
NLPIR 大数据语义智能分析平台

用户手册



自然语言处理与信息检索共享平台
Natural Language Processing & Information Retrieval Sharing Platform

<http://www.nlpir.org/>

目 录

一、NLPIR 平台简介	3
二、文件下载与说明	7
2.1 文件下载	7
2.2 文件说明	7
三、各个功能操作指南	9
3.1 精准采集	11
3.2 文档转换	16
3.3 新词、关键词提取	18
3.3.1 新词发现	19
3.3.2 关键词提取	21
3.3.3 可视化展示	22
3.4 批量分词	24
3.5 语言统计	28
3.6 文本聚类	32
3.7 文本分类	34
3.8 摘要实体	40
3.9 智能过滤	42
3.10 情感分析	47
3.11 文档去重	51
3.12 全文检索	52
3.13 编码转换	57
四、应用示范案例	59
4.1 十九大报告语义智能分析	59
4.2 文章风格对比：方文山 VS 汪峰	62
4.3 《红楼梦》作者前后同一性识别	64
五、联系我们	66
六、附录	67
6.1 下载途径	67
6.2 Github 下载	68
6.3 百度网盘下载	71

一、NLPIR 平台简介

NLPIR 大数据语义智能分析平台，针对大数据内容处理的需要，融合了网络精准采集、自然语言理解、文本挖掘和网络搜索的技术，提供客户端工具、云服务、二次开发接口。平台先后历时十八年，服务了全球四十万家机构用户，是大数据时代语义智能分析的一大利器。

开发平台由多个中间件组成，各个中间件 API 可以无缝地融合到客户的各类复杂应用系统之中，可兼容 Windows, Linux, Android, Maemo5, FreeBSD 等不同操作系统平台，可以供 Java, C, C#等各类开发语言使用。



图 1.1 NLPIR 大数据语义智能分析平台简介

NLPIR 大数据语义智能分析平台的十三大功能：

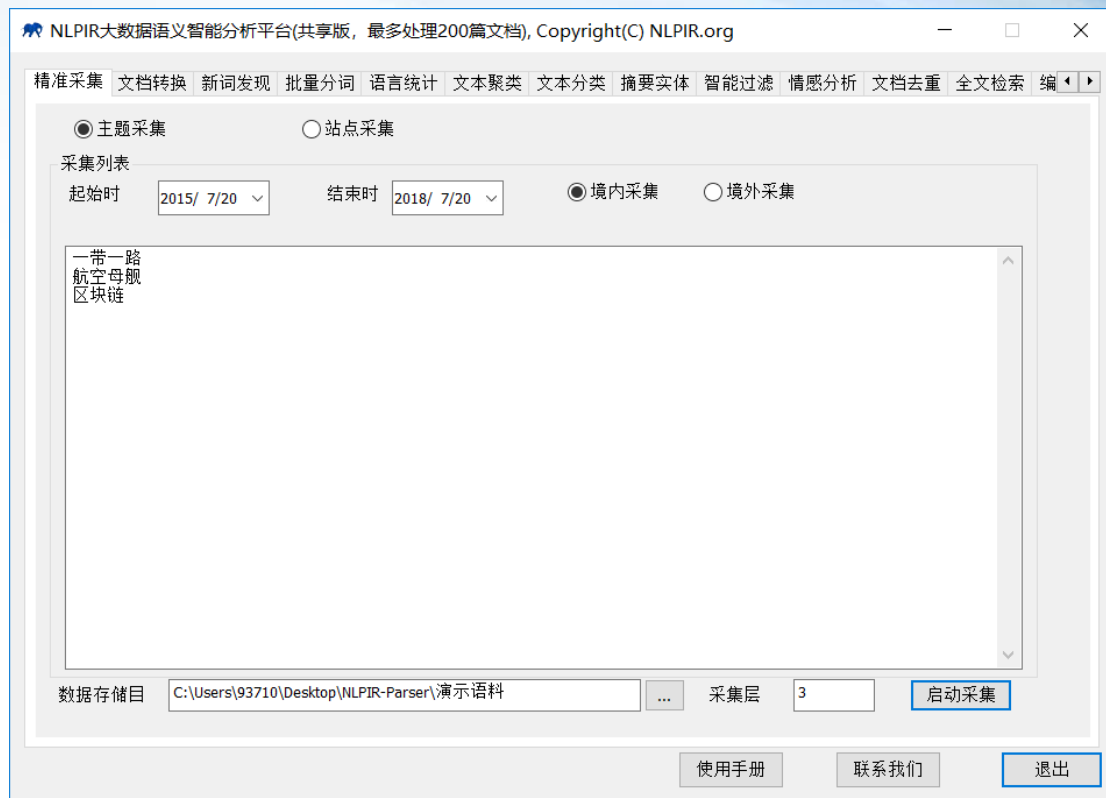


图 1.2 NLPIR 大数据语义智能分析平台客户端

1. 精准采集

对境内外互联网海量信息实时精准采集，有主题采集（按照信息需求的主题采集）与站点采集两种模式（给定网址列表的站内定点采集功能）。可帮助用户快速获取海量信息。

2. 文档转换

对 doc、excel、pdf 与 ppt 等多种主流文档格式，进行文本信息格式转换，信息抽取准确率极高，效率达到大数据处理的要求。

3. 新词发现

新词发现能从文本中挖掘出具有内涵新词、新概念，用户可以用于专业词典的编撰，还可以进一步编辑标注，导入分词词典中，提高分词系统的准确度，并适应新的语言变化。

关键词提取能够对单篇文章或文章集合，提取出若干个代表文

章中心思想的词汇或短语，可用于精化阅读、语义查询和快速匹配等。

4. 批量分词

对原始语料进行分词、自动识别人名地名机构名等未登录词、新词标注以及词性标注。可在分析过程中，导入用户定义的词典。

5. 语言统计

针对切分标注结果，系统可以自动地进行一元词频统计、二元词语转移概率统计（统计两个词左右连接的频次即概率）。针对常用的术语，会自动给出相应的英文解释。

6. 文本聚类

能够从大规模数据中自动分析出热点事件，并提供事件话题的关键特征描述。同时适用于长文本和短信、微博等短文本的热点分析。

7. 文本分类

针对事先指定的规则和示例样本，系统自动从海量文档中识别并训练分类。NLPIR 深度文本分类，可以用于新闻分类、简历分类、邮件分类、办公文档分类、区域分类等诸多方面。

8. 摘要实体

自动摘要能够对单篇或多篇文章，自动提炼出内容的精华，方便用户快速浏览文本内容。实体提取能够对单篇或多篇文章，自动提炼出内容摘要，抽取人名、地名、机构名、时间及主题关键词；方便用户快速浏览文本内容。

9. 智能过滤

对文本内容的语义智能过滤审查，内置国内最全词库，智能识别多种变种：形变、音变、繁简等多种变形，语义精准排歧。

10. 情感分析

情感分析，针对事先指定的分析对象，系统自动分析海量文档的情感倾向：情感极性 & 情感值测量，并在原文中给出正负面的得分和句子样例。

11. 文档去重

能够快速准确地判断文件集合或数据库中是否存在相同或相似内容的记录，同时找出所有的重复记录。

12. 全文检索

JZSearch 全文精准检索支持文本、数字、日期、字符串等各种数据类型，多字段的高效搜索，支持 AND/OR/NOT 以及 NEAR 邻近等查询语法，支持维语、藏语、蒙语、阿拉伯、韩语等多种少数民族语言的检索。可以无缝地与现有文本处理系统与数据库系统融合。

13. 编码转换

自动识别文档内容的编码，并进行自动转换，目前支持 Unicode/BIG5/UTF-8 等编码自动转换为简体的 GBK，同时将繁体 BIG5 和繁体 GBK 进行繁简转化。

二、文件下载与说明

2.1 文件下载

GitHub 下载地址: <https://github.com/NLPIR-team/NLPIR/tree/master/NLPIR-Parser>

下载教程参见附录。【有可能国内访问国外网址受限】

注：用户在 github 上下载 NLPIR-Parser 文件时需要专门的下载工具，建议使用 svn 工具下载文件。

2.2 文件说明

NLPIR-Parser 文件目录如下：

■ Data	Update NLPIR-Parser and 授权
■ bin-win32	Update NLPIR-Parser
■ bin-win64	Update NLPIR-Parser and 授权
■ doc	Update NLPIR-Parser and manual
■ 不良内容测试文件	Update NLPIR-Parser
■ 演示语料	update NLPIR-Parser
■ 编码转换测试文本	update NLPIR-Parser
■ 训练分类用文本	update NLPIR-Parser
■ Readme.txt	update NLPIR-Parser
■ 清理临时文件.bat	Update NLPIR-Parser and manual

图 2.1 文件目录

文件说明：

└─bin-win32	Windows 32bit 环境下的可执行程序 and 库文件，也可运行于 Win64；点击 NLPIR-Parser.exe 即可运行。
└─output	运行结果存放路径
└─bin-win64	Windows 64bit 环境下的可执行程序 and 库文件；点击 NLPIR-Parser.exe 即可运行。
└─output	运行结果存放路径



