响个股走势的其他3类因素。

随着社交媒体的发展,尤其是微博、微信等UGC(用户原创内容)的兴起,投资者在线参与股评的意愿越发强烈。网络股评具有传播快、范围广、影响力大、实时性强等特点^[2],不仅包括专业人士对股市的意见,还包括普通股民的看法,因此能够反映整个社会对股市的看法。

情感分析,是指利用文本挖掘技术,对网络评论进行语义分析,旨在识别用户的情感趋向是“高兴”还是“伤悲”,或判断用户观点是“赞同”还是“反对”^[3]。情感分析有助于提取和量化市场情感。关于情感对股市影响的研究中,通常是选取间接指标来反映股市情感(如封闭式基金折价率等),而直接从股评中提取市场情感的研究较少。本文收集新浪股吧中的股评,对其进行情感分析,在文本预处理之后提取情感指数,并利用情感指数作为外生变量建立GARCH模型,研究股评中情感与股市走势之间的关系。

2 文献综述

股评是股评者对于某一或某些股票走势、风险等的看法,集中反映了股评者对该股票的预期,较多体现了投资者情感。在股评影响股市的研究中,股评价值的理论依据包括:①证券组合和有效市场理论,重点关注股评的信息价值;②行为金融学理论,着重研究股评对股民心理的影响,进而影响股民行为和股票市场。随着理性人和有效市场假设受到质疑,第①种理论暴露出局限性,第②种理论体现出其实际价值,投资者的情感及行为成为研究焦点。

随着互联网的发展,社交媒体逐渐渗透入股评行业。股民改变了在传统媒体中的角色,不再仅是信息的使用者,也成为了信息的发布者^[2],在获取信息的同时,根据自己掌握的公众或私人信息发表自己对股市的看法^[4]。由此,股评信息更多的包含

了公众的意见和观点,为提取股民情感提供了丰富资料。同时,情感分析技术的发展为提取股民情感提供了技术基础。

2.1 股评情感对股市走势的影响

关于股评有用性研究中,学者们已经证实,在股票市场的分析中需要考虑投资者情感因素。

早期研究主要依据经典资本资产定价理论^[5,6],而随着行为金融学的发展,股评对市场影响的研究角度从股评促进市场有效性逐步转变为股评所包含的投资者情感对市场行为的影响^[7,8]。投资者情感对股市走势影响的研究中,大多选取间接指标来衡量投资者情感,如封闭式基金折价、交易量、流动性方面指标^[9];或是选取用直接指标,即将对投资者进行调查和收集的数据作为投资者情感,如美国个体投资者协会指数^[10]、投资者智能指数^[11]等,国内则主要采用央视看盘、机构看盘、华鼎多空民意调查等指标^[12]。

直接情感指数通常以问卷调查的形式来获得^[13]。中国股民情感研究起步较晚,直接情感指数多反映机构对股盘的看法,而面向广大股民的情感统计则较少,为数不多的几个也因股民参与积极性少或不了解该指数的存在等问题,造成参与者少无法代表市场的情况。此类直接指数数据难以收集,且可靠性可能较低。

间接情感指数选取间接指标来反映投资者情感,难以评价其是否能够真实反映投资者情感。另外间接情感指数依赖于相关的金融数据,但相关数据的范围却无法准确定义,例如有研究者认为天气、日照等也属于相关数据^[14],因此间接情感指数的准确性也很难得到保证。

情感分析技术能够较好的解决直接情感指数和间接情感指数的问题。情感分析通过挖掘商品评论的文本内容,识别消费者对该商品的褒贬态度^[3]。随着网络的发展,公众观点、网络留言等主观性文本越来越容易获得,因此学者们开始尝试利用情感分析技术从股评中提取投资者情感,进而研究其对股市走势的影响。

Pilar Corredor 分析了投资者情感对欧洲4个股票市场的影响,发现情感对股票市场影响显著^[15]。Porshnev A 收集 twitter 上的股评数据,提取股评者情感,并将情感分成8类,利用SVM分类,并利用神经网络与遗传算法预测道琼斯工业指数^[16]。Antweiler 和 Frank 利用 Yahoo! Finance 和 Raging

Bull 的网络股评,对其进行情感分类,最终得到景气指数,通过对股票走势和景气指数的相关性研究,认为股评不只是噪声,有助于解释股票交易量和股票的波动性^[7]。

2.2 情感分析技术

情感分析涉及自然语言处理、文本分类、文本挖掘等领域。包括3个任务:主客观文本分类、情感极性判别、情感强度判别。其中,前者属于文本预处理,后两者才涉及情感分类^[17]。本文研究对象为网络股评,属于主观文本,因此主客观文本分类不是本文研究内容。

文本情感分类有两种方法:基于统计自然语言处理的方法(统计方法)和基于情感词汇语义特性的方法(语义方法)。前者的相关研究较多,分类方法较成熟,适用于包括中文在内的多种语言;后者则由于语言的复杂性等因素应用范围较小^[18]。基于统计自然语言处理的方法,利用机器学习算法对统计语言模型进行训练,然后采用训练好的分类器对新文本进行识别^[19]。步骤包括:文本预处理、文本表示、利用分类器判断情感极性、情感强度等。

(1) 文本预处理

1) 去除文本中无关内容。包括 html 标记、无法识别的字符、广告等内容。

2) 特殊词的处理,包括停用词去除、词根还原、错字修正。

3) 词语识别及统计。将非结构化的文本表示为字、词、短语等特征形式,经过统计,得到某个字、词、短语等的数量或权重。

上述预处理涉及分词、词性标注、短语识别,目前已有较为成熟的词语识别程序,如汉语词法分析系统 ICTCLAS,其具备中文分词、词性标注、命名实体识别、支持用户词典的新词识别和关键词采集等多种功能。

(2) 文本表示:向量空间模型(VSM)

向量空间模型是统计自然语言中较有效的文本表示方法^[18]。步骤如下:①将文本分为多个特征项,这些项可以是字、词,或是短语,互异的特征项可表示为 t_k ($k=1,2,\dots,n$)。那么此文本可表示为 $\{t_1, t_2, t_3, \dots, t_k, \dots, t_n\}$ 。其中 t_k 无先后顺序关系。②依据特定原则,将相应权重 ω_k ,赋予特征项 t_k ,该文本可表示为 $\{\omega_1, t_1, \omega_2, t_2, \omega_3, t_3, \dots, \omega_k, t_k, \dots, \omega_n, t_n\}$ 。③当有多个文本时,以该文本集内每条文本为纵轴、以文本集特征项为横轴,形成矩阵 D , D

