编者按:近年来,随着信息技术的不断发展,文本分析在经济学、管理学研究中的应用日益广泛,在一定程度上提供了研究新视角。我们也注意到,文本分析在现实中存在应用落后于技术发展、文本分析提取指标与经济学理论匹配度不高、缺乏新的应用场景等问题。为此,《财贸经济》编辑部邀请中南财经政法大学金融学院李春涛老师和上海财经大学公共管理学院冯旭南老师撰写研究指引,这同时也是一篇内容丰富的文献综述。编辑部希望通过系统回顾文本分析的发展路径及其在经济学研究中的应用进展,为读者提供更加完整的认识框架,并推动文本分析更加科学合理的应用。

文本分析发展路径及在经济学研究中的应用

李春涛¹ 冯旭南² (1.中南财经政法大学金融学院;2.上海财经大学公共管理学院)

近年来,随着计算机技术的发展,文本、音频、视频等非结构化数据逐渐受到学术界关注。相对于音频、视频等难以处理的非结构化数据,基于文本数据的分析在经济学研究中得到了广泛的应用,原来一些难以测度的指标,比如金融科技发展水平、公司诚信等因文本分析而得到客观的量化,文本分析也从情感和可读性等维度让我们进一步深入认识客观世界。

一、语调分析

简单地说, 文本分析就是数单词, 例如, 我们阅读一份报告, 发现某些词汇出现的次数较多, 并基于判断报告的主题, 这就是一种朴素的文本分析方法。这种朴素的方法据说可以追溯到中世纪时期的欧洲; 20世纪初期, Mendenhall 也曾经用文本分析的方法来判断莎士比亚的一些作品是不是别人写的; 也有学者曾经研究过《红楼梦》, 并根据古汉语虚词出现的频率认定《红楼梦》可能是多个人合作撰写(Hu等, 2014); 两次世界大战中都曾经有军方通过截获的电文中字母代码出现的频数来破译电文。

在金融与会计领域,我们通过文本分析来判断文本的语调。最初,学者花费大量的精力人工阅读大量的文本,并对这些文本进行主观判断,以确定文本的语调,比如 Callahan 和 Smith(2004)研究了 1993—2001 年 71 家公司的管理层讨论与分析,用正面的项目数量减去负面项目数量构建了披露衡量指数(Composite Disclosure Score),发现高透明度公司的市值也高。Davis 等(2010)通过人工阅读筛选、基于 Diction 词典的机器分词等办法,发现年报、MD&A、新闻稿中的语气与企业未来 ROA 呈正相关。You 等(2017)对中国 8 份报纸的 20 万份财经新闻报道进行分类,然后发现基于市场化的媒体更具有对上市公司的监督作用。人工阅读的优势和劣势都很明显,优势在于人类能读懂一句话背后的潜台词,对于语义有更准确的解读,从而能够给出准确的判断;劣势在于可重复性太差、不能处理海量的文本、如果多人一起工作则难以统一标准等等。

大数据时代也带来了海量的文本,人工判断语调的方法显然难以为继,因此出现了计算机数单词的方法,它基于预设的情感词典对文本语调进行判别,是一种基于简单规则的分析方法。具体地说,我们首先构建两个字典,正面词汇字典和负面词汇字典,然后对每一篇文本(比如年报中的MD&A、业绩说明会文本、分析师报告)数单词,得到正面词汇和负面词汇分别出现的次数,如果正面词汇较多,则这篇文本的语调是正面的,否则就是负面的。

