

# 金融学文本大数据挖掘方法与研究进展<sup>\*</sup>

姚加权 张锬澎 罗平

**摘要:**在金融学领域的传统实证研究中,所用数据多局限于财务报表和股票市场数据等结构化数据。而在大数据时代,计算机技术的进步使得数据类型不断丰富,研究者开始将非结构化的文本大数据引入到金融学领域的研究中,其主要包括上市公司披露文本、财经媒体报道、社交网络文本、网络搜索指数以及 P2P 网络借贷文本等,并对文本的可读性、语气语调、相似性以及语义特征展开研究。本文首先介绍了金融学领域文本大数据挖掘步骤和方法,描述了语料获取、预处理过程、文档表示以及文档的特征抽取;然后根据不同的文本信息来源,梳理了金融学文本大数据的研究进展;最后对未来金融学文本大数据的研究方法和研究内容进行了展望。

**关键词:**文本大数据 文本分析 机器学习 深度学习 数据挖掘

## 一、引言

在金融学领域的传统实证研究文献中,研究数据多局限于财务报表数据、股票市场数据等结构化数据(structured data)。而在大数据时代,计算机技术的不断提高使得数据类型更加丰富,文本大数据已经成为计算机可以解读和分析的数据,并能够对非传统领域的经济现象展开研究(Loughran & McDonald, 2016; Teoh, 2018)。这种非结构化数据(unstructured data)在公司对外披露以及股票市场中所占的比重较大,传递形式和表达方式更为多样化,尤其是在中国这种“听话听音,听锣听声”的高语境传播环境中(林乐和谢德仁, 2016),文本大数据在金融学领域中拥有较高的研究价值。文本分析(textual analysis)是指,运用特定的方法挖掘文本信息内容,从而对文本的可读性、情绪语调、语义特征以及相似性等文本特征进行分析。通过对上市公司披露文本、财经媒体报道、社交网络文本、网络搜索指数以及 P2P 网络借贷文本等文本大数据进行挖掘和分析,研究者能够从文本的披露行为(Loughran et al, 2009; 曾庆生等, 2018)、文本的情绪和语调(Li, 2010b; Loughran & McDonald, 2011; Garcia, 2013; Davis et al, 2015; 汪昌云和武佳薇, 2015)以及文本信息的市场反应(Antweiler & Frank, 2004; 游家兴和吴静, 2012)等方面展开研究,从而为金融学领域提供更丰富的研究内容和研究视角。

文本分析型研究有较长的历史, Jones & Shoemaker(1994)以及 Cole & Jones(2005)分别对会计文本内容以及管理层讨论与分析(management discussion and analysis, MD&A)的相关文献进行了综述。随后, Li(2010a)着重于计算机语言学、自然语言处理以及统计学的大样本文本分析,按照不同主题调查了企业披露文本的相关研究。再之后, Loughran & McDonald(2016)对国外会计和金融领域中的文本分析文献以及相关方法进行了调查和描述。Guo et al(2016)总结了机器学习方法

<sup>\*</sup> 姚加权、张锬澎,暨南大学管理学院,邮政编码:510632,电子邮箱:jiaquanyao@gmail.com, zkpworking0729@163.com; 罗平,中国科学院智能信息处理重点实验室、中国科学院计算所和中国科学院大学,邮政编码:100190,电子邮箱:luop@ict.ac.cn。本文受国家自然科学基金项目(71502152, U1811461),国家社科基金重大项目(18ZDA092)和国家重点研发计划课题(2017YFB1002104)的资助。感谢冯绪、梁平汉、唐国豪、王靖一、杨海生和赵轶星的宝贵意见。本文曾在 2019 海峡两岸暨港澳金融科技青年学者论坛、2019 中国数字金融研究联盟学术年会、第十八期香樟经济学 Seminar(广州)和 2019 中国“金融科技学”年会报告,感谢与会专家的评论,感谢审稿人的修改意见,文责自负。

在财务文本大数据分析中的应用。Gentzkow et al(2019)则描述了文本大数据的分析方法以及在经济学中的应用。Cong et al(2019)描述了金融市场中典型的英文文本来源,并讨论了神经网络模型与生成统计模型在文本分析领域的应用。在文本分析研究综述方面,国内研究者如唐国豪等(2016)整理了国内外基于文本情绪分析的行为金融研究进展,并总结了主要的文本分析方法。沈艳等(2018)综述了英文文本大数据分析在经济学和金融学领域中的应用,辅助以中文文本文献。张学勇和吴雨玲(2018)以国外文献为主,从网络新闻数据、搜索引擎数据、社交网络数据以及网络论坛数据四个方面,梳理了资产定价领域中运用网络大数据挖掘技术分析投资者心理和行为的研究内容。

综合来看,以往文献着重于介绍文本分析的主要方法,但缺乏对文本大数据挖掘步骤和方法的详细介绍。本文主要在以下方面拓展了以往的研究:首先,详细介绍了文本大数据挖掘步骤和方法,描述了文本的语料获取、预处理过程、文档表示以及文档的特征抽取。其次,介绍了国内外金融学文本大数据挖掘的主要文本信息来源,并根据不同的文本信息来源,梳理了金融学文本大数据的国内外研究进展,以便把握文本大数据目前在金融学领域中的研究方向和重点领域。最后,提出了未来的研究展望,期望有助于国内研究者进一步拓展文本大数据在金融学 and 经济学领域的应用。

## 二、金融学文本大数据挖掘步骤和方法

在金融学文本大数据研究中,早期的部分学者采用人工阅读的方式识别文本信息。但随着文本数量的增大,该方法不仅耗时耗力,且提取信息的精度也由于阅读者理解能力的差别而受到制约。因此,多数学者开始将计算机处理技术引入文本大数据的分析中,该过程主要包括:语料获取、文本的预处理、文档表示以及文档的特征抽取。然后,研究者再根据需要将抽取的文档特征应用到具体的相关性分析或因果分析中,具体如图1所示。

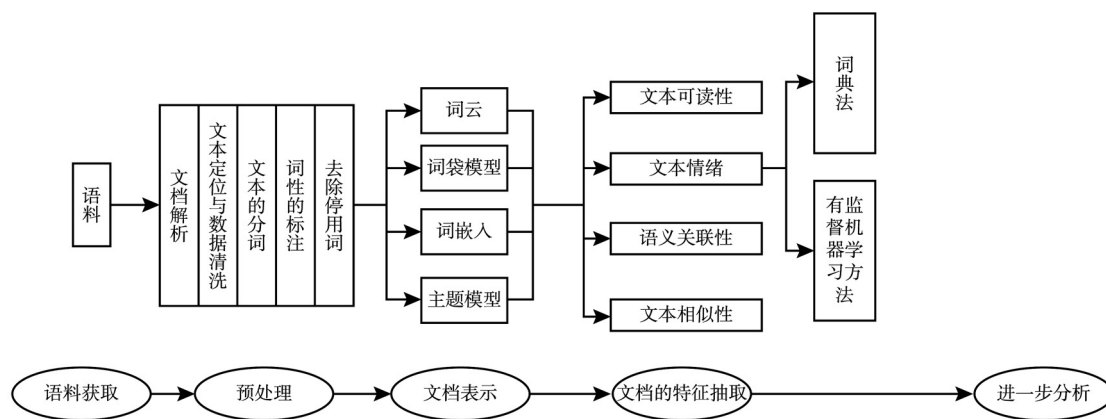


图1 文本大数据分析流程图

### (一) 语料获取

语料获取的方法主要包括:(1)手工收集。然而,该过程需要消耗大量的时间和人力成本。(2)网络抓取。由于文本量的提高以及文本大数据获取困难,大多数学者选择运用编程语言直接从网络中爬取文本大数据(Loughran & McDonald, 2014; Blankespoor et al, 2014)。该方法一方面能够及时地获取文本信息,另一方面还可以通过编程语言对文本格式和内容等进行整理,以便进行下一步分析。

### (二) 文本的预处理

在语料获取后,研究者需要对文本进行预处理,该过程主要包括文档解析、文本定位与数据清洗、文本的分词标注、词性的标注(part-of-speech tagging)以及停用词去除(stop words)五个步骤。

1. 文档解析。在信息披露监管制度下,企业需要以电子文档的形式定期或非定期地公开发布相关信息的文档。然而,这些文档仅仅实现了电子化存储,方便读者在电子设备上阅读,但这并不意

意味着机器可以自动处理,实现“机器可阅读”。在计算机领域,信息披露的电子化文档被统称为富格式文档(richly formatted documents)。这些文档包含文本段落、表格、图表等多种内容模态,通常会组织为层次化的目录结构,并经过美化的排版和格式处理以呈现给读者阅读。从文档的格式来看,绝大多数金融市场要求的信息披露文档是 PDF 格式。因此,解析富格式文档经常是进行文本预处理的第一步,即获取里面的信息内容。

在进行 PDF 文档解析的过程中需要注意两个方面:一方面,PDF 文档的生成不是一个可逆的过程,当我们使用 Word 或 Excel 的编辑器将文档导出为 PDF 文档后,虽然文档的排版格式等视觉呈现可得以保证,但文档内部的结构信息被部分或完全地丢失了。另一方面,解析后的文档是文本分析的基础,不精确的 PDF 解析可能会给后续的文本分析带来严重的影响。因此,针对金融学文本大数据,需要慎重选择文档结构的解析工具。<sup>①</sup>

2. 文本定位与数据清洗。一方面,研究者需要运用计算机程序对文本信息进行定位。例如,MD&A 部分是较多学者的研究对象,研究人员可以运用正则表达式来定位财务报告正文中 MD&A 部分的开头和结尾,进而将该部分内容提取出来。另一方面,研究者还需要对文本中被视为噪音的内容进行清洗和删除(Jiang et al, 2019),主要包括广告、超文本标记语言(HTML)、直译式脚本语言(JavaScript)等代码以及图片等。

3. 文本的分词。在英文文本中,单词被空格分开就自动完成了分词。另外,还可以通过词形还原(lemmatization)和词干提取(stemming)对单词进一步地划分。但中文文字之间没有空格切分,而且词语才是能够独立运用的最小语言单位。因此,研究者需要对中文文本进行专门的分词处理。目前,大多数学者采用了 Python 开源“jieba”中文分词模块来对企业财务报告、年度业绩说明会以及股票论坛帖子进行分词。中文文本分词存在三个难点,即切分颗粒度、歧义词的识别和新词的识别。切分颗粒度太小,容易破坏词语的意思。例如,容易将“机器学习”切分成“机器”和“学习”。针对歧义词,则应该选择合适的分词模式。例如,在使用“jieba”分词时,为了提高分词的精确度,应当选择精准分词模式。针对新词(如公司名称、产品名以及关键人物姓名),用户可以自定义词典以便帮助分词软件对新词进行识别。

4. 词性的标注。词性是识别语义信息的重要语法特征,例如名词、动词、连接词等,词性标注就是对切分后词语的词性做标记。通过词性标注,计算机能够识别词语的种类、消除词语歧义,进而能识别语法结构,降低计算机语义分析的难度。中英文在词性标注方面具有较大差异,英文单词在词性划分方面较为严谨,能通过词尾变换来揭示词性的变化,例如“-ing”、“-ness”和“-ment”等均对确认词性给予了具体的提示。但是,中文词语未对词性做出明确规范,主要靠语法和语义来识别词性,即“英语重形合、汉语重意合”。

5. 停用词去除。为了提高文本挖掘信息的精度,还需要对文本中的停用词进行剔除。停用词是指对句子语法结构很重要但本身传达意义较少的词语,其增加了文本数据的维度,提高了文本分析的成本。在英文文本中,停用词主要包括冠词(the, a)、连词(and, or)以及动词“to be”等(Gentzkow et al, 2019)。但在中文文本中,应当根据中文的语言习惯确定停用词,除了标点符号和特殊符号以外,还包括表示逻辑关系的连接词(和、然而等)以及俚语等。另外,停用词还需要根据研究的内容来决定。例如,当研究文本情感时,保留语气词以及特定的标点符号,均有利于衡量文本的情感程度。

### (三)文档表示

文本数据属于稀疏的高维度数据,计算机处理存在困难,因此对文本数据进行预处理后,还需要将文档中的信息以特定的方式表示出来,从而便于研究人员或者计算机进行下一步的分析。文档表示方法主要包括:词云(word cloud)、词袋模型(bag of words, BOW)、词嵌入(word embedding)和主题模型(topic model)。

<sup>①</sup>我们详细举例说明了选择精准 PDF 解析工具的重要性,详见本文的在线附录,有兴趣的读者可以联系作者获取。

1. 词云。词云是一种文本大数据的可视化技术。文本可视化是指将文本中比较复杂的内容和规律转化为视觉符号表达出来,进而能够使人们利用与生俱来的视觉感知快速获取文本中所蕴含的关键信息。词云技术能够描述词语在文本中出现的频率,当词语出现频率较高时,会以较大且醒目的形式呈现。

2. 词袋模型。词袋模型是一种建立在文字词组语序不重要的假设之上,将文本看作是若干个词语的集合,只计算每个词语出现次数的一种文本向量化的表示方法。该模型主要包括:独热表示法(one-hot representation)以及词频-逆文档频率法(term frequency-inverse document frequency, TF-IDF)。独热表示法操作简单。假设有两个文档“文本大数据在经济学中的应用”和“文本大数据在金融学中的应用”,基于这两个文本文档可以构建如下词表:〔“文本”,“大数据”,“在”,“经济学”,“金融学”,“中”,“的”,“应用”〕,按照该顺序进行词袋化后得到两个文档的词袋向量分别为:〔1,1,1,1,0,1,1,1〕和〔1,1,1,0,1,1,1,1〕,其中,“1”和“0”分别表示文档中有无出现这个词。然而,文档中并非每个词语均有相同的机会出现,大多数文本中只有极少数词语被经常使用,绝大多数词语很少被使用。因此,需要对每一个词语赋予其权重来更好地表示每个词语在文档中的作用。Loughran & McDonald(2011)运用 TF-IDF 方法计算了特定词语在文档中的权重。其基本公式如下:

$$idf_i = \log \frac{N}{df_i} \quad (1)$$

$$tf-idf_{i,j} = \begin{cases} \frac{(1 + \log(tf_{i,j}))}{(1 + \log(a_j))} \log \frac{N}{df_i} & \text{若 } tf_{i,j} \geq 1 \\ 0 & \text{其他} \end{cases} \quad (2)$$

式(1)中  $df_i$  定义为包含词语  $i$  的文档数量,  $N$  表示文档集中的文档总数,  $idf_i$  为逆文档频率。式(2)中  $tf_{i,j}$  为第  $j$  个文档中词语  $i$  出现的总次数,  $a_j$  为第  $j$  个文档中包含的词语数,  $tf-idf_{i,j}$  则为第  $j$  个文档中词语  $i$  的权重。但是,词袋模型存在以下问题:第一,忽略了文档中词语顺序和词语之间的语义关系,可能会产生歧义;第二,向量的维度取决于文档中词语的数量,当文档中词语数量过多时,很可能产生维度灾难。

3. 词嵌入。词嵌入是将维数为所有词的数量的高维空间嵌入到低维连续向量空间中的技术。通过词嵌入技术,可以将词语映射成低维连续向量空间中的向量,通过向量之间的距离和位置来表示文档中词语的上下文、语法和语义的相似性以及与其他词语的关系。在金融学文本分析中, Word2Vec 技术是常用的词嵌入技术,包括 CBOW(continuous bag of words)和 Skip-Gram 神经网络模型,可以通过训练使神经网络捕捉到更多词语之间的上下文信息,从而将每个词语映射成更低维度、稠密且包含更多语义信息的向量(Mikolov et al, 2013)。在 Word2Vec 技术中,词嵌入向量能得到不同词之间的类比关系,最经典的例子就是“king-queen=man-woman”,如图 2 所示。

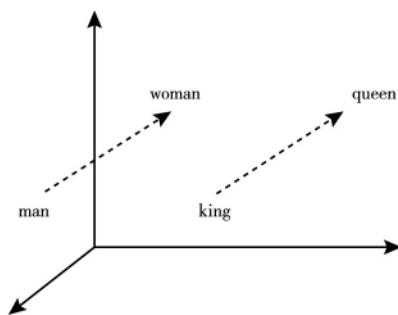


图 2 词嵌入模型实例图

4. 主题模型。最常用的主题模型是 LDA(latent dirichlet allocation)模型(Blei et al, 2003)。LDA 模型是在大规模语料集中提取主题信息的无监督机器学习方法,它假设文档生成包括两个步