(n, m)就是特征项在第 n 个文本中的权重值。称 D 为文本集的向量空间模型^[20]。

(3) 文本分类:支持向量机(SVM)

常用的文本分类算法有:Rocchio 算法、k-最近邻(KNN)、朴素贝叶斯(NB)、支持向量机(SVM)、线性最小平方拟合(LLSF)和神经网络法(NNet)等^[21]。已有研究显示,Vapnik 等提出的 SVM 表现出色。Pang 等的研究表明,SVM 分类器的表现比 ME 和 NB 都好^[22]。叶强、张紫琼以旅游博客的评论作为语料库,对 NB、SVM 的分类效果进行比较,实验显示,SVM 优于 NB^[19]。唐慧丰等对基于监督学习的中文情感分类算法进行比较,选取中心向量法、KNN、Winnow、NB 和 SVM 进行对比,实验结果表明,SVM 与其他算法相比能达到较高的精度^[21]。

2.3 股票市场的 GARCH 模型

GARCH(广义自回归条件异方差)模型在股市波动研究中倍受关注。Bollerslev 通过扩展 ARCH(自回归条件异方差)模型得到 GARCH 模型。该模型是专门针对金融数据而量身订做的回归模型,可以对误差的方差进一步建模^[23],且能较好的模拟金融数据的尖峰厚尾性,适用于异方差时间序列,因此特别适用于金融数据波动性的分析和预测。GARCH 模型可以与 ARMA(自回归滑动平均)模型结合,增加均值方程的变换形式。

本文研究情感对股市的影响,将情感作为外生变量,引入到 GARCH 模型。有两种方法:①将外生变量带入 GARCH 模型的方差方程中,构成 ARMA-GARCHX 模型;②将外生变量用于建立 ARMA 模型作为均值方程,再建立 GARCH 模型模拟其方差。其中,第②种较为常见。彭潇熟等采用 ARMA-GARCH 模型将石油价格和美元指数作为影响黄金价格走势的外生变量,该模型的整体精度好于未引入外生变量的模型^[24]。

2.3.1 ARMA 模型

ARMA 模型由 Jenkins 和 Box 所创立,借助时间序列的随机特性,运用时间序列过去值、当前值以及滞后随机扰动项的加权来建立模型,从而解释并预测时间序列的变化发展规律^[25]。

一般来说,股票的收益不仅受于目前及过去各种随机因素的影响,也取决于本身的过去值。因此 ARMA(p, q)模型表示为:

$$y_t = \omega_0 + \sum_{i=1}^p \omega_i y_{t-i} + \sum_{j=1}^q \delta_j \mu_{t-j}, (\omega_i, \delta_j \neq 0) \quad (1)$$

其中, p 和 q 为模型的自回归阶数和移动平均阶数, μ_t 为独立误差项, y_t 为平稳序列。

2.3.2 GARCH 模型

Engle 提出 ARCH 模型,较好地反映金融资产序列尖峰厚尾、波动聚集性、杠杆效应等特点。在此基础上, Bollerslev 提出 GARCH 模型,假定随机误差项的条件方差依赖于误差项前期值平方以及误差项条件方差的前期值^[23]。

GARCH 模型的一般形式为:

$$y_t = u_t + \varepsilon_t \quad (2)$$

$$\varepsilon_t | \Psi_{t-1} \sim N(0, \sigma_t^2) \quad (3)$$

$$\left. \begin{aligned} y_t &= \omega_0 + \sum_{i=1}^p \omega_i y_{t-i} + \sum_{j=1}^q \delta_j \mu_{t-j} + \gamma_t X_t, (\gamma_t \neq 0, \omega_i, \delta_j \neq 0) \\ \sigma_t^2 &= \alpha_0 + \sum_{i=1}^p \alpha_i \sigma_{t-i}^2 + \sum_{j=1}^q \beta_j \varepsilon_{t-j}^2, (\alpha_i, \beta_j \geq 0) \end{aligned} \right\} \quad (5)$$

2.3.4 ARMA - GARCHX 模型

GARCH 模型表明, ε_t 是外部输入的函数,它在一定程度上影响着金融波动。有研究认为公式(4)中的参数 δ_j 表述了在市场上,新的金融信息对股票收益造成的影响^[26]。

考虑到这些因素,将金融信息作为 GARCH 模

$$\left. \begin{aligned} y_t &= u_t + \varepsilon_t \\ \varepsilon_t | \Psi_{t-1} &\sim N(0, \sigma_t^2) \\ \sigma_t^2 &= \alpha_0 + \sum_{i=1}^p \alpha_i \sigma_{t-i}^2 + \sum_{j=1}^q \beta_j \varepsilon_{t-j}^2 + \gamma_t X_t, (\alpha_0 > 0, \alpha_i, \beta_j \geq 0) \end{aligned} \right\} \quad (7)$$

其中, p, q 代表 2 个时间滞后, $\alpha_i, \beta_j, \gamma_t$ 代表不确定的非线性相关性^[27]。

(1) 在上述公式中, y_t 为股票价格的日收益率, 即:

$$y_t = \frac{v_t}{v_{t-1}} \quad (8)$$

其中, v_t 表示针对某具体股票或指数, 第 t 天的收盘价。

(2) 在公式(5)和公式(7)中, 投资者情感 X_t 通过网络股评进行收集, 并计算得出情感指数。

ARMA-GARCHX 模型的外生变量多为虚拟变量, 常用于解决股市的星期效应。Li 等认为, 情感可以作为外生变量加入 GARCH 模型, 但是该想法并没有直接实践, 而是以此为据, 将情感作为影响因

$$\sigma_t^2 = \alpha_0 + \sum_{i=1}^p \alpha_i \sigma_{t-i}^2 + \sum_{j=1}^q \beta_j \varepsilon_{t-j}^2, (\alpha_i, \beta_j \geq 0) \quad (4)$$

其中, y_t 表示日收益率, u_t 表示确定的平均回报率, ε_t 表示随机项, 即预测残差, Ψ_{t-1} 代表在时间 t 可获得的信息集。

2.3.3 ARMAX - GARCH 模型

股票价格存在序列自相关性, 也存在尖峰后尾性, 因此结合 ARMA 模型和 GARCH 模型对股票走势进行分析。同时, 互联网的股评反映出投资者对股市走势的看法。因此, 可以在 ARMA 方程中加入外生变量情感指数, 来探究情感与股市走势的关系。