这里有几个细节:首先,我们分析的文本,本身可能是分段的,每段有多个句子,但 是数单词工作的对象是词汇,我们需要把文本拆成词汇。英文的词与词之间有空格或者标 点符号,但是本质上需要分词,比如 DCF 和 Discounted Cash Flow 本质上是同一个单词, NPV 和 Net Present Value 也是一个意思, interest rate 不能当作两个单词, CAPM 和 Beta 可能 代表同样的意思,我们需要把这些词汇统一起来,比如将 Discounted Cash Flow 全部换成 DCF。其次,中文的分词可能更加麻烦,一个主要的原因在于中文词的边界更模糊,选择 不同的切分方法对语义的影响很大。另外、中文分词常用的工具是Jieba、但是这个分词工 具并不是给金融会计的研究者做的,有很多专用词汇并没有包含在 Jieba 的字典里,因此我 们需要扩充该词典。最后、分词以后的一篇文章被拆成一大堆词组或单词、但是里面有很 多没有意义的虚词,比如标点符号、英文中的冠词、中文中的介词等,这些词本身没有意 思、但是频率还挺高、我们把这些词和标点统称为停词(stop word),需要把他们去掉。 剩下的工作就是对照着我们的字典,数一下每个单词出现的频率就行了。但是数单词的结 果,严重地依赖于我们的字典和分词的算法,字典我们可以自己去定义,并力求完善,而 分词的算法,涉及词汇的优先顺序、从左到右分还是从右向左分都会影响我们的结果。语 调分析的应用研究很多,金融会计领域较多利用年报中管理层分析与讨论(MDA)章节、 业绩说明会等文本文件构建管理层语调指标、研究管理层语调对外部市场的反应(Feldman, 2010; Bochkay 和 Levine, 2013; 林乐和谢德仁, 2017; 朱朝晖等, 2018) 以及对企业自身 的影响(Brockman 等,2013;谢德仁、林乐, 2015; 甘丽凝等2019; 原东良等2021), 同时 也有少部分研究者,如 Loughran 和 McDonald (2011) 从整个年报的语调入手,研究其对公 司股价的影响,曾庆生等(2018)和赵宇亮(2020)也采用年报的语调观察年报披露以后 企业内部交易和企业融资约束的影响。读者可以阅读如下几篇综述性文章: Feng Li (2011) 介绍了有关文本分析的研究方向与方法论的产生、Tim 等(2016)讲述了文本分析方法论 中存在的缺陷,并对未来学术发展做出展望。

二、信息含量

与语调分析类似,文献中有很多研究考察了文本信息披露是否具有信息含量,比如年报 MD&A 中是否蕴含了经理人对企业未来发展的前瞻性信息披露,是否对这些信息有市场反应(李子健等,2022),再比如分析师报告逻辑是否一致,能否促进上市公司信息融入股价(马黎珺等,2022)。关于信息含量的研究,本质上还是数单词,以前瞻性信息披露为例,研究者首先构建一个刻画前瞻性的词库,比如"未来五年""将来""预期"等,并计算出这些前瞻性信息所占的比例,或者包含前瞻性信息的句子所占的比例,从而分析前瞻性信息的市场反应。例如 Bryan(1997)通过 MD&A 中强制披露的增量信息,发现资本支出披露与即期和远期的收益相关,Li(2010)发现年报中前瞻性信息能够预测公司未来业绩,Muslu等(2015)发现年报 MD&A 中前瞻性披露有助于提升资本市场信息效率,Bozanic等(2018)研究盈余公告中前瞻性信息披露后的市场反应与分析师盈余预测准确性,马黎珺等(2019)发现证券分析师的前瞻性语句能够向市场传递增量信息。

三、文本可读性

另一个相关的领域在于文本的可读性。这一类文献最初的思路是上市公司在信息披露的时候,可能会使用晦涩难懂的语言来隐藏信息。最初的研究者把可读性定义为需要受到多少年教育的人才能读懂,比如一篇文本,小学毕业生就能读懂,可以将其可读性定义为6,如果要大学毕业才能读懂,则可读性定义为18。剩下的问题就是拿出一篇文章,分别请几个小学生、几个初中生、几个高中生或大学生去读,如果大家都能读懂,则可读性为6,如果只有大学生读得懂,则可读性为18。但是这样做显然是难以操作的,一方面成本

太高、一方面研究助手在阅读大量文本的过程中其阅读能力必然能得到提升,比如几千个文本阅读下来,小学生的阅读能力提升到了大学水平,从而使得通过人工识别进行可读性评级的方法会因为标准的渐变而失去客观性。

没有人工智能的机器可以克服这一困难。Li (2008)提出用Fog指数测度可读性,Fog指数本质上还是数单词,它包括两个维度,即每句话平均有几个单词和每个单词平均有几个音节。其基本思想是,句子越长,比如英语中的从句越多,阅读难度越大,从而可读性越差;其次,单词越长,可能越不常见,从而会降低可读性。那么好了,我们把一篇文本拿出来,先拆成一个个句子,数一下每句话多少个单词,再看看每个单词有几个音节,Fog指数就出来了。Fog指数作为最早被Li (2008)引入衡量大样本的年报可读性,随后得到了广泛的应用。比如Biddle等(2009)发现财报可读性与资本投资效率呈高度正相关;Miller(2010)发现小型投资者在年报披露前后,会减少对年报可读性差公司的投资;Lo等(2017)则发现盈余管理水平与MD&A的可读性呈负相关。

