

摘 要

随着互联网的广泛普及和大数据技术的蓬勃发展,我国政务服务平台积极与互联网融合,不断健全吸纳民意、汇集民智的工作机制。当前,在线政务留言平台逐步完善、参与人数增多,涌现出了海量的政务留言文本数据。快速准确分析留言信息反映的社情民意和用户关切,对推动社会治理具有重要决策价值。人工识别政务留言中的民众治理期望效率低下,因此,开发用户期望自动识别的算法模型和工具平台具有很强的研究与实际应用意义。

本文以我国最大的政务留言平台“领导留言板”为例,开展政务留言期望自动识别研究。目前,在线政务留言平台获取用户诉求主要依靠用户手工标注,这种选择方式主观、低效、误标率高,使得很多真正需要被及时回应的留言诉求被忽视或淹没。在研究面临的挑战中,留言期望识别面临着上下文信息不足、留言文本结构不平衡、留言文本不规范、留言数据噪声大等难点。针对以上挑战,本文将从提高政府回应效率视角,将留言期望识别问题建模为一个多分类任务,即通过分析用户留言信息,预测回应紧急程度(急需解决、建言、感谢)。本文提出了一个全新的基于深度强化学习的留言标签更正与留言期望识别的算法框架 **PecidRL**,对政务留言数据集中的噪声标签进行更正处理,并构建多种留言期望识别模型,来得到更好的留言期望识别精度。本文的主要研究工作与贡献如下:

(1) 本文提取了多视角的留言短文本特征,包括单词级别与文档级别的语义特征、基于图表示学习的短文本结构特征、情感分析特征以及语义特征与结构特征整合的融合特征。

(2) 引入强化学习模块解决留言标签噪声问题,对错误标签和模糊标签进行更正。实验结果表明,本文提出的标签更正强化学习算法模块有着出色的噪声修正能力,各项实验指标平均提升了 8.3%,最高的提升率达到 14.2%。

(3) 将多视角文本特征与经典的机器学习、深度学习、图神经网络模型融合,扩展了 19 种用户期望自动识别模型,并对比了标签更正前后的分类性能。研究结果表明,基于标签更正后的数据获得了更好的分类效果。其中,基于 **BERT** 语义特征的 **SVM** 模型具有 93.66% 的分类精度,于是将其作为最终的用户预期识别模型。

(4) 本文开发了在线服务平台 PecidRL (<http://www.csbg-jlu.info/PecidRL/>), 并提供了相关数据与代码, 旨在为政府工作人员以及相关研究人员提供工作和研究便利。

关键词:

留言期望, 多视角文本特征, 强化学习, 文本图表示, 更正与识别

Abstract

With the widespread popularity of the Internet and the vigorous development of big data technology, our government service platforms actively integrate with the Internet and constantly improve the working mechanism of absorbing public opinion and gathering public wisdom. At present, the online government petition platform has been gradually improved, the number of participants has increased, and a huge amount of petition text data has emerged. Rapid and accurate analysis of social opinion and public concerns reflected by petition has important decision-making value for promoting social governance. Manual identification of petition expectations is inefficient, so it is of strong research and practical application significance to develop algorithm and online platform for automatic petition expectation identification.

This paper takes the largest online petition platform "Message board for leaders" as an example, and carries out the research on the automatic identification of petition expectations. The online petition platform mainly relies on users' manual labeling to obtain requests, which is subjective, inefficient and highly mislabeled, making many requests that really need to be responded to in a timely manner be ignored or submerged. The research of petition expectation identification faces challenges such as insufficient contextual information, unbalanced text structure, non-standardized petition text, and noisy labeled data. To address the above challenges, this paper models the petition expectation identification problem as a multi-classification task from the perspective of improving the efficiency of government response, i.e., by analyzing petition textual information and predicting the urgency level of response (Urgency, Suggestion, and Gratitude). In this paper, we propose a novel framework PécidRL for petition correction and expectation identification based on deep reinforcement learning to correct noisy labels in petition datasets and build multiple petition expectation identification models to obtain better petition expectation identification accuracy. The main research work and contributions of this paper are as follows:

- (1) This paper extracts multi-view petition short text features, including

word-level and document-level semantic features, short text structure features based on graph representation learning, sentiment analysis features, and fusion features integrating semantic features and structure features.

