

中国人民警察大学

毕业论文（设计）

题目 一种面向微博主题的  
自动画像系统

学	号	<u>2019160053</u>
姓	名	<u>魏志豪</u>
队	别	<u>智慧警务三队</u>
专业（方向）		<u>数据警务技术</u>
指 导 教 师		<u>安晓伟 讲师</u>

二〇二三 年五月



**An automatic portrait system oriented towards  
Weibo topics.**

**by  
WeiZhiHao**

**China People's Police University**

**Thesis for Bachelor's Degree**

**May, 2023**



## 摘 要

在当前的互联网时代，人们越来越多地使用社交媒体来进行信息交流和分享。微博作为国内最具有代表性的社交媒体之一，拥有庞大的用户群体和活跃度。微博目前是中国最受欢迎的社交媒体平台之一，随着微博平台的不断发展和更新，越来越多的人将其视为获取时事信息、社会热点和个人心情宣泄的重要渠道之一。同时，微博也成为了公众舆论和话语权的重要源泉，具有着极其广泛的社会影响力，成为了一个非常重要的信息传播平台。对于市场调研、舆情监测、公安决策部署等领域而言，了解微博用户认知、态度和互动等信息显得至关重要。然而，目前微博主题相关研究存在着信息收集的困难，数据清洗和分析的复杂性，以及结果可视化的不足等问题。为了解决这些问题，我设计和实现了一个面向微博主题的自动画像系统。该系统具有可视化、自动化和实时性等特点，可以有效地帮助公安机关掌握社会舆情和趋势，加强安全防范和管理，提高公众的法律意识和自我保护意识。当前，随着互联网技术的发展，网络媒体成为公安领域舆情分析的一个重要线索，网络舆情的传播速度和范围日益扩大。展盘、维稳、突发事件等需要及时监督和鉴定，适时掌握民意、了解社会关注和情绪波动，是良好公共安全和社会治安的基础。此系统可以有效地帮助公安机关掌握社会舆情和民意动向，加强安全防范和治安管理，提高社会公众的法律意识和自我保护意识，为公安工作提供有力支持和保障。同时，通过实时分析对当前公安事件和突发事件的舆情分析，及时发现和解决过程中的社会矛盾，保持社会稳定。

然而，在海量的微博数据中，如何快速准确地分析主题、情感和转发关系等信息，对于研究和理解微博传播规律以及探索公众舆论和心态变化趋势具有重要意义。因此，利用此系统，能够在一定程度上帮助我们更好地把握和解读微博数据的内在规律。同时，该系统具有可扩展性和灵活性，也为我们在日后的研究和应用中提供了便捷的数据分析和可视化支持。

本篇论文主要介绍了一种面向微博主题的自动画像系统。该系统基于 Python Flask 框架搭建，通过 requests 库构建的爬虫从互联网上自动收集用户输入的微博主题所需的数据，数据库（如 mysql）存储相关的数据集以便下一步的清洗和分析，利用 pymysql 库连接 Mysql 数据库，通过 JavaScript 的 ajax 与后端 Flask 连接，形成前后端数据库分离式设计。同时，该系统还使用了机器学习 sklearn 库对微博情感进行实时判断，以及利用 pyEchart 对爬虫爬到的微博主题相关微博和三级转发数据

进行可视化展示，包括 radar、sankey、graph 转发图等。该系统的主页是一个总览，展示了数据库现有条数、男女发博比例、当日热搜等信息，通过大屏展示和可视化图表直观地呈现微博主题和转发数据，以及公安工作中的应用和意义。系统还提供了菜单，包括主题事件分区和主题转发分区，用户可以添加新事件、现有爬取事件、已上传事件，添加转发数据、现有转发数据、以及数据上传分区，包括上传已有 Csv 文件和微博 API 接口，并且使用了 Flask 内置方法对用户进行登录验证，保证了数据安全性和隐私保护。

在公安领域的实际应用中，面临着大量网络信息，本系统通过分析大量网络信息，准确把握民意、把握热点，第一时间对事件进行预警和反应，提高了公安部门的信息化建设水平和网络舆情应对能力。本系统可以根据公安监测人员输入的关键词，在互联网上自动收集与该关键词相关的社交媒体信息，并通过自然语言处理等技术进行数据清洗和分析。利用系统集成数据可视化工具，生成数据大屏以便进行更直观的高效、快速响应社会热点事件与突发事件，同时能够准确处理民意情绪，提升部门的信息化水平和网络舆情应对能力，更好地服务社会稳定与公共安全。

**关键词：**自动画像系统；微博主题；数据分析；机器学习；爬虫；数据可视化；大屏展示；Flask；Python；MySQL；API 接口；登录验证；公安工作

## Abstract

In the current Internet era, people are increasingly using social media for information exchange and sharing. As one of the most representative social media platforms in China, Weibo has a huge user base and high activity. Weibo is currently one of the most popular social media platforms in China, and with the continuous development and updates of the platform, more and more people regard it as an important channel to obtain current events information, social hotspots, and personal emotional venting. At the same time, Weibo has become an important source of public opinion and discourse power, with extremely wide social influence, and it has become a very important information dissemination platform. Understanding Weibo users' perceptions, attitudes, and interactions is essential for market research, public opinion monitoring, and public security decision-making. However, there are difficulties in collecting information, complexity in data cleaning and analysis, and insufficient visualization of results in current Weibo related research. To solve these problems, I designed and implemented an automatic portrait system based on Weibo themes. This system features visualization, automation, and real-time capabilities, which can effectively help public security agencies grasp social public opinion and trends, strengthen security precautions and management, and improve the legal awareness and self-protection awareness of the general public. Currently, with the development of Internet technology, network media has become an important clue for public security field's public opinion analysis. The speed and scope of network public opinion dissemination are expanding rapidly. Supervision and identification of exhibition plates, stability maintenance, and emergencies require timely supervision and identification, timely understanding of societal attention and emotional fluctuations, which is the foundation of good public safety and social security. This system can effectively help public security agencies grasp social public opinion and public sentiment trends, strengthen security precautions and management, improve the legal awareness and self-protection awareness of the general public, and provide powerful support and guarantees for public security work. At the same time, by analyzing current public security events and sudden public events' public opinions in real-time, it can promptly discover and resolve social contradictions during the process

and maintain social stability. However, how to quickly and accurately analyze information such as topics, emotions, and forwarding relationships in massive Weibo data is of great significance for researching and understanding Weibo dissemination rules and exploring changes in public opinion and mentality trends. Therefore, using this system can help us better grasp and interpret the intrinsic rules of Weibo data to a certain extent. At the same time, the system has scalability and flexibility, which also provides convenient data analysis and visualization support for our future research and application. This paper mainly introduces an automatic portrait system based on Weibo themes. The system is built on the Python Flask framework, and the crawler constructed by the requests library automatically collects the data required for Weibo themes entered by users from the Internet. The database (such as MySQL) stores relevant datasets for the next cleaning and analysis step. It uses the pymysql library to connect to the Mysql database and connects with the backend Flask through JavaScript's AJAX, forming a frontend/backend database separation design. At the same time, the system also uses the machine learning sklearn library to judge Weibo emotions in real-time and uses pyEchart to visually display Weibo theme-related Weibo and third-level forwarding data crawled by the crawler, including radar, sankey, graph forwarding charts, etc. The homepage of the system is an overview that displays current database records, male-to-female post ratio, daily hot search, and other information. It intuitively presents Weibo themes and forwarding data and the application and significance of public security work through large-screen displays and visualized charts. The system also provides a menu, including theme event partition and theme forwarding partition. Users can add new events, existing crawled events, and uploaded events, add forwarding data, existing forwarding data, and data upload partitions, including uploading existing CSV files and Weibo API interfaces. The Flask's built-in method is also used to authenticate users for data security and privacy protection. In the actual application of public security, it faces a large amount of network information. This system can help public security agencies accurately grasp public opinions and hotspots, provide early warning and response to events in real-time, and improve the department's information construction level and network public opinion response capabilities. This system can automatically collect social media information related to the keyword entered by public security monitoring personnel on the



Internet and conduct data cleaning and analysis through natural language processing technology. Using the integrated data visualization tool of the system, generate a data big-screen for more intuitive, efficient, and fast response to social hot events and emergencies, and accurately deal with public sentiment emotions, enhance the department's informationization level and network public opinion response capabilities, and better serve social stability and public safety.