件 件	└─Data	整个系统运行需要的数据文件
	└─Cluster	聚类系统运行需要的数据文件
	└─┬─Data	
	└─DeepClassifier	机器学习分类运行需要的数据文件
	└─English	英语处理需要的数据文件
	└─JZSearch	JZSearch 精准语义搜索引擎处理需要的数据文件
	└─KeyScanner	JZSearch 精准语义搜索引擎处理需要的数据文件
	└─RedupRemover	去重需要的数据文件
	└─SentimentNew	情感分析需要的数据文件
	└─┬─Data	
	└─┬─English	
	└─doc	NLPIRParser 使用手册与各模块接口文档文件
	└─┬─大数据组件接口文档	
	└─┬─Classifier	
	└─┬─Cluster	
	└─┬─DocExtractor	
	└─┬─DupRemove	
	└─┬─JZSearch	
	└─┬─KeyExtract	
	└─┬─LJSentimentAnalysis	
	└─┬─Summary	
	└─┬─WordFreq	
	└─演示语料	NLPIR-Parser 提供
	└─的测试语料，可以自行替换	
	└─┬─编码转换测试文本	NLPIR-Parser 提供的编
	└─┬─码转换测试语料，可以自行替换	
	└─┬─训练分类用文本	NLPIRParser 提供
	└─┬─的分类训练语料，可以自行替换	
	└─┬─交通	
	└─┬─体育	
	└─┬─军事	
	└─┬─政治	
	└─┬─教育	
	└─┬─经济	
	└─┬─艺术	

1. NLPIR-Parser.exe 可执行文件，本版本为共享版本（只能处理 200 个文件，总量不超过 500KB 纯文本），大规模语料处理需要购买正式版

2. 演示语料，用户可替换，必须为文本文件，如果为 GBK 以

外的编码，必须先进行编码识别与转换后方可进行其他操作。

3. 各种 dll 为各组件的调用接口，本演示程序全部基于已有的调用接口实现；

三、各个功能操作指南

首先，启动程序。

用户需要点击 C:\Users\Administrator\Desktop\NLPIR-paser/bin-win64/ 路径下的 NLPIR-Parser.exe 程序，即可打开软件，平台界面如下：

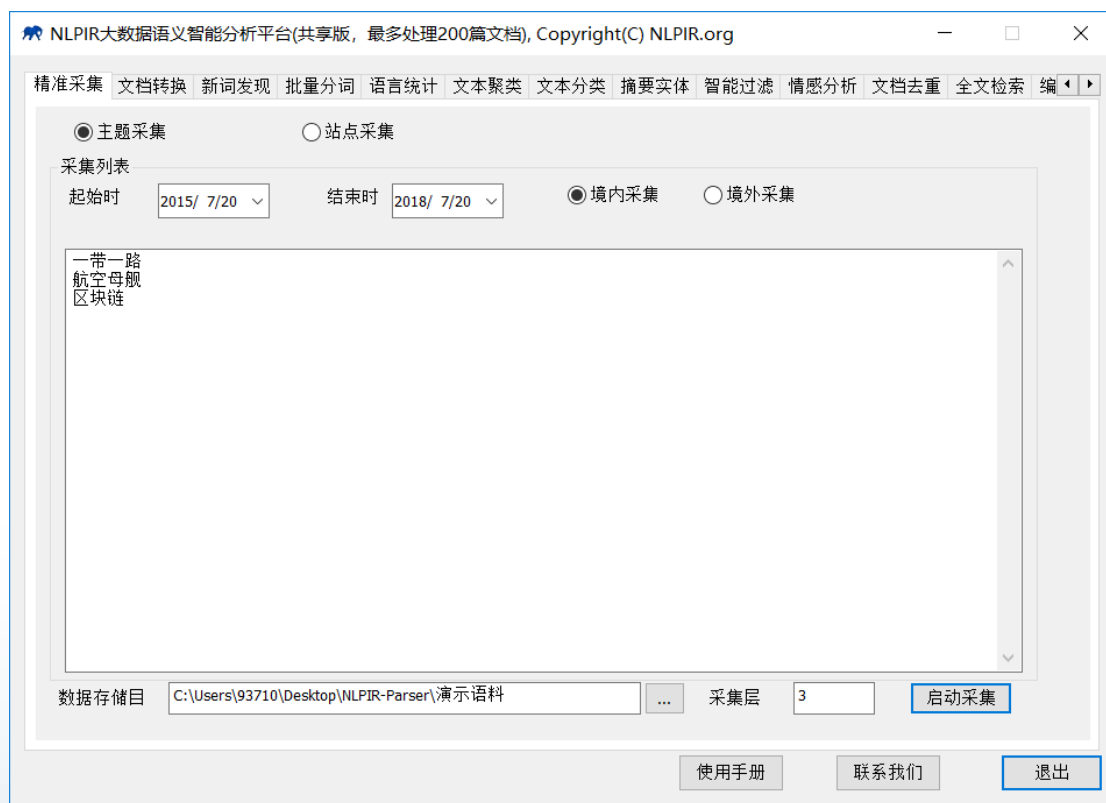


图 3.1 NLPIR 大数据语义智能分析平台界面

然后，平台界面介绍。

平台包括三大模块：“功能导航”（点击功能名称即可完成功

能切换）、“功能操作”和“基础功能”（使用手册、联系我们与退出）。

平台的十三大功能（由左至右）：精准采集，文档转换、新词发现、批量分词、语言统计、文本聚类、文本分类、摘要实体、智能过滤、情感分析、文档去重、全文检索和编码转换。用户可根据需要选择使用。

点击“使用手册”，即可打开平台使用手册文档，帮助用户了解平台，指导用户进行各项功能操作。

NLPIR 大数据语义智能分析平台

用户手册



图 3.2 使用手册

点击“联系我们”，联系信息框弹出，用户可查看咨询。



图 3.3 联系我们

注：平台内置测试语料，但用户仍可定义自己的语料（新建文件夹放入自己的语料）。

3.1 精准采集

用户点击“精准采集”（第一个功能模块），进入精准采集模块。

精准采集功能可实现对境内外互联网海量信息的实时精准采集。精准采集包括主题采集（按照信息需求的主题采集）与站点采集两种模式（给定网址列表的站内定点采集功能）。可帮助用户快速获取海量信息。用户可自定义采集模式、采集时间区域、采集主题站点与采集存储。

➤ 主题采集

按照给定的关键词或主题词进行信息采集。

Step1: 定义主题词。选择“主题采集”，在采集模块输入关键词，例如“一带一路”、“航空母舰”与“区块链”等三个主题，系统将按此关键词进行主题采集，获取主题相关的主流新闻报道、BBS 与博客等内容。

Step2: 采集设置。用户可自定义采集时间（系统默认采集时段为近 3 年，用户可在此时间段内自定义自己的采集时间）。选择采集区域“境内采集”（或境外采集，需要启动翻墙措施方可使用）。

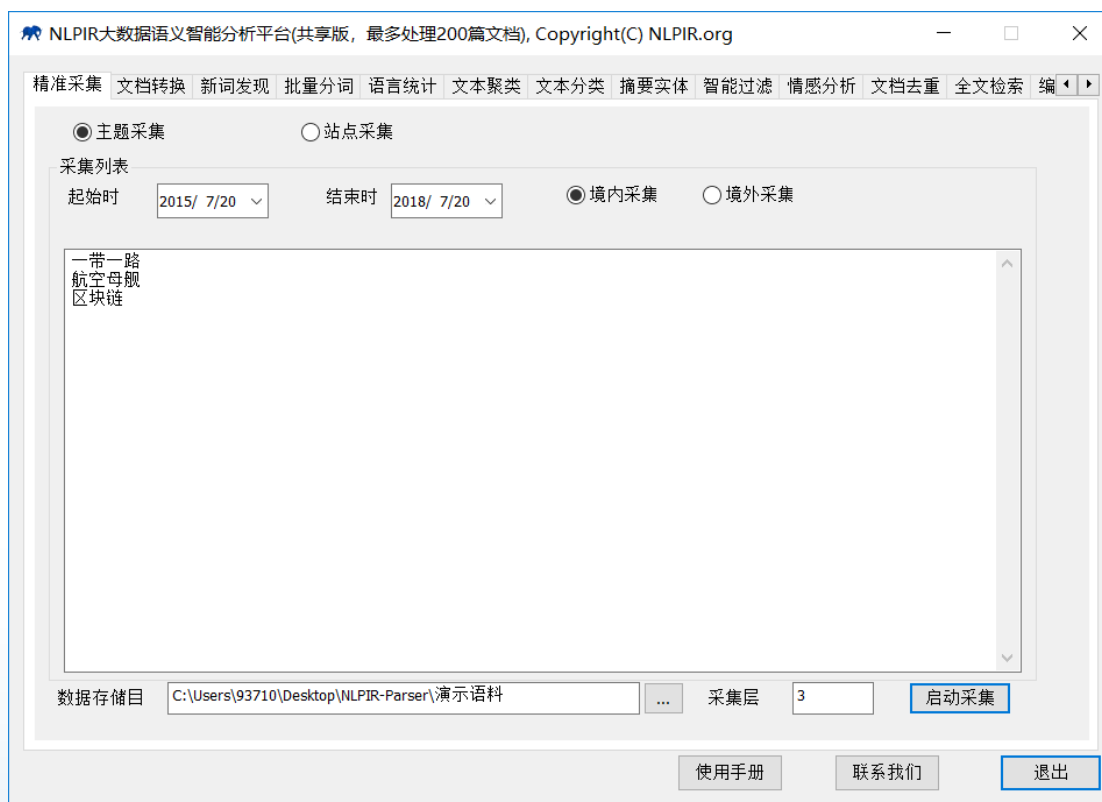


图 3.4 主题采集

Step3: 定义采集存储。选择语料存放路径（默认路径：NLPIR-Parser\演示语料）。点击“启动采集”，系统弹出信息采集窗口，开始采集信息，用户了解信息采集过程与详情。