(2) The reinforcement learning module is introduced to address the noisy label problem correcting the wrong and fuzzy labels. The experimental results show that the reinforcement learning module for label correction proposed in this paper has excellent noise correction performance with an average improvement of each experimental metric of 8.3% and the highest improvement rate reaching 14.2%.

(3) By fusing multi-view text features with traditional machine learning models, classical deep learning models, and graph neural network, 19 petition expectation identification models are extended and the classification performance before and after label correction is compared. The results of the experiments show that better classification performance are obtained based on the label-corrected data. Among them, the Peti-SVM-bert has the highest identification accuracy of 93.66% , which is decided as the final petition expectation identification model.

(4) An online web-server is developed for PeciRL with source code and dataset used in this paper aiming to maximize the convenience of government staff as well as related researchers to use the tool online.

Keywords:

petition expectation, multi-view text features, reinforcement learning, text graph representation, correction and identification

目 录

第 1 章 绪论	1
1.1 研究背景及意义	1
1.2 研究现状	2
1.2.1 文本分类	3
1.2.2 在线请愿分析	5
1.2.3 处理噪声标签	6
1.3 政府留言中期望识别问题面临的挑战	8
1.4 主要研究内容	9
1.5 论文结构	9
第 2 章 基于强化学习的留言标签更正算法	11
2.1 PccidRL 算法的整体介绍	11
2.2 数据集的构建	13
2.3 多视角文本特征的构建	15
2.4 基于强化学习的留言标签更正算法	18
2.4.1 标签选择器	19
2.4.2 标签判别器	23
2.5 本章小结	23
第 3 章 基于更正标签的留言期望识别模型	25
3.1 拓展的留言期望识别模型	25
3.2 评估指标	28
3.3 实验结果及分析	29

3.3.1 实验环境介绍	29
3.3.2 实验结果	30
3.3.3 个案研究	34
3.3 本章小结	37
第 4 章 PucidRL 在线服务平台开发	39
第 5 章 总结与展望	42
参考文献	44
作者简介及在学期间取得的科研成果	50
致谢	52

第 1 章 绪论

1.1 研究背景及意义

随着互联网的大规模普及、大数据新技术的不断涌现，我国的政务服务也向着与互联网融合的方向发展转型，在致力于实现政务服务信息化的过程中，建设了许多在线的一体化服务引擎，旨在加强不同服务部门之间的数据互通，及时解决各类民生问题。其中，政务留言在线平台的建设使得民众信访的方式从原来单一的线下信访，转变为成本低廉且便捷的在线留言进行网络信访的方式。

网络信访是公民表达诉求、寻求政府帮助以保障其利益的一种公民参与机制。作为公众与政府部门之间的互动交流平台，政务留言平台在促进政府管理效率、解决公民实际需求、锁定和解决可能造成不良社会矛盾和社会影响的事务方面做出了巨大贡献。详细了解民众的意见与建议，及时掌握民众的诉求，是提高国家治理水平、提升人民生活满意度的基础与重点工作。随着政务留言在线平台的广泛使用，大量的信访留言数据呈爆炸式增长。

领导留言板¹是人民日报社人民网于 2006 年创建的网络政务留言平台，供民众在线表达诉求、反映问题、提出意见建议，并由政府相关部门的工作人员对留言内容进行回复与处理。在上传留言时，留言者可以选择“求助”、“投诉”、“咨询”、“建言”、“感谢”五种标签来表达留言者对于留言的期望。根据留言者手动标注的留言期望，政府工作人员可以选择不同的回复策略和响应水平来更有效的处理留言。然而，手动标注留言期望面临两个不可控因素：第一，留言者在手动选择标签时会出现误选和漏选的操作；第二，留言者在上传留言时倾向于放大留言的紧急程度以期望获得更高的关注。这些由误标、误选和漏选导致的错误标签，更有一些留言者会发布一些无意义的留言，这会导致政府的低效率和回复延迟。图 1.1 展示了来自领导留言板的一条留言实例。该条留言选择的标签为“建言”，然而从文字内容中可以明显看出，该条留言描述的诉求与留言者选择的标签并不一致，或许标注为“求助”或者“投诉”更为合适。