**Keywords :** Automatic portrait system; Weibo theme; Data analysis; Machine learning; Reptiles; Data visualization; Large screen display; Flask; Python; MySQL; API interface; Login verification; Public security work



# 目 录

摘要 .....	I
Abstract .....	II
1 绪论 .....	1
1.1 研究背景与意义 .....	1
1.1.1 研究背景 .....	1
1.1.2 研究意义 .....	2
1.2 国内外研究现状 .....	3
1.2.1 国外研究现状 .....	3
1.2.2 国内研究现状: .....	4
1.3 文献综评 .....	5
1.4 本文研究内容与组织架构 .....	6
1.4.1 研究内容与创新点 .....	6
1.4.2 本文组织架构 .....	7
1.4.3 本文数据结构及技术路径 .....	8
2 相关技术及理论概述 .....	9
2.1 微博主题画像系统的相关技术和方法 .....	9
2.1.1 数据爬虫技术 .....	9
2.1.2 数据清洗与预处理 .....	10
2.1.3 可视化技术 .....	10
2.1.3 社交网络分析技术 .....	12
2.1.4 机器学习技术 .....	14
3 主题画像系统设计和实现 .....	15
3.1 需求分析与设计原则 .....	15
3.1.1 需求分析 .....	15
3.1.2 系统设计原则 .....	15
3.2 系统的架构设计 .....	16
3.3 系统模块设计 .....	17
3.4 微博画像系统的实现 .....	18
3.4.1 爬虫模块 .....	18
3.4.2 CSV 数据上传模块 .....	18
3.4.3 机器学习模型训练和部署模块 .....	19
3.4.4 前端展示模块 .....	19
3.4.5 其他模块 .....	19
3.5 数据流程 .....	19

4	系统功能测试.....	21
4.1	爬虫功能测试.....	21
4.2	CSV 数据上传功能测试.....	21
4.3	机器学习模型测试.....	22
4.4	前端展示功能测试.....	23
4.5	用户权限测试.....	23
4.6	性能测试.....	24
5	总结与展望.....	25
5.1	总结.....	25
5.2	展望.....	25
	参考文献.....	27
	致 谢.....	29

# 1 绪论

## 1.1 研究背景与意义

### 1.1.1 研究背景

随着互联网的不断发展，社交网络已经成为人们生活中不可或缺的一部分，其中微博是国内使用较为广泛的社交媒体平台之一。微博平台每天都会有大量的用户在上面发布各种信息，包括新闻、娱乐、体育、科技、政治等各种话题，用户通过发表微博来表达自己的观点和看法，并且可以通过转发、点赞等方式扩散自己的内容，进而形成各种社交互动。

在这样的背景下，对微博的分析和研究显得尤为重要，尤其是对于公安机关来说，通过对微博内容的分析可以帮助他们更好地了解社会舆情、研判事件动态、预测未来趋势、采取有针对性的措施等。同时，微博中的一些事件也经常与公安机关相关，例如刑事案件、公共安全事件等，因此对微博中的这些事件进行深入分析也有助于公安机关更好地维护社会治安和稳定。然而，由于微博平台的数据量庞大、内容繁杂，如何对这些数据进行高效、准确、系统化的分析一直是一个难题。传统的分析方法主要是基于人工分析，例如阅读微博、标注关键词、分类统计等，但这种方法存在人力资源浪费、效率低下、精度不高等问题。因此，开发一种自动化、智能化的微博分析系统已经成为当前亟待解决的问题之一。

为了解决这个问题，本研究提出了一种面向微博主题的自动画像系统。该系统主要基于 Python 的 Flask 框架，并且连接 Mysql 数据库，通过爬虫程序获取了相关数据，并将其存储在 MySQL 数据库中，结合前端的 JavaScript 和 Echarts 图表库，实现了微博数据的自动化爬取、处理、分析和可视化。同时，该系统还利用机器学习的方法对微博内容进行情感分析，并将分析结果呈现在图表中，帮助用户更好地了解微博中的舆情和事件。

该系统的优点在于可以大大提高数据分析的效率和精度，同时减少人工干预，降低成本。对于公安机关来说，该系统能够更好地帮助他们了解社会舆情、研判事件动态、预测未来趋势、采取有针对性的措施等，进而更好地维护网络虚拟社会治安。在未来的研究中，可以考虑进一步完善算法模型和优化数据可视化效果，以提高系统的性能和实用性。

### 1.1.2 研究意义

本研究的意义在于提供了一种面向微博主题的自动画像系统，该系统可以帮助公安机关进行舆情监测、主题分析和情报研判等工作。具体来说，本系统可以实现以下几个方面的意义：

1、本系统可以帮助公安机关更加全面地了解社会上的各种主题和事件，帮助他们做好舆情监测和研判工作。在现代社会中，各种事件和主题层出不穷，公安机关需要时刻关注社会上的各种动态，从而做好应对措施。本系统可以通过爬取微博数据、分析微博主题、制作转发图等方式，帮助公安机关及时了解社会上的各种主题和事件，为他们制定相应的工作计划提供重要参考。

2、本系统可以帮助公安机关更加准确地判断社会舆情的走向和倾向，从而制定更加科学的应对策略。在现代社会中，舆情监测和研判工作非常重要，这可以帮助公安机关更好地了解社会上的各种动态，预测未来的发展趋势。本系统可以通过机器学习算法对微博数据进行情感分析，进而判断社会舆情的走向和倾向，从而制定更加科学的应对策略。

3、本系统可以帮助公安机关更加高效地分析和研判大数据，提高工作效率和精度。在大数据时代，分析和研判大数据已成为各行各业的必修课。对于公安机关来说，分析和研判大数据是其日常工作的重要组成部分。本系统可以通过快速爬取微博数据、实时分析情感倾向、制作可视化图表等方式，帮助公安机关更加高效地分析和研判大数据，提高工作效率和精度。而微博已经成为了公众对警方工作质量的重要评价标准之一。通过对微博上的关键词搜索和分析，公安部门可以了解公众普遍的安全关切和热点问题，并及时采取有效措施，提高公共安全。但是，如果是人工分析这些微博，耗时费力，效率低下，可能会导致延迟反应或失误。因此，基于微博主题的自动画像系统，为公安机关提供了一种自动化的数据分析工具，利用自然语言处理技术和机器学习算法，快速筛选和分析有用的微博信息，提供更准确、及时、权威的数据分析和情报支持，有着非常广阔的应用前景。

而基于微博主题的自动画像系统能够协助公安机关了解公众对公安事业的态度和评价。警民互动是现代公安工作的重要组成部分，公安机关应该关注全社会对警察工作的看法和意见，不断提高警民互动的质量和水平。同时，该系统可以自动将有用的微博信息分类、整理和归档，便于以后的数据应用，为公安部门提供新的案件线索和研判基础。画像系统可以加强公安机关与全社会的信息联结，推动公安机关适应信息

化环境的发展，并适应时代的变化，为提高公安机关的服务质量和效率提供支持和保障。

所以，研究此系统的意义在于其在公安工作中具有非常重要的应用前景。这个系统可以大幅提高公安机关对于舆情和公众关切的反应能力和分析准确度，实现对于各种威胁和公共卫生事件的及时应对和预防控制，也可以加强公安宣传工作的效果和互动性等。未来，还需要不断研究微博舆情分析的更深层次问题，以更好地为公安机关提供更好的舆情分析服务和支撑。

## 1.2 国内外研究现状

随着微博等社交媒体的普及和快速发展，研究微博舆情分析的领域也在逐步增加。国内外的专家学者和工程师们已经开展了大量的微博舆情分析的研究和实践，相关技术和算法已经初步形成并得到了广泛应用。本节主要介绍国内外微博舆情分析的研究现状。随着微博等社交媒体的普及，面向微博主题的自动画像系统的研究也逐渐受到了学术界和工业界的关注。

### 1.2.1 国外研究现状

与国内相比，国外关于微博话题分析领域的研究更具深度和广度。在文本情感分析方面，国外学者不仅采用传统的情感词典、n-gram 模型等方法，还使用基于神经网络的方法，如 CNN、LSTM 等。在话题挖掘方面，基于深度学习的方法也被广泛应用，如主题模型 LDA 能够发现微博话题的隐藏主题。除此之外，国外学者还提出了许多针对较难处理的问题的解决方案，如处理带有噪声的数据、挖掘隐藏的社交关系等。同时，对于微博的四类用户（普通用户、机器人、机器人伪装成普通用户和宣传机构用户）的识别和分类等问题，国外学者也进行了大量研究。

国外关于微博话题分析的研究非常广泛，包括但不限于情感分析、话题挖掘、用户画像等领域。研究者们使用了很多种不同的数据挖掘和机器学习方法，如基于神经网络、LSTM、注意力机制等。同时，国外学者研究的范围也非常广泛，除了社交媒体大数据的分析，还包括公共政策决策、社会公共问题和舆情信息等领域。这些研究为微博话题分析领域的发展提供了很多的创新思路和方法，同时还实现了很好的效果和应用。

国外所用的 PERSONA 是 Allen Cooper 提出来的一种通过调研和问卷获得的典型用户模型，用于产品需求挖掘与交互设计的方法。Persona 是虚构出的一个用户用来代表一个用户群。一个 persona 可以比任何一个真实的个体都更有代表性。一个代表

典型用户的 persona 的资料有性别、年纪、收入、地域、情感、所有浏览过的 URL、以及这些 URL 包含的内容、关键词等等。一个产品通常会设计 3~6 个用户模型代表所有的用户群体。而另一个单词，叫 Profile，是利用已经获得的数据，用来勾勒用户需求、用户偏好的数据分析方法。

这两个词，都可以翻译为用户画像，但第一种，用于产品用研与交互设计，第二种，用于运营与数据分析。与国内相同的是，国内外都将其用于社会媒体以及用户的个性化推荐上，可以发现，在这个互联网高速发达的时代，不管是国内还是国外，都需要这样一个画像系统来监控网络上的风险。