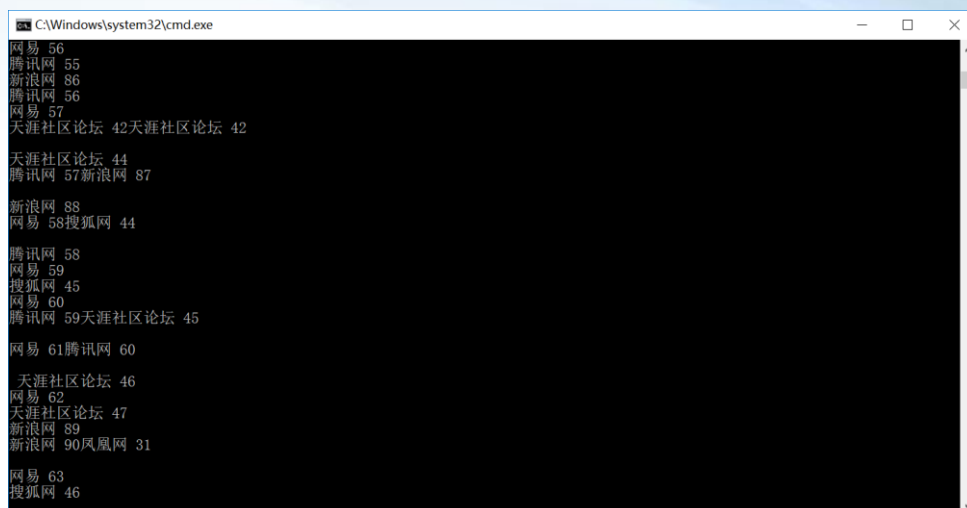


图 3.5 信息采集过程

系统提示采集全部结束，用户可关闭此窗口。

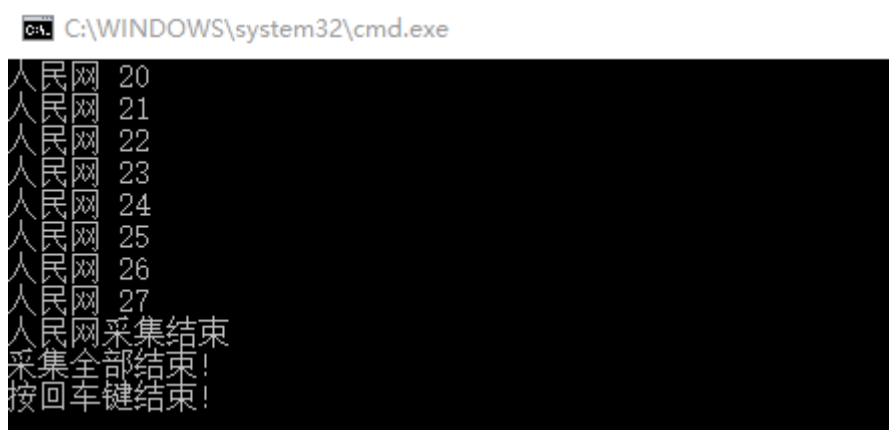


图 3.6 采集结束

采集完成以后，用户可查看采集结果（默认：\NLPIR-Parser\演示语料），采集结果文件夹包括：境内新闻、境外新闻与 bbs 以及通用采集。其中的子目录中的数字指的是文章发布的日期，如 境内新闻 20180301：指的是 2018 年 3 月 1 日的境内新闻。

NLPIR-Parser > 演示语料		
名称	修改日期	类型
bbs	2018/3/6 18:44	文件夹
境内新闻	2018/3/6 18:46	文件夹

NLPIR-Parser > 演示语料 > 境内新闻		
名称	修改日期	类型
境内新闻(20180209)	2018/3/6 18:46	文件夹
境内新闻(20180208)	2018/3/6 18:46	文件夹
境内新闻(20180213)	2018/3/6 18:46	文件夹
境内新闻(20180214)	2018/3/6 18:46	文件夹
境内新闻(20180220)	2018/3/6 18:46	文件夹
境内新闻(20180223)	2018/3/6 18:46	文件夹
境内新闻(20180302)	2018/3/6 18:45	文件夹
境内新闻(20180301)	2018/3/6 18:45	文件夹
境内新闻(20180305)	2018/3/6 18:45	文件夹
境内新闻(20180228)	2018/3/6 18:44	文件夹