步骤:第一步,假定每个文档均有对应的主题分布,在文档的主题分布中抽取一个主题;第二步,假定每个主题都有对应的词语分布,从上一步抽取的主题所对应的词语分布中抽取一个词语。通过将这两步迭代拟合到文档中的每个词语,即可得出每个文档的主题分布和每个主题的词语分布。

Huang et al(2018)指出 LDA 模型具有以下优势:首先,该模型克服了手动编码的局限性,能够对大量文本文档进行分类;其次,LDA 模型能够提供可靠且具有可复制性的文本主题分类,排除了人工文本分类的主观性;最后,LDA 模型不需要研究者预先指定相应的规则和关键词。然而,该模型的局限性在于预设主题个数的方式中加入了人的主观因素,这会影响主题个数的选择,进而影响主题的生成和文本的主题归类<sup>①</sup>。

#### (四)文档的特征抽取

目前,在金融学领域中针对文档的特征抽取方面主要包括四个方面:文本可读性(textual readability)、文本情绪(textual sentiment)、语义关联性(textual relatedness)以及文本相似性(textual similarity)。

1. 文本可读性。文本的可读性反映了读者理解文本信息的难易程度,文本可读性较低时,投资者会难以理解文本编辑者所传达的信息,进而会影响到投资者的投资行为。Li(2008)将迷雾指数(fog)应用到了文本分析中,指出迷雾指数越小时,年报的可读性越强。另外,学者们还利用了年报中的字数(You & Zhang,2009)和年报电子文档的大小(Loughran & McDonald,2014)来衡量年报的可读性。以往多数研究采用迷雾指数来衡量文本可读性(Li,2010b;Lehavy et al,2011),但这种方式仍具有一些问题。例如,如果将文本中每个句子的词语随机排序,那么文章将完全无法理解,但是原来的文本与随机排序后的文本所计算的迷雾指数完全相同(Jones & Shoemaker,1994)。此外,Loughran & McDonald(2014)研究指出,迷雾指数在测量商业文本可读性时具有局限性。正如 Loughran & McDonald(2016)提出,当衡量企业信息披露可读性时,将公司复杂性和年报可读性分开是困难的。当公司具有多种业务时,内部业务结构复杂的公司很可能会因为业务复杂性而披露难以阅读和理解的年报。因此,在衡量公司年报文本可读性时,应当考虑将企业的业务复杂性因素剔除。

2. 文本情绪。目前,文本情绪的提取方法主要包括词典法(dictionary-based approach)和有监督机器学习方法(supervised machine learning)。

词典法是指运用情绪词典来研究文本情绪或者语气语调的方法。一旦确定了词典,就可以使相关研究具有可复制性。另外,在构建词典时,金融领域的知识也尤为重要,只有将金融学知识应用于词典的构建过程中才能使词典更适用于金融文本的分析。针对英文文本大数据,国外已经形成了多部具有影响力的词典,例如 Henry 词典(Henry Word Lists,2008)、LM 词典(Loughran and McDonald Word Lists,2011)、哈佛大学通用调查词典(Harvard General Inquirer Word Lists, GI)、文辞乐观与悲观词典(Diction Optimism and Pessimism Word Lists)等。结合以上词典,已有文献能够对媒体报道情绪(Tetlock,2007;Tetlock et al,2008;Solomon et al,2014)、电话会议文本语气语调(Price et al,2012)以及财务报告的语气语调(Feldman et al,2010;Merkley,2014)进行分析。针对中文文本大数据,大多数学者在参考英文词典及其他词库的基础上构建自己的词典展开研究(汪昌云和武佳薇,2015;曾庆生等,2018)。另外,姚加权等(2019)通过词典重组和深度学习算法构建了针对金融领域正式文本和非正式文本的中文情绪词典。该文指出基于该词典构建的上市公司年报语调指标和社交媒体情绪指标能够有效地预测上市公司股票的收益率、成交量等市场因素以及股价崩盘风险。

Jegadeesh & Wu(2013)指出,在词典法中选择合适的加权方法至少跟选择准确的词典一样重

<sup>①</sup>为了避免该局限性,学者们研究了一系列方法,例如复杂度得分(perplexity score)(Blei et al,2003)和层次狄利克雷过程(hierarchical Dirichlet process)(Teh et al,2006)。

要。在金融学领域中,多数学者采用了简单比例加总权重法衡量文本情绪,具体公式为:

$$Pos = \text{积极、正面词的个数} / \text{文本总词数}, \quad (3)$$

$$Neg = \text{消极、负面词的个数} / \text{文本总词数} \quad (4)$$

$$Tone = \frac{(Pos - Neg)}{(Pos + Neg)} \quad (5)$$

公式(5)中  $Tone$  为管理层净正面语调指标,  $-1 \leq Tone \leq 1$ , 当  $Pos$  大于  $Neg$  时,  $Tone$  越大, 从而说明管理层语调更加正面。此外, 林乐和谢德仁(2016)在稳健性检验中进一步采用了 TF-IDF 方法来衡量词语权重。然而, 词典法的局限性在于, 一方面, 构建针对特定文本的词典时, 需要相关领域的专业知识, 这样就导致构建出来的特定词典无法简单地应用于其他文本。另一方面, 词典法仅关注特定的关键词, 从而会忽略文档的上下文关系。

除了词典法, 学者们还运用有监督机器学习方法对文本情绪进行了分类。其中, 有监督机器学习方法是指将有标签的数据集分为训练集和测试集, 利用训练集来训练模型, 然后将训练的模型应用到测试集中, 使用测试集的预测结果来评估模型。在金融学文本情绪分类中, 学者常用的有监督机器学习方法为朴素贝叶斯和支持向量机。

朴素贝叶斯(naïve Bayesian)是一种基于贝叶斯理论的有监督机器学习算法。首先, 输入训练集学习文档词语归类关系, 得到文档归类的先验概率以及条件概率分布。其次, 根据贝叶斯条件概率公式计算已知文档属于不同文档类别的条件概率。最后, 基于最大后验假设把该文档归为具有最大后验概率的一类。Loughran & McDonald(2016)指出了朴素贝叶斯方法的三个优点: 首先, 该方法是文本分析中最古老且成熟的方法之一; 其次, 该方法基于机器学习来阅读文本, 使得研究大量文本信息成为可能; 最后, 文本度量规则的确定排除了研究人员的主观性。然而, 该方法建立在词语集合中词语属性相互独立的假设之上, 但实际上该假设很难成立。在金融学领域, 朴素贝叶斯方法已经广泛地运用到了文本分析中(Antweiler & Frank, 2004; Das & Chen, 2007; Li, 2010b; Jegadeesh & Wu, 2013)。

支持向量机(support vector machine, SVM)是一种基于统计学习理论和结构风险最小原理的有监督机器学习算法。其基本思想为, 将语料库中的文档通过核函数映射为高维度特征空间中的一个样本点, 然后根据训练集, 在特征空间中找到最优分类的超平面, 使得它能尽可能多地将两类数据点正确分开, 同时使分开的两类数据点距离超平面最远。Guo et al(2016)指出, 与朴素贝叶斯方法相比, 支持向量机能够实现更好的样本外预测精度。朴素贝叶斯方法在学习涵盖了所有文档, 引入了较多噪音, 而该方法只关注不同类型惩罚函数选择的支持向量, 从而能够避免过度拟合的问题。而且, 该方法可利用核函数解决线性不可分问题。该方法的弊端在于: 分类结果对核函数的选择比较敏感, 难以针对具体问题选择出最佳的核函数。Antweiler & Frank(2004)运用该方法对互联网股票留言板中的文本情绪进行了分类。

需要指出的是, 在运用有监督机器学习方法时, 检验和评价模型的分类效果至关重要。其主要方法为交叉验证法, 该方法的基本思想为将原始数据切分为较小子集, 并将其随机重新组合为训练集和测试集, 在此基础上对模型反复进行训练、测试, 并根据结果对模型进行选择。此外, 学者还可以通过一系列指标对模型的分类水平进行评估, 例如正确率、精度、召回率、F1 值等指标。

Henry & Leone(2016)指出, 与词典法相比, 有监督机器学习方法的局限性在于该方法需要人工编码提供训练集, 这比用词典法进行词频测量所需的工作量更大。此外, 在对训练集进行人工编码的过程中, 编码结果容易受到研究人员的主观行为影响, 进而影响机器学习的分类效果, 并且研究的可复制性较差。尽管存在以上缺陷, 但对没有预定词典的文本来说, 有监督机器学习方法仍是一种合理且有效的分类工具, 特别是在分类精确度比可复制性更重要的情况下。

3. 语义关联性。语义关联性就是根据某一类词语去识别文本语义特征的过程。具体而言, 首

先依照某一类关键词构建词表,然后计算词表中词语在文档中的词频,进而识别出文本中与关键词语义相关的语义特征。例如,Loughran et al(2009)根据“corporate responsibility”“social responsibility”“socially responsible”等关键词识别了公司对道德相关术语的使用情况。此外,学者还可以运用词嵌入技术,根据空间中词向量之间的距离(即语法和语义的相似性)来处理词语语义关联性问题。例如,Li et al(2020)运用 Word2Vec 技术扩展了不同类型企业文化的关键词。

4. 文本相似性。目前,较多学者运用了余弦相似度指标来衡量财务报告的相似性(Hoberg & Phillips,2010,2016;王雄元等,2018)和专利文本的相似程度(Kelly et al,2018)。假设文本  $d_1$  和  $d_2$  对应的文本向量分别为  $a=(w_{a1},w_{a2},\cdots,w_{an})$ , $b=(w_{b1},w_{b2},\cdots,w_{bm})$ ,则文本  $d_1$  和  $d_2$  的余弦相似度计算公式如下:

$$\cosine\ similarity(d_1,d_2)=\frac{a\times b}{\|a\|\times\|b\|}=\frac{\sum_1^n w_{ai}\times w_{bi}}{\sqrt{\sum_1^n w_{ai}^2}\times\sqrt{\sum_1^n w_{bi}^2}}\quad (6)$$

其中, $n$  为特征个数, $w_{ai},w_{bi}$  为特征  $i$  在两个文本中的权重。该公式取值在 0 和 1 之间,数值越大表明文档相似度越大。

通过总结金融学文本大数据挖掘步骤和方法可以发现,中英文文本大数据分析的主要区别在于文本的预处理过程。相对于英文文本大数据,中文文本大数据的预处理过程更加复杂。在对中文文本大数据进行文本特征抽取前,研究者需要进行分词处理。而且,中文词语未对词性做出明确规范,主要依靠文本的语法和语义来识别词性。在文本可读性方面,尽管多数学者参考迷雾指数构建中文文本的可读性指标,但是他们运用的常用词或复杂词的词典不同,从而导致构建的可读性指标也不尽相同。

### 三、金融学文本大数据分析

本文接下来将根据不同文本信息来源介绍金融学文本大数据的研究内容,国外研究的文本类型主要包括:上市公司披露的文本信息(如财务报告、电话会议、招股说明书等),财经媒体报道,社交网络文本以及搜索指数等。针对以上文本来源,国外文献运用文本挖掘方法抽取文本特征(例如可读性、语气语调、相似性和语义特征等),并探讨了这些文本特征与企业表现、股票市场行为之间的关系。

#### (一)对上市公司披露文本信息的研究

上市公司披露的文本信息能够反映企业披露行为,也可用于衡量企业财务和经营状况以及向股票市场传达出公司管理层对未来发展的信心程度。其中,针对财务报告、电话会议文本以及招股说明书的研究比较广泛,研究重点在于文本的可读性、语气语调以及相似性。

1. 上市公司披露文本的可读性研究。可读性较高的公司披露文本能够更好地向投资者传达公司信息。Li(2008)运用迷雾指数衡量了财务报告的可读性,发现财务报告可读性较高的企业利润更加持久。Lehavy et al(2011)指出提高企业财务报告的可读性能够降低分析师盈余预测的离散度,提高分析师盈余预测的准确度。Guay et al(2016)发现,财务报告可读性较差的公司会通过自愿披露的方式来减轻可读性的负面影响。Lo et al(2017)发现管理者会策略性地操控财务报告可读性,以误导或影响投资者对企业业绩的评价。Bushee et al(2018)将电话会议文本语言的复杂性区分为信息成分和混淆成分,发现信息成分能够降低信息不对称程度,而混淆成分会加剧信息不对称。

2. 上市公司披露文本的语气语调研究。在语气语调的研究方面,学者们认为公司披露文本的语气语调能够用于预测企业表现和股票市场的变化。Mayew et al(2015)研究发现财务报告中 MD&A 部分的语气具有增量信息,能够预测企业的破产概率和持续经营能力。Davis et al(2015)指出,管理者的乐观情绪会影响电话会议中管理者的语气。Allee & Deangelis(2015)考察了电话会议

中语气分散程度,发现管理者语气分散程度与企业未来业绩和管理者的决策相关,而且语气分散程度还会影响分析师和投资者对信息的感知。Bochkay et al(2020)开发了一个极端语气词典,研究发现,管理者在电话会议中使用极端词汇后,企业的股票交易量会显著增加,股价反应会更加强烈。Jiang et al(2019)基于公司财务报告以及电话会议文本构建了经理情绪指标,指出该经理情绪指标能够有效预测股票收益,且该指标的预测能力超过常用的宏观经济变量及投资者情绪。

3. 上市公司披露文本的相似性研究。相似性也是公司披露文本信息的重要特征,一方面,基于企业之间财务报告内容的相似性可以研究不同企业之间的关系。例如,Hoberg & Phillip(2010)发现两家公司财务报告中产品描述部分越相似,那么两者之间发生并购的可能性越高,且并购绩效越好。Hoberg & Phillips(2016)根据不同公司财务报告产品描述部分的相似性构建了公司特有的产品市场竞争对手集合,并基于此形成了时变的行业划分标准。另一方面,企业之间和企业不同时期财务报告的相似性也为研究企业披露行为模板化提供了契机(Brown & Tucker, 2011; Lang & Stice-Lawrence, 2015)。

4. 上市公司披露文本的语义特征研究。还有一系列文献针对公司披露文本的语义特征展开研究。例如,Buehlmaier & Whited(2018)基于 MD&A 部分运用朴素贝叶斯方法构建了融资约束指标,发现受融资约束的公司的股票收益较高。Bochkay et al(2019)根据电话会议研究了企业 CEO 信息披露的风格,指出 CEO 的前瞻性信息披露行为及乐观态度会在其任期内下降,外部雇用和经验不足的 CEO 更倾向于披露前瞻性信息,而年轻的 CEO 在信息披露方面表现出更大程度的乐观态度。Hanley & Hoberg(2019)将 LDA 模型和 Word2Vec 技术相结合,从银行年报中提取与风险相关的语义主题,并结合投资者的交易模式研究发现,金融行业新显露的风险信号有助于监管金融市场的稳定性。

## (二)对财经媒体报道的研究

财经媒体报道为研究股票市场提供了丰富的文本大数据,常见的英文媒体报道有《华尔街日报》、《纽约时报》以及《金融时报》等。

1. 媒体关注和媒体情绪研究。部分文献根据财经媒体发布新闻的数量及正负面新闻研究媒体报道对股票市场的影响。Hillert et al(2014)基于美国 45 家报纸在 1989—2010 年间发表的约 220 万篇文章进行研究,发现媒体报道会加剧投资者的偏见,收益的可预测性在媒体关注度高的公司中更强。Baloria & Heese(2018)指出受到 FNC(Fox News Channel)媒体倾向性报道威胁的民主党关联企业在大选前会隐瞒负面消息,而在大选后释放负面消息,证明企业会为了名誉资本而避免负面媒体报道。Frank & Sanati(2018)研究发现,股票市场存在对好消息反应过度、对坏消息反应不足的现象,正面新闻冲击后股价会出现反转,负面新闻冲击会引起股价漂移。此外,还有文献研究了媒体情绪对未来房价的预测作用(Soo, 2018)。

2. 经济政策不确定性研究。媒体报道中还包含着经济政策的信息,Baker et al(2016)根据多个主要经济体中具有代表性的媒体报道,运用文本挖掘技术构建了经济政策不确定性指数(economic policy uncertainty, EPU),该指数能够连续且定量地描述经济政策的不确定性。随后,Gulen & Ion(2016)运用该指数研究了经济政策不确定性对企业投资的影响,发现宏观经济政策不确定性会影响微观企业的财务决策,抑制企业投资。Bonaime et al(2018)基于该指数研究了经济政策不确定性对企业并购的影响,发现经济政策不确定性的上升会减少并购交易的价值和数量。