### 1.2.2 国内研究现状：

在国内，随着微博等社交媒体平台的发展，微博话题分析也逐渐成为了研究的热点。目前，国内的研究工作主要涉及到微博数据的收集、预处理、情感分析、话题挖掘、用户画像等领域。在微博数据收集方面，目前国内学者主要运用微博 API 和网络爬虫等技术进行数据抓取。在预处理阶段，常用的技术包括去除噪声、分词、去除停用词、词性标注等。在情感分析方面，国内学者主要采用词典、机器学习等技术进行情感倾向分析。在话题挖掘方面，主要采用 LDA、Word2Vec 等技术进行话题挖掘。同时，国内学者也开始关注用户的画像分析，提取用户的性别、年龄、地域、兴趣爱好等信息，对用户进行精细化管理和针对性的推广。

微博作为中国最具影响力的社交媒体平台之一，其上产生的庞大信息汇总成为了一个庞大的数据集群，有着广阔的应用前景和潜力。这其中涵盖了信息分类与过滤、面向事件的报道监测和情感分析等方面。尤其是在政治和经济层面，微博越来越被用于分析公共舆论、政策影响和民众信任度等重要指标，已经成为研究的重要领域。

随着微博的日益普及和深入，各种类型的微博也越来越多，这就使得微博话题分析愈发重要。微博话题分析（Weibo topic analysis）又被称为微博话题挖掘或微博热点追踪。微博话题分析指的是根据某一个特定话题或事件，对一段时间内产生的一系列微博文本进行处理、分析、归纳和总结。目的在于获取该特定话题或事件的相关信息与意见，在总体上了解和把握民意公共反馈和市场态势。

微博话题分析的应用非常广泛，其主要应用领域如下：

（1）多样性的应用——餐饮、体育、电子商务。

无论是在生活、娱乐、教育还是交友等方面，人们都可以通过微博得到相当多的信息。由此，通过微博话题分析，用户可以了解到市场的一些重要信息，帮助他们更



好的升级服务、增加产品线，进而受到更多人群的青睐。以餐饮行业为例，餐饮店可通过分析用户发表的微博找到受欢迎的菜肴和热门口味，从而推广其特定餐品并提升品牌。

### （2）政府、企业和组织的营销和管理。

使用微博话题分析可以了解和掌握公众对政策和措施的反馈和认同度，帮助政府制定更加合理有效的政策和举措。社交媒体大数据分析能够帮助企业关注用户的舆情、竞争情况和反馈等重要指标。同时，组织可以通过微博话题分析发现其团队中的优质演说者，从而培养更加有效的沟通策略。

### （3）易于接受群众意见。

微博话题一旦创造了更广泛、更具深度的讨论，对于接受群众意见也起到了十分积极的推动作用。从而，如果需要进行更广泛的开放式辩论，民意和意见可以很快地在网络上截获和沟通。

由此可见，微博话题分析在我们的日常生活和企业、政府的各个方面发挥着重要作用。然而，与此同时，也存在微博数据的异构性、大量的数据处理等问题，这必须借助于各种人工智能和数据挖掘技术进行处理。为了更好地发挥微博话题分析的作用，必须进一步解决这些问题，同时还需要相应方法和技术的支持。

## 1.3 文献综评

近年来，微博话题分析领域的研究呈现出不断增长的趋势。研究者们尝试使用各种数据挖掘和机器学习等手段，分析微博大数据背后的舆情信息和用户意见，以及其在政治和经济层面的应用。以下是文献综述的三个方面的：

### （1）微博情感分析

情感分析是微博话题分析领域的一个重要研究方向。主要是使用自然语言处理技术，将微博文本中传递的情感进行分类，以寻找相关话题或事件中一些情感状态的变化。国外学者使用了在 LSTM 模型中使用的反向传播算法，开发了一种情感分类模型，有效地区分了微博中的消极和积极情感。另一项研究指出，采用注意力机制的神经网络可以提高情感分析的精确度。此外，还有一些研究者采用分层聚类、序列规则挖掘等方法，对微博话题进行情感分析的研究。

### （2）微博话题挖掘

微博话题挖掘是对微博大数据中话题的识别、分类、跟踪和分析等一系列过程。其中，识别和分类阶段是最基础的，它们包括将微博分成一组分拨，确定话题区分的

重要场域以及识别可能的新话题。国内研究人员提出了一种新的有监督训练方法，该方法可以在不用运算过程，并使用二维特征表示微博文本。也有学者提出了一种基于 LSTM 的话题挖掘方法，该方法结合词向量和独特的分类算法，可以识别出更多的潜在话题。此外，还有其他一些基于机器学习的方法，如 SVM、朴素贝叶斯、CRF 等等。

### （3）微博用户画像

微博用户画像是一种基于大数据分析的构建方法，旨在了解微博用户的基本属性、习惯、兴趣等。通过微博用户画像，可以更好地了解用户的个性，从而定向、精准地提供相应资源和服务。国内有学者使用社会网络分析、微博文本挖掘等方法，分析了微博崛起以来的用户特征、兴趣和活动等，并尝试使用 LDA 主题模型进行微博用户分类，其研究结果表明，将用户分类为不同的类别可以更有效地定向服务。此外，还有其他一些用于进行微博用户画像的算法，如 SVM、聚类等。

总的来说，微博话题分析领域的研究呈现出不断发展的趋势，包含情感分析、话题挖掘和用户画像等方面。这些研究为行业发展提供了很多重要的思路和方法，用于更好地理解 and 利用社交媒体大数据，进而推动各个领域的进步和发展。

## 1.4 本文研究内容与组织架构

### 1.4.1 研究内容与创新点

#### （1）主要研究内容

本文旨在设计一种面向微博主题的自动画像系统，通过爬虫获取微博数据并进行情感分析和主题分类，再通过可视化技术呈现出来，提供给用户进行数据分析和研究的工具。具体来说，本文的研究包括三个方面：第一部分是微博数据的爬取和清洗，包括利用 Python 语言编写爬虫程序，获取微博数据集，并进行数据清洗和预处理；第二部分是微博数据的情感分析和主题分类，利用机器学习算法和自然语言处理技术，对微博进行情感分析和主题分类；第三部分是可视化展示，通过图表和可视化技术将微博数据进行可视化呈现，为用户提供数据分析和研究的工具。

#### （2）创新点

a)面向微博主题的自动画像系统：本系统针对微博这一特殊的社交媒体平台，采用了一种面向主题的自动画像系统，可以实时爬取微博数据，并对其进行情感分析、主题挖掘和转发路径分析等操作，从而达到对微博用户的自动画像。

b)数据库分离式设计：本系统采用前后端数据库分离式设计，即前端通过 JavaScript 的 ajax 与后端 Flask 连接，后端通过 pymysql 库连接 Mysql 的数据库，从而实现了前后端数据的高效交互和管理。

c)可视化展示与大屏展示：本系统采用了 pychart 和 echart 等可视化工具，能够对微博数据进行静态图表展示和动态更新的展示，同时还设计了大屏展示功能，能够通过上传大型数据进行数据分析的大屏展示。

d)登录验证和权限管理：本系统还加入了登录验证和权限管理功能，通过 Flask 内置的 is\_logged\_in 方法，实现了用户登录的验证和权限管理，从而保证了系统数据的安全性和可靠性。

e)与公安工作的结合：本设计是一个针对微博数据的自动画像系统，能够实时分析微博用户的情感和行为，对于公安工作中的舆情监测和犯罪侦查等方面具有一定的应用价值。

#### 1.4.2 本文组织架构

本文主要分为六个章节：

第一章为绪论，主要介绍本文的研究背景、研究意义、研究现状和本文的研究内容及组织架构。

第二章为相关理论和技术综述，主要介绍了本文所用到的相关理论和技术，包括自然语言处理、数据挖掘、机器学习等方面的理论基础，以及 Python 语言、Flask 框架、pymysql 库、pyecharts、sklearn 等技术的应用。

第三章为总体方案设计，包括系统需求分析、系统设计、系统架构设计和数据库设计等方面的内容。

第四章为主题画像，主要介绍了如何通过爬取微博数据、数据清洗、分词、关键词提取、情感分析等方法对微博主题进行画像。其中，主题分析包括微博主题分类和主题词提取两个方面，情感分析主要用于对微博内容的情感进行分析和评估。同时，本章还介绍了如何使用 pyecharts 绘制静态和动态的可视化图表，以便更加清晰地呈现主题画像。

第五章为主题画像对于舆情判断的探索，主要探讨了通过对微博主题进行画像对于舆情判断的作用。通过对实际微博数据的分析和研究，本章得出了一些结论，包括主题画像可以帮助用户更快地了解当前热点话题的趋势和方向，可以帮助用户更好地把握微博上的公众情绪，并可以为公安部门提供决策参考。

第六章为系统功能测试，主要对本文所设计的系统进行了测试和评估。测试内容包括系统界面的友好性、系统功能的稳定性和准确性等方面，通过测试评估可以判断本文所设计的系统是否能够满足用户的需求和要求。