图 3.7 采集结果文件

➤ 站点采集

站点采集指的是按照给定的网址，在该网址内部垂直采集。

Step 1: 选择“站点采集”，输入站点地址，例如：

<http://news.sina.com.cn/>。

Step 2: 定义采集时间、区域与采集结果存放路径，点击“启动采集”，系统开始采集任务。

系统将站点采集的结果保存在“通用采集”文件夹中，文件目录：`\NLPIR-Parser\演示语料`。

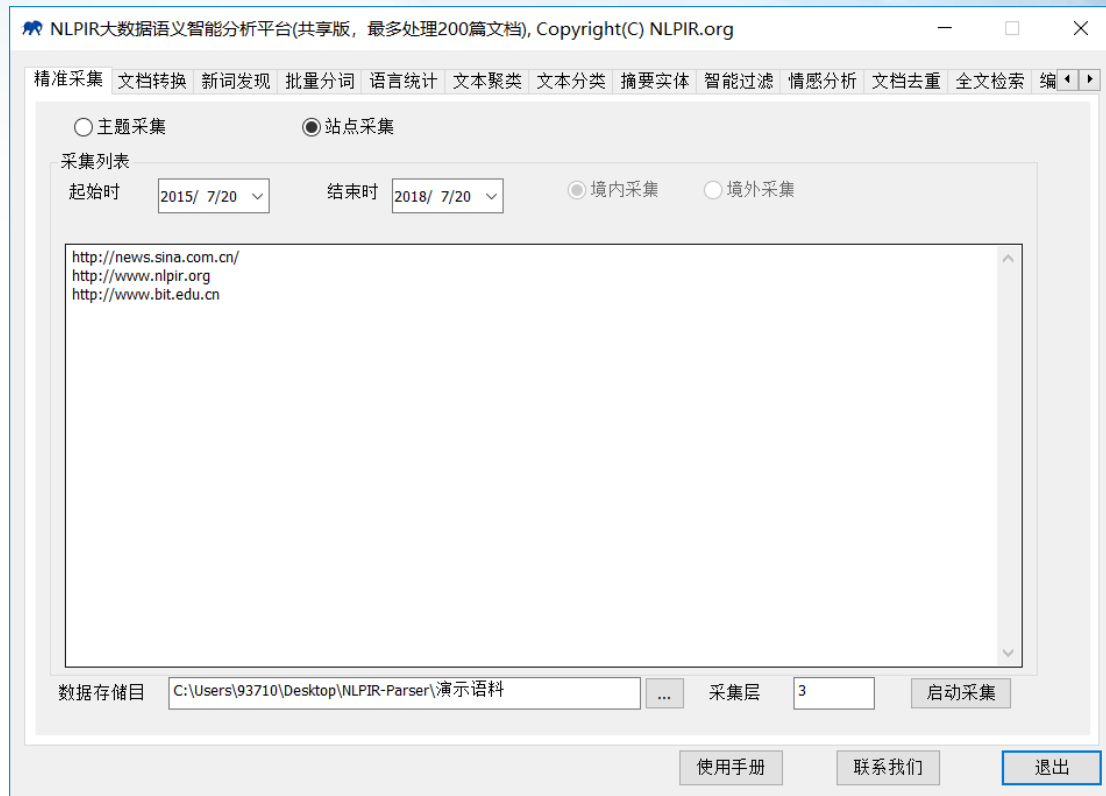


图 3.8 站点采集

站点采集程序启动后，显示采集过程与详情，采集完毕后窗口自动关闭。

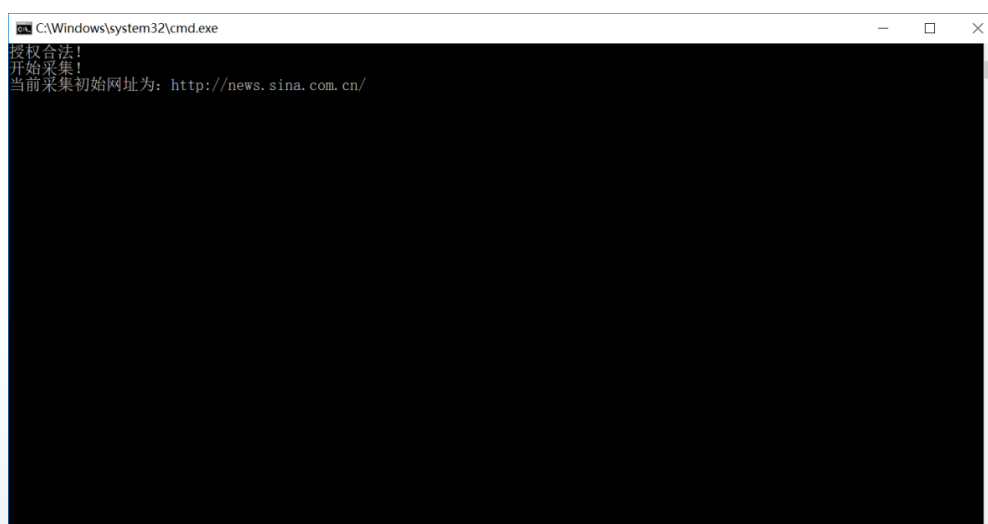


图 3.9 站点采集

PIR-Parser > 演示语料 > 站点采集

名称	修改日期	类型	大小
"高铁一姐"曾被多人诈骗 她挖的坑一直...	2018/7/20 9:55	XML 文档	7 KB
"货拉拉"当网约车载客? 平台-司机载客...	2018/7/20 9:55	XML 文档	10 KB
"货拉拉"载客视频网上热传 发生事故谁...	2018/7/20 9:55	XML 文档	8 KB
"特普会"发言被骂惨 特朗普急回应: 昨...	2018/7/20 9:55	XML 文档	3 KB
"特普会"刚结束 美国就拘捕一名俄女特...	2018/7/20 9:55	XML 文档	2 KB
《九州》拍网剧, "天空城"里一片新面孔...	2018/7/20 9:55	XML 文档	5 KB
2年前交1.5亿彩礼又退婚 "新郎"起诉: ...	2018/7/20 9:55	XML 文档	13 KB
3万6就能买到京牌和车是真的? 如何参...	2018/7/20 9:55	XML 文档	8 KB
4人集资10.8亿被判刑: 千余人参与 实际...	2018/7/20 9:55	XML 文档	3 KB
7岁女儿目睹爸爸救起落水母子 自制奖状...	2018/7/20 9:55	XML 文档	4 KB
14岁女孩从瀑布高处跳水玩耍 被急流卷...	2018/7/20 9:55	XML 文档	2 KB
16岁上大学的他将成福建最年轻市长(图-...	2018/7/20 9:55	XML 文档	6 KB
24岁小伙因还不清高利贷喝农药自杀 抢...	2018/7/20 9:55	XML 文档	3 KB
31名老赖子女就读高收费私立学校 法院...	2018/7/20 9:55	XML 文档	7 KB
31岁富太称能拿紧俏房源 亲友老公都被...	2018/7/20 9:55	XML 文档	7 KB
72岁英国男资助非洲贫困家庭 性虐当地...	2018/7/20 9:55	XML 文档	4 KB

图 3.10 站点采集结果文件

3.2 文档转换

用户点击功能导航栏“文档转换”，系统进入“文档转换”模块。

文档转换功能对 doc、excel、pdf 与 ppt 等多种主流文档格式，进行文本信息抽取，信息抽取准确率极高，达到大数据处理的要求。

Step1: 选择待处理文件。在“文档所在路径”输入框中输入或选择需要抽取的文档文件（用户可选择电脑中的任何文档），例如:\NLPIR-Parser\文档转换。

Step2: 定义结果存储。在“结果存放路径”选择文档转换完成文件存放的地址路径，例如:\NLPIR-Parser\文档转换。

Step3: 点击“文档解析抽取”，系统弹出文档转换处理窗口，开始文档转换。

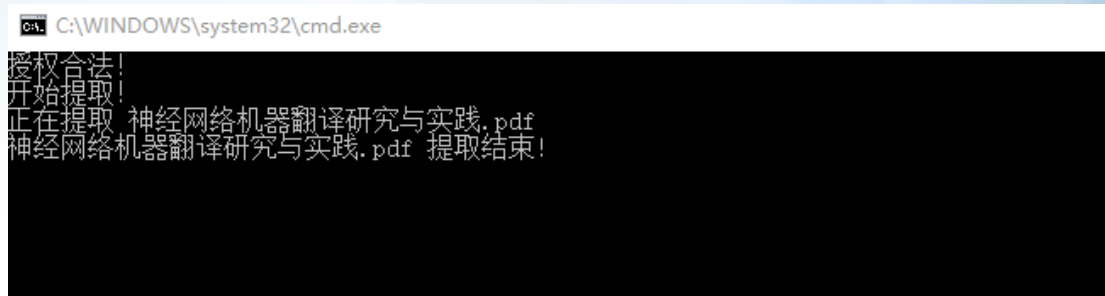


图 3.11 文档转换

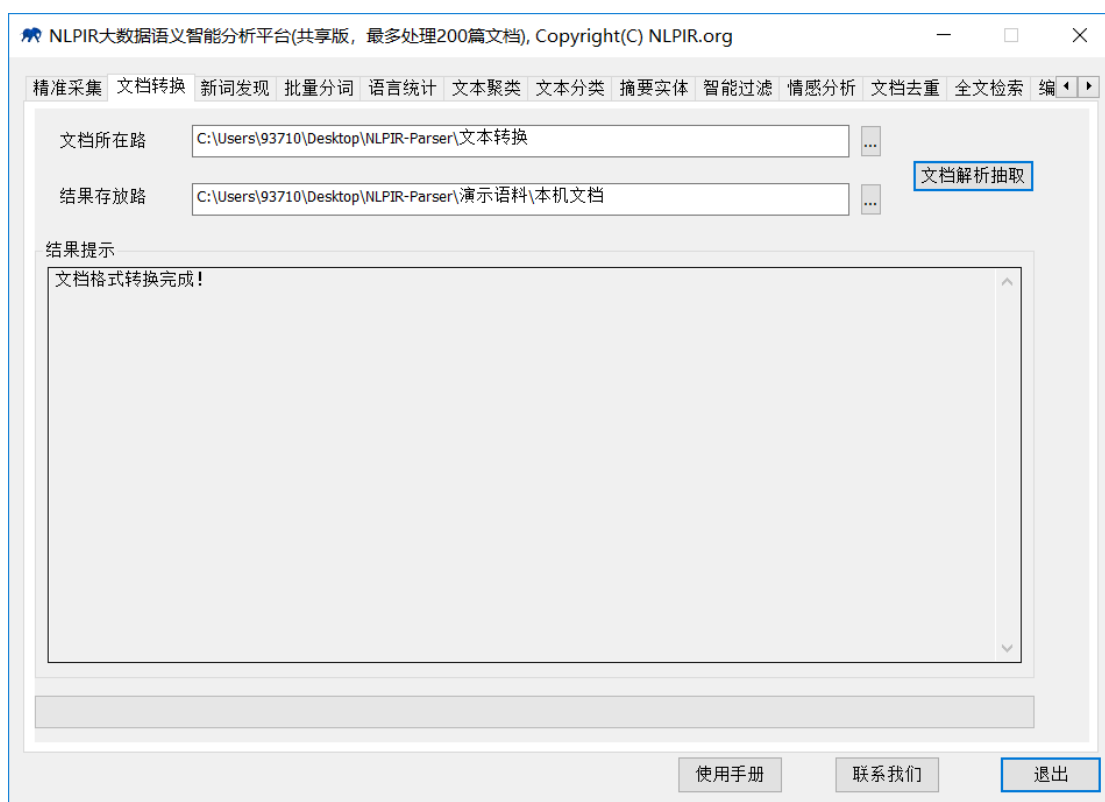


图 3.12 文档转换完成

文档转换结果文件会自动打开（用户也可打开文档转换结果存储目录查看结果文件），抽取完成的文档以文本文件的格式保存。通过结果文件与文件原文的对比，可发现文件抽取具有非常高的准确率。







NLPIR-Parser > 文档抽取			
名称	修改日期	类型	
 精准营销中用户画像挖掘.ppt	2017/1/17 15:04	Microsoft Po	
 精准营销中用户画像挖掘.ppt.txt	2018/3/1 14:55	TXT 文件	
 神经网络机器翻译研究与实践.pdf	2017/1/17 15:37	WPS PDF 文档	
 神经网络机器翻译研究与实践.pdf.txt	2018/3/1 14:52	TXT 文件	
 一带一路.docx	2018/3/1 14:23	Microsoft Wc	
 一带一路.docx.txt	2018/3/1 14:52	TXT 文件	