结合公式(1)ARMAX-GARCH 模型:

型方差方程的外生变量, 即:

$$\sigma_t^2 = \gamma_0 + \sum_{i=1}^p \gamma_i \sigma_{t-i}^2 + \sum_{j=1}^q \delta_j \varepsilon_{t-j}^2 + \omega_t X_t, (\omega \neq 0, \gamma_i, \delta_j \geq 0) \quad (6)$$

其中, X_t 是第 t 天的网络金融信息量, 此处认为是情感因素。因此加入外生变量的 GARCH 模型可以表示为

素加入 SVM 模型的矩阵, 通过由情感指数、前 n 天股价、前 n 天股票交易量构成的 SVM 模型矩阵, 并分出训练集和测试集, 利用 SVM 进行分类和预测, 效果良好^[26]。

已有研究多在 ARMA 方程中加入外生变量而较少在 GARCH 模型方差方程中加入, 对情感因素对股市走势的影响上也缺乏实证研究。而外生变量加入方差方程和均值方程哪个更优, 过往研究没有给出答案, 仍需探索。

2.4 研究评述

投资者情感影响股市走势的研究中, 大多选取直接或间接指标来衡量投资者情感, 并未利用情感分析技术从股评中提取投资者情感进行分析。

GARCH 模型适用于异方差、尖峰厚尾的金融时间序列,且大多数的股指走势符合 GARCH(1,1) 模型。但将情感作为外生变量的 GARCH 模型的拟合效果仍待探索。为此,本文希望:

(1) 针对直接及间接情感指标难以准确衡量反映市场情感的问题,本文拟利用情感分析的方法提取股评中情感值。搜集新浪股吧中的股评,利用情感分析技术提取情感。选取 SVM 分类器,训练和预测每条股评的情感值,并利用情感指数公式整合每天的情感指数。

(2) 针对网络股票评论的具体预测方法的不足,利用 ARMA-GARCH 模型,将提取的情感指数作为外生变量分别加入模型的均值方程和方差方程,来拟合和预测股票价格的波动情况,GARCH(1,1) 作为基本对比模型,并且将 2 个模型进行对比分析。

3 模型构建

基于情感分析的 GARCH 模型设计如图 1 所示,包括 3 部分:

(1) 收集网络股评并计算情感值。收集信息包括股评内容,发表日期。股评的情感包括看涨(+1)和看跌(-1)两类。

(2) 根据网络股评的情感值,以日为单位,计算情感指数,获得网络股评中每日的投资者情感指数,将其作为 GARCH 模型的外生变量 X 。

(3) 将情感指数 X 分别带入 GARCH 模型中的均值方程和方差方程,建立 ARMA-GARCHX 模型和 ARMAX-GARCH 模型,对股指进行动态拟合及预测,并将 2 个模型进行比较。

3.1 情感分析

本文通过情感分析技术,将股评情感结构化,识别其看涨看跌属性。情感分类已有的研究显示,基于统计的分类方法,准确率高于基于语义的方法,因此本文采取基于统计的方法对股票评论进行情感分类^[28]。基本过程如图 2 所示:预处理、文本表示(特征项选择、特征项降维、特征项权重设置)、分类器处理,最终得到一个有关情感类别的输出。

步骤 1: 预处理

预处理阶段对标点符号的处理我们采用以下方式:虽然某些标点符号(如叹号、问号、“^ - ^”等)会对情感强度产生影响,但非标点文本足以表达绝大部分情感,加之网络用语的标点使用不规范。因此,本文将删除标点、表情等符号,以简化处理过程。

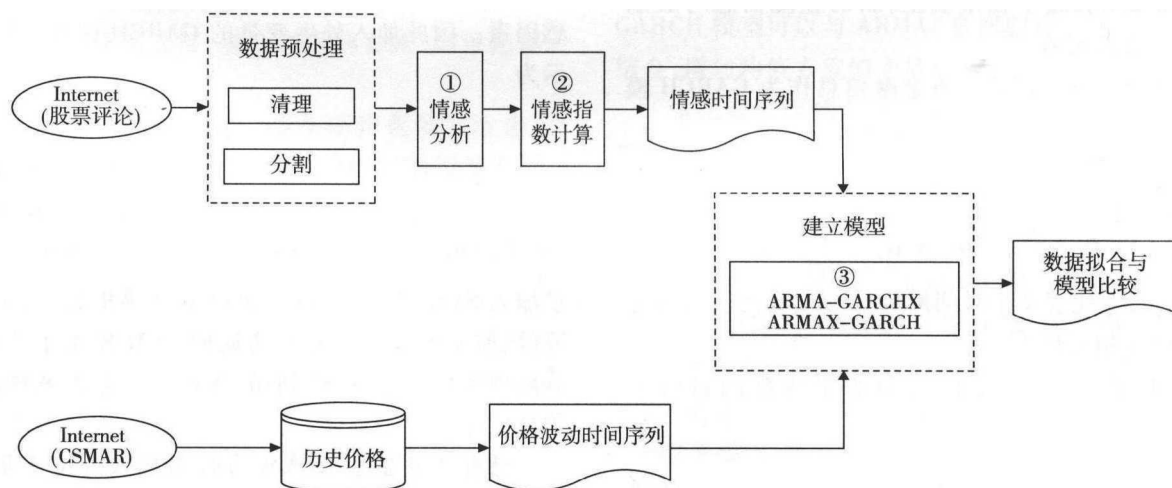


图1 基于情感分析的 GARCH 模型设计步骤

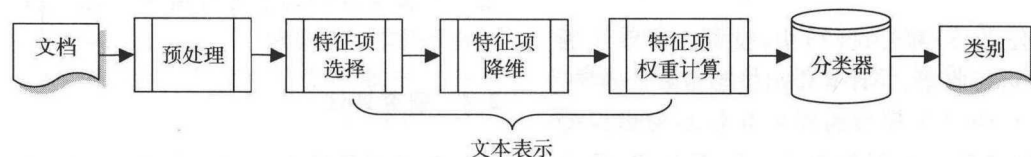


图2 情感分类的过程示意图

步骤 2：文本表示

采用向量空间模型 (Vector Space Model, VSM) 表示文本。基本过程是：对文本进行分词处理，然后根据训练样本集生成特征项的序列 $T = T(t_1, t_2, \dots, t_n)$ ，再根据 T 对训练样本集和测试样本集中的文档进行赋值，生成向量 $D = D(t_1, w_1, t_2, w_2, \dots, t_n, w_n)$ ，简记为 $D = D(w_1, w_2, \dots, w_n)$ 。其中 w_k 为特征项 t_k 的权重。VSM 模型涉及 3 个问题：特征项选择、特征项降维和特征项权重计算。