也有学者把Fog指数引用到中文语境中,句长的问题都好解决,但是中文没有所谓的双音节字,因此比如Qiu等(2013)通过统计每篇文档的汉字平均笔画、基本汉字占比、文章中句式完整的比例等信息,构造了衡量可读性的指标,发现可读性会吸引更多分析师的关注。

虽然 Fog 指数得到了广泛的应用,但是对 Fog 的质疑也从来没有停止过,特别是多音节词是否就晦涩难懂。Loughran 和 McDonald(2014)发现年报中有大量易于理解的多音节词汇,因此认为 Fog 指数测度年报可读性并不合理。Loughran 和 McDonald(2014)进一步提出利用年报大小(比如文件的字节数)测度年报的可读性,其基本的思路是,披露越多的年报,可能是经理人把关键信息淹没在信息森林的行为。沿袭这一思路,Jin 等(2018)利用年报的字数、字符数、页数度量可读性,发现可读性的提高有助于提升企业透明度,降低企业的代理成本。任宏达和王琨(2018,2019)采用年报文本大小和年报页数衡量可读性,研究了两类企业的信息披露质量:依赖社会关系获取资源的企业(关系型企业)、竞争市场中的企业。这一方法简洁、客观且可重复,但是年报大小往往也会因为公司的业务范围、披露详细程度不同而不同。这就像我们读文言文和现代汉语,显然前者更简洁,但是可读性更差。

Bonsall等(2017)提出了Bog指数,其基本的思路正如Loughran和McDonald(2014)所述,多音节单词并非不常用(比如我们常常使用 congratulations 和 populations,但是很少使用 thy 和 xeme),Bonsall等,(2017)提出可以根据常用性将词汇分级,那些使用了较多常用词汇的年报必然是更可读的。关于中文可读性的研究,更多地是采用类似于Bog指数的方法考察年报语义的复杂程度,比如孟庆斌等(2017)利用常见汉字词汇,王克敏等(2018)则利用转折词(虽然、但是、然而等)占比、会计术语密度、次常用字密度,李春涛等(2020)则利用年报常用词频构建年报可读性指标。

四、机器学习

前面提到的文本分析方法基本上都可以看作数单词,但是人类语言的复杂性,导致这种数单词的方法往往具有一定的局限性。比如,"过去一年,汽车行业复苏强劲,本公司销售量也得到了快速复苏,我们预期在未来三年左右实现扭亏为盈",这句话使用的全部是正面词汇,但是表达的是完全负面的意思,即公司未来三年将一直亏损。词频分析难以穷尽全部的可能,随着计算机科学和计算语言学的发展,学者们逐渐将机器学习引入到文本分析中。

机器学习的文本分析应用大致分为两个方向,一是利用有监督模型如朴素贝叶斯分析、 支持向量机等对文本进行分类,其中语调分析是最广泛的应用场景,它们通过学习文本表 征和已知标签关系,并基于训练好的模型对未标记的海量文本进行推理。由于非结构化的文本数据包含的信息无法用事先定义好的方式或数据模型建模,因此这类方法首先利用文本表示技术如BOW、N-gram、TF-IDF、word2vec等将文本数据转化为词向量,再根据应用场景建模判别文本类型。一般认为,这种预结构化的学习方式对文本特征的挖掘是浅层的。深层次的学习基于人工神经网络,通过学习一组非线性变换将输入文本直接映射到输出,从而将特征工程集成到模型训练过程中,常用的模型包括RNN、BERT等,Li等(2020)对此进行了详细的、介绍,感兴趣的读者可以参阅。二是利用无监督方法挖掘文本信息,如文本相似度、主题模型等。从建模形式上看,有监督方法是通用的,它们一般只需根据场景的不同更换标签,进而通过数据驱动的方式学习输入和标签的关系。相比之下,无监督方法仅需输入文本,而不需要"标准答案"指导训练。

以朴素贝叶斯在语调分析中的应用为例,假如我们不想通过穷尽各种词频组合的方式 判断其语调,而是希望借助计算机快速判断其传达的是正面还是负面语调,文本的用词 (内容、频率等)成为判断该文本所包含信息属性的重要依据。实现这一思路的算法在本 质是计算特定关键词出现时整个文本所传达信息为某种(正/负)属性的条件概率,即:

 $P[Tone | KW_1, KW_2, \dots, KW_n]$.