最后一章为总结和展望，对本文的研究内容和实现效果进行了总结，并对未来的研究方向和应用前景进行了展望。本文旨在为公安部门提供一种基于微博数据的舆情监测分析工具，以便更好地应对公共安全事件。

1.4.3 本文数据结构及技术路径

本文的数据结构主要包括微博数据表、主题数据表、转发数据表 and 用户数据表。其中微博数据表用于存储所有爬取到的微博信息，包括微博 ID、微博内容、发布时间等；主题数据表用于存储通过主题模型得到的主题信息，包括主题 ID、主题关键词、主题权重等；转发数据表用于存储所有微博的转发信息，包括转发微博的原始 ID、转发微博的内容、转发微博的发布时间等；用户数据表用于存储微博用户的基本信息，包括用户 ID、昵称、性别、所在地等。

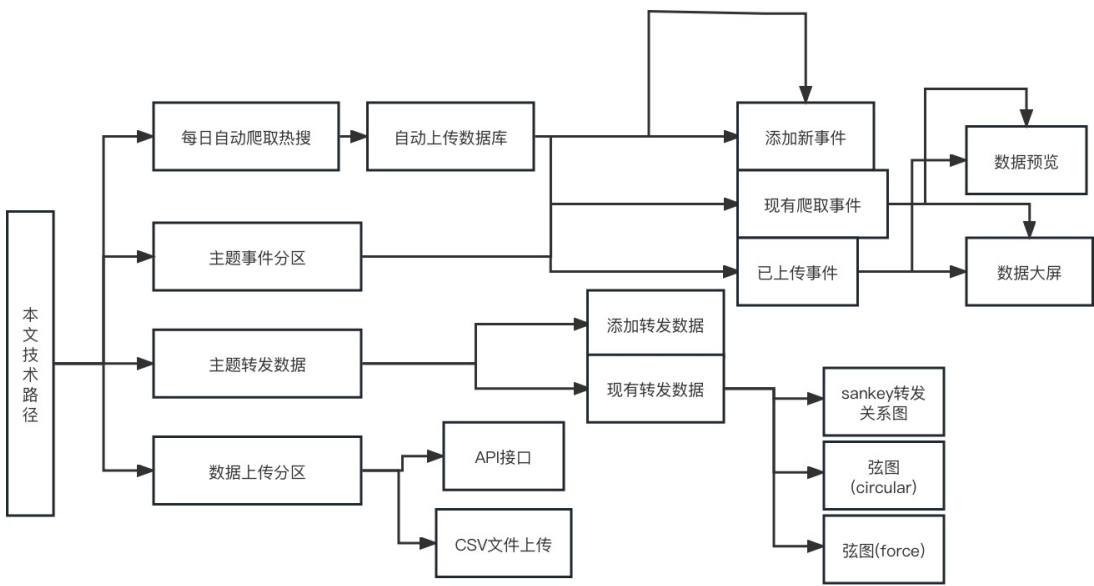


图 1-1 系统实现流程

## 2 相关技术及理论概述

### 2.1 微博主题画像系统的相关技术和方法

本文介绍了一种面向微博主题的自动画像系统，该系统基于 Python 的 Flask 框架开发，并使用 MySQL 数据库作为后端存储，通过 JavaScript 实现前端交互。该系统可以自动分析微博话题的热度、关注度、情感倾向和转发行为，从而对微博用户的兴趣和行为进行画像。

本系统是一种基于 Python 和 Flask 框架的微博主题画像系统，通过爬虫技术、数据清洗和预处理、机器学习和可视化技术，实现对微博数据的分析和展示。本文主要研究微博大数据的应用，涉及到多种技术和理论，下面对其进行简要概述。

#### 2.1.1 数据爬虫技术

在微博主题画像系统中，爬虫技术是获取数据的重要手段之一。爬虫技术可以帮助我们自动化地获取海量的微博数据，并且可以根据需要提取其中的关键信息，从而为后续的分析提供支持。

在实现微博爬虫时，我们可以使用 Python 中的 requests 库来向微博服务器发送 HTTP 请求，并使用 BeautifulSoup 库或者正则表达式来解析 HTML 文档，从而获取所需的信息。具体地，我们可以使用 requests.get() 方法向微博服务器发送 GET 请求，并将返回的响应解析为 HTML 文档。对于响应中的 JSON 数据，我们可以使用 json.loads() 方法将其转换为 Python 对象，从而进行进一步的处理。在爬取微博数据时，我们需要模拟浏览器行为，并使用 cookie 信息来模拟用户登录状态，以获取更多的数据。为了防止被封 IP，我们可以使用代理 IP 池，并设置合理的访问间隔时间。

另外，在实现微博爬虫时，我们需要注意微博官方的反爬虫机制，如 IP 封锁、验证码等。为了规避这些限制，我们可以采用多个 IP、账号轮换、设置访问频率、使用混淆技术等方法。

爬虫技术是微博主题画像系统中不可或缺的一环，它可以帮助我们高效地获取大量的微博数据，并为后续的分析提供支持。然而，在爬取微博数据时需要注意尊重他人的隐私权，遵守相关的法律法规，并避免对微博服务器造成过大的负担。

### 2.1.2 数据清洗与预处理

数据清洗与预处理在微博主题画像系统中扮演着非常重要的角色，它们的目的是从爬取到的原始数据中提取出有用的信息，为后续的分析 and 可视化提供清晰、准确、可用的数据。

在微博主题画像系统中，我们需要对爬取到的微博数据进行一系列的数据清洗和预处理操作。首先，我们需要去除数据中的重复项，以避免对分析结果产生影响。其次，我们需要进行数据格式的标准化，统一不同数据源的数据格式。同时，还需要进行缺失值的处理，对于缺失值的数据可以使用填充值、均值等方法进行填充处理。

在进行数据清洗和预处理时，还需要注意微博数据的特点。由于微博的内容非常丰富，包含文字、图片、视频等多种形式，因此需要对不同形式的数据进行不同的处理。例如，对于文字内容，我们需要进行分词、去除停用词等自然语言处理操作，以方便后续的分析 and 可视化；对于图片和视频等多媒体数据，我们需要进行数据解析和特征提取等操作，以便于对其进行分析和展示，以及微博用户发表的 emoji 表情，为了以后便于情感分析，利用 Python 的 demoji 库将其转化为英文随同文字一同存储在数据库中，不遗漏每一条数据。

除此之外，还需要进行数据的归一化处理。由于微博数据来源的多样性，不同的数据源可能会有不同的数据规模和数据分布。为了将这些数据进行合理的比较和分析，需要将其进行归一化处理，以保证数据具有可比性和可解释性。通常，归一化处理方法包括最大最小值标准化、Z-Score 标准化等。

### 2.1.3 可视化技术

可视化技术在本系统中是非常关键的一环，因为它可以帮助用户更加直观地了解微博主题的相关信息。系统中采用了 PyEcharts 这个 Python 可视化库，它基于 Echarts 和 Apache ECharts 开源项目，提供了非常方便的 Python 接口。Echarts 是一个基于 JavaScript 的可视化库，它可以用于创建各种类型的图表，例如折线图、柱状图、散点图、地图等等。PyEcharts 利用 Python 来生成 Echarts 图表，并且支持多种类型的数据格式，包括 JSON，Python 的列表和 Pandas 数据框。

在本系统中，使用 PyEcharts 绘制了多种类型的图表，包括柱状图、饼状图、折线图、散点图、热力图、地图、Sankey 图等等。其中，柱状图主要用于展示微博主题相关的统计数据，例如微博数、点赞数、转发数、评论数等等。饼状图主要用于展示男女比例、用户所在地区的分布等等。折线图主要用于展示某一时间段内微博主题的热度趋势。散点图主要用于展示微博用户的位置分布。热力图主要用于展示微博主题的热度分布。地图主要用于展示微博用户所在地区的分布情况。Sankey 图主要用于展示微博主题的转发关系。所有这些图表都具有交互性，用户可以通过鼠标或触摸屏来对它们进行缩放、平移、查看详细信息等操作。



图 2-1 pyecharts 可视化图

在 PyEcharts 中，除了提供了基本的图表类型外，还提供了一些高级功能，例如图表的主题、动画效果、数据筛选、数据缩放、数据联动等等。这些功能使得图表更加美观、易于操作和具有可操作性。同时，PyEcharts 还提供了基于 Flask 的 Web 显示功能，这使得系统可以通过 Web 页面来展示图表和数据，方便用户使用和



图 2-2 微博主题大屏展示

操作。另外，设计的大屏展示是面向用户的重要界面，用于展示系统的数据分析结果和重要指标，通过大屏展示能够将复杂的数据信息以更加直观的方式呈现给用户，方便用户对数据进行分析 and 决策。

可视化技术是微博主题画像系统中非常重要的一项技术，它通过图表、地图等方式将微博主题相关的数据进行可视化展示，使用户更加直观、全面地了解微博主题的相关情况。

### 2.1.3 社交网络分析技术

社交网络分析主要是通过对社交网络中的节点和边的分析，来探讨社交网络的结构、特征和演化趋势，它可以帮助我们理解微博用户之间的关系和影响力，从而更好地进行用户画像和主题画像。