(1) 特征项选择。NVAA 模型选取词语作为基本的特征项单元，即提取句中名词、动词、形容词、副词组成特征项单元，用以表示语句。例如表 1：

表 1 词性示例

词性	举例
名词	走势, 大盘
形容词	不错, 很好
动词	看涨, 涨停
副词	很, 也

在提取特征项的过程中，按照以下原则处理：

1) 考虑否定词影响，以便更好的反应股评真实情感。将否定词（“不”、“不是”、“没有”等）与其后紧随的词组合在一起考虑。例如“此股票近期不会大涨”，提取的特征词应为“不会大涨”。

2) 删除形容词、副词中的“的”、“地”、“得”等字眼。与情感表达无关，删除后更简单。

3) 股票市场专用词，如“回调”、“抬轿”、“庄家”、“站岗”等，分词时将词汇准确分类。

4) 在评论中，无可避免会出现一些错别字，碰到这种情况，人为地改正这些错字，因为错误是随机的，偶然的，所以必须加以纠正才能正确反映作者的原意。

表 2 给出两个特征选取的例子。

(2) 特征降维和特征权重计算

1) 布尔权重值

布尔权重是最简单的特征权重计算方式，如果某一特征项出现在文本中，则权重为 1，否则为 0。研究表明语言的褒贬倾向主要取决于某情感词语是

否出现，而不是出现的次数，因此在特征权重计算方式中，布尔值表现最为出色，SVM 达到 82.9% 的分类准确率^[21]。

表 2 特征选取示例

文件	名词	形容词	动词	副词
明天一定会涨的，涨多少要看庄的水平了。	庄		涨	一定
	水平		涨	
			要	
	两天		看	
航天军工板块都是创新高的股票，可能要回调了。	航天军工板块	新高	是	很
	股票		创	都
			要	
			回调	

2) 文档频率

文档频率 DF (Document Frequency) 是指出现某个特征项的文档的频率。通过汇总特征项和特征值出现频率，根据频率对特征项由高到低排序，再根据设定的阈值筛选出频率较高的特征项，实现降维。因为 DF 方法简单易行，可扩展性好，分类准确率相对较高，适合超大规模文本数据集的特征降维，所以本文实验中采取 DF 方法作为降维的方法^[29]。

本文利用文档频率法降维后获得出现频率较高的特征项，用这些特征项将文本表示成向量形式，并采用布尔权重法设置权重。由此，完成语料的结构化，将语料转化为机器可识别的有序数据，用于进行情感分类。

步骤 3：SVM 分类器

根据上文研究可知，SVM 算法在诸多分类器中运用最为广泛，分类准确率较高。为此使用 Chih-Chung Chang 和 Chih-Jen Lin 的 LIBSVM-A Library for Support Vector Machines 所提供的 SVM 分类器进行实验操作，操作步骤如图 3 所示。

3.2 情感指数计算

经过 3.1 的情感分析，利用 SVM 自动分类，可

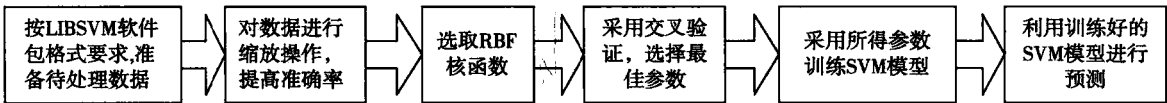


图 3 SVM 操作步骤

以获得每篇股票评论的情感倾向。本节将通过构建情感指数,从大量的股票评论量化出投资者情感值。选择“看涨指数(*bullisness index*)”,以日为单位,对每天的每篇股票评论的情感值进行整合,从而形成投资者情感值。该看涨指数由 Antweiler 和 Frank 提出,在 Antweiler 的研究和其他研究中被验证是较好的选择^[7]。

$$bullisness\ index = \ln \left[\frac{1 + M^{bull}}{1 + M^{bear}} \right]$$

其中, M^{bull} 指一天当中看涨的股票评论的数量, M^{bear} 指一天当中看跌的股票评论的数量。看涨指数若大于0,说明投资者的整体情感是看涨的,如果看涨指数小于0,则整体情感是看跌的。

3.3 GARCH 模型

将3.2节得到的情感指数带入 GARCH 模型的均值方程和方差方程[公式(5)和公式(7)],以基础的 GARCH(1,1)模型为对比模型,得出结果并进行分析。

4 实验分析

4.1 数据描述

网络股评通常集中在微博和贴吧,而股吧数据更适合本研究。首先,新浪微博有“微博看市”,但该版块多为机构股评,且针对对象较为混乱,有些专家发表大盘走势意见,有些则提供应买哪只股、何时买等操作性意见;其次,微博数量大且话题广,较难

筛选与个股相关的所有微博评论,且发言人对某股或对应的公司的言论不仅仅是针对股票投资。例如,以“中国石油”检索出的微博,相当一部分是针对油价波动的言论。而股吧中不仅有机专家,更多的是一般股民;并且股吧言论目的一致性较强,都是针对股票投资,而极少涉及其他方面。

根据 Alexa 中国官方网站的排名,以新浪财经(finance.sina.com.cn)的股票股吧为股评来源,进行数据收集和情感提取。共收集新浪股吧中评论数最多的30只组成沪深300指数的股票,既保证股票股评的数量,以便进行情感分析,提取情感指数;又因股票取自沪深300指数的组成股票,并且评论数量多说明其关注度高,能够保证股票具有一定代表性。具体股票信息请见表3:

表3中评论数目是剔除了空值、无法识别的评论(例如评论中仅有“¥、&”等符号)、零散日期评论(如某股票贴吧评论从2012年1月开始,基本每个交易日都有评论,而2011年中偶尔有个别日期有评论,则剔除2011年及其之前的不连续的评论)之后的评论数目。