由贝叶斯公式可知:

$$P[Tone | KW_1, KW_2, \dots, KW_n] = \frac{P[KW_1, KW_2, \dots, KW_n | \text{Tone}] P[Tone]}{P[KW_1, KW_2, \dots, KW_n]}.$$

我们可以通过经人工标记信息属性的样本(训练样本)计算上述等式右边部分的三个概率。以朴素贝叶斯为例,我们可以首先通过人工将其中 500 个样本信息披露归纳为正面或负面,计算现有样本中文本信息特定属性的概率,即P[Positive]、P[Negative];然后,将这 500 份文本进行分词,并假定这些关键词(这里,我们事实上把文本表征为一个只包含 0 和 1 的 OneHot 向量,其他常用的方式包括 TF-IDF 等)的使用先天是独立的,可以计算 出 $P[KW_1, KW_2, ..., KW_n|$ Positive] 和 $P[KW_1, KW_2, ..., KW_n|$ Negative];最后,计算出 $P[Tone|KW_1, KW_2, ..., KW_n]$,即经过训练的一个朴素贝叶斯分类器。最后可以通过增加训练样本的量,不断优化上述朴素贝叶斯分类器,并用于判断其余海量文本的语调。

Loughran 和 McDonald (2016) 明确指出朴素贝叶斯方法适应于阅读文本,使得研究大量文本信息成为可能。机器学习的方法发展迅猛,当前流行的还有文本相似度、支持向量机和主题模型等,其应用也越来越广泛。机器学习的方法主要用于处理文本主题分类和情感分析问题,国内方面,比如马黎珺等 (2019) 利用机器学习将文本分为历史语句和前瞻性语句,用于研究分析师报告的前瞻性描述和企业未来进步的关系; 张宗新和吴钊颖 (2021) 利用机器学习对媒体情绪进行有效识别,并得出了"分析师的预测倾向受到媒体情绪影响"的结论。不同文本之间的相似度分析代表含义不同,其中,计算不同企业年报之间的文本相似度可以用来衡量企业产品市场的竞争程度,相同企业不同研报文本相似度可以衡量研报的信息增量(Hoberg 和 Philips,2016; 刘昌阳等 2020),不仅如此,同一家公司的前后年份的 MDA 文本相似度可以有效衡量企业的形式主义程度(钱爱民和朱大鹏,2020)。Bybee 等 (2021)利用主题模型从华尔街日报中挖掘了多个与宏观经济状态相关的主题变量,他们发现这些变量和经济活动相关的数字测度密切相关,并且在解释经济动态变化方面超过了标准的数字经济指标。

五、专业文本特征词提取

利用文本挖掘技术进行专业特征文本提取,可以有效构建特征指数,比如,货币政策立场指数、互联网金融指数(李成、高智贤,2014;郭品、沈悦,2015;沈悦、郭品2015)。同时根据专业术语和定量文本的出现频率度量文本的复杂程度和精确度(李晓溪等,2019),也可以根据短视词汇、企业研发词汇、企业文化词汇、数字化转型词汇等度量企业管理层短视程度、企业研发披露程度、企业文化建设程度以及企业转型化程度等(李岩琼、姚颐,2020;胡楠等,2021;权小锋、朱宇翔,2022;黄大禹等2021)。

从当前已发表文献可以得出,尽管文本分析和文本挖掘提取非结构化数据,可以对结构化数据进行补充,为研究者提供新视角,在经济学研究领域的应用,也正处于蓬勃发展阶段。但笔者认为该方法在当前经济学研究领域应用中还存在一定的挑战和机遇。具体可以从当前经济学领域中文本分析的局限性和未来研究展望两方面进行论述。