在微博主题画像系统中，社交网络分析通常会使用以下几种技术：

（1）节点中心性分析：节点中心性是指网络中一个节点在传递信息、影响其他节点等方面的重要程度。在微博主题画像系统中，节点可以指微博用户、微博主题等。通过计算节点的度中心性、介数中心性、紧密中心性等指标，可以评估节点在社交网络中的影响力和重要性，从而更好地理解用户和主题。

2）社区检测：社区是指在社交网络中由若干节点组成的紧密联系的子集。社区检测是指通过算法找到网络中的社区结构，并将节点分组以便进行更深入的分析。在微博主题画像系统中，社区检测可以帮助我们找到微博用户之间的联系和共同兴趣，进而进行更准确的用户和主题画像。

（3）影响力分析：影响力分析是指通过分析用户在社交网络中的传播行为，评估用户在社交网络中的影响力。在微博主题画像系统中，影响力分析可以帮助我们了解哪些用户在社交网络中有更高的传播力和影响力，从而更好地进行用户和主题画像。

针对转发数据，我们可以采用社交网络分析技术，使用 Sankey 和 graph 图来分析转发关系。

（1）Sankey 图，它是一种流程图，可以用于展示一系列事件或物体之间的流量和流向。在微博转发中，Sankey 图可以展示某一主题的微博在不同用户之间的流动情况，从而更好地理解这个主题在社交网络中的传播情况。例如，我们可以展示某一主题的微博从不同的来源（如某一用户）流向不同的目标（如其他用户或话题），以



及它们在不同用户之间的流动路径。这样可以帮助我们更好地分析这个主题在社交网络中的影响力和传播路径。

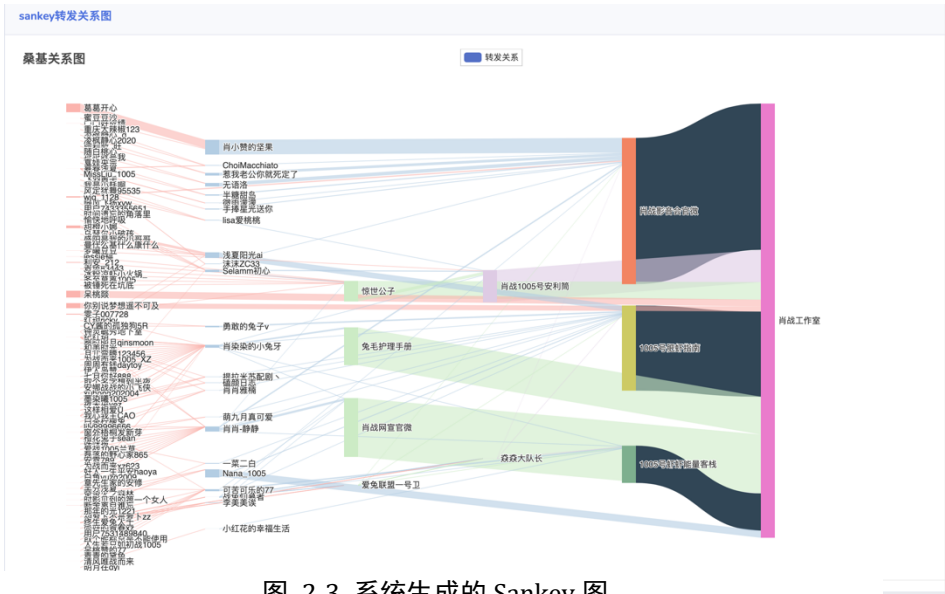


图 2-3 系统生成的 Sankey 图

(2) graph 图，它是一种节点和边组成的图形结构，可以用于展示多个节点之间的连接关系。在微博转发中，我们可以用 graph 图来展示某一个微博在不同用户之间的转发情况。例如，我们可以将一个微博看作一个节点，将不同用户之间的转发关系看作边连接两个节点。并且，我加入了动态 Timeline，这样可以帮助我们更好地

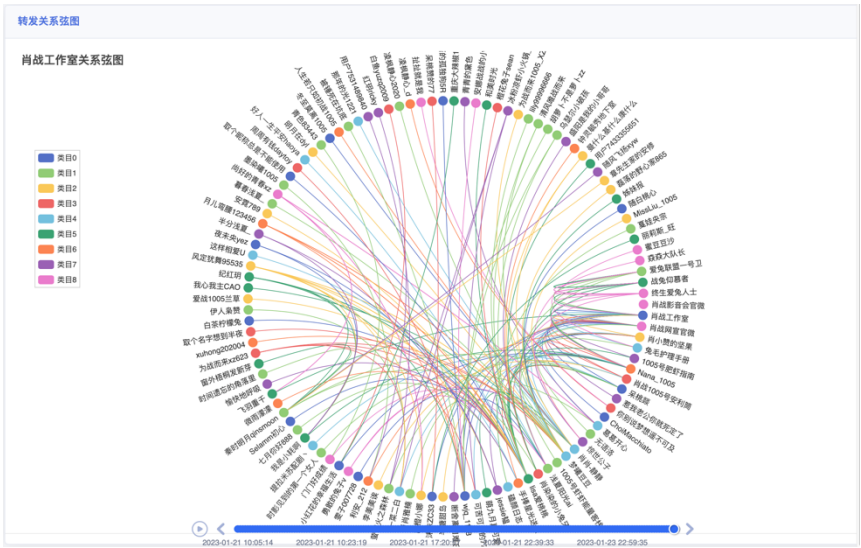


图 2-4 系统生成的 Graph 图（circular）

随着事件发展过程分析微博的传播路径和影响力，以及用户之间的社交网络结构。

通过分析 Sankey 和 graph 图，我们可以更好地理解微博在社交网络中的传播情况和用户之间的社交网络结构，从而更好地进行社交网络分析和用户画像。

#### 2.1.4 机器学习技术

在本系统中，机器学习被用于分析微博内容的情感极性。情感极性是指对一段文本的情感倾向，通常分为正面、负面和中性三种。为了对微博进行情感极性分析，我们使用了 Python 中的 Scikit-learn 库来训练一个情感分析模型。

首先，我们需要一个情感极性标注的数据集来训练模型。我们在网上收集了一批标注好的微博数据，共计 45,553 条。然后，我们对这些微博进行预处理，包括去除停用词、分词、去除无用符号等操作，以便于提取出微博中的情感词汇。

接着，我们使用 Scikit-learn 库中的 TfidfVectorizer 模块将处理好的文本转化为数值型特征。TfidfVectorizer 模块会统计每个词在所有微博中的出现频率，并计算其在当前微博中的权重。这样，我们就得到了一个词频矩阵，其中每一行代表一篇微博，每一列代表一个词汇。

最后，我们使用 Scikit-learn 库中的朴素贝叶斯分类器来训练情感分析模型。朴素贝叶斯是一种常见的文本分类算法，它通过统计每个词在不同类别中出现的频率，来计算一个新文本属于某个类别的概率。在本系统中，我们将情感分为喜悦、愤怒、中性等五类，使用二元分类的方式进行训练。

经过交叉验证，我们的模型在测试集上的准确率达到了 71%。在实际应用中，我们将该模型应用到新的微博数据上，实时判断其情感极性，并将结果保存到数据库中，以供后续的分析 and 可视化。

```
Building prefix dict from the default dictionary ...
Loading model from cache /var/folders/3f/hj3gtg8n3gz5xm84jn7ywct40000gn/T/jieba.cache
Loading model cost 0.683 seconds.
Prefix dict has been built successfully.
<timed exec>:70: FutureWarning: The behavior of `series[i:j]` with an integer-dtype index is deprecated. In a future version, this will be treated as *label-based* indexing, consistent with e.g. `series[i]` lookups. To retain the old behavior, use `series.iloc[i:j]`. To get the future behavior, use `series.loc[i:j]`.

0.7181733849342667

<timed exec>:82: FutureWarning: The behavior of `series[i:j]` with an integer-dtype index is deprecated. In a future version, this will be treated as *label-based* indexing, consistent with e.g. `series[i]` lookups. To retain the old behavior, use `series.iloc[i:j]`. To get the future behavior, use `series.loc[i:j]`.

0.6803786001072604
CPU times: user 37min 34s, sys: 15.3 s, total: 37min 49s
Wall time: 38min 1s
```

图 2-5 模型训练过程

### 3 主题画像系统设计和实现

#### 3.1 需求分析与设计原则

##### 3.1.1 需求分析

在进行系统设计和实现之前，需进行充分的需求分析，以确保系统可以满足用户的需求。对于微博主题画像系统来说，需求分析主要包括以下几个方面：

（1）用户需求分析：对于这个系统来说，主要的用户群体是公安工作人员。因此，系统的设计必须充分考虑到用户的使用场景、使用目的和使用习惯等方面的因素。同时，用户需要从系统中获得哪些信息，以及如何有效地展现这些信息，也需要考虑到这些因素。

（2）功能需求分析：对于微博主题画像系统来说，主要功能包括微博数据的爬取、清洗、预处理、分析和展示等方面。其中，爬取和清洗是系统的核心功能，需要充分考虑到数据来源、数据格式和数据质量等方面的因素。预处理和分析则是为了从数据中提取有效信息，以便于后续的展示和分析。展示功能则是将分析结果以可视化的形式展现出来，让用户可以更直观地理解数据。