图 3.13 文档转换结果文件

神经网络机器翻译研究与实践.pdf		神经网络机器翻译研究与实践.pdf.txt
摘要:		摘要:
机器翻译一直以来都是自然语言处理的一大重要方向，集成了自然语言处理		机器翻译一直以来都是自然语言处理的一大重要方向，集成了自然语言处理
各项先进的技术。尽管在过去的几十年中，传统的统计机器翻译（SMT）发展迅		各项先进的技术。尽管在过去的几十年中，传统的统计机器翻译（SMT）发展迅
速，但翻译质量仍然不能满足用户的需求。近几年来，神经网络机器翻译（NMT）		速，但翻译质量仍然不能满足用户的需求。近几年来，神经网络机器翻译（NMT）
在机器翻译这个课题上，获得了令人瞩目的成果。NMT 由单一的深度神经网络		在机器翻译这个课题上，获得了令人瞩目的成果。NMT 由单一的深度神经网络
组成，直接学习源语言到目标语言的翻译，是一个端到端的系统，在短短的两年		组成，直接学习源语言到目标语言的翻译，是一个端到端的系统，在短短的两年
间，就拥有了超越传统机器翻译的能力。Google、Baidu 等公司纷纷将线上机		间，就拥有了超越传统机器翻译的能力。Google、Baidu 等公司纷纷将线上机
器翻译系统的算法替换为 NMT，可以说，NMT 是未来机器翻译方向的大势。本		器翻译系统的算法替换为 NMT，可以说，NMT 是未来机器翻译方向的大势。本
文将简要描述机器翻译方向的研究进展，讨论几个关键难点的现有解决方案。接		文将简要描述机器翻译方向的研究进展，讨论几个关键难点的现有解决方案。接
着，对实际的模型做训练，研究其翻译效果。最后加入数据对比试验，尝试增大		着，对实际的模型做训练，研究其翻译效果。最后加入数据对比试验，尝试增大
训练数据量，以获得数据对训练结果的实际影响。		训练数据量，以获得数据对训练结果的实际影响。
一、 机器翻译研究进展综述		

图 3.14 文档转换效果对比

3.3 新词、关键词提取

用户点击“新词发现”，系统切换进入“新词发现”功能模块。

新词发现模块包括新词发现与关键词抽取两个功能。

3.3.1 新词发现

新词发现能从文本中挖掘出具有内涵新词、新概念，用户可以用专业词典的编撰，还可以进一步编辑标注，导入分词词典中，提高分词系统的准确度，并适应新的语言变化。

Step1: 选择语料源。在“语料源所在路径”输入框中输入或选择需要提取新词的语料所在路径，用户需要事先定义并存放需要处理的语料源，比如：`\NLPIR-Parser\十九大报告`。在“新词存放地址”选择结果文件的存储路径。

如果“语料源所在路径”是通过选择文件夹方式确定，则系统会自动指定“新词存放地址”为`\NLPIR-Parser\output\关键词分析\NewTermlist.txt`；如果“语料源所在路径”是由手动输入，则需要指定输出的“新词存放地址”。

Step2: 点击“新词提取”，系统开始进行发现新词任务。

新词提取结果输出到“新词存放地址”所指定的文件路径，也会输出到结果提示框中。

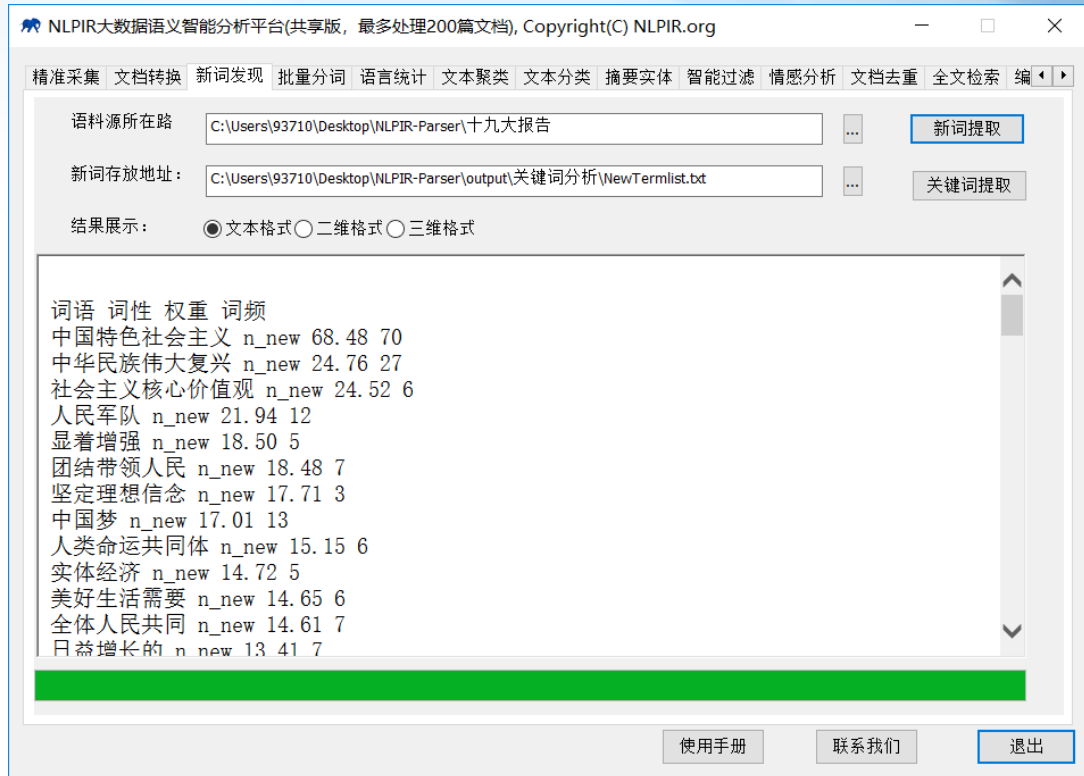


图 3.15 新词提取

新词提取完成后，系统会自动打开结果文件。

NewTermlist.txt - 记事本

文件(F) 编辑(E) 格式(O) 查看(V) 帮助(H)

Word	Part-Of-Speech	Weight	Frequency
中国特色社会主义	n_new	68.48	70
中华民族伟大复兴	n_new	24.76	27
社会主义核心价值观	n_new	24.52	6
人民军队	n_new	21.94	12
显着增强	n_new	18.50	5
团结带领人民	n_new	18.48	7
坚定理想信念	n_new	17.71	3
中国梦	n_new	17.01	13
人类命运共同体	n_new	15.15	6
实体经济	n_new	14.72	5
美好生活需要	n_new	14.65	6
全体人民共同	n_new	14.61	7
日益增长的	n_new	13.41	7

图 3.16 新词结果文件

NewTermlist 是新词提取结果文件。新词提取内容包括：词语、词性、权重和词频统计。

本步骤所得到的新词，可以作为分词标注器的用户词典导入，从

而使分词结果更加准确。对于不需要导入新词的用户，本步骤可以跳过。

3.3.2 关键词提取

关键词提取能够对单篇文章或文章集合，提取出若干个代表文章中心思想的词汇或短语，可用于精化阅读、语义查询和快速匹配等。

Step1: 选择语料源文件夹，以十九大报告为例：\NLPIR-Parser\十九大报告。

Step2: 点击“关键词提取”，系统即可开始进行关键词提取。
关键词存放路径默认为：\NLPIR-Parser\output\关键词分析\Keylist.txt。