4.2 特征值统计

收集包括万科A(sz000002)和浦发银行(sh600000)等30只股票,在新浪股吧中2011年3月至2014年9月的约200万条评论。将2014年第1季度的股评作为训练语料,约5万条,进行人工标注。其余作为测试语料,通过SVM模型进行训练,自动标注。

表3 股票名称、代码及评论数量

万科 A (sz000002) 56381 条评论	河北钢铁 (sz000709) 41730 条评论	民生银行 (sh600016) 65798 条评论	退市长油 (sh600087) 40155 条评论	贵州茅台 (sh600519) 57279 条评论
中国宝安 (sz000009) 67293 条评论	京东方 A (sz000725) 63799 条评论	中国石化 (sh600028) 46221 条评论	特变电工 (sh600089) 41691 条评论	上海石化 (sh600688) 33512 条评论
中兴通讯 (sz000063) 35275 条评论	苏宁云商 (sz002024) 93134 条评论	中信证券 (sh600030) 60795 条评论	同方股份 (sh600100) 35246 条评论	大商股份 (sh600694) 38192 条评论
TCL 集团 (sz000100) 84179 条评论	浦发银行 (sh600000) 82582 条评论	三一重工 (sh600031) 39989 条评论	中科英华 (sh600110) 45230 条评论	东方集团 (sh600811) 31277 条评论
中联重科 (sz000157) 68660 条评论	武钢股份 (sh600005) 26171 条评论	中国联通 (sh600050) 53804 条评论	重庆啤酒 (sh600132) 63422 条评论	四川长虹 (sh600839) 64191 条评论
攀钢钒钛 (sz000629) 72982 条评论	包钢股份 (sh600010) 84181 条评论	海信电器 (sh600060) 39854 条评论	广汇能源 (sh600256) 46260 条评论	梅雁吉祥 (sh600868) 46386 条评论

人工标注将股评分为4类:①看涨:对股票市场整体发展看好,预期会上涨;②看跌:对股票市场发展态势比较悲观,预计会下跌;③盘整:没有明确对股票的涨跌做出判断,但内容表达了与股市相关的信息;④其他:表达的信息与股票完全无关。

因为“盘整”和“其他”类股评的评论内容混乱,在分类过程中很可能对分类效果产生影响,降低分类的准确率。因此,剔除掉与股票无关的“其他”类型的股评,同时剔除态度不明确的“盘整”类股评,仅留下对本文研究有价值的信息,主要是看涨信息和看跌信息。

利用中科院的汉语分词系统 ICTCLAS 进行分词,通过对训练集进行人工标注剔除盘整和其他类股评,利用 DF 方法,提取看涨看跌类股评中出现频率较高的特征词。将 libsvm 作为 SVM 分类器,利用其将测试集分为2类,看涨(+1)和看跌(-1)。

采用分类准确率来评价分类结果:①分类准确率 $P = (A + D) / (A + B + C + D)$;②正类准确率 $P_p = A / (A + B)$;③负类准确率 $P_n = D / (C + D)$ 。A, B, C, D 的含义如表3所示。

表4 分类准确率

	实际为看涨的 评论数	实际为看跌的 评论数
标注为看涨的评论数	A	B
标注为看跌的评论数	C	D

表5 情感分类结果

方法	P_p	P_n	P
基于统计的方法	96.74%	94.81%	95.78%

表5结果显示,采用3.1中的情感分类方法,整体的分类准确率达95.78%,分类效果较好,可以接受的。因此,可以采用该方法对测试集数据进行分类。

4.3 实验结果和分析

建立 GARCH 模型之前,首先检验30只股票价格的平稳性,此30只股票均不存在单位根,为平稳序列。其次,检验30只股票的股指序列是否存在 ARCH 效应。经检验,除1只股票(武钢股份,sh600050)外,其余29只股票均存在不同程度的 ARCH 效应,满足设立 GARCH 模型的条件。根据

剩余29只股票数据,建立 ARMAX-GARCH(1,1)模型和 ARMA-GARCHX(1,1)模型,并与基础的 GARCH(1,1)模型相比较,得出结论。其中 X 是指股票情感指数的时间序列,ARMAX-GARCH 模型中,情感指数作为外生变量加入 GARCH 模型的均值方程,对股票价格的时间序列进行拟合,构建函数,如公式(5);而 ARMA-GARCHX 是将情感指数作为外生变量加入方差方程进行拟合,如公式(7)。

构建 GARCH 模型时,需要同时满足以下条件,否则构造的 GARCH 模型拟合无效。

(1)数据存在 ARCH 效应(条件异方差性),才有必要建立 GARCH 模型。通过检测数据 ARCH LM 的滞后10阶的结果,观察 P 值,如果 P 值小于0.1,说明数据存在 ARCH 效应,否则不存在 ARCH 效应,无法建立 GARCH 模型;

(2)GARCH 模型系数 α 和 β [见公式(4)、公式(5)、公式(7)],均大于零,且和小于1;

(3)各个参数系数值要通过检验,显著不等于0,P 值需小于0.1;

(4)建立的 GARCH 模型满足上述2个条件后,再次检查模型的 ARCH 效应,ARCH 效应消失,即 P 值大于0.1。

针对同一组数据,不同的 GARCH 模型进行比较时,将用到对数似然函数值、AIC 值和 SC 值。通常,具有较高的对数似然函数值、较小的 AIC 值和 SC 值的模型较另一个更优。

表6、表7和表8分别为模型 GARCH(1,1)、ARMAX-GARCH(1,1)和 ARMA-GARCHX(1,1)的检验结果。表中 L 表示对数似然函数值,AIC 表示赤池信息量准则值,SC 值表示施瓦茨准则值。

α, β 表示 GARCH 模型方程的系数, γ 表示情感指数系数。公式请见公式(5)和公式(7)。

情感滞后期数是指在模型拟合中用的是当期还是前 n 期的情感指数。如果是当期,则情感滞后期数为0,如果用的是前一天的情感指数,则情感滞后期为-1。

股票代码中 * 代表此股数据在模型中,相较于其他2个模型拟合效果最优。而较优的模型除符合各约束条件外,还应有较高的对数似然函数值和较小的 AIC、SC 值。以上3表中 L 代表对数似然函数值。3个模型通过比较对数似然函数值、AIC 和 SC 值来得出较优结果。

表6 GARCH(1,1)