六、文本分析的局限性

文本分析和文本挖掘首先需要对文本数据进行收集和清洗,当前经济学研究领域的文本信息披露形式不具有统一性,其中部分披露的文本文件包含图片或属于加密文件,从而加大文本提取难度,应规范文本披露形式以提高文本信息提取质量。其次,获取文本信息渠道不同,同时文本本身噪声较多,比如网络上获得的文本往往夹杂着很多 html 标签,研究者处理方式不同,会导致构建的数据指标存在差异,重复性较低,应构建权威文本数据库、研发具有不同词汇优先级的分词方法、透明的专用字典,从而规范数据处理过程,提高可复制性。再次,关于专业文本特征词提取,缺乏权威的专业词库,以至于各类指标定义不一,指标可操控性加大。最后,经济学领域更倾向于利用特征词提取文本信息,应结合自然语言领域中文本摘要、命名实体识别等较新的文本提取技术。

七、未来研究展望

文本分析的流行,一方面得益于计算机运算能力和自然语言识别能力的提升,另一方面得益于越来越多的海量文本信息,这是供给方的原因。从需求方来说,过去很多难以测度的指标,在引入文本分析以后成为可能。除了前述的语调、相似度和可读性之外,学术界一直关心经理人的诚信问题,但是在文本分析之前,诚信难以客观测度,文本分析技术为诚信的测度提供了可能,比如,Guiso等(2015)、张维迎和柯荣住(2002)、戴亦一等(2019)和林斌等(2016),通过分析员工对经理人的评价测度诚信;Kizirian等(2005)和Dikolli等(2020)则利用年度股东信文本的措辞(是否使用了为自己开拓的词汇);姜付秀等(2015)和翟胜宝等(2015)则利用企业官方公布的文化规范进行度量。

金融科技是当前研究的热点问题,但是一个地区金融科技的发展水平测度一直是一个难题,李春涛等(2020)利用百度新闻搜索每一个地区金融科技关键词获得的页面数量,构建了不同地区不同年份的金融科技指标。

未来的研究一方面是引入更先进的算法,提高文本分析指标的计算效率和提升可重复性;另一方面是引入更多的可以用来进行文本分析的语料库,比如 Larcker 和 Zakolyukina (2012)创建了针对电话会议的词表,并分析了企业高管的会议语音,认为欺骗性高管会使用更多的通俗化语言和更积极的词汇;Obaid等(2021)分析了报道中的文本和嵌入图像,发现图像在一定程度上能替代文本分析,尤其是在市场出现极端趋势时。当然还需仔细考察如"机器学习"等,以纯客观为主导的新兴技术,是否契合当前领域的研究逻辑,否则会使研究产生过多的"噪音"。

参考文献:

[1] 戴亦一,张鹏东,潘越.老赖越多,贷款越难?——来自地区诚信水平与上市公司银行借款的证据[[].金融研

究,2019(08).