3. 媒体偏向、谣言和假新闻研究。财经媒体报道还存在本地偏向现象。Gurun & Butler(2012)指出媒体在报道本地公司新闻时使用更少的否定词,出现该现象的原因在于本地公司投入了更多广告支出。他们还发现,异常的本地媒体偏向与企业股票估值密切相关。另外,市场上还存在公司谣言和虚假新闻。Ahern & Sosyura(2015)通过对公司并购的谣言进行研究,发现媒体更偏向于发布具有新闻价值公司的谣言,媒体发布的谣言会争夺投资者有限的注意力,从而导致股票价格的过度反应和随后逆转。Kogan et al(2019)发现虚假新闻会提高所涉及股票的交易量和价格波动性,当虚假新闻曝光后,假新闻所在平台上所有新闻对股票交易量和价格波动性的影



响会降低。

### (三)对社交网络文本的研究

随着社交网络的兴起,学者开始将微博和股票论坛等社交网络纳入金融学领域的研究范围内,社交网络中的文本信息与股票市场的关系是该领域的研究重点。

1. 社交网络文本情绪研究。早期,Antweiler & Frank(2004)以雅虎财经网络论坛中的帖子为研究对象,发现以帖子数量衡量的关注度指标能够有效地预测股票收益率和市场波动情况,帖子情绪分歧与同期股票交易量正相关。Chen et al(2014)以股票论坛 Seeking Alpha 上的文章和评论为对象进行研究发现,文章和评论的语气能够预测企业未来的股票收益。Cookson & Niessner(2020)通过对美国股票论坛 StockTwits 上的帖子和用户信息进行研究发现,大约一半的投资者分歧是由不同投资理念造成的,而且投资者分歧可以有效预测股票市场中的异常交易量。

2. 策略性信息披露研究。公司也可以通过社交网络平台策略性地发布信息。例如,Blankespoor et al(2014)发现公司在 Twitter 上发布新闻的链接能够降低公司股票的买卖差价,提高交易深度。研究也表明,社交网络作为信息传播的途径,能够优化投资者获取信息的能力,降低投资者信息搜寻成本。Jung et al(2018)分析了标普 1500 家公司在 Twitter 上的信息发布情况,研究发现,当存在坏消息时公司会减少 Twitter 的发文量,对于投资者成熟程度较低的公司和社交媒体受众面较广的公司而言,这种策略性信息披露行为更明显。

### (四)对搜索指数的研究

在互联网技术不断发展的背景下,网络搜索指数是衡量投资者对股票关注程度的有效指标。Da et al(2011)获取了单个股票每周的谷歌搜索指数,利用股票的搜索频率直接测量了投资者关注度。研究指出,运用搜索指数能够更及时地度量投资者关注度,搜索指数的增加能够预测未来两周股价的上涨以及一年内的股价反转。Chi & Shanthikumar(2017)通过搜索位置检验了投资者对不同地区股票的关注情况,发现投资者存在“本地偏见”,即投资者更倾向于关注本地的股票,这种“本地偏见”还会影响市场对盈余公告的反应。此外,已有文献还研究了特定关键词的搜索指数与股票市场的关系。例如,Da et al(2015)通过确定 118 个词的搜索量与同期市场收益之间的历史关系,选择 30 个最负面的词语作为特定搜索关键词,并使用其构建 FEARS 指数,发现该指数预测了股票市场短期收益逆转、波动性的暂时性增加以及共同基金从股票基金流入债券基金。

### (五)对其他英文文本的研究

还有学者针对其他文本进行了研究。例如,De Franco et al(2015)运用 2002—2009 年间分析师报告研究发现,可读性较高的分析师报告能够提高股票交易量。Hwang & Kim(2017)分析了封闭式投资公司对外披露的年度股东报告的可读性对公司价值的影响,研究发现,在相对不透明的信息环境中,当公司发布可读性较差的年度股东报告时,投资者会产生怀疑等负面情绪,导致公司发生折价交易,从而降低了公司价值。Green et al(2019)运用 Glassdoor 上员工对雇主的评级信息进行研究,发现,员工对雇主的评价与企业的销售增长和盈利能力相关,并能够预测一个季度后的未预期盈余。Huang(2018)发现亚马逊(Amazon.com)上的消费者产品评论可以预测公司股票收益。Chen et al(2019)基于美国专利文本数据研究了金融科技创新对金融行业的价值影响。Ryans(2019)运用朴素贝叶斯文本分类方法,结合公司未来发生财务重述和资产减值的情况,将问询函划分为重述问询函和非重述问询函、减值问询函和非减值问询函,研究了问询函与财务报告质量的相关性。Bandiera et al(2020)以 CEO 日志为研究文本,通过机器学习算法将 CEO 划分为“领导者”和“管理者”,研究了 CEO 行为与企业绩效的关系。

## 四、金融学文本大数据的国内研究进展

国内研究的文本信息来源主要包括:上市公司披露的文本信息(如财务报告、业绩说明会文本、招股说明书)、财经媒体报道、社交网络文本、网络搜索指数以及 P2P 网络借贷文本。国内文献通过

提取各种信息的文本特征对企业表现、股票市场以及网络借贷市场展开研究。与国外相比,随着国内 P2P 网络借贷的兴起,该类文本成为国内文本分析的研究重点之一。

1. 对上市公司披露文本信息的研究。针对上市公司年报,丘心颖等(2016)研究了上市公司年报可读性与分析师信息解读之间的关系,发现年报可读性与分析师预测质量不存在显著关系,说明中国分析师未能有效发挥专业解读信息的作用。曾庆生等(2018)研究了年报语调与年报披露后内部人交易行为之间的关系,发现企业管理者在编制年报时存在“口是心非”的现象,积极的年报语调却伴随着管理者较高的股票卖出规模。王雄元等(2018)基于年报 MD&A 中与风险相关的信息进行研究,前后两年年报中风险段落内容的相似性能够降低审计费用。针对公司披露的其他文本,林乐和谢德仁(2016)指出投资者能够识别上市公司年度业绩说明会中的管理层语调。Yan et al (2019)发现 IPO 招股说明书中不确定性或负面语调与 IPO 初期股票收益及后续股票收益的波动显著相关,而且降低了股票的长期回报率。

2. 对财经媒体报道的研究。证监会规定上市公司必须在《上海证券报》《中国证券报》《证券时报》《金融时报》《经济日报》《中国改革报》《中国日报》《证券市场周刊》“七报一刊”中公布企业重大信息。另外,我国还拥有百度新闻、新浪财经、和讯网等网络新闻媒体,这均为研究中国股票市场提供了丰富的媒体报道信息。游家兴和吴静(2012)以财经报纸为研究文本进行研究,发现媒体情绪越高涨或越低落时,资产定价偏误的现象越严重。汪昌云和武佳薇(2015)以财经媒体报道为研究对象进行研究,发现媒体负面语气的下降会提高 IPO 抑价率、IPO 超募资金比例和承销商费用占比。王靖一和黄益平(2018)还研究了网络媒体情绪对网贷市场的影响。此外,媒体报道可以划分为市场导向媒体报道和政策导向媒体报道。Piotroski et al(2017)指出中国媒体的集团化改革使得政策导向的媒体报道更加集中于政治目标,市场导向的媒体报道更加集中于商业目标。You et al(2018)从信息监督的角度研究发现,与政策导向媒体报道相比,市场导向的媒体报道能够提供更多关于企业方面的信息,且只有市场导向的媒体报道能够对公司治理产生重大影响。此外,在衡量宏观经济政策不确定性方面,Huang & Luk(2020)利用中国大陆多家报纸构建了新的并且频率更高的中国 EPU 指数,研究发现,新的中国 EPU 指数能够预测中国的股票价格、就业和产出情况。

3. 对社交网络文本的研究。何贤杰等(2016)基于新浪微博博文研究指出,公司治理水平越高,公司越倾向于开通微博并且发布更多的公司信息。东方财富股吧及雪球网等股票论坛为研究我国股票市场中投资者关注和投资者情绪提供了机会。Huang et al(2016)基于东方财富股吧发帖信息发现,中国投资者也存在“本地偏见”现象,这种偏见在欠发达地区、大型公司、非沪深 300 指数、低成交量且名称表明公司所在地的股票中尤其明显。孙书娜和孙谦(2018)研究发现,根据“雪球网”用户自选股信息构建的投资者关注在短期内会提高股票价格和股票交易量,但该影响会随着时间的推移而逐渐衰减。另外,Jiang, Liu & Yang(2019)表明股票论坛中投资者之间的交流也会对股票收益产生影响。

4. 对搜索指数的研究。在对搜索指数的研究方面,部分学者基于网络搜索指数构建了投资者个股关注指标,研究了投资者关注与资产定价的关系。例如,俞庆进和张兵(2012)以百度指数作为投资者关注度指标进行研究,发现投资者关注对股票当期收益有正向价格压力,但这种压力会在短期内实现反转。还有学者运用特定关键词的搜索指数展开研究。例如,曾建光(2015)根据“余额宝被盗”的百度搜索指数构建了投资者网络安全风险感知指标,研究发现,投资者对互联网安全风险感知越强烈,要求的风险补偿越高,并且移动互联网投资者风险感知要强于电脑端投资者。

5. 对 P2P 网络借贷文本的研究。在我国金融改革和金融创新的背景下,P2P 网络借贷(peer to peer lending)掀起了新的热潮,一些学者对 P2P 网络借贷成功率的影响因素展开研究。陈霄等(2018)发现可读性较强的借款描述能够向投资者传递积极信息,提高借款的成功率。彭红枫和林川(2018)分析了借款描述中特定词汇比重对借贷成功率的影响,研究发现,积极语气词汇和金融词汇

比重与借款成功率正相关,消极语气词汇比重、强语气词汇比重和弱语气词汇比重与借款成功率负相关。此外,还有不少学者基于 P2P 网络借贷文本对借款利率、筹资效率等方面展开研究。

除了对以上五个文本数据来源进行分析以外,还有学者针对分析师报告、年报问询函、私人会议总结报告等其他中文文本进行了广泛研究。<sup>①</sup>

## 五、金融学文本大数据研究展望

本文总结并介绍了金融学文本大数据挖掘步骤和方法,描述了文本的语料获取、预处理过程,文档表示以及文档特征的抽取。另外,本文根据不同的文本信息来源,梳理了国内外金融学文本大数据的研究内容。从已有文献来看,对文本大数据挖掘与分析正处于蓬勃发展阶段。本文认为未来针对金融学文本大数据挖掘与分析还可以从以下几个方面进一步深入探讨。

1. 丰富研究内容,开拓更多文本数据来源。文本大数据在金融学领域中的研究内容和信息来源还可以进一步细化和丰富。例如在财经媒体报道方面,研究者可以不仅对媒体报道的数量和情绪进行分析,还可以对报道事件类型进行判断,从而识别出有关企业并购、IPO、财务舞弊、高管个人新闻等方面的媒体报道,研究不同事件类型的报道对企业和股票市场的影响。在股票论坛方面,研究者可以根据投资者对企业的关注度来构建企业网络,研究不同企业之间的竞争关系。从数据来源方面,还可以开拓更多的文本大数据。例如微信公众号、政府工作报告、国务院政策文件、法院裁判文书、招聘网站、企业发布的业绩修正公告、社会责任鉴证意见、内部控制评价报告等。

2. 运用新的文本信息提取方法。目前,在金融学领域的文本分析研究中,应用较为广泛的仍是无法反映上下文含义的“词袋”方法。然而,在自然语言处理领域(natural language processing, NLP)仍有许多新的分析方法和工具,它们在金融学领域的文本分析中尚未得到足够的重视,但具有较大的潜力。例如:(1)命名实体识别(named entity recognition,NER)。NER 属于 NLP 领域的重要基础工具,能够识别待处理文本中的命名实体,从而提取出时间、地点、人名、机构、货币、百分比和日期,常用的 NER 工具有 Stanford NER。(2)关系提取(relation extraction)。通常运用有监督的机器学习方法从含有实体对的句子中提取出实体对之间的对应关系,并对它们的共现性进行分析。(3)文本摘要(summarization),即使用计算机算法压缩文本内容的过程,摘要的长度取决于压缩率。Cardinaels et al(2018)研究指出,基于算法形成的摘要比管理层披露摘要的积极程度更低,且算法摘要能够使投资者对企业股价做出较为保守的估计。

3. 将深度学习引入文本领域的学术研究。深度学习方法在 NLP 领域得到迅猛发展。深度学习模型主要包括:卷积神经网络(convolutional neural networks,CNN)模型、循环神经网络(recurrent neural networks,RNN)模型以及其变体长短期记忆网络(long short term memory networks,LSTM)模型,生成对抗网络(GAN),强化学习,以及目前在 NLP 领域流行的 BERT、XLNet 等模型。将深度学习引入文本领域将会展开更丰富的研究内容,同时还会提高文本信息提取的准确性。例如,Cao et al(2018)利用 LSTM 模型的变种在企业披露文本中检查出了不一致错误。

深度学习与传统机器学习方法在特征表示及模型参数数量等方面具有较大差异,如表 1 所示。Heaton et al(2016)指出深度学习方法在金融领域中的研究具有以下优势:第一,模型考虑了与预测问题相关的尽可能多的数据信息;第二,能够捕获输入数据间的非线性关系,提高样本内的拟合程度;第三,能够有效避免浅层结构的过度拟合问题。此外,当训练集数据量提高到一定程度时,深度学习信息提取准确度要明显高于传统机器学习方法。

<sup>①</sup>我们详细梳理了近年来关于金融学文本大数据分析的大约 170 篇国内期刊论文,详见本文的在线附录,有兴趣的读者可以联系作者获取。

表 1 传统机器学习与深度学习方法的比较

| 指标                  | 传统机器学习模型<br>(例如 logistic、SVM 等) | 深度学习模型<br>(例如 CNN、RNN 等) |
|---------------------|---------------------------------|--------------------------|
| 特征表示                | 需要手动提取特征,输入被表示为向量               | 无须手动提取特征,特征被表示自动学习       |
| 模型参数数量              | 与特征个数线性相关                       | 可达到千万、亿个参数规模             |
| 所需训练样本<br>(以保证高准确率) | 即使增大训练样本,性能也可能无法提升              | 巨大(百万、千万个观测值级别)          |
| 可处理任务               | 输出为单个随机变量                       | 输出为随机变量的结构体              |
| 计算代价                | 小                               | 大(往往需要 GPU 进行运算)         |

4. 构建具有针对性的中文情绪词典。词典法以预先设定的词典为基础来计算文本中不同类型词语的词频,并结合合适的加权方法来提取文本信息。然而,在中文文本情绪分类方面,词典法的应用目前仍处于探索阶段。多数学者选择已有的英文情绪词典以及词库作为参照来构建中文文本情绪词典,这就导致了构建的词典缺乏针对中文语境的问题。另外,不同来源及类别的文本信息在语言使用方面具有不同特点。例如,企业年报中专业术语较多,社交网络媒体中俚语和表情符号使用较多(姚加权等,2019)。因此,应当针对不同的中文文本内容构建具有针对性的中文文本情绪词典,并在未来的研究中持续验证且更新词典内容。

5. 改进文本可读性指标。目前在文本可读性的衡量方面,多数学者参考迷雾指数展开分析。但是,语言的语序和逻辑关系是影响可读性的重要因素,如果仅考虑句子的长度和复杂字词的比例而忽略语序和逻辑,仍无法准确衡量读者对文本的理解程度。目前,已有新的指标对文本的可读性进行衡量。例如,StyleWriter 软件包中的 Bog 指数(Bonsall IV et al,2017),捕获了语言学家强调的且又简单的英文文本特征。任宏达和王琨(2018)指出,运用机器学习方法衡量的可读性指标具有全面性和综合性,还可以克服自然语言的障碍。因此可以预见的是,在未来的研究中,更多学者会考虑采用新的指标展开研究,并使用机器学习甚至是深度学习模型,构建更加综合和准确的可读性指标。

另外,基于文本中表格内容的统计信息可能会构造出更有效的可读性指标。在以往的研究中,一般会将文档中的表格删除,只分析文本段落中的内容。然而,表格内容包含的数字化信息往往比文字信息更加客观和容易理解。如前所述,利用更精准的文档结构识别技术识别出文档中的所有表格后,可以计算每页平均的表格数量,文档中数字和文字的相对比例以及数字在表格和文本段落中出现的比例等也可以视为可读性指标。

6. 提高研究的可复制性。文本大数据的非结构化特征使得将其转化为结构化数据的过程比较复杂,其转变方法会影响研究的可复制性,即他人能否按照文章描述的研究思路和方法得出一致的结论。已经有学者在论文中以附录的形式详细介绍研究的文本分析方法与步骤(Huang et al,2018;姚加权等,2019)。因此,在未来的文本分析中,为了提高研究的可复制性,作者应当详细记录文档的预处理过程、文档表示以及特征的抽取方式。无论是使用词典法还是使用较为复杂的机器学习和深度学习方法,研究者都应当详细地揭示影响研究结果的关键词、词典以及具体思路和算法等。