股票代码	L	AIC	SC	α	B
sz000002	2348.107	-5.12032	-5.10454	0.033663	0.928622
sz000009	1935.11	-4.28406	-4.26808	0.054369	0.897877
sz000063	2164.481	-4.63126	-4.61051	0.359264	0.115291
sz000100	2023.396	-4.83348	-4.81652	0.081952	0.849175
sz000157	2438.287	-5.23503	-5.21423	0.029347	0.969242
sz000629	2088.091	-4.69615	-4.67997	0.018794	0.975066
sz000709*	2710.853	-5.53549	-5.51551	0.094451	0.847551
sz000725	2508.788	-5.45112	-5.43012	0.063511	0.90148
sz002024	2214.278	-4.71564	-4.6898	0.032901	0.958549
sh600000	3354.695	-7.15747	-7.13161	0.061999	0.930534
sh600005	2842.913	-6.0189	-5.99319	0.22501	0.349275
sh600010	2051.952	-4.48893	-4.46253	0.396643	0.55475
sh600016	2499.13	-5.18407	-5.1689	0.146509	0.818599
sh600028	2682.682	-5.83155	-5.80003	0.091772	0.646357
sh600030	2258.554	-4.79183	-4.77122	0.071706	0.175501
sh600031	2382.636	-4.96275	-4.94752	0.064937	8.70E-01
sh600060	2109.142	-4.61367	-4.59785	0.055235	0.870215
sh600087	1369.934	-4.80963	-4.77905	0.07568	0.495472
sh600089	1303.687	-4.81887	-4.77908	0.239806	0.408314
sh600100	2313.132	-4.78784	-4.77269	0.027973	9.62E-01
sh600110	1565.131	-4.57432	-4.55444	0.065176	0.827835
sh600132	1895.34	-4.77611	-4.7525	0.444311	0.51575
sh600256	2109.578	-4.60454	-4.58874	0.059253	0.791723
sh600519	1820.088	-5.02376	-4.9984	0.15507	0.448754
sh600688	2198.591	-4.82418	-4.79242	0.639719	0.262039
sh600694	2216.115	-5.02183	-5.00012	0.032847	0.962067
sh600811	2145.395	-4.81754	-4.79599	0.42385	0.122418
sh600839	2261.8	-5.06457	-5.04845	0.050427	0.938859
sh600868	2399.624	-4.96088	-4.93059	0.3479	0.622762

表 7 ARMAX - GARCH

股票代码	L	AIC	SC	α	β	γ	情感指数滞后期数
sz000002 *	2398.895	-5.22685	-5.20054	4.59E-02	0.904989	0.008676	0
sz000009 *	1954.053	-4.32384	-4.30254	0.053305	0.909583	0.024132	0
sz000063 *	2168.501	-4.65841	-4.62195	-0.00947	0.134874	-0.00268	-4
sz000100 *	2096.689	-5.00404	-4.97576	0.102466	0.861443	0.018602	0
sz000157 *	2448.278	-5.25436	-5.22837	0.030096	0.968321	0.008593	0
sz000629 *	2103.797	-4.72927	-4.7077	0.019661	0.974331	0.020531	0
sz000709	2697.858	-5.55229	-5.52714	0.091305	0.843763	0.004982	0
sz000725 *	2511.169	-5.45412	-5.42788	0.063336	0.899647	0.003904	0
sz002024 *	2225.88	-4.73827	-4.70726	0.032796	0.958837	0.014432	0
sh600000 *	3367.141	-7.18193	-7.15089	0.065608	0.920075	0.004848	0
sh600005 *	2844.721	-6.02062	-5.98976	0.220808	0.349577	0.002449	0
sh600010 *	2071.41	-4.52941	-4.49773	0.353756	0.570767	0.018605	0
sh600016 *	2527.756	-5.24145	-5.22122	0.182352	0.807699	0.017734	0
sh600028	2684.859	-5.81097	-5.77957	1.80E-01	0.540047	0.001588	0
sh600030	无法建立模型,经格兰杰因果关系检验,情感指数不是该股票走势的格兰杰原因						
sh600031 *	2435.928	-5.08031	-5.0549	0.023852	0.971245	-0.00317	-3
sh600060 *	0.522886	-4.65426	-4.63316	0.042665	0.918961	0.021163	0
sh600087 *	1366.68	-4.79465	-4.75643	0.064402	0.532319	0.002876	0
sh600089	无法建立模型,情感指数系数不显著,说明情感对该股票拟合影响不大						
sh600100 *	2324.994	-4.81035	-4.79016	0.027343	0.963364	0.011718	0
sh600110 *	1602.339	-4.67742	-4.64428	0.064131	0.84496	0.008507	0
sh600132 *	1901.21	-4.78841	-4.7589	0.399199	0.499241	0.008987	0
sh600256 *	2149.939	-4.69866	-4.67228	0.050428	0.84606	0.00498	-3
sh600519 *	1839.587	-5.07216	-5.03413	-0.0244	0.628155	0.003178	-1
sh600688 *	2291.681	-5.00697	-4.97531	0.386763	0.383158	-0.00107	0
sh600694 *	2221.87	-5.03262	-5.00549	0.034179	0.961927	0.006874	0
sh600811 *	2199.194	-4.93407	-4.90174	0.03486	0.7958	0.002402	-1
sh600839 *	2273.266	-5.08804	-5.06654	0.050476	0.938068	0.013476	0
sh600868 *	2409.467	-4.97921	-4.94386	0.312105	0.617209	0.009832	0

表 8 ARMA - GARCHX

股票代码	L	AIC	SC	α	β	γ	情感指数滞后期数
sz000002	2384.734	-5.213	-5.18662	0.056648	0.871032	-0.0000912	-4
sz000009	1935.968	-4.28374	-4.26244	0.052128	0.896945	5.51E-05	0
sz000063	2159.839	-4.6369	-4.60568	0.050349	0.29342	-0.000182	-3
sz000100	2024.645	-4.83408	-4.81145	0.072479	0.870315	6.19E-05	0