- [2] 甘丽凝,陈思,胡珉,王俊秋.管理层语调与权益资本成本——基于创业板上市公司业绩说明会的经验证据[J].会计研究,2019(06):27-34.
- [3] 郭品,沈悦.互联网金融对商业银行风险承担的影响:理论解读与实证检验[[].财贸经济,2015(10):102-116.
- [4] 郭品,沈悦.互联网金融加重了商业银行的风险承担吗?——来自中国银行业的经验证据[J].南开经济研究,2015(04):80-97.
- [5] 胡楠,薛付婧,王昊楠.管理者短视主义影响企业长期投资吗?——基于文本分析和机器学习[J].管理世界,2021,37(05):139-156+11+19-21.
- [6] 黄大禹,谢获宝,孟祥瑜,张秋艳.数字化转型与企业价值——基于文本分析方法的经验证据[J].经济学家,2021(12):41-51.
- [7] 姜付秀,石贝贝,李行天."诚信"的企业诚信吗?——基于盈余管理的经验证据[[].会计研究, 2015(8).
- [8] 李成,高智贤.货币政策立场与银行信贷的异质性反应——基于信贷传导渠道的理论解读与实证检验[J]. 财贸经济,2014(12):51-63.
- [9] 李春涛,张计宝,张璇. 年报可读性与企业创新[]]. 经济管理, 2020, 42(10):156-173.
- [10] 李晓溪,杨国超,饶品贵.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[J].经济研究,2019,54(05):181-198.
- [11] 李岩琼,姚颐.研发文本信息:真的多说无益吗?——基于分析师预测的文本分析[J].会计研究,2020(02):26-42.
- [12] 李子健, 李春涛, 冯旭南. 非财务信息披露与资本市场定价效率[J]. 财贸经济, 2022, 43(9):38-52.
- [13] 林斌,陈颖,舒伟,郑颖.社会信任与公司违规[[].中国会计评论, 2016(3).
- [14] 刘昌阳,刘亚辉,尹玉刚.上市公司产品竞争与分析师研究报告文本信息[[].世界经济,2020,43(02):122-146.
- [15] 马黎珺, 吴雅倩, 伊志宏, 刘嫣然. 分析师报告的逻辑性特征研究:问题,成因与经济后果[J]. 管理世界, 2022, 38(8): 217-231.
- [16] 马黎珺, 伊志宏, 张澈. 廉价交谈还是言之有据?——分析师报告文本的信息含量研究[J]. 管理世界, 2019, 35(7):182-200.
- [17] 孟庆斌,杨俊华,鲁冰.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,2017(12):132-150.
- [18] 钱爱民,朱大鹏.财务报告文本相似度与违规处罚——基于文本分析的经验证据[J].会计研究,2020(09):44-58
- [19] 权小锋,朱宇翔. "员工关爱"文化、成本粘性与公司绩效[J].财贸经济,2022,43(07):118-133.
- [20] 任宏达,王琨.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[J].会计研究,2019(03):32-39.
- [21] 任宏达,王琨.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[J].南开管理评论,2018,21(05):128-138.
- [22] 王克敏,王华杰,李栋栋,戴杏云.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[J].管理世界,2018,34(12):120-132+194.
- [23] 谢德仁,林乐. 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[J]. 会计研究, 2015,(2):20-27.
- [24] 张维迎,柯荣住.信任及其解释:来自中国的跨省调查分析[J].经济研究, 2002(10).
- [25] 张宗新,吴钊颖.媒体情绪传染与分析师乐观偏差——基于机器学习文本分析方法的经验证据[J].管理世界,2021,37(01):170-185+11+20-22.
- [26] 赵宇亮.年报净语调对企业债权融资的影响研究[J].经济管理,2020,42(07):176-191.
- [27] 周建,原东良,马雨飞.MD&A 语调会影响企业履行社会责任吗?——基于信息增量与印象管理的视角[J]. 管理学刊,2021,34(06):88-107.