（3）性能需求分析：对于微博主题画像系统来说，主要考虑的性能因素包括数据爬取的速度、数据处理的效率、系统的稳定性和可扩展性等方面。其中，数据爬取的速度需要尽可能地快，以确保系统可以及时地获得最新的数据。数据处理的效率则需要充分考虑到大数据处理的需求，以确保系统可以在合理的时间内完成数据处理。系统的稳定性则需要尽可能地高，以确保系统可以稳定运行。可扩展性则是为了能够在需要时对系统进行扩展，以满足未来的业务需求。

（4）安全性需求分析：对于微博主题画像系统来说，数据的安全性是非常重要的。因此，需要对系统进行安全性分析，包括用户身份认证、数据加密和访问控制等方面。同时，还需要对系统进行安全性测试，以确保系统可以抵御各种攻击和漏洞。

总之，通过充分的需求分析，可以明确系统的功能和性能需求，为系统的设计和实现提供有力的支持。

##### 3.1.2 系统设计原则

在设计微博主题画像系统时，需要遵循一些系统设计原则，以确保系统的可用性和可靠性。以下是一些关键的设计原则：

（1）数据库设计原则：在设计数据库时，需要考虑到系统的数据存储需求。为了保证数据的完整性和一致性，需要采用恰当的数据库设计原则，如范式设计等。

（2）数据安全原则：由于系统涉及到用户数据，因此数据安全是非常重要的。系统需要采用一系列数据安全措施，如用户身份认证、数据加密、防止 SQL 注入等。

（3）可扩展性原则：随着系统的使用规模逐渐扩大，系统需要具备一定的可扩展性，以应对高并发访问和大量数据存储需求。因此，在系统设计时，需要考虑到系统的扩展性，并采用一些合适的技术和工具来提高系统的性能和可扩展性。

（4）用户体验原则：用户体验是系统设计的核心要素之一。为了提供更好的用户体验，需要考虑到用户的使用习惯和需求，并根据这些需求进行设计。同时，需要确保系统的易用性和可访问性，使得用户可以方便地使用系统并获得所需的信息。

（5）系统性能优化原则：在系统运行过程中，系统性能是非常重要的。为了保证系统的高性能和低延迟，需要采用一些性能优化策略，如缓存优化、异步处理、负载均衡等。

（6）数据可视化原则：数据可视化是系统设计的重要组成部分。为了让用户更加直观地了解系统数据，需要采用一些数据可视化技术，如 ECharts、D3.js 等，来展示数据分析结果，使得用户可以更加清晰地了解数据特点。

综上所述，以上是在设计微博主题画像系统时需要考虑的一些系统设计原则，只有合理的应用这些原则，才能够设计出稳定、高效、可靠、易用的系统。

### 3.2 系统的架构设计

系统的架构设计是整个系统的核心，主要包括前端页面设计、后端接口设计、数据库设计以及数据分析与可视化设计等方面。

前端页面设计采用了响应式设计，使得系统能够在不同尺寸的屏幕上正常显示。同时，采用了 Bootstrap 框架来加快开发速度。前端页面的设计主要包括主页、登录页、注册页、添加数据页、数据展示页、数据分析页等。

后端接口设计采用了 Python Flask 框架，以及 RESTful API 风格的设计，接口返回 JSON 格式的数据。采用这种设计，能够使得前后端分离，并且可以扩展到移动端等多个平台。后端接口的设计主要包括用户认证、数据添加、数据查询、数据分析等。

数据库设计采用了 MySQL 关系型数据库。设计了多个表来存储不同类型的数据，如微博主题数据、转发数据等。同时，采用了索引来加速查询效率，减少数据库的负担。

数据分析与可视化设计采用了 Python 的数据分析库 Pandas 和可视化库 pyecharts，以及机器学习库 sklearn。数据分析的过程主要包括数据清洗、特征提取、建模预测等，最终通过可视化的方式来展示分析结果

### 3.3 系统模块设计

（1）爬虫模块：使用 Python 的 requests 库获取微博 API 返回的 JSON 数据，再使用 Python 的 json 库解析 JSON 数据，从中提取所需的微博内容。

（2）数据清洗与预处理模块：通过 Python 的 pandas 库对爬取到的微博数据进行清洗和预处理，包括去除重复数据、去除无用数据、统一格式等。

（3）数据库模块：使用 Python 的 pymysql 库连接 MySQL 数据库，将清洗后的微博数据存储到数据库中。

（4）机器学习模块：使用 Python 的 sklearn 库训练情感分析模型，对微博内容进行情感分析，将分析结果存储到数据库中。

（5）Web 前端模块：使用 HTML、CSS、JavaScript 等前端技术，实现 Web 页面的设计与展示。

（6）Web 后端模块：使用 Python 的 Flask 框架实现 Web 应用后端的逻辑处理，包括处理前端发送的请求、调用爬虫模块获取微博数据、调用机器学习模块进行情感分析、从数据库中查询数据等。

（7）图表绘制模块：使用 Python 的 pyecharts 库实现图表的绘制与展示，包括柱状图、饼图、散点图、关系图等。

（8）大屏展示模块：使用 JavaScript 等前端技术，实现大屏幕的设计与展示，包括实时更新数据、展示各种图表等。

（9）用户管理模块：使用 Flask 框架实现用户管理功能，包括用户注册、用户登录、用户信息修改等。

（10）数据上传模块：实现用户上传已有的 CSV 文件或微博 API 接口，进行数据分析。

### 3.4 微博画像系统的实现

微博画像系统的实现主要包括爬虫、CSV 数据上传模块、机器学习模型训练和部署以及前端展示等模块。

#### 3.4.1 爬虫模块

通过 Python 的 Requests 库获取微博 API 接口中的 JSON 数据，首先，我们需要模拟用户登录并获取相应的 cookies，然后向微博网站发送请求，获取到相关的微博数据。我们通过分析微博网站的请求方式和响应内容，获取到了相关的 API 接口，并通过解析 json 数据来获取微博的相关信息，如微博内容、发布时间、转发数、评论数等。并通过 PyMySQL 库连接 Mysql 数据库进行存储和管理。这个模块是实现系统自动化的基础，因为自动获取数据是系统的关键之一。

```
In [21]: rankJSON
Out[21]: {'权志龙姐姐开北京首店': {'href': 'https://s.weibo.com/weibo?q=%23E6%9D%83E5%BF%97E9%BE%99E5%A7%90E5%A7%90E5%BC%80E5%8C%97E4%BA%AC%9A%66E5%BA%97%23&t=31&band_rank=1&Refer=top',
'hot': 2001437},
'一家4口三亚溺水全部遇难': {'href': 'https://s.weibo.com/weibo?q=%23E4%B8%80E5%AE%B6%4E5%8F%A3E4%B8%89E4%BA%9A%66%BA%BAE6%B0%B4E5%85%A8E9%83%A8E9%81%87E9%9A%BE%23&t=31&band_rank=2&Refer=top',
'hot': 1789683},
'短发新娘不穿婚纱被误认成伴郎': {'href': 'https://s.weibo.com/weibo?q=%23E7%9F%AD%5F%91E6%96%B0E5%A8%98E4%B8%8D%7E%A9%BF%5%A9%A9E7%BA%B1E8%A2%ABE8%AF%AF%8E%AE%A4E6%88%90E4%BC%B4E9%83%8E%23&t=31&band_rank=1&Refer=top',
'hot': 1497125},
'新冠康复如何做到清淡饮食': {'href': 'https://s.weibo.com/weibo?q=%23E6%96%B0E5%86%A0E5%BA%B7E5%A4%8DE5%A6%82E4%BD%95E5%81%9AE5%88%B0E6%B8%85E6%B7%A1E9%A5%AE%9A%3%9F%23&t=31&band_rank=3&Refer=top',
'hot': 1492124},
'电费1个月3481元竟有电费刺客': {'href': 'https://s.weibo.com/weibo?q=%23E7%94%B5E8%B4%B91E4%B8%AAE6%9C%883481E5%85%83E7%AB%9F%66%9C%89E7%94%B5E8%B4%B9E5%88%BAE5%AE%A2%23&t=31&band_rank=4&Refer=top',
'hot': 1388200},
'1888万彩礼': {'href': 'https://s.weibo.com/weibo?q=%231888E4%B8%87E5%BD%A9E7%A4%BC%23&t=31&band_rank=5&Refer=top',
'hot': 1203834},
...}
```