续表

股票代码	L	AIC	SC	α	β	γ	情感指数滞后期数
sz000157	2372.278	-5.09092	-5.06492	0.077618	0.50304	-0.000373	0
sz000629	2089.817	-4.71371	-4.69208	0.004968	0.987644	5.36E-05	-4
sz000709	2642.204	-5.43399	-5.4139	0.096623	0.807988	4.66E-05	-2
sz000725	2440.148	-5.30174	-5.28074	0.112267	0.840176	2.00E-05	0
sz002024	2170.14	-4.62357	-4.60289	0.034929	9.52E-01	2.78E-05	-1
sh600000	无法消除 ARCH 效应,说明还有其他因子影响着股票市场的波动特征						
sh600005	2808.084	-5.94504	-5.91933	0.191594	0.310672	5.84E-05	0
sh600010	2051.821	-4.49138	-4.45967	0.45112	0.238275	0.00096	-3
sh600016	2508.05	-5.21134	-5.19108	0.093695	0.871063	-9.72E-05	-3
sh600028*	2752.688	-5.95811	-5.9267	0.053133	0.80977	-3.96E-05	0
sh600030	无法建立模型,经格兰杰因果关系检验,情感指数不是该股票走势的格兰杰原因						
sh600031	2392.78	-4.98181	-4.96152	0.004101	0.989035	4.05E-05	-1
sh600060	2109.964	-4.61328	-4.59218	0.03026	0.954446	4.26E-05	-3
sh600087	无法建立模型,情感指数系数不显著,说明情感对该股票拟合影响不大						
sh600089	1309.132	-4.83908	-4.79929	0.277535	0.166301	-0.000586	0
sh600100	2308.496	-4.79105	-4.77081	0.027592	0.958999	2.99E-05	-4
sh600110	1596.224	-4.68007	-4.64682	0.067697	0.78743	0.000213	-4
sh600132	1830.811	-4.61316	-4.58955	0.0597	0.918919	6.21E-05	0
sh600256	2107.669	-4.60826	-4.58715	0.062376	0.753495	-8.84E-05	-3
sh600519	1835.753	-5.06432	-5.03263	0.058847	0.44234	-0.000246	0
sh600688	2019.973	-4.43118	-4.39942	0.58447	0.211261	-0.000149	-3
sh600694	2182.867	-4.9576	-4.93586	0.02776	0.972185	-1.56E-05	-3
sh600811	1944.218	-4.36761	-4.34064	0.087754	0.455916	-0.000741	-3
sh600839	无法建立模型,无法消除 ARCH 效应,从而说明还有其他因子影响着股票市场的波动特征						
sh600868	2321.72	-4.8044	-4.77408	0.218512	0.432786	-0.000149	-3

上述3个模型中,GARCH(p, q)参数 p, q 均取1,情感指数 X 从当期 X_0 (情感指数滞后期数为0)到前5期 X_{-5} (情感指数滞后期数为-5),分别带入模型拟合,选取最优模型结果。

上述模型均经过ARCH检验,消除ARCH效应。且模型系数 α 和 β 满足系数约束条件 $\alpha > 0$, $\beta > 0$, $0 < \alpha + \beta < 1$ 。其中“无法建立模型”即模型系数不满足约束条件或最终无法消除ARCH效应。

结果显示,29只股票均可建立的GARCH(1,1)模型,但有少数几只股票无法建立ARMAX-GARCH或ARMA-GARCHX模型。原因主要有如下3类:

(1)sh600030:该股票在加入情感指数的两个

模型中均无法建立模型,经过格兰杰因果检验发现情感不是股票走势的原因。表8为sh600030格兰杰因果检验结果:

表8 sh600030 格兰杰因果检验

股票代码	原假设	F值	P值
sh600030	H0: 情感指数不是股票走势的格兰杰原因	3.00164	0.0107
	H1: 股票走势不是情感指数的格兰杰原因	1.68268	0.1361

(2)sh600000和sh600839:在ARMA-GARCHX

表 9 GARCH(1,1) 与 ARMA-GARCHX(1,1) 比较表

GARCH(1,1) 较优的股票	sz000002, sz000009, sz000100, sz000629, sh600010, sh600016, sh600028, sh600031, sh600060, sh600089, sh600110, sh600519	共 12 只股票
ARMA-GARCHX(1,1) 较优的股票	sz000063, sz000157, sz000709, sz000725, sz002024, sh600005, sh600100, sh600132, sh600256, sh600688, sh600694, sh600811, sh600868	共 13 只股票

模型建立后无法消除 ARCH 效应,应该还有其他因素影响股票波动。

(3)sh600089 和 sh600087: 情感指数系数不显著,说明情感对股票走势影响不大。

虽然 GARCH(1,1)模型应用最为广泛,但拟合效果相比于表 6 中 ARMAX-GARCH(1,1)模型较差。表 6 中显示,有 2 只股票无法进行 ARMAX-GARCH(1,1)模型的拟合,但其余 27 只股票中有 25 只股票的对数似然函数值、AIC 和 SIC 值均优于其他 2 个模型。说明该模型的拟合效果较其他 2 个模型更好。而通过对比 GARCH(1,1)和 ARMA-GARCHX 模型,在 2 个模型均可以建立的情况下,后者没有明显优于前者。表 9 为两模型比较结果。

股评的情感作用于不同的时期,有些股评反映作者当前的情感态度,有些则是对未来走势的预测。从情感指数的滞后期数来看,ARMAX-GARCH(1,1)模型虽拟合效果较好,但是情感滞后期数基本为 0,即该模型更多的利用股评情感里反映当期股票走势的部分,而较少利用反映未来走势的部分。而 ARMA-GARCHX 模型中,情感指数的滞后阶数集中于 -3 和 -4,即在该模型中,情感指数能够较好的预测未来 3~4 天的走势。因此在预测方面 ARCH-GARCHX 模型具有优越性。

5 结 论

本文将情感分析和机器学习方法相结合,利用改进的综合情感指数,提取并计算基于网络股评的投资者情感,并作为 GARCH 模型的输入项来预测股票价格的波动。研究结果显示,

(1)通过情感分析技术,从网络股评信息中提取投资者情感是可行的。且该方法具有数据样本量大、时效性强、更加真实、更接近投资者的真实情感等优点。

(2)采用提取出的投资者情感和 GARCH 模型,将投资者情感作为一个输入项来预测股票价格的波动(金融风险),这是可行的。通过实验发现,将情感

因素作为外生变量加入均值方程,比加入方差方程拟合效果更好,并且优于对比模型 GARCH(1,1)。

(3)利用 0 到 -5 期情感指数分别建立模型并选取最优者,发现最优的 ARMAX-GARCH 模型情感指数多集中于当期,即在该模型中情感指数多用来反映当前股票价格。而最优的 ARMA-GARCHX 模型的情感指数多集中于前 3、4 期,即该模型中情感指数中反映未来股票价格的部分被较多的应用于此模型。对此现象的解释为,组成情感值的情感可以根据其作用时间分为 2 类:第 1 类是对未来一段时期的预测展望,第 2 类是对已发生事情的评价,可以是当天的也可以是之前某天的。结果显示,当期情感作用于均值方程时,拟合效果更好,这可能是因为当期情感中对过去股票的评价反映了当前股票价格,但是并不能说明当期股评情感是影响股票价格的因素,说明第 1 类情感属于噪声信息。而第 2 类情感包含更多有价值的预测信息,更多的作用于方差方程,改变股价的波动幅度,更能有效预测未来股票的价格走势。