#### 参考文献:

- 陈霄 叶德珠 邓洁,2018:《借款描述的可读性能够提高网络借款成功率吗》,《中国工业经济》第 3 期。
- 何贤杰 等,2016:《上市公司网络新媒体信息披露研究:基于微博的实证分析》,《财经研究》第 3 期。
- 林乐 谢德仁,2016:《投资者会听话听音吗?——基于管理层语调视角的实证研究》,《财经研究》第 7 期。
- 彭红枫 林川,2018:《言之有物:网络借贷中语言有用吗?——来自人人贷借款描述的经验证据》,《金融研究》第 11 期。
- 丘心颖 郑小翠 邓可斌,2016:《分析师能有效发挥专业解读信息的作用吗?——基于汉字年报复杂性指标的研究》,《经济学(季刊)》第 4 期。
- 任宏达 王琨,2018:《社会关系与企业信息披露质量——基于中国上市公司年报的文本分析》,《南开管理评论》第 5 期。

- 沈艳 陈赞 黄卓,2018:《文本大数据分析在经济学和金融学中的应用》,北京大学工作论文。
- 孙书娜 孙谦,2018:《投资者关注和股市表现——基于雪球关注度的研究》,《管理科学学报》第6期。
- 唐国豪 姜富伟 张定胜,2016:《金融市场文本情绪研究进展》,《经济学动态》第11期。
- 汪昌云 武佳薇,2015:《媒体语气、投资者情绪与IPO定价》,《金融研究》第9期。
- 王靖一 黄益平,2018:《金融科技媒体情绪的刻画与对网贷市场的影响》,《经济学(季刊)》第4期。
- 王雄元 高曦 何捷,2018:《年报风险信息披露与审计费用——基于文本余弦相似度视角》,《审计研究》第5期。
- 姚加权等,2019:《语调、情绪及其市场影响——基于金融领域中文情绪词典的研究》,暨南大学工作论文。
- 游家兴 吴静,2012:《沉默的螺旋:媒体情绪与资产误定价》,《经济研究》第7期。
- 俞庆进 张兵,2012:《投资者有限关注与股票收益——以百度指数作为关注度的一项实证研究》,《金融研究》第8期。
- 曾建光,2015:《网络安全风险感知与互联网金融的资产定价》,《经济研究》第7期。
- 曾庆生等,2018:《年报语调与内部人交易:“表里如一”还是“口是心非”?》,《管理世界》第9期。
- 张学勇 吴雨玲,2018:《基于网络大数据挖掘的实证资产定价研究进展》,《经济学动态》第6期。
- Allee, K. D. & M. D. Deangelis(2015), “The structure of voluntary disclosure narratives: Evidence from tone dispersion”, *Journal of Accounting Research* 53(2):241—274.
- Ahern, K. R. & D. Sosyura(2015), “Rumor has it: Sensationalism in financial media”, *Review of Financial Studies* 28(7):2050—2093.
- Antweiler, W. & M. Frank(2004), “Is all that talk just noise? The information content of internet stock message boards”, *Journal of Finance* 59(3):1259—1294.
- Baker, S. R. et al(2016), “Measuring economic policy uncertainty”, *Quarterly Journal of Economics* 131(4):1593—1636.
- Baloria, V. P. & J. Heese(2018), “The effects of media slant on firm behavior”, *Journal of Financial Economics* 129(1):184—202.
- Bandiera, O. et al(2020), “CEO behavior and firm performance”, *Journal of Political Economy*, forthcoming.
- Blankespoor, E. et al(2014), “The role of dissemination in market liquidity: Evidence from firms’ use of twitter”, *Accounting Review* 89(1):79—112.
- Blei, D. M. et al(2003), “Latent Dirichlet allocation”, *Journal of Machine Learning Research* 3: 993—1022.
- Brown, S. V. & J. W. Tucker(2011), “Large-sample evidence on firm’s year-over-year MD&A modifications”, *Journal of Accounting Research* 49(2):309—346.
- Bochkay, K. et al(2020), “Hyperbole or reality? Investor responses to extreme language in earnings conference calls”, *Accounting Review* 95(2):31—60.
- Bochkay, K. et al(2019), “Dynamics of CEO disclosure style”, *Accounting Review* 94(4): 103—140.
- Bonsall IV, S. et al(2017), “A plain English measure of financial reporting readability”, *Journal of Accounting and Economics* 63(2—3):329—357.
- Bonaime, A. et al(2018), “Does policy uncertainty affect mergers and acquisitions?”, *Journal of Financial Economics* 129(3):531—558.
- Buehlmaier, M. M. & T. M. Whited(2018), “Are financial constraints priced? Evidence from textual analysis”, *Review of Financial Studies* 31(7):2693—2728.
- Bushee, B. J., I. D. Gow & D. J. Taylor(2018), “Linguistic complexity in firm disclosures: Obfuscation or information?”, *Journal of Accounting Research* 56(1):85—121.
- Cao, Y. et al(2018), “Towards automatic numerical cross-checking: Extracting formulas from text”, In: *Proceedings of the 2018 World Wide Web Conference*, pp. 1795—1804.
- Cardinaels, E. et al(2018), “Automatic summaries of earnings releases: Attributes and effects on investors’ judgments”, Tilburg University Working Paper, <https://ssrn.com/abstract=2904384>.
- Chen, H. et al(2014), “Wisdom of crowds: The value of stock opinions transmitted through social media”, *Review of Financial Studies* 27(5):1367—1403.
- Chen, M. A., Q. Wu & B. Yang(2019), “How valuable is FinTech innovation?”, *Review of Financial Studies* 35(5):2062—2106.
- Chi, S. & D. M. Shanthikumar(2017), “Local bias in google search and the market response around earnings announcements”, *Accounting Review* 92(4):115—143.

- Cole, C. J. & C. L. Jones(2005), "Management discussion and analysis: A review and implication for future research", *Journal of Accounting Literature* 24:135—174.
- Cookson, J. A. & M. Niessner(2020), "Why don't we agree? Evidence from a social network of investors", *Journal of Finance* 75(1):173—228.
- Cong, L. W., T. Liang & X. Zhang(2019), "Analyzing textual information at scale", Cornell University Working Paper, <https://ssrn.com/abstract=3449822>.
- Da, Z., J. Engelberg & P. Gao(2011), "In search of attention", *Journal of Finance* 66(5): 1461—1499.
- Da, Z., J. Engelberg & P. Gao(2015), "The sum of all fears investor sentiment and asset prices", *Review of Financial Studies* 28(1):1—32.
- Das, S. R. & M. Y. Chen(2007), "Yahoo! for amazon: Sentiment extraction from small talk on the web", *Management Science* 53(9):1375—1388.
- Davis, A. K. et al(2015), "The effect of manager-specific optimism on the tone of earnings conference calls", *Review of Accounting Studies* 20(2):639—673.
- De Franco, G. et al(2015), "Analyst report readability", *Contemporary Accounting Research* 32(1): 76—104.
- Feldman, R. et al(2010), "Management's tone change, post earnings announcement drift and accruals", *Review of Accounting Studies* 15(4):915—953.
- Frank, M. Z. & A. Sanati(2018), "How does the stock market absorb shocks?", *Journal of Financial Economics* 129(1):136—153.
- Garcia, D. (2013), "Sentiment during recessions", *Journal of Finance* 68(3):1267—1300.
- Gentzkow, M. et al(2019), "Text as data", *Journal of Economic Literature* 57(3):535—574.
- Green, T. C. et al(2019), "Crowdsourced employer reviews and stock returns", *Journal of Financial Economics* 134(1):236—251.
- Guay, W. et al(2016), "Guiding through the fog: Financial statement complexity and voluntary disclosure", *Journal of Accounting and Economics* 62(2—3):234—269.
- Gulen, H. & M. Ion(2016), "Policy uncertainty and corporate investment", *Review of Financial Studies* 29(3):523—564.
- Guo, L., F. Shi & J. Tu(2016), "Textual analysis and machine learning: Crack unstructured data in finance and accounting", *Journal of Finance and Data Science* 2(3):153—170.
- Gurun, U. G. & A. W. Butler(2012), "Don't believe the hype: Local media slant, local advertising, and firm value", *Journal of Finance* 67(2):561—598.
- Hanley, K. W. & G. Hoberg(2019), "Dynamic interpretation of emerging risks in the financial sector", *Review of Financial Studies* 32(12):4543—4603.
- Heaton, J. B. et al(2016), "Deep learning in finance: Deep portfolios", *Applied Stochastic Models in Business and Industry* 33(1):3—12.
- Henry, E. & A. J. Leone(2016), "Measuring qualitative information in capital markets research: Comparison of alternative methodologies to measure disclosure tone", *Accounting Review* 91(1): 153—178.
- Hillert, A. et al(2014), "Media makes momentum", *Review of Financial Studies* 27(12): 3467—3501.
- Hoberg, G. & G. Phillips(2010), "Product market synergies and competition in mergers and acquisitions: A text-based analysis", *Review of Financial Studies* 23(10):3773—3811.
- Hoberg, G. & G. Phillips(2016), "Text-based network industries and endogenous product differentiation", *Journal of Political Economy* 124(5):1423—1465.
- Huang, Y., H. Qiu & Z. Wu(2016), "Local bias in investor attention: Evidence from China's internet stock message boards", *Journal of Empirical Finance* 38:338—354.
- Huang, J. (2018), "The customer knows best: The investment value of consumer opinions", *Journal of Financial Economics* 128(1):164—182.
- Huang, A. H. et al(2018), "Analyst information discovery and interpretation roles: A topic modeling approach", *Management Science* 64(4):2833—2855.
- Huang, Y. & P. Luk (2020), "Measuring economic policy uncertainty in China", *China Economic Review* 59, No. 101367.

- Hwang, B. H. & H. H. Kim(2017), "It pays to write well", *Journal of Financial Economics* 124(2): 373—394.
- Jegadeesh, N. & D. Wu(2013), "Word power: A new approach for content analysis", *Journal of Financial Economics* 110(3):712—729.
- Jiang, F. et al(2019), "Manager sentiment and stock returns", *Journal of Financial Economics* 132(1):126—149.
- Jiang, L. , J. Liu & B. Yang(2019), "Communication and comovement: Evidence from online stock forums", *Financial Management* 48(3):805—847.
- Jones, M. J. & P. A. Shoemaker(1994), "Accounting narratives: A review of empirical studies of content and readability", *Journal of Accounting Literature* 13:142—184.
- Jung, M. J. et al(2018), "Do firms strategically disseminate? Evidence from corporate use of social media", *Accounting Review* 93(4):225—252.
- Kelly, B. et al(2018), "Measuring technological innovation over the long run", NBER Working Paper, No. w25266.
- Kogan, S. et al(2019), "Fake news: Evidence from financial markets", MIT Sloan School of Management Working Paper, <https://ssrn.com/abstract=3237763>.
- Lang, W. & L. Stice-Lawrence(2015), "Textual analysis and international financial reporting: Large sample evidence", *Journal of Accounting and Economics* 60(2—3):110—135.
- Lehavy, R. et al(2011), "The effect of annual report readability on analyst following and the properties of their earnings forecasts", *Accounting Review* 86(3):1087—1115.
- Li, F. (2008), "Annual report readability, current earnings, and earnings persistence", *Journal of Accounting and Economics* 45(2—3): 221—247.
- Li, F. (2010a), "Textual analysis of corporate disclosure: A survey of the literature", *Journal of Accounting Literature* 29:143—165.
- Li, F. (2010b), "The information content of forward-looking statement in corporate filings: A naïve Bayesian machines learning approach", *Journal of Accounting Research* 48(5):1049—1102.
- Li, K. et al(2020), "Measuring corporate culture using machine learning", University of British Columbia Working Paper, <https://ssrn.com/abstract=3256608>.
- Lo, K. et al(2017), "Earnings management and annual report readability", *Journal of Accounting and Economics* 63(1):1—25.
- Loughran, T. & B. McDonald(2011), "When is a liability not a liability? Textual analysis, dictionaries and 10—Ks", *Journal of Finance* 66(1):35—65.
- Loughran, T. & B. McDonald(2014), "Measuring readability in financial disclosures", *Journal of Finance* 69(4): 1643—1671.
- Loughran, T. & B. McDonald(2016), "Textual analysis in accounting and finance: A survey", *Journal of Accounting Research* 54(4):1187—1230.
- Loughran, T. , B. McDonald & H. Yun(2009), "A wolf in sheep's clothing: The use of ethics-related terms in 10—K reports", *Journal of Business Ethics* 89(1):39—49.
- Mayew, W. J. et al(2015), "MD&A disclosure and the firm's ability to continue as a going concern", *Accounting Review* 90(4):1621—1651.
- Merkley, K. J. (2014), "Narrative disclosure and earnings performance: Evidence from R&D disclosures", *Accounting Review* 89(2):725—757.
- Mikolov, T. et al(2013), "Distributed representations of words and phrases and their compositionality", In: *Advances in Neural Information Processing Systems* 26 (NIPS 2013), pp. 3111—3119.
- Piotroski, J. D. , T. J. Wong & T. Zhang(2017), "Political bias in corporate news: The role of conglomeration reform in China", *Journal of Law and Economics* 60(1):173—207.
- Price, S. M. et al(2012), "Earnings conference calls and stock returns: The incremental informativeness of textual tone", *Journal of Banking & Finance* 36(4):992—1011.
- Ryans, J. P. (2019), "Textual classification of SEC comment letters", *Review of Accounting Studies*, forthcoming.
- Solomon, D. H. et al(2014), "Winners in the spotlight: Media coverage of fund holdings as a driver of flows", *Journal of Financial Economics* 113(1):53—72.

- Soo, C. K (2018), "Quantifying sentiment with news media across local housing markets", *Review of Financial Studies* 31(10):3689—3719.
- Teh, Y. W. et al(2006), "Sharing clusters among related groups: Hierarchical Dirichlet processes", *Journal of the American Statistical Association* 101(476):1566—1581.
- Tetlock, P. C. (2007), "Giving content to investor sentiment: The role of media in the stock market", *Journal of Finance* 62(3):1139—1168.
- Tetlock, P. C. et al(2008), "More than words: Quantifying language to measure firms' fundamentals", *Journal of Finance* 63(3):1437—1467.
- Teoh, S. H. (2018), "The promise and challenges of news datasets for accounting research", *Accounting, Organizations and Society* 68—69:109—117.
- Yan, Y. et al(2019), "Uncertainty and IPO initial returns: Evidence from the tone analysis of China's IPO prospectuses", *Pacific-Basin Finance Journal* 57, No. 101075.
- You, H. & X. Zhang(2009), "Financial reporting complexity and investor underreaction to 10-K information", *Review of Accounting Studies* 14(4):559—586.
- You, J., B. Zhang & L. Zhang(2018), "Who captures the power of the pen?", *Review of Financial Studies* 31(1): 43—96.

### Text Mining in Financial Big Data and Its Research Progress

YAO Jiaquan<sup>1</sup> ZHANG Kunpeng<sup>1</sup> LUO Ping<sup>2,3</sup>

(1. Jinan University, Guangzhou, China; 2. Chinese Academy of Sciences, Beijing, China;

3. University of Chinese Academy of Sciences, Beijing, China)

**Abstract :** Traditional empirical studies in the field of finance usually rely on structured data such as financial statements and stock market trading data. In the era of big data, data types have enriched with the improvement of computer technology and researchers have begun to introduce textual big data into the field of finance, mainly including the disclosure documents of listed companies, financial media reports, social network texts, internet search index, P2P online lending texts, and have examined the readability, tone, similarity and semantic characteristics of the text. This paper first introduces the steps and methods of textual big data mining in the field of finance, describing the corpus acquisition, preprocessing, document representation and the extraction process of document features. In addition, according to different sources of textual information, this paper introduces the research progress in financial textual big data. Finally, this paper provides a comprehensive research prospect on the research methods and topics of financial textual big data.