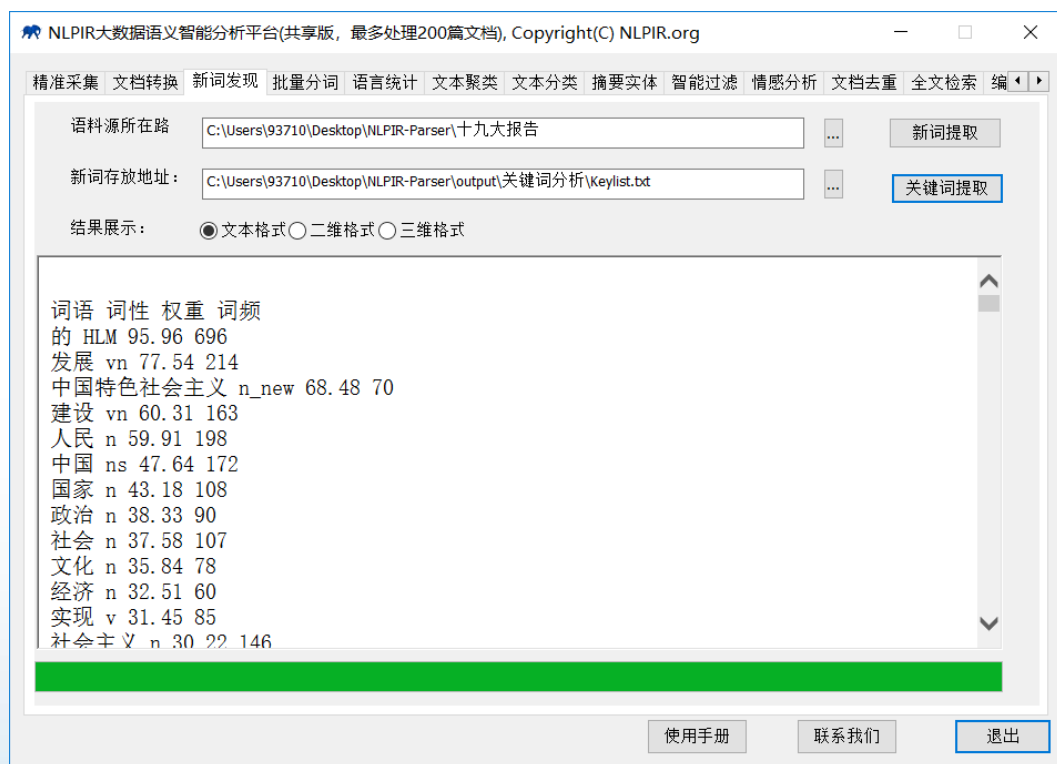


图 3.17 关键词提取

关键词提取完成以后，系统自动打开结果文件 keylist。



Word	Part-Of-Speech	Weight	Frequency
的	HLM	95.96	696
发展	vn	77.54	214
中国特色社会主义	n_new	68.48	70
建设	vn	60.31	163
人民	n	59.91	198
中国	ns	47.64	172
国家	n	43.18	108
政治	n	38.33	90
社会	n	37.58	107
文化	n	35.84	78
经济	n	32.51	60
实现	v	31.45	85
社会主义	n	30.22	146

图 3.18 keylist

关键词分析内容包括：词语、词性、权重和词频统计。系统默认词汇以权重值高低排序。

3.3.3 可视化展示

系统可实现对于新词、关键词提取结果的高维可视化展示，可视化形式有三种：文本格式、二维格式与三维格式。用户可根据需要直接使用，无需再次设计美化。

- 文本格式：以文本的形式展示提取结果

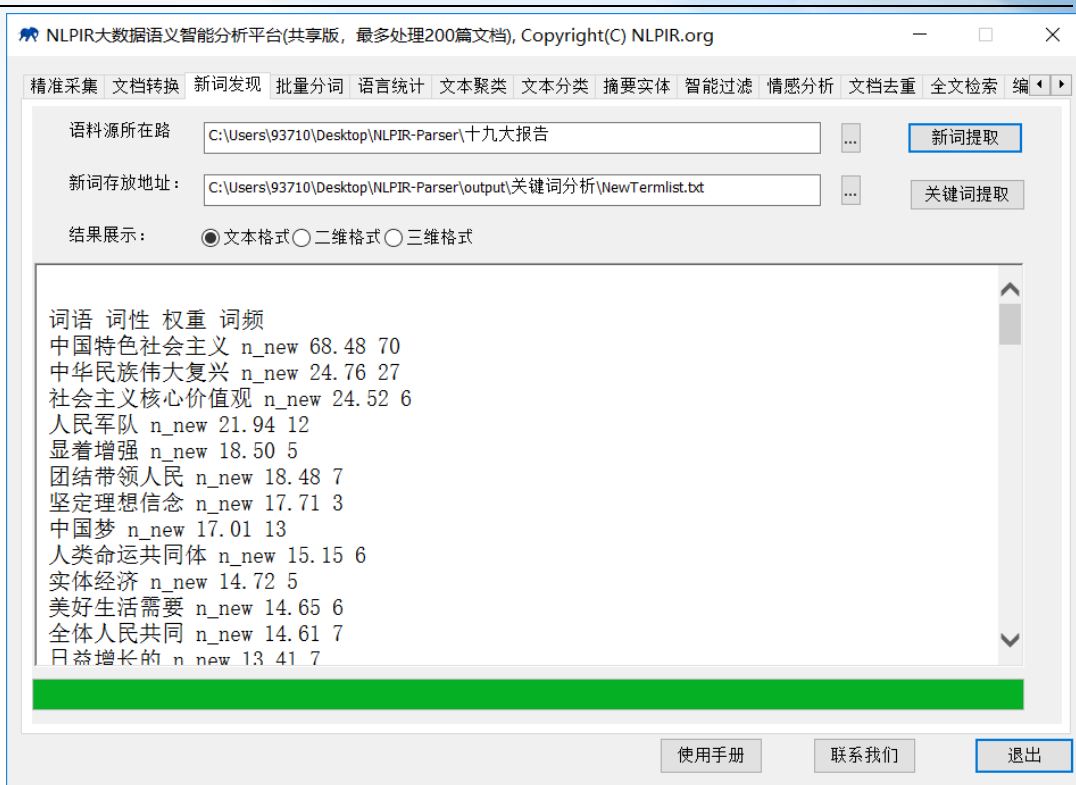


图 3.19 文本格式

➤ 二维格式: top42 词汇的词云形式展示效果, 非常直观。



图 3.20 二维格式

户的自定义词典。例如，将十九大报告提取新词作为用户新词导入。

Step1: 在“新词存放地点”指定新词文件，选择新词提取 new termlist 文件（默认）。点击“编辑”，系统弹出词典文件，用户可对新词文件进行编辑（注：每行一个用户词与词性，系统给出的标注默认为 newword，用户可以根据实际情况进行校对，词性可以标注为任意字符串，系统不做限制）。编辑完成后保存新词文件并关闭。