改进的方向包括:①改进情感分析方法,提高股评情感分类的准确率;②分析情感指数的构成,将股评情感分为对当期股价的看法和对未来股价走势的看法,可以通过股评语句中特定的时间词语来对股评进行此种分类,分析 2 种股评对 GARCH 模型均值方程和方差方程的影响;③本文的 GARCH 模型 2 个参数 p, q 均为 1,参数的限制和情感指数中的噪声可能是造成 ARMA-GARCHX 模型结果不够理想的原因。将来需排除股评中对已发生事情的评价,分离出股评中对未来的评价,研究此类情感作为外生变量对 GARCH 模型拟合效果的影响。

参 考 文 献

[1] 宋敏晶. 基于情感分析的股票预测模型研究[D]. 哈尔滨工业大学,2013.
[2] Dan Gillmor. We the media: Grassroots Journalism by the people, for the people [OL]. [2004-06-30]. <http://www.authorama.com/we-the-media.html>.
[3] Pang B, Lee L, Vaithyanathan S. Thumbs up. Sentiment

- classification using machine learning techniques [C]// Proceedings of the 2002 Conference on Empirical Methods in Natural Language Processing, Morristown, NY, USA: ACM Press, 2002: 79-86.
- [4] Logue D E, Tuttle D L. Brokerage house investment advice [J]. Financial Review, 1973, 28(1): 38-54.
- [5] Baker M P, Wurgler J. Investor sentiment in the stock market [J]. Journal of Economic Perspectives, 2007, 21(2): 129-151.
- [6] 马春林, 倪苏云, 吴冲锋. 股评家关注股票基本面因素: 股市信息供应量的影响因素分析 [J]. 上海经济研究, 2002(5): 47-51.
- [7] Antweiler W, Frank M. Is all that talk just noise? The information content of internet stock message boards [J]. Journal of Finance, 2004, 59(3): 1259-1295.
- [8] 池丽旭, 庄新田. 我国投资者情绪对股票收益影响——基于面板数据的研究 [J]. 管理评论, 2011(6): 41-48.
- [9] Baker M, Wurgler J. Investor sentiment and the cross-section of stock returns [J]. The Journal of Finance, 2006, 61(4): 1645-1680.
- [10] Brown G. Volatility, sentiment and noise traders [J]. Financial Analysts Journal, 1999, 55: 82-90.
- [11] Roger C G, Meir S. Bullish and bearish [J]. Financial Analysts Journal, 1998, 54(6): 63-72.
- [12] 季美惠, 宋顺林, 王思琪. 投资者情绪与会计信息价值相关性——基于中国上市公司的实证分析 [J]. 中大管理研究, 2014(2): 1-15.
- [13] 黄少安, 刘达. 投资者情绪理论与中国封闭式基金折价 [J]. 南开经济评论, 2009, 12(1): 96-101.
- [14] Saunders E M. Stock prices and Wall Street weather [J]. The American Economic Review, 1993, 83(5): 1337-1345.
- [15] Corredor P, Ferrer E, Santamaria R. Investor sentiment effect in stock markets: Stock characteristics or country-specific factors? [J]. International Review of Economics and Finance, 2013, 6(27): 572-591.
- [16] Porshnev A, Redkin I, Shevchenko A. Machine learning in prediction of stock market indicators based on historical data and data from Twitter sentiment analysis [C]// 2013 IEEE 13th International Conference on Data Mining Workshops. NY, USA: IEEE Press, 2013: 440-444.
- [17] 郑丽娟. 中文在线评论的用户情感分析及应用 [D]. 同济大学, 2014.
- [18] 金聪, 金平. 网络环境下中文情感倾向的分类方法 [J]. 语言文字应用, 2008, 5(2): 139-144.
- [19] 张紫琼, 叶强, 李一军. 互联网商品评论情感分析研究综述 [J]. 管理科学学报, 2010, 13(6): 84-96.
- [20] Tsatsaronis G, Panagiotopoulou V. A Generalized Vector Space Model for text retrieval based on semantic relatedness [C]// 12th Conference of the European Chapter of the Association for Computational Linguistics. NY, USA: ACM Press, 2009: 70-78.
- [21] 唐慧丰, 谭松波, 程学旗. 基于监督学习的中文情感分类技术比较研究 [J]. 中文信息学报, 2007, 21(6): 55-94.
- [22] Pang B, Lee L. Opinion mining and sentiment analysis [J]. Foundations and Trends in Information Retrieval, 2008, 2(1-2): 1-135.
- [23] Freisleben B, Ripper K. Volatility estimation with a neural network [C]// Proceedings of the IEEE/IAFE on Computational Intelligence for Financial Engineering. NY, USA: IEEE Press, 1997: 177-181.
- [24] 彭潇熟, 张德生. 国际黄金价格具有外生变量的 GARCH 预测模型 [J]. 黄金, 2011(1): 10-14.
- [25] Box G E P, Jenkins G M, Reinsel G C. Time series analysis: forecasting and control [J]. Journal of Marketing Research, 1977, 68(3): 343-344.
- [26] Li N, Liang X, Li X, et al. Network environment and financial risk using machine learning and sentiment analysis [J]. Human and Ecological Risk Assessment, 2009, 15(2): 227-252.
- [27] Ye Q, Lin B, Li Y J. Sentiment classification for Chinese reviews: A comparison between SVM and semantic approaches [C]// Proceedings of 2005 International Conference on Machine Learning & Cybernetics. NY, USA: IEEE Press, 2005: 2341-2346.
- [28] 王洪伟, 郑丽娟, 刘仲英. 中文网络评论的情感特征项选择研究 [J]. 信息系统学报, 2012(10): 76-86.
- [29] Azam N, Yao J T. Comparison of term frequency and document frequency based feature selection metrics in text categorization [J]. Expert Systems with Applications, 2012, 39(5): 4760-4768.

(责任编辑 魏瑞斌)