**Keywords :** Textual Big Data; Textual Analysis; Machine Learning; Deep Learning; Data Mining

(责任编辑:刘洪愧)

(校对:李仁贵)



# 金融学文本大数据挖掘方法与研究进展——在线附录

姚加权\* 张锟澎† 罗平‡

本部分内容是《经济学动态》2020年第4期中《金融学文本大数据挖掘方法与研究进展》的在线附录，共包括两个部分：“附录 I 金融学文本大数据国内研究”和“附录 II 文档解析”。

“附录 I 金融学文本大数据国内研究”是正文部分“四、金融学文本大数据的国内研究进展”的拓展内容。附录 I 详细梳理了近年来对上市公司披露文本信息、财经媒体报道、社交网络文本、搜索指数、P2P 网络借贷文本等文本大数据进行研究的大约 170 篇国内期刊论文。

“附录 II 文档解析”是正文部分“二、金融学文本大数据挖掘步骤与方法”中“二、文本的预处理”中“1. 文档解析”的拓展内容。附录 II 详细举例说明了选择精准 PDF 解析工具的重要性。

## 附录 I 金融学文本大数据国内研究

国内研究的文本信息来源主要包括：上市公司披露的文本信息（如财务报告、业绩说明会文本、招股说明书）、财经媒体报道、社交网络文本、网络搜索指数以及 P2P 网络借贷文本。国内文献通过提取各种信息的文本特征对企业表现、股票市场以及网络借贷市场展开研究。与国外相比，随着国内 P2P 网络借贷的兴起，该类文本成为国内文本分析的研究重点之一。

### （一）对上市公司披露文本信息的研究

1. 对上市公司年报的研究。上市公司年报作为缓解市场信息不对称的重要文件，为信息披露研究提供了丰富的文本信息。目前，学者主要从年报的可读性，年报披露行为模板化，年报信息操纵以及年报信息披露的市场反应等方面展开研究。

早期，阎达五和孙蔓莉（2002）对我国深市 B 股上市公司的年报可读性进行了分析，他们指出我国 B 股上市公司年报处于较为难读和非常难读的水平之间。随后，丘心颖等（2016）研究了上市公司年报可读性与分析师信息解读之间的关系，发现年报可读性与分析师预测质量不存在显著关系，说明中国分析师未能有效发挥专业解读信息的作用。任宏达和王琨（2018）研究发现依赖社会关系获取资源的企业年报的可读性较低，在股权分置改革后这类企业年报的可读性得到了显著的提升。王克敏等（2018）研究指出管理者会出于自利动机操纵年报信息的复杂性，该操纵行为对数字信息的操纵具有替代作用。

投资者或分析师可以及时跟进年报中的文本信息来判断企业的发展状况，但如果年报中文本信息内容更新较慢且重复率高，存在披露行为模板化，就会严重影响年报的价值含量。对此，沈洪涛和苏亮德（2012）以我国重污染行业上市公司年报中的环境信息披露为研究对象，发现企业环境信息披露存在同形性和模仿行为。蒋艳辉等（2014）以管理层讨论与分析

\*暨南大学管理学院，电子邮箱：jiaquanyao@gmail.com。

†暨南大学管理学院，电子邮箱：zkpworking0729@163.com。

‡中国科学院智能信息处理重点实验室、中国科学院计算所和中国科学院大学，电子邮箱：luop@ict.ac.cn。

(Management Discussion and Analysis, MD&A) 部分为研究文本,发现同一内容信息在不同时期报告中重复出现会显著增加公司股权资本成本。赵子夜等(2019)发现,本公司和其他公司同期的 MD&A 部分文本相似性,以及在高财务风险条件下公司当期与上期 MD&A 部分文本相似性会导致信息披露不足,进而产生负面的经济后果。

公司内部人员很可能会对年报中文本信息的语气语调进行操纵,从而使得披露的信息更有利于公司的利益。曾庆生等(2018)研究了年报语调与年报披露后内部人交易行为之间的关系,发现企业管理者在编制年报时存在“口是心非”的现象,积极的年报语调却伴随着管理者较高的卖出股票规模。王华杰和王克敏(2018)研究发现企业管理层会通过操纵年报文本信息的语调来补充公司对数字信息的操纵行为。

对年报信息含量及市场反应的研究也是年报文本分析的重要内容。针对年报 MD&A 中的风险信息披露,已有研究表明,年报风险信息披露能够提高分析师的预测准确度(王雄元等,2017),降低银行贷款利率(王雄元和曾敬,2019),降低企业权益资本成本(王雄元和高曦,2018)。王雄元等(2018)基于年报 MD&A 中与风险相关的信息进行研究发现,前后两年年报中风险段落内容的相似性能够降低审计费用。其次,针对年报中 MD&A 部分,薛爽等(2010)指出 MD&A 中的信息有利于投资者预测企业未来的盈利能力。孟庆斌等(2017)研究发现 MD&A 的信息含量越高,未来股价崩盘风险越低。陈艺云(2019)研究表明 MD&A 部分传递的管理层语调为预测财务困境提供了新信息,但该信息并没有在市场交易价格中得到充分反映。除年报中的 MD&A 部分外,学者还对年报其他内容进行了分析。王艳艳等(2018)以年报中审计报告的关键审计事项段为研究对象,发现关键审计事项段的可读性、语气语调以及精确度会影响审计报告的沟通价值。吴璇等(2019)以年报附注中披露的经营范围信息为研究对象,发现了公司股价与该公司经营业务相似公司股价的收益联动效应。李哲(2018)以沪深重污染行业上市公司环境信息词频数衡量企业环境信息披露的详细程度,研究了市场监管者和交易者对“多言寡行”环境信息披露模式的行为反应,发现环境信息披露结构失序影响了信息传递效率和资产定价效率,使得市场监管者和交易者摒弃了这种披露模式。

还有文献根据年报文本内容构建了相关指标展开研究。孙蔓莉等(2012)基于年报 MD&A 部分内容,构建了业绩自利性归因指标,研究发现我国企业存在业绩自利性归因行为。孙蔓莉等(2013)进一步指出,该行为会对业绩差的公司的股价起到支撑作用。姜付秀等(2015)基于年报、内控评价报告等文本信息识别了企业的诚信文化,他们指出以诚信作为企业文化的企业盈余管理水平较低,且诚信对正负向盈余管理水平均能发挥一定程度的抑制作用。吴建祖和肖书锋(2016)通过对年报中董事会报告内容进行文本分析,构建了企业高管团队创新注意力转移指标,研究了高管团队创新注意力转移对企业研发投入跳跃<sup>1</sup>的影响。姜付秀等(2017)基于中国上市公司年报采用文本分析方法构建了上市公司的融资约束指标,研究指出多个大股东的公司具有更低的融资约束水平。柴才等(2017)挖掘了年报中关于竞争战略特点的文字描述,对公司竞争战略进行了界定。潘健平等(2019)运用年报中董事会报告和公司官网文本信息构建了企业合作文化指标,研究了企业合作文化与企业创新的关系。任宏达和王琨(2019)基于年报中董事会报告内容,构建了产品市场竞争指标。

<sup>1</sup> 企业研发投入跳跃是研发投入变化的一种极端情况,是指研发投入在一段时间内脱离历史趋势显著、紧凑的最大变化(吴建祖和肖书锋,2016)。

2. 对业绩说明会文本的研究。我国上市公司举行的业绩说明会(Earnings Communication Conference)也提供了丰富的文本信息。在业绩说明会上,公司内部人员对参会者提出的问题进行实时回答,他们回答问题的语气语调会向市场传达出具有价值含量的信息。基于此,我国学者从企业业绩、投资者以及分析师三个角度,对业绩说明会文本中企业管理层语调展开了研究。早期,谢德仁和林乐(2015)研究了业绩说明会文本中管理层语调对企业未来业绩的预测作用,发现企业业绩说明会上管理层语调能够提供预测企业未来业绩的增量信息。从投资者角度,林乐和谢德仁(2016)指出投资者能够识别上市公司年度业绩说明会中的管理层语调。林乐和谢德仁(2017)进一步指出,分析师在一定程度上能够利用业绩说明会上管理层语调中的信息来调整自身荐股行为。

3. 对招股说明书的研究。招股说明书是我国股份公司公开发行股票时就募股事宜发布的公告,主要包括公司状况、经营计划以及股票发行情况等信息。对此,我国学者开始用文本分析方法衡量招股说明书的隐含风险信息,并对其展开研究。郝项超和苏之翔(2014)将招股说明书中的重大风险提示信息分为标准风险提示信息和特有风险提示信息,研究了这两类风险提示信息对公司 IPO 抑价的影响。他们发现主板上市公司中特有风险提示信息可以显著降低 IPO 抑价,而标准风险提示信息对 IPO 抑价无显著影响。随后,姚颐和赵梅(2016)研究了招股说明书中风险信息的披露意愿和市场反应。一方面,未来低业绩的公司更倾向于披露风险信息;另一方面,披露较多的总风险、经营风险和财务风险信息能够显著降低公司 IPO 抑价,提高首日上市的流动性。Yan et al (2019)发现 IPO 招股说明书中不确定性或负面语调与 IPO 初期股票收益及后续股票收益的波动显著相关,而且降低了股票的长期回报率。

4. 对社会责任报告的研究。近年来,投资者、社会公众以及政府对企业社会责任履行情况的关注度不断提高,企业发布的社会责任报告成为企业信息披露中的重要部分。企业社会责任报告描述了劳工的实践、企业慈善以及环境保护等内容,向利益相关者传达了企业可持续发展的重要文本信息。段钊等(2017)针对我国 2009-2015 年上市公司发布的社会责任报告,运用文本挖掘技术提出了一种新的评价企业社会责任信息披露主客观性的方法。他们发现,我国上市公司社会责任报告主观性得分总体为正态分布,呈历年上升趋势,由于行业异质性和政策调控的存在,不同行业和年份之间的得分呈现显著差异。张继勋等(2019)研究了社会责任报告语调与投资者感知社会责任之间的关系。研究发现,在财务信息披露诚信度比较高的情况下,公司采用积极语调披露社会责任信息能够提高投资者感知的社会责任。

5. 对上市公司披露的其他文本信息的研究。除以上文本信息,国内学者还对上市公司披露的其他文本进行了研究。吴冬梅和刘运国(2012)以独立董事辞职公告为研究文本,从辞职人数和原因两个角度研究了独立董事辞职公告与董事会决议公告、股东大会决议公告和年报等其他公告的捆绑披露行为。张宁和刘春林(2012)以上市公司针对市场负面传闻发布的澄清公告为研究文本,检验了内容详细的澄清公告在股价恢复中的作用。刘春林和张宁(2012)发现公司的澄清方式、公司声誉、停牌等因素均会影响澄清公告的澄清效果。贾明等(2014)研究了公司发布澄清公告对正负面传闻的辟谣作用以及传闻来源和澄清公告特征对该辟谣作用的影响。Liu et al (2017)以上市公司访客沟通记录为研究文本,研究了我国上市公司与共同基金经理之间的沟通对共同基金交易行为的影响。

## （二）对财经媒体报道的研究

财经媒体报道通过提供及时且较为容易理解的信息来影响市场投资者（张纯和吴明明，2015）。证监会规定上市公司必须在《上海证券报》、《中国证券报》、《证券时报》、《金融时报》、《经济日报》、《中国改革报》、《中国日报》、《证券市场周刊》“七报一刊”中公布企业重大信息。另外，我国还拥有百度新闻、新浪财经、和讯网等网络新闻媒体，这均为研究中国股票市场提供了丰富的媒体报道信息。

1. 对媒体关注的研究。多数学者运用媒体报道数量来衡量媒体关注度，研究了媒体报道与资产定价（饶育蕾等，2010；黄俊和陈信元，2013；黄俊和郭照蕊，2014；Kim et al, 2019）、公司行为（徐莉萍等，2011；孔东民等，2013；应千伟等，2017）、公司治理（罗进辉，2012；周开国等，2016）、高管薪酬契约有效性（罗进辉，2018）、公司超额现金持有水平（罗进辉等，2018）、分析师盈余预测（周开国等，2014；谭松涛等，2015）、审计师的专业判断（吕敏康和冉明东，2012）、政府审计功能（池国华等，2018）以及政府决算披露质量（张琦和郑瑶，2018）的关系。

2. 对媒体情绪的研究。已有文献将媒体报道区分为正面媒体报道和负面媒体报道，对媒体情绪进行研究。首先，我国学者从不同的角度研究了媒体情绪与资产定价的关系。游家兴和吴静（2012）以财经报纸为研究文本进行研究发现，媒体情绪越高涨或越低落时，资产误定价的现象越严重。在媒体情绪与 IPO 定价方面，游家兴和郑建鑫（2013）发现媒体报道情绪越乐观，新股发行的抑价程度越大。汪昌云和武佳薇（2015）以财经媒体报道为研究对象进行研究发现，媒体负面语气的下降会提高 IPO 抑价率、IPO 超募资金比例和承销商费用占比。邵新建等（2015）指出，IPO 公司会利用媒体公关活动增加媒体正面报道数量，上调证券发行价格，但从长期来看，证券交易价格会逐渐调整至实际价值。

其次，我国学者还研究了媒体情绪对审计意见和网络借贷的影响。在审计意见方面，张龙平和吕敏康（2014）指出，当媒体对上市公司的评价越高时，上市公司被出具标准无保留意见的可能性越高。吕敏康和冯丽丽（2017）发现，媒体对上市公司的较高评价会降低审计质量。在网络借贷方面，王靖一和黄益平（2018）研究了网络媒体情绪对网贷市场的影响。张皓星和黄益平（2018）研究了互联网金融情绪对网络借贷违约率的影响，发现网络借贷存在“反向挤兑”现象，即行业情绪变差时，新增订单和正在还款的订单的违约率均会上升。

再次，媒体情绪与企业表现和企业行为也具有密切联系。我国的媒体报道可以作为公司治理的替代渠道（于忠泊等，2011；Borochin & Cu, 2017），它在约束利益相关者的行为和完善公司治理方面发挥着积极的作用。已有研究表明，媒体的正面报道能够预测创业企业的发展前景（罗炜等，2017），促进污染型企业履行社会责任（李百兴等，2018）。媒体的负面报道能够促进企业下一期的业绩改善（郑志刚等，2011），减少控股股东的掏空行为（Ye et al, 2015），提高企业的环保投资（王云等，2017），促进重污染企业绿色并购行为（潘爱玲等，2019），抑制公司激进性避税行为（刘笑霞和李明辉，2018），抑制企业创新（杨道广等，2017），遏制上市公司财务重述行为的发生（戴亦一等，2011），促使上市公司更换高质量的审计师（戴亦一等，2013），提高审计定价（刘笑霞等，2017），以及在上市公司重大资产重组过程中发挥事中监督和公司治理的作用（陈泽艺等，2017）。另外，杨洁和郭立宏（2017）区分了国有和民营企业，研究发现负面报道后，民营企业比国有企业更偏好进行印象管理。然而，

媒体也可能会和公司发生合谋行为。才国伟等（2015）研究了企业在股权再融资事件中与媒体的合谋行为，发现企业股权再融资期间，媒体正面报道倾向显著增强，进而吸引投资者投资行为。易志高等（2017）基于纸质媒体和互联网媒体发现，高管及家属减持期间，公司会利用媒体炒作来抬高股票价格以实现高管财富转移。金字超等（2018）进一步指出，在公司和媒体关系由亲近到敌意的过程中，资本市场的媒体对公司的行为远非“监督”与“合谋”两种可能，还可能包括媒体利用影响力对公司进行“威胁”的行为。

3. 对市场导向和政策导向媒体报道的研究。结合我国具体情况，媒体报道可以划分为市场导向媒体报道和政策导向媒体报道。李培功和沈艺峰（2010）指出，相对于政策导向媒体报道，市场导向媒体报道具有更加积极的公司治理作用。杨玉龙等（2018）基于股价同步性和知情交易概率两个指标，发现两类媒体对资本市场的影响存在差异，与市场导向媒体报道相比，政策导向媒体报道能够更好地提升资本市场的信息效率，实现信息的挖掘和传递。Piotroski et al（2017）指出中国媒体的集团化改革使得政策导向的媒体报道更加集中于政治目标，市场导向的媒体报道更加集中于商业目标。You et al（2018）从信息监督的角度研究发现，与政策导向媒体报道相比，市场导向的媒体报道能够提供更多关于企业方面的信息，且只有市场导向的媒体报道能够对公司治理产生重大影响。