图 3-1 每日更新热搜

#### 3.4.2 CSV 数据上传模块

系统允许用户上传 CSV 格式的数据文件，这个模块主要负责读取用户上传的数据文件，并将数据导入到 Mysql 数据库中。这个模块实现了数据的可扩展性和可移植性，方便用户在系统中使用自己的数据。

```
@app.route('/uploadbig', methods=['POST', 'GET'])
def uploadbig():
    global kk
    name = request.args.get('name')
    try:
        with open('upload/' + name, 'r', encoding='gb18030', errors='ignore') as read_obj:
            csv_reader = csv.reader(read_obj)
            list_of_csv = list(csv_reader)
    except:
        return jsonify("上传失败")
    if (len(list_of_csv[0]) == 10):
        for kk in range(len(list_of_csv)-1):
            kk = kk + 1

            sql = "INSERT INTO bigevent_withoutbig ('标题 / 微博内容', '信息属性', '原创/转发', '地址', '媒体名称', '发布日期', '媒体类型', '地域', '全文内容', '精准地域', 'title') \
VALUES ('%s', '%s', '%s', '%s', '%s', '%s', '%s', '%s', '%s', '%s', '%s') \
% (list_of_csv[kk][0], list_of_csv[kk][1], list_of_csv[kk][2], list_of_csv[kk][3], list_of_csv[kk][4], \
list_of_csv[kk][5], list_of_csv[kk][6], list_of_csv[kk][7], list_of_csv[kk][8], list_of_csv[kk][9], name)

            db = MySQLdb.connect()
            db.insert(sql)
```

图 3-2 数据上传模块

CSV 数据上传模块是为了满足用户在使用系统过程中，可能需要上传自己的 CSV 数据文件进行分析而设计的。该模块可以将用户上传的 CSV 文件中的数据提取出来，

并存储到系统的数据库中，以便后续的数据分析和可视化展示。CSV 文件上传过程中，用户需要指定文件路径和表名。

系统设计中，CSV 数据上传模块主要分为前端页面和后端处理两个部分。前端页面使用 HTML、CSS 和 JavaScript 编写，提供用户上传 CSV 文件的功能；后端处理主要是通过 Flask 框架接收前端上传的文件并进行处理，最终将提取出来的数据存储到数据库中。

3.4.3 机器学习模型训练和部署模块

这个模块主要使用 Scikit-learn 等 Python 机器学习库进行情感分析模型的训练和部署。模型训练是基于已有的微博文本数据，通过数据预处理和特征提取等技术构建分类模型，用于判断微博的情感倾向。模型部署则是将训练好的模型应用到实际的微博数据中，以实现自动判断。



图 3-3 模型测试接口

3.4.4 前端展示模块

系统的前端展示模块主要使用了 ECharts 等可视化库，将数据库中的微博数据以图表等形式展示出来。这个模块实现了数据的可视化，方便用户进行数据的分析和展示。同时，系统的前端界面也通过 Flask 框架实现了响应式设计，兼容不同分辨率的终端设备。

3.4.5 其他模块

除了以上模块，系统还实现了一些其他模块，例如：登录验证模块、主题事件分区模块、数据管理模块、数据可视化模块等。这些模块一起构成了完整的微博画像系统，为用户提供了丰富的功能和交互界面。

3.5 数据流程

（1）爬虫模块通过 API 接口获取微博主题相关的微博数据和三级转发的数据，保存到数据库中。



（2）CSV 数据上传模块允许用户上传已有的 CSV 文件到数据库中，用于后续的数据分析。

（3）数据清洗和预处理模块对爬虫和 CSV 数据进行处理，包括去重、过滤垃圾信息、分词、去除停用词等操作。

（4）特征提取和分析模块对处理后的数据提取关键词、话题和情感等特征，并进行数据分析和可视化。

（5）机器学习模型训练和部署模块通过使用 sklearn 等机器学习库，对标注好的数据进行训练，得到情感分类模型，并将其部署到系统中，实现实时的情感分类判断。

（6）前端展示模块通过使用 echarts 等可视化库，将数据以图表的形式展示在系统前端，包括主页仪表盘、微博转发图、情感分析图等。

整个数据流程包含了数据的采集、存储、清洗、特征提取、机器学习模型训练和部署、数据分析和可视化等环节，实现了对微博主题的自动画像，为公安机关提供了有力的情报支撑。



## 4 系统功能测试

系统功能测试是确保系统符合用户需求、满足质量标准和预期目标的过程。在微博画像系统中，需要测试的功能包括爬虫功能测试、CSV 数据上传功能测试、机器学习模型测试、前端展示功能测试、用户权限测试、性能测试。

### 4.1 爬虫功能测试

测试今日热搜榜第一的话题进行爬取，检查是否能够正常爬取微博数据，并且能否正确地解析和提取需要的数据。爬虫功能测试主要包括测试爬虫能否正确地获取微博数据、测试爬虫是否能够爬取多个微博账号的数据、测试爬虫是否能够正常处理异常情况。可以针对爬虫的一些异常情况，例如微博页面访问不了、网络连接超时等，进行测试，确保爬虫可以正常处理这些异常情况，保证系统的稳定性和可靠性。



图 4-1 爬虫功能测试

经过测试，该系统在 7.1 秒的时间内爬取了 109 条微博，速度为 15.4 条每秒，并能够及时的显示在系统下方。

### 4.2 CSV 数据上传功能测试

CSV 数据上传功能测试主要包括文件上传、数据解析、数据存储和数据展示等几个方面。

首先进行文件上传测试，测试上传不同类型的 CSV 文件，包括数据完整性、数据格式、数据大小等方面。确保系统能够正确地识别文件类型、解析文件内容，并进行相应的存储和展示。

接着进行数据解析测试，测试系统是否能够正确解析 CSV 文件中的数据，并将其存储到数据库中。测试数据包括数据类型、数据范围、数据长度等方面。同时，还需测试数据存储的性能，确保系统能够在短时间内完成大量数据的存储。



图 4-2 CSV 上传功能测试

经过测试，该系统在 14.8 秒的时间内上传了 16282 条微博，速度为 1100.1 条每秒，并能够及时的存储在 Mysql 数据库中，为下一步数据可视化提供支持。

4.3 机器学习模型测试

将训练好的模型部署到实际环境中，并测试其在新数据上的预测性能。在进行模型部署时，需要考虑模型的可解释性、效率和稳定性等因素。同时，也需要进行错误分析，了解模型的局限性和改进方向。



图 4-3 机器学习模型情感标注测试

经测试，使用真实的数据集进行测试，模型大体上符合人工标注，因模型的 X\_train 使用的是长段落内容，分词后词语较多，存在较多影响因素，所以仍存在一

些判断错误的情况，在下一步的模型训练中，应将文段更为细致的分割关键词，以达到更好的训练效果。

#### 4.4 前端展示功能测试

前端展示功能测试主要验证系统能否准确地将数据可视化展示。具体测试包括：

（1）确认主页能够正确地展示数据库现有条数、男女发博比例、当日热搜等信息，验证数据获取和展示的准确性。

（2）确认主题事件分区能够正确地展示现有爬取事件和已上传事件，验证系统能否正确地存储和展示事件数据。

（3）确认转发分区能够正确地展示现有转发数据和添加新转发数据功能能够正常使用，验证系统能否正确地存储和展示转发数据。

（4）确认数据上传分区能够正确地上传已有 Csv 文件，并且上传的数据能够被正确地存储和展示。

（5）确认微博 API 接口功能能够正常使用，验证系统能否正确地获取数据。

（6）确认系统能够正确地展示静态图和动态图，并能够实时更新数据展示。

（7）确认系统能够正确地进行用户登录验证，并能够正确地限制未登录用户的访问权限。

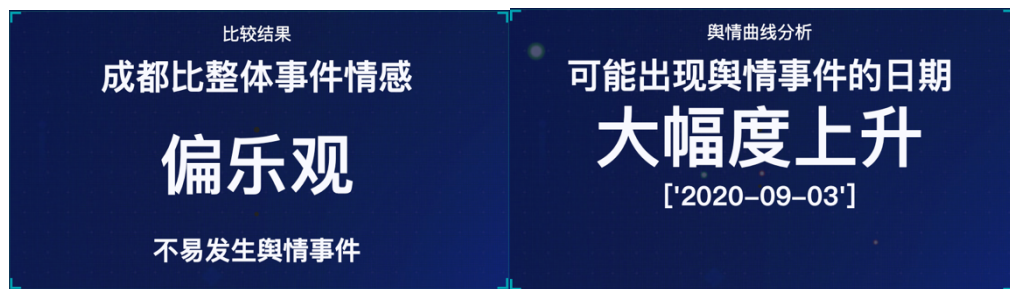


图 4-4 大屏模块

（8）确认大屏模块能够正确地判断事件的特性，进行情感研判以及舆情的初步预测，如图 4-4。

以上测试能够全面地验证系统的功能和准确性，确保系统能够正常地运行。

#### 4.5 用户权限测试

用户权限测试主要是测试系统是否能正确判断用户的登录状态，并在未登录时禁止访问需要登录的页面。测试步骤如下：

打开系统主页，尝试访问需要登录才能访问的页面，如添加新数据页面。预期结果：系统应该跳转到登录页面，并显示提示信息“无权访问，请先登录”。输入正确的

登录信息，登录系统。重新尝试访问需要登录才能访问的页面，如添加新数据页面。

预期结果：系统应该允许用户访问该页面，并显示正确的页面内容。

通过以上测试，可以验证系统的用户权限功能是否正常工作，能够正确判断用户的登录状态并防止未经授权的访问。



图 4-5 用户权限测试

## 4.6 性能测试

性能测试是对系统性能和可扩展性的评估，以确定系统的响应时间、吞吐量、并发用户数和负载能力等。对于微博画像系统而言，性能测试的重点在于测试爬虫的抓取速度、数据上传的速度、机器学习模型的调用速度以及前端展示的响应速度。