- [28] 朱朝晖, 许文瀚. 上市公司年报语调操纵、非效率投资与盈余管理[J]. 审计与经济研究, 2018,(3):63-72.
- [29] 曾庆生,周波,张程等. 年报语调与内部人交易:"表里如一"还是"口是心非"? [J]. 管理世界, 2018,(9):143-160.
- [30] Biddle, G.,G. Hilary, and R. Verdi, How Does Financial Reporting Quality Relate to Investment Efficiency? Journal of Accounting and Economics 2009, 48: 112 131.
- [31] Bochkay K, Levine C. Using MD& A to improve earnings forecasts[]]. SSRN Electronic Journal, 2013.
- [32] Bonsall, S.B., Leone, A.J., Miller, B.P., & Rennekamp, K. A Plain English Measure of Financial Reporting Readability. Journal of Accounting and Economics, 2017,63:329-357.
- [33] Bozanic Z, Roulstone, D. T., and Van Buskirk, A., Management earnings forecasts and other forward-looking statements[J]. Journal of Accounting and Economics, 2018, 65(1): 1-20.
- [34] Bybee, L., Kelly, B.T., Manela, A.and Xiu, D. Business News and Business Cycles. Working
- [35] Paper 29344. National Bureau of Economic Research, 2021.
- [36] Brockman, P., Li,X.,and Price, S. M., Do managers put their money where their mouths are? Evidence from insider trading after conference calls[J]. SSRN Electronic Journal, 2013.
- [37] Bryan, S. H., Incremental Information Content of Required Disclosures Contained in Management Discussion and Analysis. The Accounting Review, 1997,72: 285 301.
- [38] Callahan, C. M., and R. E. Smith, "Firm Performance and Management's Discussion and Analysis Disclosures: An Industry Approach." Working paper, University of Arkansas Fayetteville, 2004.
- [39] Davis, A.K., & Tama-Sweet, I. Managers' Use of Language Across Alternative Disclosure Outlets: Earnings Press Releases Versus MD&A. Financial Accounting Journal, 2011.
- [40] Davis, A.K., Piger, J., and Sedor, L.M.. Beyond the Numbers: Measuring the Information Content of Earnings Press Release Language. Usefulness of Quantitative & Qualitative Information in Earnings Announcements, 2011.
- [41] Dikolli, S.S., Keusch, T., Mayew, W. J., and Steffen, T. D., "CEO Behavioral Integrity, Auditor Responses, and Firm Outcomes", The Accounting Review, 2020, 95(2), pp. 61-88.
- [42] Feldman, R., S. Govindaraj, J. Livnat, and Segal, B., Management's Tone Change, Post Earnings
 Announcement Drift and Accruals. Review of Accounting Studies. Volume 15, Number 4, 2010, 39: 915-953.
- [43] Guiso, L., Sapienza, P. and Zingales, L., "The Value of Corporate Culture", Journal of Financial Economics, 2015, 117(1): 60-76.
- [44] Hagen, L., Harrison, T.M., Uzuner, Ö., May, W., Fake, T., and Katragadda, S., E-petition popularity: Do linguistic and semantic factors matter? Gov. Inf. Q., 2016,33: 783-795.
- [45] Hoberg, G., Discussion of Using Unstructured and Qualitative Disclosures to Explain Accruals. Journal of Accounting and Economics, 2016, 62:228-233.
- [46] Hu, X., Wang, Y., and Wu, Q., Multiple Authors Detection: a Quantitative Analysis of Dream of the Red Chamber. Adv. Data Sci. Adapt. Anal., 2014,6.
- [47] Kizirian, T. G., Mayhew, B. W. and Sneathen, J., The Impact of Management Integrity on Audit Planning and Evidence. Auditing: A Journal of Practice & Theory, 2005, 24(2):49-67.
- [48] Larcker, D. F., and Zakolyukina, A. A., Detecting Deceptive Discussions in Conference Calls. Journal of Accounting Research, 2012, 50: 495-540.
- [49] Li, F., The information content of forward looking statements in corporate filings—A naïve Bayesian machine learning approach[J]. Journal of Accounting Research, 2010, 48(5): 1049-1102.
- [50] Li, F., Textual Analysis of Corporate Disclosures: A Survey of the Literature, 2011.
- [51] LI, F., Annual Report Readability, Current Earnings, and Earnings Persistence. Journal of Accounting and Economics, 2008, 45: 221 47.

- [52] Lo, K., Ramos, F., and Rogo, R., Earnings Management and Annual Report Readability. CGN: Disclosure & Accounting Decisions (Topic), 2016.
- [53] Loughran ,T., Mcdonald B, Textual Analysis in Accounting and Finance: A Survey[J]. Journal of Accounting Research, 2016, 54.
- [54] Loughran ,T., McDonald B. When is a Liability not a Liability? Textual Analysis, Dictionaries, and 10-Ks[J]. The Journal of Finance, 2011,66(1):35-65
- [55] Loughran, T., and B. McDonald. ,Measuring Readability in Financial Disclosures. The Journal of Finance, 2014,69(4): 1643-1671.
- [56] Miller, B. P., "The Effects of Reporting Complexity on Small and Large Investor Trading." The Accounting Review ,2010,85: 2107-2043.
- [57] Muslu, V., Radhakrishnan, S., Subramanyam, K. R., et al. Forward-looking MD&A Disclosures and the Information Environment[]]. Management Science, 2015, 61(5): 931-948.
- [58] Muslu, V., Radhakrishnan S, Subramanyam K R, et al. Forward-looking MD & A Disclosures and the Information Environment[J].Management Science,2015,61(5): 931 948.
- [59] Obaid, K., and Pukthuanthong, K., A Picture is Worth a Thousand Words: Measuring Investor Sentiment by Combining Machine Learning and Photos from News. FEN: Behavioral Finance (Topic), 2021.
- [60] Qiu, X., Jiang, S., and Deng, K., Automatic Assessment of Information Disclosure Quality in Chinese Annual Reports. NLPCC,2013.
- [61] You, J., Zhang, B., and Zhang, L., Who Captures the Power of the Pen? Governance, 2017.