4. 对媒体报道的其他研究。在对媒体报道的其他研究方面：游家兴等（2018）研究了我国财经媒体报道在不同经济发展水平、制度建设以及社会信任程度方面的有偏性。朱恩伟等（2019）以百度新闻为语料来源，运用新闻文本共现性衡量了银企关系强度，研究了银企关系影响企业信贷融资能力的规律。Huang & Luk（2020）利用中国大陆的多家报纸构建了新的并且频率更高的中国经济政策不确定性（Economic Policy Uncertainty, EPU）指数，研究发现，新的 EPU 指数能够预测中国的股票价格、就业和产出情况。此外，学者还针对媒体中不同的信息展开研究：①股评家文章。宋军和吴冲锋（2003a）研究了股评家对于大盘预测信息的准确性及影响预测的影响因素。宋军和吴冲锋（2003b）研究了股评家在预测大盘时的羊群行为。②媒体荐股栏目。朱宝宪和王怡凯（2001）以《上海证券报》中《为您选股》栏目为研究文本，检验了媒体推荐股票的效果。徐永新和陈婵（2009）以《上海证券报》中《实力机构周末荐股精选》栏目为研究文本，研究了媒体股票推荐行为引起的市场反应。李志生等（2017）基于《中国证券报》、《上海证券报》以及《证券时报》上发布的荐股信息，研究了媒体荐股的市场反应、效率和媒体荐股对财富流动的影响。③媒体中的传闻和谣言。赵静梅等（2010）基于股市谣言发现，股票市场中的传谣会对股价产生异常冲击，传谣和辟谣会在短期显著提升涉谣股票交易的活跃度和流动性。彭情和唐雪松（2019）基于股市传闻发现，股市传闻会显著降低会计盈余的价值相关性，扰乱资本市场的内在逻辑。④公司丑闻。Yu et al（2015）指出，企业丑闻存在行业内溢出效应，治理水平良好的公司可以减少行业丑闻的影响。Hung et al（2015）将公司丑闻划分为政治丑闻、市场丑闻和混合丑闻，研究了不同类型丑闻对股票市场的影响。

### （三）对社交网络文本的研究

博客、股票论坛等社交网络为市场参与者提供了一个重要的信息沟通和交流的场所，帮助投资者之间实现“零距离”接触，进而能够有效提高投资者的信息获取和解读能力，优化市场效率，降低股价崩盘风险，缓解市场信息不对称（丁慧等，2018a，2018b）。我国学者

提取了社交网络中的文本信息，展开了社交网络与股票市场关系的研究。

1. 对微博文本信息的研究。从微博的开通意愿和信息披露内容方面，何贤杰等（2016）基于新浪微博博文研究指出，公司治理水平越高，公司越倾向于开通微博并且发布更多的公司信息。从市场反应方面，徐巍和陈冬华（2016）指出，当微博对已公告信息进行传播时，会带来强烈的市场反应。已有文献表明，微博披露信息能够提高当日公司股票的超额回报和超额交易量（徐巍和陈冬华，2016），帮助分析师理解企业信息，提高分析师的盈余预测精度（胡军和王甄，2015；胡军等，2016）。此外，开通微博的公司的股价同步性较低（胡军和王甄，2015），而且微博信息中经营活动及策略类信息占比越高，公司股价同步性越低（何贤杰等，2018）。也有学者根据微博信息构建了特定指标。例如，刘海飞等（2017）基于微博信息的关注度、信赖度和更新频率三个角度，构建了微博信息质量指标，研究了微博信息质量与股价同步性的内在关联关系。

2. 对股票论坛的研究。在我国互联网中还拥有东方财富股吧、和讯网以及雪球网等股票论坛，学者们采用股票论坛发帖量以及股票关注人数衡量了投资者关注度，进而对投资者关注与股票市场的关系展开研究。Huang et al（2016）基于东方财富股吧发帖信息发现，中国投资者存在“本地偏见”现象，这种偏见在欠发达地区、大型公司、非沪深 300 指数、低成交量且名称表明公司所在地的股票中尤其明显。金字超等（2017）运用和讯网个股关注度研究了资本市场中注意力的配置规律，发现公司在发布财务公告时能够获得更多的市场关注，而其他公司同时发布公告会分散市场注意力。孙书娜和孙谦（2018）研究发现，根据“雪球网”用户自选股信息构建的投资者关注在短期内会提高股票价格和股票交易量，但该影响会随着时间的推移而逐渐衰减。

我国学者还采用文本分析方法从股票论坛中提炼了投资者情绪指标和投资者情绪分歧指标。杨晓兰等（2016）指出，股票论坛中投资者的本地关注与股票收益率之间的关系受到投资者情绪的影响。部慧等（2018）发现投资者情绪对股票收益率和交易量存在显著的当期影响。尹海员和吴兴颖（2019）构建了日内高频投资者情绪指标，研究发现日内投资者情绪能够正向预测股票市场运行，这种作用在交易日下午表现更加明显。此外，段江娇等（2017）基于东方财富股吧中个股帖子，构建了投资者情绪分歧指标，研究指出投资者当日情绪分歧越大，股票市场未来两天的交易量也越大。另外，Jiang, Liu & Yang（2019）表明股票论坛中投资者之间的交流也会对股票收益产生影响。

#### （四）对搜索指数的研究

1. 基于搜索指数的投资者关注度研究。通过网络搜索指数构建投资者个股关注指标，进而研究投资者关注对资产定价的影响已然成为研究热点。张永杰等（2011）研究发现互联网中的开源信息蕴含着成交量中未包含的有效信息内容，能够解释个股的异常日收益率。宋双杰等（2011）利用谷歌趋势提供的搜索量数据，研究了投资者关注对 IPO 异象的影响，结果指出 IPO 前投资者关注对于市场热销程度、首日超额收益和长期表现有更好的解释力和预测力。俞庆进和张兵（2012）以百度指数作为投资者关注度指标进行研究发现，投资者关注对股票当期收益有正向价格压力，但这种压力会在短期内实现反转。赵龙凯等（2013）研究了投资者关注度与同时期股票表现的关系，指出投资者关注度与同时期股票收益率存在正相关关系。冯旭南（2014）通过对“业绩预告”事件前后百度超额搜索量研究发现，中国

投资者具有信息获取能力，能够在一定程度上提前对“业绩预告”的信息作出反应。张谊浩等（2014）发现尽管股市在某种程度上可以影响网络搜索行为，但是网络搜索能更大程度上预测股票市场表现。

2. 基于关键词搜索的指标研究。学者对搜索指数的研究并不局限于构建投资者关注指标。例如，曾建光（2015）根据“余额宝被盗”的百度搜索指数构建了投资者网络安全风险感知指标，研究发现，投资者对互联网安全风险感知越强烈，要求的风险补偿越高，并且移动互联网投资者风险感知要强于电脑端投资者。沈悦和郭品（2015）以支付结算、资源配置、风险管理和网络渠道四个维度中关键词的百度搜索指数为基础构建了互联网金融指数，研究表明互联网金融的发展可以提升商业银行的全要素生产率，并且该影响在不同类型的商业银行之间存在差异。李春涛等（2020）通过金融科技关键词的百度新闻高级检索，构建了地区金融科技发展水平指标，发现地区金融科技发展水平显著提高了当地的企业创新。

### （五）对 P2P 网络借贷文本的研究

在我国金融改革和金融创新的背景下，P2P 网络借贷（Peer to Peer Lending）掀起了新的热潮。P2P 网络借贷是指借贷双方不借助金融机构，而是直接在网络借贷平台上进行交易的无抵押贷款方式。P2P 网络借贷属于新生事物，存在着信息不对称和地域歧视问题（廖理等，2014a），因此，对 P2P 网络借贷进行研究具有现实意义。

1. 对 P2P 网络借贷成功率影响因素的研究。多数学者对 P2P 网络借贷成功率的影响因素展开了广泛的研究。王博等（2017）指出，在借贷市场中信息披露有助于抑制信息扭曲，提高信息透明度。已有文献表明，借款人提供更多的描述性信息（李焰等，2014），而且语言长度越长时（廖理等，2015a），其借款成功率越高。在借款者个人信息方面，借款人所在行业和头衔（张海洋和蔡航，2018；胡金焱和李建文，2019）以及个人信用评级认证、社会资本（李悦雷等，2013；王会娟和廖理，2014）均会对借款成功率产生显著影响。从借款描述的可读性方面，陈霄等（2018）发现可读性较强的借款描述能够向投资者传递积极信息，提高借款的成功率。叶德珠和陈霄（2017）进一步指出，借款描述中的标点和字数能够产生增量信息，对借贷成功率产生影响。从借款描述的语言风格方面，彭红枫和林川（2018）分析了借款描述中特定词汇比重对借贷成功率的影响，研究发现，积极语气词汇和金融词汇比重与借款成功率正相关，消极语气词汇比重、强语气词汇比重和弱语气词汇比重与借款成功率负相关。在交易方式方面，丁杰等（2018）研究了人人贷中的双向交易者交易行为，研究发现双向交易者具有双重信息价值，他们招标和参与投标的借款项目成功率较高。

2. 对 P2P 网络借贷其他因素的研究。学者们还利用 P2P 网络借贷文本对借款违约率、借贷利率、投资者行为、筹资效率以及问题平台识别展开研究。在借款违约率的研究方面，廖理等（2015b）基于人人贷文本数据指出，高等教育年限提高了借款人的自我约束能力，高学历借款者如约还款的概率更高。胡金焱等（2018）发现由于提供的可认证信息不足，农民在获得借款后发生违约的可能性较高。陈林等（2019）进一步指出，如果借款描述字数过多并且存在重复语句，则违约风险较大。在借款利率的研究方面，已有文献指出在我国非完全市场化和信息不对称的背景下所形成的利率有效性较低，仍有相当高比例的违约风险未包含在利率之中（廖理等，2014b），并且借款利率也没有起到明显的违约风险信号的作用（张晓玫等，2016）。吴雨等（2018）还研究了房价与 P2P 网络借款利率的关系，研究表明房价

上涨显著提高了借款利率。在投资者行为的研究方面,已有研究指出 P2P 网络借贷中的投资者具有一定的风险意识(邓东升和陈钊,2019),他们会表现出对收益的追求及风险的规避(胡金焱和宋唯实,2017)。另外,已有研究也发现 P2P 市场中投资者存在羊群行为,投资者更倾向于对订单完成进度较高的订单投资(廖理等,2015c)。廖理等(2018a)发现了投资者的群体智慧,即受投资者追捧的借款项目违约率更低。廖理等(2018b)发现了投资者的学习效应,即有借款者经历投资者的投资业绩会显著提升。在筹资效率的研究方面,周雄伟等(2017)基于拍拍贷的“拍活宝”的借贷文本信息指出,一定程度的平台参与投资能够促进借款人筹资效率的提高,但是过度的平台参与会降低留给普通投资者的投资份额,抑制普通投资者的积极性,从而降低平台的筹资效率。最后,在问题平台的识别方面,基于我国“网贷之家”和“网贷天眼”等 P2P 平台,研究表明利率奇高是问题平台的首要标志(叶青等,2016;向虹宇等,2019)。此外,平台公司管理层的信息披露程度、公司治理水平,是否有第三方资金存管等因素也能反映出 P2P 平台的风险(王修华等,2016;何光辉等,2017;李苍舒和沈艳,2018)。

#### (六) 对其他中文文本的研究

除了对以上五个文本数据来源进行分析以外,学者还针对其他中文文本进行了研究,包括分析师报告、年报问询函、私人会议总结报告、管理层业绩预告、法院判决书以及网络众筹项目文本。

①分析师报告。马黎珺等(2019)发现分析师报告中前瞻性语句向投资者传递了增量信息,并成功预测了企业未来基本面的变化。伊志宏等(2019)指出分析师报告中公司特质性信息越多,股价同步性越低。

②年报问询函。陈运森等(2018a)发现证券交易所发放的年报问询函具有信息含量。已有文献指出,年报问询函能够提高审计质量(陈运森等,2018b),降低公司股价崩盘风险(张俊生等,2018),识别并抑制企业盈余管理行为(陈运森等,2019;刘柏和卢家锐,2019),缓解企业并购重组过程中的信息不对称(李晓溪等,2019a)。

③私人会议总结报告。Bowen et al(2018)收集了我国 2012-2014 年深交所上市公司内部人员与投资者、分析师进行私人会议的总结报告,研究发现私人会议中的信息会影响参会者的股票买卖行为和交易时间,并且报告的语调与市场随后盈余公告的反应以及股票未来表现密切相关。

④管理层业绩预告。李晓溪等(2019b)研究了年报问询函对管理层业绩预告的影响,发现收到年报问询函的公司会提高业绩预告的积极性和预测精确度,降低业绩预告文本内容的复杂性,提高业绩预告文本内容的详细程度,从而减轻问询函的负面影响。

⑤法院判决书。丁骋骋和邱瑾(2016)收集了法院判决书以及各大新闻网站中非法集资的报道,考察了非法集资主角的性别差异和其他微观个体特征。梁平汉和江鸿泽(2020)基于裁判文书网中的传销刑事判决书研究发现,网络传销显著降低了当地社会的信任程度,金融可得性的提高能够有效抑制网络传销的发展。

⑥网络众筹项目文本。彭红枫和米雁翔(2017)以京东东家私募股权融资平台上融资成功的创投板项目为研究对象,研究了股权众筹质量的有效信号以及项目不确定性对信号质量的调节作用。顾乃康和赵坤霞(2019)以京东众筹、淘宝众筹和众筹网的信息为研究对象,探索了产品众筹的动态过程。



## 附录 II 文档解析

在计算机领域，精确地识别文档中的表格、段落等内容块（Physical Objects）以及内容块之间的逻辑关联（如层次目录结构）等研究工作，统称为“文档结构识别”（Mao et al, 2003）。

现在，大量的 PDF 转 Word 的工具或者 PDF 的解析工具效果都不佳。常见的问题主要包括：一方面，PDF 文档中呈现出的一个文本段落，转化后被换行符切割为多个部分。例如，附录图 1 为某上市公司年报的一部分，采用粗糙的 PDF 解析工具解析后，段落格式产生错乱，如附录图 2 所示。另一方面，PDF 文档中的表格，特别是无线框的表格，转化后出现严重格式混乱。例如，附录图 3 为某上市公司年报中的无线框表格，采用粗糙 PDF 解析工具解析后，表格出现错位问题，如附录图 4 所示。这些问题本质上就是 PDF 文档内部的存储内容的限制。然而，精准的 PDF 解析工具能够识别完整的段落，并准确的记录表格信息，例如附录图 5 和附录图 6 分别是采用精准 PDF 解析工具对附录图 1 和附录图 3 的解析结果。因此，我们在这里强调，PDF 的生成并不是一个可逆的过程，需要精准的 PDF 解析工具才能识别完整的文本内容。

### 41. FINANCIAL RISK MANAGEMENT OBJECTIVES AND POLICIES

(continued)

The Group manages its capital structure and makes adjustments to it in light of changes in economic conditions and the risk characteristics of the underlying assets. To maintain or adjust the capital structure, the Group may adjust the dividend payment to shareholders, return capital to shareholders or issue new shares. The Group is not subject to any externally imposed capital requirements. No changes were made in the objectives, policies or processes for managing capital during the years ended 31 December 2017 and 31 December 2016.

### 41. 金融風險管理目標及政策(續)

本集團因應經濟狀況變化及相關資產的風險特徵管理其資本結構及作出調整。為保持或調整資本結構，本集團可能調整向股東派發之股息、向股東退還資本或發行新股。本集團無須符合任何外部施加的資本要求。於截至2017年12月31日及截至2016年12月31日止年度內，管理資本的目標、政策或程序並無變動。

附录图 1 某上市公司年报的一部分

### 41. FINANCIAL RISK MANAGEMENT 41. 金融风险管理目标及政策(续) OBJECTIVES AND POLICIES (continued)

The Group manages its capital structure and makes adjustments to it in light of changes in economic conditions and the risk characteristics of the underlying assets. To maintain or adjust the capital structure, the Group may adjust the dividend payment to shareholders, return capital to shareholders or issue new shares. The Group is not subject to any externally imposed capital requirements. No changes were made in the objectives, policies or processes for managing capital during the years ended 31 December 2017 and 31 December 2016.