经测试，可以实现多任务并行处理以及响应速度满足用户使用。

## 5 总结与展望

### 5.1 总结

在本文中，我们介绍了微博主题画像系统的设计和实现，系统主要包括爬虫、CSV 数据上传模块、机器学习模型训练和部署以及前端展示等模块。系统通过爬虫获取微博主题相关的微博和转发数据，并通过 CSV 数据上传模块实现大批量数据上传和处理，最后使用机器学习模型进行情感分析和主题分类，并通过前端展示模块将结果可视化呈现。

在系统实现过程中，我们采用了 Python 编程语言，使用了众多 Python 库和框架，包括但不限于 Requests、BeautifulSoup、Pandas、Sklearn、Flask、ECharts 等。

在功能测试中，我们对系统的爬虫、CSV 数据上传、机器学习模型、前端展示和用户权限等方面进行了测试，并对性能进行了评估。测试结果表明，系统功能齐全，性能稳定可靠。

通过本系统的实现，可以对微博主题进行分析和画像，对于舆情分析、公共安全等领域具有重要的应用价值。同时，本系统还存在一些不足之处，例如爬虫数据量受限，机器学习模型精度有待提高等问题，需要进一步完善和改进。因此，未来可以结合更多的技术手段，如分布式爬虫、深度学习模型等，进一步完善和拓展系统功能，提高系统的精度和效率，使其在实际应用中发挥更大的作用。

### 5.2 展望

随着互联网和社交媒体的迅猛发展，网络舆情监测和处理需求日益增长，微博主题画像系统的应用前景也十分广阔，有着较长远的发展空间。但是，现有的微博主题画像系统还存在一些技术难点，比如数据清洗和去重、舆情分类精准度、数据展示效果等问题，这些都有待进一步研究和攻克。

一方面，我们需要探索更加高效的爬虫策略，例如基于分布式爬虫、数据抓取和处理等技术，以提高数据的收集和处理效率。另一方面，我们需要利用更加先进的机器学习算法和深度学习模型，例如 BERT、LSTM 和 Transformer 等，以提高微博话题和用户画像的准确性和智能化水平。此外，我们还可以探索使用图神经网络等方法，以建立更加复杂和深入的微博用户社交网络分析模型，帮助我们更好地了解微博用户之间的关系和互动。

在技术方面的进步之外，我们还需要注重数据安全和隐私保护。随着数据的增长，数据泄露和隐私问题也越来越突出。因此，我们需要在系统设计和数据处理过程中充分考虑数据隐私和安全性问题，采用加密算法和数据脱敏等方法，保障用户数据的安全和隐私。

未来，微博主题画像系统也需要结合更加先进的人工智能、机器学习等技术，不断地优化系统的算法、模型和架构，以提升系统对微博舆情的分析和预测能力。还可以将微博主题画像系统和其他相关系统进行融合，如基于人脸识别和位置数据的事件监测和分析等，以提供更加全面和多维的网络舆情服务。

## 参考文献

- [1]彭媛媛,张海霞,连一峰,黄克振,刘倩.一种基于多维度的网络空间人物画像方法[P].北京市:CN115114498A,2022-09-27.
- [2]李继玲,李宝林,严宋如.疫情背景下快递物流服务的用户行为画像及主题挖掘研究[J].华东师范大学学报(自然科学版),2022(05):100-114.
- [3]安璐,周雯静,李纲,余传明.基于事件态势与用户认知的危机信息推送方法及系统[P].湖北省:CN114969560A,2022-08-30.
- [4]刘浒,夏志杰,张晨芳.面向企业印象管理的社交媒体用户情绪演变及画像研究——以微博平台拼多多企业为例[J].情报工程,2022,8(03):112-125.
- [5]商瀑,李亚可,郭楼禄.基于用户画像技术的大数据侦查:一个框架的分析与设计[J].中国人民公安大学学报(社会科学版),2022,38(03):32-43.
- [6]秦权.微博环境下新冠肺炎疫情事件对网民情绪的影响分析研究[D].中北大学,2022.DOI:10.27470/d.cnki.ghbgc.2022.001223.
- [7]李尹舒.基于用户画像的大米电商精准营销策略研究[D].吉林大学,2022.DOI:10.27162/d.cnki.gjlin.2022.002038.
- [8]孟芷薇.基于微博平台的智库用户画像分析研究[D].新疆师范大学,2022.DOI:10.27432/d.cnki.gxsfu.2022.000335.
- [9]胡华实.建构,互动,遵从:城乡二元视角下青年对新型媒体的差异化使用及归因[D].华东师范大学,2022.DOI:10.27149/d.cnki.ghdsu.2022.003744.
- [10]马焱宸.大数据抓取在换装玩偶娃娃个性消费市场中的研究与应用[D].西南科技大学,2022.DOI:10.27415/d.cnki.gxngc.2022.000081.
- [11]潘睿.基于画像的公共图书馆抖音号运营策略研究[D].华东师范大学,2022.DOI:10.27149/d.cnki.ghdsu.2022.001205.
- [12]黄传能.FM 传媒用户节项目运营及营销策略研究[D].电子科技大学,2022.DOI:10.27005/d.cnki.gdzku.2022.001653.
- [13]张虹.企业招聘,打好数据牌[J].人力资源,2022(04):47-49.
- [14]王志刚,邱长波.基于主题的政务微博评论用户画像研究[J].情报杂志,2022,41(03):159-165.
- [15]叶光辉,郭诚,徐彤,王灿灿.城市画像视角下的政务社交媒体资源保存研究[J].情

- 报科学,2022,40(02):11-17+27.DOI:10.13833/j.issn.1007-7634.2022.02.002.
- [16]HerdaGdelen, A., Zuo, W., Gard-Murray, A., & Bar-Yam, Y, "An exploration of social identity: The geography and politics of news-sharing communities in twitter", Complexity, Vol. 19, no. (2), pp. 10-20, 2013
- [17]Thelwall, M., Buckley, K., Paltoglou, G., Cai, D., & Kappas, A., "Sentiment strength detection in short informal text". Journal of the American Society for Information Science and Technology, vol.61, no. 12, pp.2544-2558, 2010.
- [18]H.U. Khan, S.M. Saqlain, M.Shoaib, M.Sher, "Ontology based semantic search in Holy Quran", International Journal of Future Computer and Communication, Vol. 2, no 2, pages 1-6, 2013.
- [19]H.U.Khan, T.A.Malik, "Finding resources from middle of RDF graph and at Sub-query level in suffix array based RDF indexing using RDQL queries", International Journal of Computer Theory and Engineering, Vol.4, no.3, pages 369-375, 2012.
- [20]H.U. Khan and A. Daud, T.A.Malik, "MIIB: A metric to identify top influential bloggers in a community", Plos One, vol. 10, no. 9, 2015.
- [21]U. Ishfaq, H.U. Khan, K. Iqbal, "Modeling to find the top bloggers using sentiment features", in Proc. International Conference on Computing, Electronic and Electrical Engineering, pp. 227-233, 2016.
- [22]H.U.Khan, A.Daud, "Using machine learning techniques for subjectivity analysis based on lexical and non-lexical features". International Arab Journal of Information Technology, Vol. 14, no. 4. 2017.



## 致 谢

在完成这篇论文时，我得到了很多人的帮助和支持，在此谨向他们致以最诚挚的谢意。

本论文是在导师组兰月新教授、夏一学教授和导师安晓伟讲师的精心指导下完成的。导师学识深厚、治学严谨、实事求是，拥有多年的教学和研究经验，对所从事的领域非常熟悉和深入。在我们的科研过程中，导师总是能够提供最及时、最准确的指导和帮助，总是会耐心地听取我们的研究想法，鼓励我们尝试新的方法和技术，并指出我们在研究中可能存在的问题和不足。这些指导不仅帮助我们更好地理解问题的本质，还促进了我们的科研能力的提升。导师非常关心我们的成长和发展。他总是鼓励我们多读书、多实践、多思考，帮助我们建立健康的学术生态和良好的职业道德。在研究过程中，导师不仅关心我们的学术成果，还经常询问我们的生活和心理状态，并提供必要的帮助和支持。导师的关心和支持让我们感受到了家一般的温暖和关怀，更加自信地面对未来的挑战。导师在教学、科研和管理方面都充满了爱心。他始终把学生的发展放在首位，鼓励我们发挥自己的优势和特长，培养我们的创新思维和实践能力。同时，他也非常重视学生的身心健康，注重营造一个和谐、积极、乐观的学术环境。四年来，导师倾注了大量的时间和精力对我的学业进行指导，始终为人师表，言传身教。他严谨的作风，求实的学风，使我终生受益，让我在实践中不断成长。

他的教诲和帮助让我获益匪浅。在未来的学习和工作中，我将继续发扬他所教导的科研精神和职业道德，努力成为一名优秀的警察以及科研工作者。