¹ <http://liuyan.people.com.cn/>



图 1.1 领导留言板中的留言实例

目前，关注信访留言的价值及影响的研究越来越多，然而针对信访留言中的期望预测，以及分析政府对于不同紧急程度的留言期望应采用的响应水平的研究相对少。有研究利用主题模型分析信访留言中的主题及词频^[1,2]，以求分析民众提出意见的重点关注领域。然而，这些研究结果只能从宏观上反映公民的意见，政府管理者仍然需要大量精力分辨对于不同的留言期望应做出何种回应。在线的政务留言平台，例如领导留言板，在创建时为留言者提供了可以手动选择留言期望标签的机制。但是，对于留言期望的不当识别或者错误选择不可避免地会产生大量的噪声数据，这需要大量人力对其进行进一步检测，这无疑是在耗时费力且出错率高的。有研究表明，在极端情况下，如果事态紧急的留言诉求没有得到及时有效地回应，有可能会引发一系列社会问题，造成一定负面影响^[3]。因此，及时发现急需处理的紧急诉求对于留言者、政府以及社会而言都是至关重要的。

综上所述，手动对海量的留言文本数据进行处理与回复无疑是耗时耗力且低效的，然而，简单的自动识别模型无法对留言数据中大量的误标注噪声数据进行处理。因此，如何构建机器学习模型自动地更正误标注标签并对留言期望进行高效识别从而提高政府的工作效率，无论从理论还是实践应用方面来说，都是极具价值的。这使得政府可以快速识别留言紧急程度，进而快速确定政府管理者的不同响应水平，针对不同的响应级别从而采取不同的策略，并相应的安排及时且具有针对性的服务资源。

1.2 研究现状

政府留言中的期望识别问题是一个计算机与社会科学交叉领域下的研究问

题。在计算机领域，留言期望识别问题被视为机器学习中的文本分类问题；在社科领域内，留言期望识别问题可以从留言内容分析与留言影响探究两个方向着手。此外，本研究面对的重要挑战是如何在含有噪声标签的不规范的短文本数据集中学习到可靠信息。本小节将从上述三个不同研究目标的角度阐述留言期望识别问题的研究现状。

1.2.1 文本分类

文本分类问题一直是自然语言处理领域的热点话题，并被广泛应用于许多热门的应用领域，例如垃圾邮件检测^[4]、情感分析^[5]、新闻自动分类^[6]以及查询意图识别^[7]等。一般来说，文本分类任务可以分为两步来进行。第一步，从文档中提取文本特征，将文本以一种可以输入机器学习模型的有效数据形式表现出来。第二步，根据提取的特征构建一个合适的机器学习分类器^[8]。

对于第一步，即提取文本特征，较为常用的文本特征提取方法是独热编码。这一方法将一个词映射为一个一维向量，该向量只有一个元素是 1，其余都是 0，向量的维度即语料库中的词数。独热编码本质上是一个词袋模型，词袋模型使用一组无序的词来表示文档，不考虑文本中词的语法和顺序。词袋模型的主要思想是用文档中的单词出现的次数来表示文档，因为它只考虑了词的频率，而忽略了文档中上下文的关系，从而会失去文本的部分语义，因此这一模型有很大的局限性。TF-IDF 模型^[9]能够弥补词袋模型的上述弱点。TF-IDF 模型的核心是评估单词对文档的重要性。一个词的重要性随其在文档中出现的次数成比例增加，但同时又随其在语料库中的频率成反比减少。2013 年，Mikolov 等人提出了一种全新的词嵌入技术 word2vec^[10]，这是自然语言处理领域中的一项伟大成就。word2vec 模型本质上是一个浅层神经网络，它利用独热编码器将单词从高维空间映射到低维空间，拥有 Skip-gram 和 CBOW 两种形式。2018 年，谷歌发布了另一个优秀的语言表示预训练模型 BERT^[11]。BERT 使用多个双向 Transformer^[12]的 Encoder 作为主要框架，能够更深层地捕捉到语句中的双向关系。BERT 利用海量语料进行自监督训练，来学习优秀的语义特征作为词嵌入。基于 BERT 的词嵌入可以成功被应用于许多下游任务^[13,14]。与 word2vec 相比，基于 BERT 的词嵌入更专注于句子层面，其能挖掘到的语义信息也更加深入与丰富。BERT 具有强大的语