附录图 2 粗糙 PDF 解析工具的解析结果——段落

|                         | 2017 年 12 月 31 日 |                | 2016 年 12 月 31 日 |                |
|-------------------------|------------------|----------------|------------------|----------------|
|                         | 账面价值             | 公允价值           | 账面价值             | 公允价值           |
| 金融资产                    |                  |                |                  |                |
| 持有至到期投资 <sup>(i)</sup>  | 226,095          | 221,529        | 208,431          | 213,596        |
| 应收款项类投资 <sup>(ii)</sup> | <u>391,399</u>   | <u>390,484</u> | <u>250,141</u>   | <u>249,868</u> |
| 金融负债                    |                  |                |                  |                |
| 应付债券 <sup>(iii)</sup>   | (398,340)        | (392,321)      | (301,765)        | (299,957)      |
| 本行                      |                  |                |                  |                |

附录图 3 某公司年报中无线框表格

|             | 2017 年 12 月 31 日 |           | 2016 年 12 月 31 日 |           |
|-------------|------------------|-----------|------------------|-----------|
|             | 账面价值             | 公允价值      | 账面价值             | 公允价值      |
| 金融资产        |                  |           |                  |           |
| 持有至到期投资(i)  | 226,095          | 221,529   | 208,431          | 213,596   |
| 应收款项类投资(ii) | 391,399          | 390,484   | 250,141          | 249,868   |
| 金融负债        |                  |           |                  |           |
| 应付债券(iii)   | (398,340)        | (392,321) | (301,765)        | (299,957) |
| 本行          |                  |           |                  |           |

附录图 4 粗糙 PDF 解析工具解析结果——无线框表格

"text": "The Group manages its capital structure and makes adjustments to it in light of changes in economic conditions and the risk characteristics of the underlying assets. To maintain or adjust the capital structure, the Group may adjust the dividend payment to shareholders, return capital to shareholders or issue new shares. The Group is not subject to any externally imposed capital requirements. No changes were made in the objectives, policies or processes for managing capital during the years ended 31 December 2017 and 31 December 2016.",

"text": "本集團因應經濟狀況變化及相關資產的風險特徵管理其資本結構及作出調整。為保持或調整資本結構，本集團可能調整向股東派發之股息、向股東退還資本或發行新股。本集團無須符合任何外部施加的資本要求。於截至2017年12月31日及截至2016年12月31日止年度內，管理資本的目標、政策或程序並無變動。",

附录图 5 精准 PDF 解析工具解析结果——段落

```

"0_1": { "align": "left", "value": "2017 年 12 月 31 日" },
"0_3": { "align": "center", "value": "2016 年 12 月 31 日" },
"1_0": { "align": "unknown", "value": "" },
"1_1": { "align": "center", "value": "账面价值" },
"1_2": { "align": "center", "value": "公允价值" },
"1_3": { "align": "center", "value": "账面价值" },
"1_4": { "align": "center", "value": "公允价值" },
"2_0": { "align": "left", "value": "金融资产" },
"2_1": { "align": "unknown", "value": "" },
"2_2": { "align": "unknown", "value": "" },
"2_3": { "align": "unknown", "value": "" },
"2_4": { "align": "unknown", "value": "" },
"3_0": { "align": "left", "value": "持有至到期投资(i)" },
"3_1": { "align": "center", "value": "226,095" },
"3_2": { "align": "center", "value": "221,529" },
"3_3": { "align": "center", "value": "208,431" },
"3_4": { "align": "center", "value": "213,596" },
"4_0": { "align": "left", "value": "应收款项类投资(ii)" },
"4_1": { "align": "center", "value": "391,399" },
"4_2": { "align": "center", "value": "390,484" },
"4_3": { "align": "center", "value": "250,141" },
"4_4": { "align": "center", "value": "249,868" },
"5_0": { "align": "left", "value": "金融负债" },
"5_1": { "align": "unknown", "value": "" },
"5_2": { "align": "unknown", "value": "" },
"5_3": { "align": "unknown", "value": "" },
"5_4": { "align": "unknown", "value": "" },
"6_0": { "align": "left", "value": "应付债券(iii)" },
"6_1": { "align": "center", "value": "(398,340)" },
"6_2": { "align": "center", "value": "(392,321)" },
"6_3": { "align": "center", "value": "(301,765)" },
"6_4": { "align": "center", "value": "(299,957)" }

```

附录图 6 精准 PDF 解析工具解析结果——无线框表格

随着近几年深度学习的广泛应用，文档结构识别的相关研究工作取得了飞速的进展。基本的方法是，以深度学习中的目标检测（Object Detection）和页面分割（Image or Page Segmentation）模型为基础，在大量带标注的页面数据支撑下，实现对页面中的内容块的自动检测，并进一步判断内容块之间的从属逻辑关系（He et al, 2017; Viana & Oliveria, 2017; Yang et al, 2017）。在公开的免费工具方面，PDFlux<sup>1</sup>提供客户端和网页端的工具，识别输入的 PDF 文档中的表格和文本段落。在该网站上，PDFlux 对 PDF 文档中的有线框和无线框表格的识别准确率分别为 99.9% 和 97%，已达到工业应用的精准级别。与同类的软件比较，如 PDFTables<sup>2</sup>，我们使用了中国上市公司年报 PDF 为测试样本，PDFlux 不仅能识别表格，还可以识别文本段落，且识别的精准度（如表格外线、表格内部结构）有大幅度的优势。但是，PDFTables 的识别速度比 PDFlux 要快一个数量级。同时，这两种免费工具并不支持对 PDF 文件的批量处理。

如前所述，粗糙的 PDF 解析工具并不能正确识别文本段落和表格，因此计算出的句子平均长度可能有很大的偏差。经过初步实验，我们对比了粗糙和精准的 PDF 解析工具对计算文档中句子平均长度的影响：我们以 100 份随机抽样的中国 PDF 格式的年报为例，发现

<sup>1</sup> 详见：<http://www.pdfbox.com/>。

<sup>2</sup> 详见：<https://pdftables.com/>。

计算出的两组句子平均长度并不存在统计相关性。因此，我们在此强调：针对金融学领域的文档，需要慎重选择文档结构的解析工具。

## 参考文献（在线附录部分）

- (1) 部慧 解峥 李佳鸿 吴俊杰, 2018:《基于股评的投资者情绪对股票市场的影响》,《管理科学学报》第4期。
- (2) 才国伟 邵志浩 徐信忠, 2015:《企业和媒体存在合谋行为吗?——来自中国上市公司媒体报道的间接证据》,《管理世界》第7期。
- (3) 柴才 世忠 钦华, 2017:《管薪酬激励与公司业绩——基于三种薪酬激励视角下的经验研究》,《会计研究》第6期。
- (4) 陈林 谢彦妮 李平 李强, 2019:《借款陈述文字中的违约信号——基于P2P网络借贷的实证研究》,《中国管理科学》第4期。
- (5) 陈霄 叶德珠 邓洁, 2018:《借款描述的可读性能够提高网络借款成功率吗》,《中国工业经济》第3期。
- (6) 陈艺云, 2019:《基于信息披露文本的上市公司财务困境预测:以中文年报管理层讨论与分析为样本的研究》,《中国管理科学》第7期。
- (7) 陈运森 邓祎璐 李哲, 2018a:《非处罚性监管具有信息含量吗?——基于问询函的证据》,《金融研究》第4期。
- (8) 陈运森 邓祎璐 李哲, 2018b:《非行政处罚性监管能改进审计质量吗?——基于财务报告问询函的证据》,《审计研究》第5期。
- (9) 陈运森 邓祎璐 李哲, 2019:《证券交易所一线监管的有效性研究:基于财务报告问询函的证据》,《管理世界》第3期。
- (10) 陈泽艺 李常青 魏志华, 2017:《媒体负面报道影响并购成败吗——来自上市公司重大资产重组的经验证据》,《南开管理评论》第1期。
- (11) 池国华 杨金 谷峰, 2018:《媒体关注是否提升了政府审计功能?——基于中国省级面板数据的实证研究》,《会计研究》第1期。
- (12) 戴亦一 潘越 陈芬, 2013:《媒体监督、政府质量与审计师变更》,《会计研究》第10期。
- (13) 戴亦一 潘越 刘思超, 2011:《媒体监督、政府干预与公司治理:来自中国上市公司财务重述视角的证据》,《世界经济》第11期。
- (14) 邓东升 陈钊, 2019:《互联网金融风险与投资者风险意识——来自网贷平台交易数据的证据》,《财贸经济》第2期。
- (15) 丁骋骋 邱瑾, 2016:《性别与信用:非法集资主角的微观个体特征——基于网络数据挖掘的分析》,《财贸经济》第3期。
- (16) 丁慧 吕长江 陈运佳, 2018a:《投资者信息能力:意见分歧与股价崩盘风险——来自社交媒体“上证e互动”的证据》,《管理世界》第9期。
- (17) 丁慧 吕长江 黄海杰, 2018b:《社交媒体、投资者信息获取和解读能力与盈余预期——来自“上证e互动”平台的证据》,《经济研究》第1期。
- (18) 丁杰 李悦雷 曾燕 李仲飞, 2018:《P2P网贷中双向交易者的双重信息价值及信息传递》,《南开管理评论》第2期。
- (19) 段江娇 刘红忠 曾剑平, 2017:《中国股票网络论坛的信息含量分析》,《金融研究》第10期。
- (20) 段钊 何雅娟 钟原, 2017:《企业社会责任信息披露是否客观——基于文本挖掘的我国上市公司实证研究》,《南开管理评论》第4期。
- (21) 冯旭南, 2014:《中国投资者具有信息获取能力吗?——来自“业绩预告”效应的证据》,《经济学(季刊)》第3期。

- (22) 顾乃康 赵坤霞, 2019:《实时的社会信息与互联网产品众筹的动态性——基于大数据的采集与挖掘研究》,《金融研究》第 1 期。
- (23) 郝项超 苏之翔, 2014:《重大风险提示可以降低 IPO 抑价吗? ——基于文本分析法的经验证据》,《财经研究》第 5 期。
- (24) 何光辉 杨咸月 蒲嘉杰, 2017:《中国 P2P 网络借贷平台风险及其决定因素研究》,《数量经济技术经济研究》第 11 期。
- (25) 何贤杰 王孝钰 赵海龙 陈信元, 2016:《上市公司网络新媒体信息披露研究: 基于微博的实证分析》,《财经研究》第 3 期。
- (26) 何贤杰 王孝钰 孙淑伟 朱红军, 2018:《网络新媒体信息披露的经济后果研究——基于股价同步性的视角》,《管理科学学报》第 6 期。
- (27) 胡金焱 宋唯实, 2017:《借贷中投资者的理性意识与权衡行为——基于“人人贷”数据的实证分析》,《金融研究》第 7 期。
- (28) 胡金焱 李建文, 2019:《信贷市场存在行业歧视吗? ——以 P2P 网络借贷为例的研究》,《财贸经济》第 7 期。
- (29) 胡金焱 李建文 张博, 2018:《P2P 网络借贷是否实现了普惠金融目标》,《世界经济》第 11 期。
- (30) 胡军 王甄, 2015:《微博、特质性信息披露与股价同步性》,《金融研究》第 11 期。
- (31) 胡军 王甄 陶莹 邹隽奇, 2016:《微博、信息披露与分析师盈余预测》,《财经研究》第 5 期。
- (32) 黄俊 陈信元, 2013:《媒体报道与 IPO 抑价——来自创业板的经验证据》,《管理科学学报》第 2 期。
- (33) 黄俊 郭照蕊, 2014:《新闻媒体报道与资本市场定价效率——基于股价同步性的分析》,《管理世界》第 5 期。
- (34) 贾明 阮宏飞 张喆, 2014:《上市公司澄清公告的辟谣效果研究》,《管理科学》第 2 期。
- (35) 姜付秀 石贝贝 李行天, 2015:《“诚信”的企业诚信吗? ——基于盈余管理的经验证据》,《会计研究》第 8 期。
- (36) 姜付秀 王运通 田园 吴恺, 2017:《多个大股东与企业融资约束——基于文本分析的经验证据》,《管理世界》第 12 期。
- (37) 蒋艳辉 马超群 熊希希, 2014:《创业板上市公司文本惯性披露、信息相似度与资产定价——基于 Fama-French 改进模型的经验分析》,《中国管理科学》第 8 期。
- (38) 金宇超 靳庆鲁 李晓雪, 2017:《资本市场注意力总量是稀缺资源吗? 》,《金融研究》第 10 期。
- (39) 金宇超 靳庆鲁 严青蕾, 2018:《合谋与胁迫:作为经济主体的媒体行为——基于新闻敲诈曝光的事件研究》,《管理科学学报》第 3 期。
- (40) 孔东民 刘莎莎 应千伟, 2013:《公司行为中的媒体角色: 激浊扬清还是推波助澜? 》,《管理世界》第 7 期。
- (41) 梁平汉 江鸿泽, 2020:《金融可得性与互联网金融风险防范——基于网络传销案件的实证分析》,《中国工业经济》第 4 期。
- (42) 李百兴 王博 卿小权, 2018:《企业社会责任履行、媒体监督与财务绩效研究——基于 A 股重污染行业的经验数据》,《会计研究》第 7 期。
- (43) 李苍舒 沈艳, 2018:《风险传染的信息识别——基于网络借贷市场的实证》,《金融研究》第 11 期。
- (44) 李春涛 闫续文 宋敏 杨威, 2020:《金融科技与企业创新——新三板上市公司的证据》,

《中国工业经济》第 1 期。

(45) 李培功 沈艺峰, 2010:《媒体的公司治理作用:中国的经验证据》,《经济研究》第 4 期。

(46) 李晓溪 杨国超 饶品贵, 2019a:《交易所问询函有监管作用吗?——基于并购重组报告书的文本分析》,《经济研究》第 5 期。

(47) 李晓溪 饶品贵 岳衡, 2019b:《年报问询函与管理层业绩预告》,《管理世界》第 8 期。

(48) 李焰 高弋君 李珍妮 才子豪 王冰婷 杨宇轩, 2014:《借款人描述性信息对投资人决策的影响——基于 P2P 网络借贷平台的分析》,《经济研究》第 S1 期。