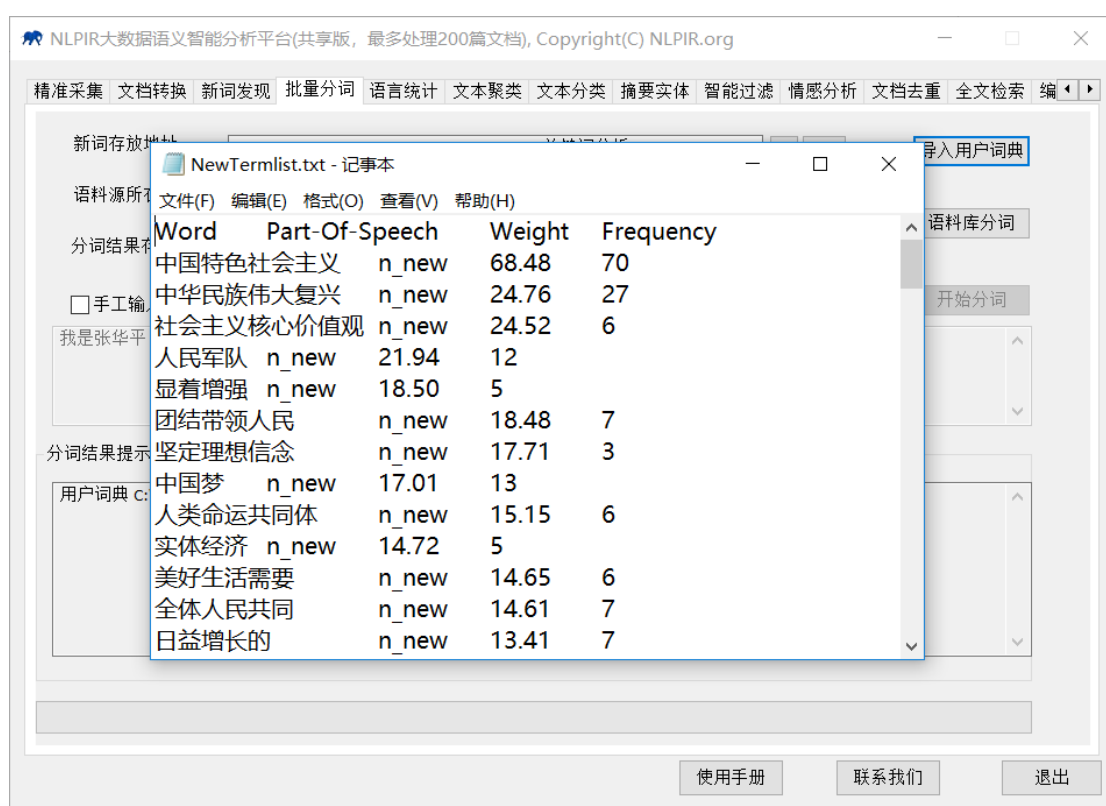


图 3.22 编辑用户词典

Step2: 点击“导入用户词典”，系统开始导入用户词典，并在结果提示框中会显示是否导入成功。对于不需要导入新词的用户，本步骤可以跳过。

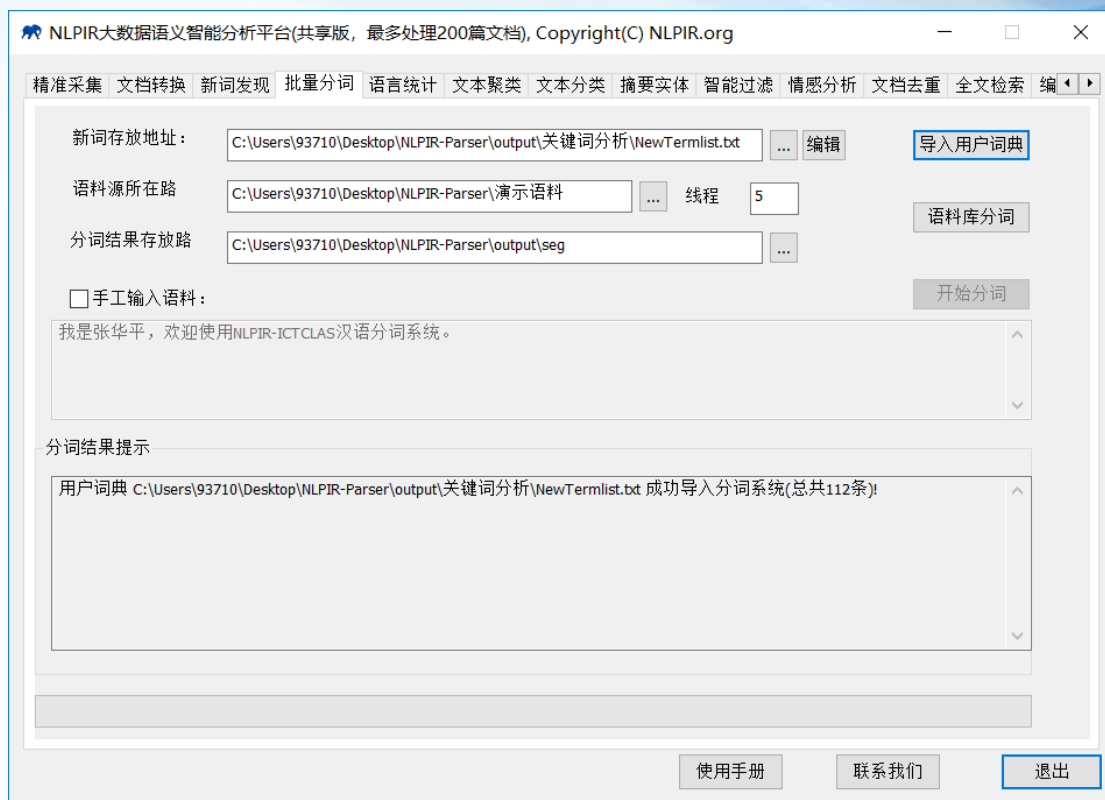


图 3.23 导入用户词典

2) 批量分词

Step1: 选择待分词文件，定义“语料源所在路径”（以十九大报告为例），文件路径：**NLPIR-Parser\十九大报告**。该目录下的语料可以与新词发现中所使用的语料相同，也可以不同，根据用户需求确定。

选择语料源所在路径后，系统会指定默认的“分词结果存放路径”为：**NLPIR-Parser\bin-win64\output\seg**。用户也可以指定其它输出路径。分词及词性标注结果以 **txt** 格式文件存放，文件名与源语料中的文件名一致。