言表征和特征提取能力，成为近年来自然语言处理领域最具突破性的技术之一，目前，BERT 在 11 个不同的自然语言处理任务中达到了 SOTA 的表现。

对于第二步，即构建机器学习模型，目前有很多不同的机器学习模型被应用于文本分类。最初主要是传统的机器学习模型被应用于文本分类，例如支持向量机 (SVM)、逻辑回归 (LR)、随机森林 (RF) 和朴素贝叶斯 (NB) 等。随着大量深度学习模型的出现，利用深度学习方法解决文本分类问题成为一个新的研究热点。2015 年，Yoon Kim 提出了为文本分类设计的卷积神经网络 (CNN) 模型 TextCNN^[15]，它利用 n-grams 信息作为句子中的局部特征，通过 CNN 卷积来获得。在此之前，卷积神经网络多被用于处理图像问题，TextCNN 模型可以说是将卷积神经网络应用于文本问题的开山之作。在 TextCNN 的基础上，Yubo Chen 等人提出了一种动态多池卷积方法 DMCNN^[16]，将特征图从卷积层分割到池层，以提高子采样的效率。2017 年，Rie Johnson 和 Tong Zhang 提出了 DPCNN^[17]，通过增加网络深度与等长卷积和下采样来提取长距离的文本依赖关系。除了 CNN 之外，另一个经典的深度学习模型 RNN^[18]也被应用于文本分类问题，RNN 在处理序列数据方面有着天然的优势。AttBLSTM 模型^[19]在 RNN 的基础上引入了注意力机制，在 LSTM 层之后加入注意力层，将词级特征合并为句级特征。RCNN 模型^[20]结合 CNN 和 RNN，在 Bi-LSTM 的基础上获得上下文信息，用以进行文本分类。

近年来，图表示学习成为了新的研究热点，这为文本分提供了另一个全新的研究视角。将自然语言文本视为序列数据可以捕捉到更多的上下文语义信息，但文本同时包含着类似图的结构信息，将文本转化为图可以从新的角度来表示句法结构和语义信息。图的节点可以用来表示文本单位，边可以用来表示节点之间不同类型的关系，如词义、语义关系和上下文重叠等^[21]。TextGCN^[22]利用一个简单的两层图卷积网络进行文本分类，利用文档-词的关系和全局词共现信息，将整个语料库构建为一个图，以捕捉更深层次和更丰富的文本结构和语义信息。虽然 TextGCN 在文本分类问题上有着优秀的表现，但它很难在新样本上运行，而且全局图的构建不够灵活，内存消耗很大。为了解决上述的局限性，Lianzhe Huang 等人提出了一种用于文本分类的文本级图神经网络模型^[23]，为每个文档构建一个独立的图，并设计一个消息传递机制来在图神经网络中传达上下文信息。另一个

为每个文档构建独立图的文本分类模型是 TextING^[24], TextING 利用门控图神经网络来学习单词节点的嵌入, 节点通过来自其邻居和自身表示的信息来进行更新, 并汇总成文档的图级表示。Linmei Hu 等人提出了用于短文分类的异质图注意网络 HGAT^[25], 将文本转化为异质图结构, 在节点层和类型层使用两层注意机制, 以学习更丰富的语义与异构图结构信息。

1.2.2 在线请愿分析

研究人员认为, 网上请愿是政府部门和公民之间沟通的有效渠道之一。网上请愿平台旨在增加政府工作的透明度、提高政府的管理效率, 特别是政府机构可以对公民的不同诉求采取不同的回应策略。随着网上请愿平台的普及、用户的增多, 请愿信息也随之大量涌现。然而, 这些请愿信息大多是非结构化的文本数据, 同时带有不可控的误标、误选、误填等高噪声信息, 这不可避免地给请愿分析工作带来巨大的困难和挑战。