(49) 李悦雷 郭阳 张维, 2013:《中国 P2P 小额贷款市场借贷成功率影响因素分析》,《金融研究》第 7 期。

(50) 李哲, 2018:《“多言寡行”的环境披露模式是否会被信息使用者摒弃》,《世界经济》第 12 期。

(51) 李志生 李好 刘淳 张霆, 2017:《天使还是魔鬼?——分析师媒体荐股的市场效应》,《管理科学学报》第 5 期。

(52) 廖理 李梦然 王正位, 2014a:《中国互联网金融的地域歧视研究》,《数量经济技术经济研究》第 5 期。

(53) 廖理 李梦然 王正位, 2014b:《聪明的投资者:非完全市场化利率与风险识别——来自 P2P 网络借贷的证据》,《经济研究》第 7 期。

(54) 廖理 吉霖 张伟强, 2015a:《语言可信吗? 借贷市场上语言的作用——来自 P2P 平台的证据》,《清华大学学报(自然科学版)》第 4 期。

(55) 廖理 吉霖 张伟强, 2015b:《借贷市场能准确识别学历的价值吗?——来自 P2P 平台的经验证据》,《金融研究》第 3 期。

(56) 廖理、李梦然、王正位、贺裴菲:《观察中学习: P2P 网络投资中信息传递与羊群行为》,《清华大学学报(哲学社会科学版)》, 2015c 年第 1 期。

(57) 廖理 向佳 王正位, 2018a:《P2P 借贷投资者的群体智慧》,《中国管理科学》第 10 期。

(58) 廖理 向佳 王正位, 2018b:《网络借贷的角色转换与投资者学习效应》,《中国工业经济》第 9 期。

(59) 林乐 谢德仁, 2016:《投资者会听话听音吗?——基于管理层语调视角的实证研究》,《财经研究》第 7 期。

(60) 林乐 谢德仁, 2017:《分析师荐股更新利用管理层语调吗?——基于业绩说明会的文本分析》,《管理世界》第 11 期。

(61) 刘柏 卢家锐, 2019:《交易所一线监管能甄别资本市场风险吗?——基于年报问询函的证据》,《财经研究》第 7 期。

(62) 刘春林 张宁, 2012:《上市公司传闻的澄清效果研究——来自中国证券市场的证据》,《管理科学学报》第 5 期。

(63) 刘锋 叶强 李一军, 2014:《媒体关注与投资者关注对股票收益的交互作用: 基于中国金融股的实证研究》,《管理科学学报》第 1 期。

(64) 刘海飞 许金涛 柏巍 李心丹, 2017:《社交网络、投资者关注与股价同步性》,《管理科学学报》第 2 期。

(65) 刘笑霞 李明辉 孙蕾, 2017:《媒体负面报道、审计定价与审计延迟》,《会计研究》第 4 期。

(66) 刘笑霞 李明辉, 2018:《媒体负面报道、分析师跟踪与税收激进度》,《会计研究》第

9 期。

(67) 罗进辉, 2012:《媒体报道的公司治理作用——双重代理成本视角》,《金融研究》第 10 期。

(68) 罗进辉, 2018:《媒体报道与高管薪酬契约有效性》,《金融研究》第 3 期。

(69) 罗进辉 李小荣 向元高, 2018:《媒体报道与公司的超额现金持有水平》,《管理科学学报》第 7 期。

(70) 罗炜 何顶 洪莉莎 常国珍, 2017:《媒体报道可以预测创业企业的发展前景吗?》,《金融研究》第 8 期。

(71) 吕敏康 冯丽丽, 2017:《媒体报道、职业能力异质性与审计质量》,《审计研究》第 3 期。

(72) 吕敏康 冉明东, 2012:《媒体报道影响审计师专业判断吗?——基于盈余管理风险判断视角的实证分析》,《审计研究》第 6 期。

(73) 马黎珺 伊志宏 张澈, 2019:《廉价交谈还是言之有据?——分析师报告文本的信息含量研究》,《管理世界》第 7 期。

(74) 潘健平 潘越 马奕涵, 2019:《以“合”为贵? 合作文化与企业创新》,《金融研究》第 1 期。

(75) 潘爱玲 刘昕 邱金龙 申宇, 2019:《媒体压力下的绿色并购能否促使重污染企业实现实质性转型》,《中国工业经济》第 2 期。

(76) 孟庆斌 杨俊华 鲁冰, 2017:《管理层讨论与分析披露的信息含量与股价崩盘风险——基于文本向量化方法的研究》,《中国工业经济》第 12 期。

(77) 彭红枫 米雁翔, 2017:《信息不对称、信号质量与股权众筹融资绩效》,《财贸经济》第 5 期。

(78) 彭红枫 林川, 2018:《言之有物:网络借贷中语言有用吗?——来自人人贷借款描述的经验证据》,《金融研究》第 11 期。

(79) 彭情 唐雪松, 2019:《流言招来的“是非”:股市传闻与盈余价值相关性》,《管理世界》第 3 期。

(80) 丘心颖 郑小翠 邓可斌, 2016:《分析师能有效发挥专业解读信息的作用吗?——基于汉字年报复杂性指标的研究》,《经济学(季刊)》第 4 期。

(81) 饶育蕾 彭叠峰 成大超, 2010:《媒体注意力会引起股票的异常收益吗?——来自中国股票市场的经验证据》,《系统工程理论与实践》第 2 期。

(82) 任宏达 王琨, 2019:《产品市场竞争与信息披露质量——基于上市公司年报文本分析的新证据》,《会计研究》第 3 期。

(83) 任宏达 王琨, 2018:《社会关系与企业信息披露质量——基于中国上市公司年报的文本分析》,《南开管理评论》第 5 期。

(84) 邵新建 何明燕 江萍 薛熠 廖静池, 2015:《媒体公关、投资者情绪与证券发行定价》,《金融研究》第 9 期。

(85) 沈洪涛 苏亮德, 2012:《企业信息披露中的模仿行为研究——基于制度理论的分析》,《南开管理评论》第 3 期。

(86) 沈悦 郭品, 2015:《互联网金融、技术溢出与商业银行全要素生产率》,《金融研究》第 3 期。

(87) 宋军 吴冲锋, 2003a:《中国股评家的羊群行为研究》,《管理科学学报》第 1 期。

(88) 宋军 吴冲锋, 2003b:《中国股评家预测行为的实证研究》,《数理统计与管理》第 3 期。

(89) 宋双杰 曹晖 杨坤, 2011:《投资者关注与 IPO 异象——来自网络搜索量的经验证据》,



《经济研究》第 S1 期。

(90) 孙蔓莉 王竹君 蒋艳霞, 2012:《代理问题、公司治理模式与业绩自利性归因倾向——基于美、中、日三国的数据比较》,《会计研究》第 1 期。

(91) 孙蔓莉 蒋璐 孙健, 2013:《业绩归因的自利性披露及市场反应研究——汇率单边升值情境下的纺织业表现》,《会计研究》第 4 期。

(92) 孙书娜 孙谦, 2018:《投资者关注和股市表现——基于雪球关注度的研究》,《管理科学学报》第 6 期。

(93) 谭松涛 甘顺利 阚铎, 2015:《媒体报道能够降低分析师预测偏差吗》,《金融研究》第 5 期。

(94) 汪昌云 武佳薇, 2015:《媒体语气、投资者情绪与 IPO 定价》,《金融研究》第 9 期。

(95) 王博 张晓玫 卢露, 2017:《网络借贷是实现普惠金融的有效途径吗——来自“人人贷”的微观借贷证据》,《中国工业经济》第 2 期。

(96) 王华杰 王克敏, 2018:《应计操纵与年报文本信息语气操纵研究》,《会计研究》第 4 期。

(97) 王会娟 廖理, 2014:《中国 P2P 网络借贷平台信用认证机制研究——来自“人人贷”的经验证据》,《中国工业经济》第 4 期。

(98) 王靖一 黄益平, 2018:《金融科技媒体情绪的刻画与对网贷市场的影响》,《经济学(季刊)》第 4 期。

(99) 王克敏 王华杰 李栋栋 戴杏云, 2018:《年报文本信息复杂性与管理者自利——来自中国上市公司的证据》,《管理世界》第 12 期。

(100) 王雄元 李岩琼 肖恣, 2017:《年报风险信息披露有助于提高分析师预测准确度吗?》,《会计研究》第 10 期。

(101) 王雄元 高曦 何捷, 2018:《年报风险信息披露与审计费用——基于文本余弦相似度视角》,《审计研究》第 5 期。

(102) 王雄元 高曦, 2018:《年报风险信息披露与权益资本成本》,《金融研究》第 1 期。

(103) 王雄元 曾敬, 2019:《年报风险信息披露与银行贷款利率》,《金融研究》第 1 期。

(104) 王修华 孟路 欧阳辉, 2016:《P2P 网络借贷问题平台特征分析及投资者识别——来自 222 家平台的证据》,《财贸经济》第 12 期。

(105) 王艳艳 许锐 王成龙 于李胜, 2018:《关键审计事项段能够提高审计报告的沟通价值吗?》,《会计研究》第 6 期。

(106) 王云 李延喜 马壮 宋金波, 2017:《媒体关注、环境规制与企业环保投资》,《南开管理评论》第 6 期。

(107) 吴冬梅 刘运国, 2012:《捆绑披露是隐藏坏消息吗?——来自独立董事辞职公告的证据》,《会计研究》第 12 期。

(108) 吴建祖 肖书锋, 2016:《创新注意力转移、研发投入跳跃与企业绩效——来自中国 A 股上市公司的经验证据》,《南开管理评论》第 2 期。

(109) 吴璇 田高良 李玥婷 薛宇婷, 2019:《经营信息披露与股票收益联动——基于财务报告文本附注的分析》,《南开管理评论》第 3 期。

(110) 吴雨 李洁 尹志超, 2018:《房价上涨对 P2P 网络借贷成本的影响分析》,《金融研究》第 11 期。

(111) 向虹宇 王正位 江静琳 廖理, 2019:《网贷平台的利率究竟代表了什么?》,《经济研究》第 5 期。

(112) 谢德仁 林乐, 2015:《管理层语调能预示公司未来业绩吗?——基于我国上市公司年度业绩说明会的文本分析》,《会计研究》第 2 期。

- (113) 徐莉萍 辛宇 祝继高, 2011:《媒体关注与上市公司社会责任之履行——基于汶川地震捐款的实证研究》,《管理世界》第3期。
- (114) 徐巍 陈冬华, 2016:《自媒体披露的信息作用——来自新浪微博的实证证据》,《金融研究》第3期。
- (115) 徐永新 陈婵, 2009:《媒体荐股市场反应的动因分析》,《管理世界》第11期。
- (116) 薛爽 肖泽忠 潘妙丽, 2010:《管理层讨论与分析是否提供了有用信息?——基于亏损上市公司的实证探索》,《管理世界》第5期。
- (117) 阎达五 孙蔓莉, 2002:《深市B股发行公司年度报告可读性特征研究》,《会计研究》第5期。
- (118) 杨道广 陈汉文 刘启亮, 2007:《媒体压力与企业创新》,《经济研究》第8期。
- (119) 杨洁 郭立宏, 2017:《声明还是缄默:负面报道后国企和民企印象管理行为差异研究》,《南开管理评论》第1期。
- (120) 杨晓兰 沈翰彬 祝宇, 2016:《本地偏好、投资者情绪与股票收益率:来自网络论坛的经验证据》,《金融研究》第12期。
- (121) 杨玉龙 吴文 高永靖 张倩男, 2018:《新闻媒体、资讯特征与资本市场信息效率》,《财经研究》第6期。
- (122) 姚颐 赵梅, 2016:《中国式风险披露、披露水平与市场反应》,《经济研究》第7期。
- (123) 叶德珠 陈霄, 2017:《标点与字数会影响网络借贷吗——来自人人贷的经验证据》,《财贸经济》第5期。
- (124) 叶青 李增泉 徐伟航, 2016:《P2P网络借贷平台的风险识别研究》,《会计研究》第6期。
- (125) 伊志宏 杨圣之 陈钦源, 2019:《分析师能降低股价同步性吗——基于研究报告文本分析的实证研究》,《中国工业经济》第1期。
- (126) 易志高 潘子成 茅宁 李心丹, 2017:《策略性媒体披露与财富转移——来自公司高管减持期间的证据》,《经济研究》第4期。
- (127) 尹海员 吴兴颖, 2019:《投资者高频情绪对股票日内收益率的预测作用》,《中国工业经济》第8期。
- (128) 应千伟 吕昊婧 邓可斌, 2017:《媒体关注的市场压力效应及其传导机制》,《管理科学学报》第4期。
- (129) 游家兴 吴静, 2012:《沉默的螺旋:媒体情绪与资产误定价》,《经济研究》第7期。
- (130) 游家兴 陈志锋 肖曾昱 薛小琳, 2018:《财经媒体地域偏见实证研究》,《经济研究》第4期。
- (131) 游家兴 郑建鑫, 2013:《媒体情绪、框架依赖偏差与IPO异象——基于议程设置理论的研究视角》,《投资研究》第12期。
- (132) 于忠泊 田高良 齐保垒 张皓, 2011:《媒体关注的公司治理机制——基于盈余管理视角的考察》,《管理世界》第9期。
- (133) 俞庆进 张兵, 2012:《投资者有限关注与股票收益——以百度指数作为关注度的一项实证研究》,《金融研究》第8期。
- (134) 曾建光, 2015:《网络安全风险感知与互联网金融的资产定价》,《经济研究》第7期。
- (135) 曾庆生 周波 张程 陈信元, 2018:《年报语调与内部人交易:“表里如一”还是“口是心非”?》,《管理世界》第9期。
- (136) 张纯 吴明明, 2015:《媒体在资本市场中的角色:信息解释还是信息挖掘?》,《财经研究》第12期。

- (137) 张海洋 蔡航, 2018:《头衔的价值——来自网络借贷的证据》,《经济学(季刊)》第4期。
- (138) 张皓星 黄益平, 2018:《情绪、违约率与反向挤兑——来自某互金企业的证据》,《经济学(季刊)》第4期。
- (139) 张继勋 蔡闫东 倪古强, 2019:《社会责任披露语调、财务信息诚信与投资者感知——一项实验研究》,《南开管理评论》第1期。
- (140) 张俊生 汤晓建 李广众, 2018:《预防性监管能够抑制股价崩盘风险吗?——基于交易所年报问询函的研究》,《管理科学学报》第10期。
- (141) 张龙平 吕敏康, 2014:《媒体意见对审计判断的作用机制及影响——基于新闻传播学理论的解释》,《审计研究》第1期。
- (142) 张宁 刘春林, 2012:《传闻澄清的市场反应研究——澄清公告详细性的作用》,《经贸经济》第3期。
- (143) 张琦 郑瑶, 2018:《媒体报道能影响政府决算披露质量吗?》,《会计研究》第1期。
- (144) 张晓玫 王博 周玉琴, 2016:《非完全利率市场化下网络借贷的利率定价有效吗——来自“人人贷”的微观借贷证据》,《南开管理评论》第4期。
- (145) 张谊浩 李元 苏中锋 张泽林, 2014:《网络搜索能预测股票市场吗?》,《金融研究》第2期。
- (146) 张永杰 张维 金曦 熊熊, 2011:《互联网知道的更多么?——网络开源信息对资产定价的影响》,《系统工程理论与实践》第4期。
- (147) 赵静梅 何欣 吴风云, 2010:《中国股市谣言研究: 传谣、辟谣及其对股价的冲击》,《管理世界》第11期。
- (148) 赵龙凯 陆子昱 王致远, 2013:《众里寻“股”千百度——股票收益率与百度搜索量关系的实证探究》,《金融研究》第4期。
- (149) 赵子夜 杨庆 杨楠, 2019:《言多必失? 管理层报告的样板化及其经济后果》,《管理科学学报》第3期。
- (150) 郑志刚 丁冬 汪昌云, 2011:《媒体的负面报道、经理人声誉与企业业绩改善——来自我国上市公司的证据》,《金融研究》第12期。
- (151) 周开国 应千伟 陈晓娴, 2014:《媒体关注度、分析师关注度与盈余预测准确度》,《金融研究》第2期。
- (152) 周开国 应千伟 钟畅, 2016:《媒体监督能够起到外部治理的作用吗?——来自中国上市公司违规的证据》,《金融研究》第6期。
- (153) 周雄伟 朱恒先 李世刚, 2017:《“平台参与投资”与 P2P 筹资效率——基于拍拍贷平台“拍活宝”数据的经验研究》,《中国工业经济》第4期。
- (154) 朱宝宪 王怡凯, 2001:《证券媒体选股建议效果的实证分析》,《经济研究》第4期。
- (155) 朱恩伟 吴璟 刘洪玉, 2019:《基于新闻文本共现性的银企关系分析——以房地产上市公司为例》,《金融研究》第2期。
- (156) Borochin, P. & W. Cu(2017), “Alternative Corporate Governance: Domestic Media Coverage of Mergers and Acquisitions in China”, *Journal of Banking and Finance* 87: 1-25.
- (157) Bowen, R. M., S. Dutta, S. Tang & P. Zhu(2018), “Inside the ‘Black Box’ of Private In-House Meetings”, *Review of Accounting Studies* 23(2): 487-527.
- (158) He, D., S. Cohen, B. Price, D. Kifer & C. L. Giles (2017), “Multi-scale multi-task fcn for semantic page segmentation and table detection”, *In International Conference on Document Analysis and Recognition (ICDAR)*.
- (159) Huang, Y., H. Qiu & Z. Wu(2016), “Local Bias in Investor Attention: Evidence from

China's Internet Stock Message Boards", *Journal of Empirical Finance* 38: 338-354.

(160) Huang, Y. & P. Luk(2020), "Measuring economic policy uncertainty in China", *China Economic Review* 59, 101367.

(161) Hung, M., T. J. Wong & F. Zhang(2015), "The Value of Political Ties Versus Market Credibility: Evidence from Corporate Scandals in China", *Contemporary Accounting Research* 32: 1641-1675.

(162) Jiang, L., J. Liu & B. Yang(2019), "Communication and comovement: Evidence from online stock forums", *Financial Management* 48(3): 805-847.

(163) Kim, J. B., L. Li, Z. Yu & H. Zhang(2019), "Local Versus Non-local Effects of Chinese Media and Post-earnings Announcement Drift", *Journal of Banking and Finance* 106: 82-92.

(164) Liu, S., Y. Dai & D. Kong(2017), "Does It Pay to Communicate with Firms? Evidence from Firm Site Visits of Mutual Funds", *Journal of Business Finance & Accounting* 44: 611-645.

(165) Mao, S., A. Rosenfeld & T. Kanungo(2003), "Document structure analysis algorithms: A literature survey", *Proc Spie Electronic Imaging* 5010: 197-207.

(166) Piotroski, J.D., T.J.Wong & T.Zhang(2017), "Political bias in corporate news: The role of conglomeration reform in China", *Journal of Law and Economics* 60(1):173-207.

(167) Viana, M. P. & D. A. B. Oliveria(2017), "Fast CNN-based document layout analysis", *In Proceedings of the IEEE International Conference on Computer Vision (ICCV)*.

(168) Yan, Y., X. Xiong, J. G. Meng & G. Zou(2019), "Uncertainty and IPO initial returns: Evidence from the Tone Analysis of China's IPO Prospectuses", *Pacific-Basin Finance Journal* 57: 101075.

(169) Yang, X., E. Yumer, P. Asente, M. Kralej, D. Kifer & C. L. Giles(2017), "Learning to extract semantic structure from documents using multimodal fully convolutional neural network", *In Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*.

(170) Ye, Y., L. Huang & M. Li(2015), "Negative Media Coverage, Law Environment and Tunneling of Controlling shareholder: Evidence from Chinese Listed Companies", *China Finance Review International* 5: 3-18.

(171) You, J., B. Zhang & L. Zhang(2018), "Who Captures the Power of the Pen?", *Review of Financial Studies* 31: 43-96.

(172) Yu, X., P. Zhang & Y. Zheng(2015), "Corporate governance, Political Connections, and Intra-industry Effects: Evidence from corporate Scandals in China", *Financial Management* 44: 49~80.