目前已经有大量针对请愿内容和请愿产生的影响进行分析的研究工作。请愿内容分析这一任务主要被认为是一个文本分类问题, 有许多经典的文本分类方法已经被应用于请愿识别任务。对于非结构化的文本数据, 主题模型可以有效地发现文本中的核心内容, 从而减少个体差异的影响。Loni Hagen 提出了一个基于 LDA 主题模型^[26]从请愿内容中提取出现的主题, 通过训练决定最佳的 k 个主题后, 这些主题进一步通过人工分析或计算机进行评估和分析^[27]。Narang Kim 和 Soongoo Hong 结合监督学习和非监督学习, 提出了一个请愿自动识别模型, 该模型在 LDA 提取的主题基础上对请愿数据进行聚类, 然后将带有新类别的数据送入谷歌预训练的 CNN 分类器, 以实现最佳的请愿识别效果^[1]。类似的, Woo Yun Hui 和 Kim Hyon Hee 将 LDA 与 K-means 聚类方法结合用以确定请愿数据中的关键性主题, 并提取多种不同的文本特征输入 LSTM 模型进行最后的请愿预测^[28]。

针对请愿影响的分析, 目前已有许多计算模型被提出来衡量请愿在政府政策制定中造成影响的能力。对于政府来说, 了解公民的关注和期望, 对实现改善民生环境这一目标非常重要^[29]。Logistic 回归作为经典的传统机器学习方法, 被应用于研究政府对请愿内容的反应中话语策略的影响^[30]。情感分析也被用来捕捉请愿中所表达的主观情感, 预测对请愿期望的情感倾向^[31]。除了机器学习方法外,

许多数学模型也被应用于请愿影响的分析。争议性度量被用来探索请愿数据中潜在主题的相关性，旨在研究争议性问题对社会经济指标的潜在影响^[32]。此外，基于大数据分析的方法也被应用于请愿影响，利用常用的在线搜索引擎百度和谷歌构建不同类型的请愿流行指数，并结合格兰杰因果分析法进一步探讨对于请愿信息中的风险预知在经济、环境、公共生活等方面的影响^[33]。

作为文本内容分析的有效工具，LDA 主题模型也常被用于辅助请愿影响的分析研究。LDA 被用来挖掘领导人政策性公开演讲及公众对演讲的回复等文本内容的主题，用以探寻公众及领导人对经济、政治和社会等方面政策的支持程度^[31]，或预见经济形势以制定相应的政策^[32]。此外，Logistic 回归算法被用于预测政府针对请愿留言不同的回复策略而产生的不同影响^[30]。Junxiang Wang 等人提出了一个有不确定度预估（Uncertainty Estimation）单元的链式多任务学习框架来识别请愿数据集中有潜力成功的请愿内容^[34]。当一条请愿收集到足够多的签名时，它可以引起对该问题负责的决策者的注意，而请愿的成功意味着相应决策者决定采取行动来解决请愿中的问题，与上述研究不同的是，该研究的数据是时序的，其学习过程是一个动态的增量过程。

1.2.3 处理噪声标签

标签具有噪声是留言识别模型构建中面临的一个重要挑战，这同时也是机器学习领域近年来一个研究热点。为了探究深度神经网络（DNNs）泛化能力的影响，Chiyuan Zhang 等人进行了系统的对比实验，对数据标签进行了不同程度的污染，包括干扰部分正确标签以及将所有标签替换为随机类别的标签^[43]，实验结果证明，无论何种程度的污染，深度神经网络都很容易拟合噪声数据。这一实验结果足以说明噪声数据对机器学习模型的负面影响之大。训练数据并不总是与真实样本分布一致（ground-truth）的情况被归类为弱监督学习中的不精确监督（Inaccurate Supervision）^[44]。为了防止模型对噪声过拟合并提升模型的泛化能力，模型中常常加入正则化模块，常见的正则化模块包括数据增强（data augmentation）^[45]、权重衰减（weight decay）^[46]、dropout^[47] 和批标准化（batch normalization）^[48]。然而，实验表明，即使上述正则化模块在模型中全部被激活，其在真实的数据集上训练的模型表现与在含有噪声的数据集上训练的模型表现

以上内容仅为本文档的试下载部分，为可阅读页数的一半内容。如要下载或阅读全文，请访问：<https://d.book118.com/068023132106006043>