Step2: 点击“语料库分词”，系统开始分词与词性标注。系统会在完成时自动为用户打开分词结果目录。

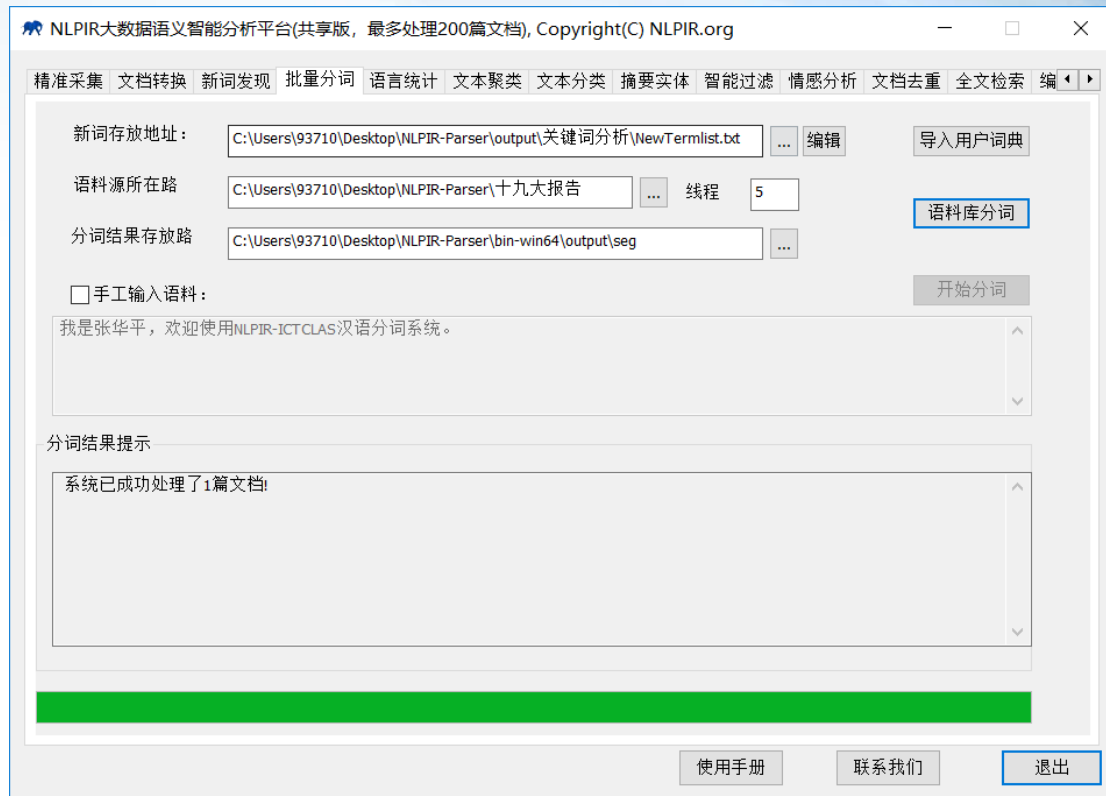


图 3.24 分词成功

分词结果文件地址：NLPIR-Parser\bin-win64\output\seg。

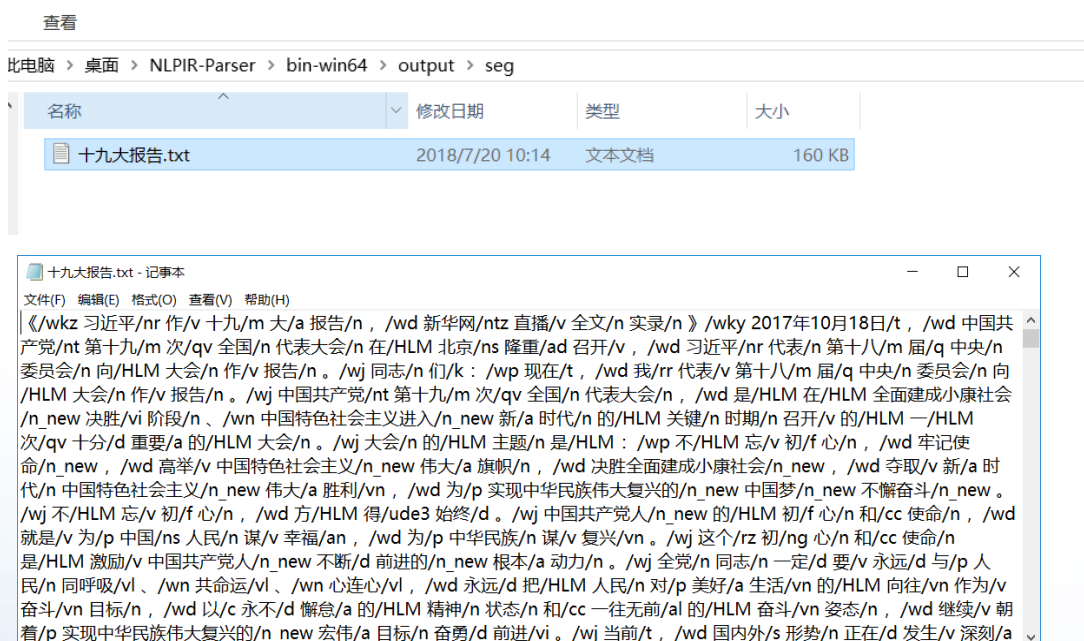


图 3.25 分词结果文件

3) 手动输入语料分词

系统支持用户手动输入语料进行分词。

选择“手动输入语料”，输入语料，点击“开始分词”，系统进行分词。分词结果会呈现在分词结果提示框中。

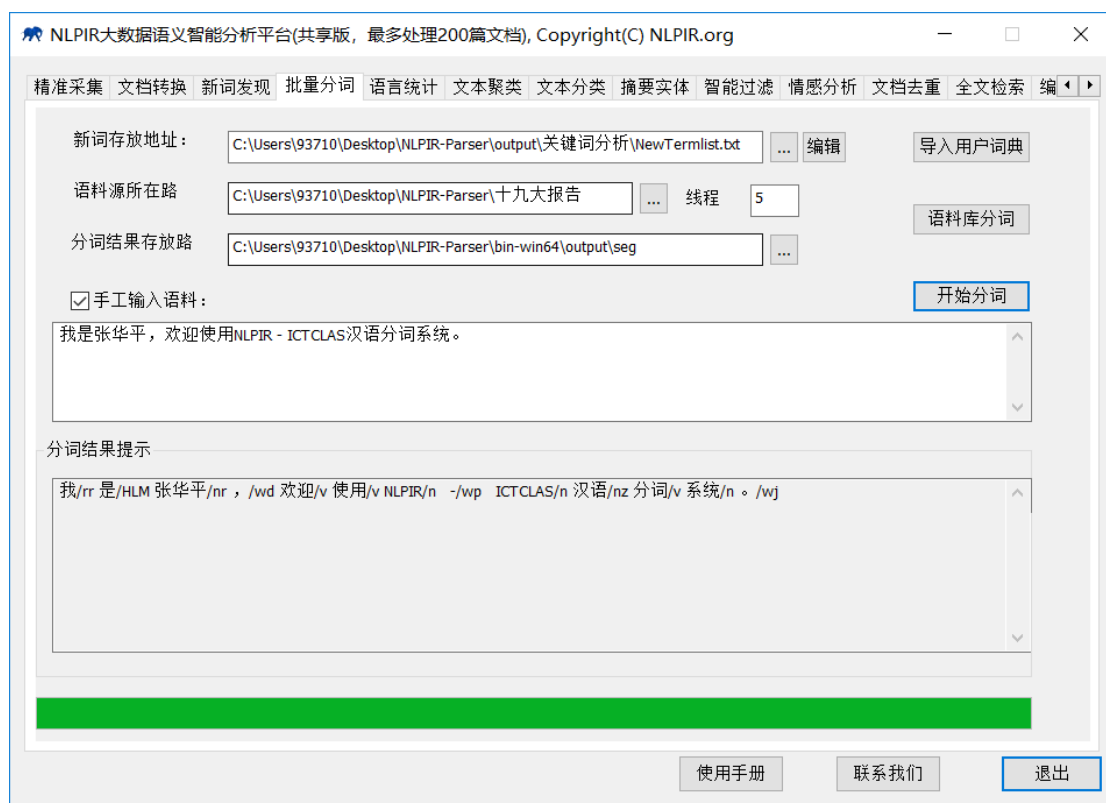


图 3.26 手动输入分词

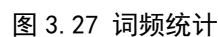
“我是张华平，欢迎使用 NLPIR - ICTCLAS 汉语分词系统”的分词结果为：我/rr 是/vshi 张华平/nr ， /wd 欢迎/v 使用/v NLPIR/n -/wp ICTCLAS/n 汉语/nz 分词/v 系统/n 。 /wj

3.5 语言统计

语言统计功能针对切分标注结果，系统可以自动地进行一元词频统计、二元词语转移概率统计(统计两个词左右连接的频次即概率)。针对常用的术语，会自动给出相应的英文解释。

用户点击“语言统计”，进入系统语言统计功能模块。

Step2: 点击“词频统计与翻译”，系统开始统计词频、共现词对频率等信息。



- ✧ 按字典排序的词频统计文件 "E:\NLPIR-Parser\bin-win64\output\FreqSortByWord.xls"

	A	B	C	D	E	F	G	H	I	J
1	总词数为: 3139, 所有词的平均频率为: 4.355527									
2	词语	词性	词频	一元概率	信息熵					
3	2017年10月18日	t	1	0.000073	0.000697					
4	戮	n	1	0.000073	0.000697					
5	爱	v	5	0.000366	0.002894					
6	爱	vn	1	0.000073	0.000697					
7	爱国	a	3	0.000219	0.001849					
8	爱国人士	n	1	0.000073	0.000697					
9	爱国统一战线	n_new	3	0.000219	0.001849					
10	爱国者	n	1	0.000073	0.000697					
11	爱国主义	n	2	0.000146	0.001292					
12	爱护	v	1	0.000073	0.000697					
13	爱护	vn	1	0.000073	0.000697					
14	安邦定	nr	1	0.000073	0.000697					
15	安定团结	nl	1	0.000073	0.000697					
16	安居乐业	vl	1	0.000073	0.000697					
17	安康	a	1	0.000073	0.000697					
18	安排	v	1	0.000073	0.000697					
19	安排	vn	2	0.000146	0.001292					
20	安全	a	2	0.000146	0.001292					
21	安全	an	49	0.003584	0.020182					
22	安全感	n	1	0.000073	0.000697					
23	安全观	n	2	0.000146	0.001292					

图 3.28 FreqSortByWord

按字典排序词频统计结果包括：词频统计结果（总词数与平均频率）、词语、词性、词频、一元概率与信心熵。其中，一元概率指的是单个词独立出现的概率，信息熵指的是该词包含的信息广

度，其公式为：
$$H(X) = -\sum_{i=1}^n P(X) \log P(X)$$
。

✧ 按词频排序的统计结果文件 "E:\NLPIR-Parser\bin-win64\output\FreqTrans.xls"

按词频排序的统计内容如下，包括：词语、词性、词频、一元概率、信息熵与译文。

1	总词数为: 3139, 所有词的平均频率为: 4.355527						
2	词语	词性	词频	一元概率	信息熵	译文	
3	的	HLM	625	0.045714	0.141043	target; bull's-eye 有~放矢 shoot the arrow at the target; have a definite object in	
4	党	n	195	0.014263	0.060618	① (政党) political party; party ② (指中国共产党) the Party (the Communist Party)	
5	人民	n	151	0.011044	0.049764	the people; popular (adj.) 世界各国~ peoples of the world ~之间的联系和交流 people	
6	是	HLM	145	0.010606	0.048217	① (对, 正确) correct; right ② (表示答应) yes; right ~, 我就来。 Yes, I'm coming	
7	建设	vn	144	0.010532	0.047957	build; construct; construction (n.) 社会主义~ socialist construction ~有中国特色的	
8	坚持	v	131	0.009582	0.044535	persist in; persevere in; uphold; insist on; stick to; adhere to ~原则 adhere to pr	
9	国家	n	105	0.00768	0.037395	country; state; nation 发展中~ developing countries 中等发达~ moderately developed	
10	发展	v	101	0.007387	0.036257	① (变化) develop; expand; grow; development(n.) ~生产力 development of production	
11	在	HLM	94	0.006875	0.034238	① (存在, 生存) exist; be living ② (表示位置) at 在120 公里处 at 120 kilometers 他	
12	社会	n	93	0.006802	0.033947	society; social (adj.) 工业~ industrial society 农业~ agricultural society 社会主	
13	新	a	92	0.006729	0.033654	① (跟“老”或“旧”相对) new; fresh; up-to-date ~发明 a new invention ~技术 new	
14	发展	vn	91	0.006656	0.033361	① (变化) develop; expand; grow; development(n.) ~生产力 development of production	
15	政治	n	90	0.006583	0.033067	politics; political affairs	
16	要	v	90	0.006583	0.033067	① (重要) important; essential ~事 an important matter ② (希望得到) want; ask fo	
17	制度	n	89	0.00651	0.032773	① (规章) rules; regulations 税收~ tax rules and regulations ② (体系) system; in	
18	推进	vi	81	0.005925	0.030385	① (推动前进) push on; carry forward; advance; give impetus to ~国民经济信息化 try	
19	中国	ns	78	0.005705	0.029475	China; Chinese (adj.)	
20	体系	n	77	0.005632	0.02917	system; setup 经济~ economic system 思想~ ideological system	
21	实现	v	74	0.005413	0.028248	realize; achieve; bring about ~工业化和经济的社会化、市场化、现代化 accomplish indu	

图 3.29 FreqTrans.xls

“党”的译文：① (政党) political party; party ② (指中国共产党) the Party (the Communist Party of China) 入~ join the Party 整~ Party consolidation ③ (集团) clique; faction; gang 死~ sworn follower ④ (偏袒) be partial to; take sides with ⑤ (亲族) kinsfolk; relatives 父~ father's kinsfolk。

✧ Bigrams 输出文件

"E:\NLPIR-Parser\bin-win64\output\Bigrams.xls"

Bigrams 结果包括：二元词对总数、前一个词、后一个词、共现频次与二元词对信息熵。共现频次指的是两个词以前后顺序同时出现的频率，二元词对信息熵指的是这两个词包含的信息广度。如下：“党”和“的”以“党的”共现形式出现了 84 词，频率为 0.430769，其信息熵值为 0.031287。

1	二元词对总数为: 11958					
2	前一个词	后一个词	共现频次	二元概率	二元词对信息熵	
3	党	的	84	0.430769	0.031287	
4	。	(50	inf	0.02052	
5	,	坚持	43	inf	0.018122	
6	。	要	38	inf	0.016358	
7	,	是	38	inf	0.016358	
8	新	时代	35	0.380435	0.015277	
9	建设	,	34	0.236111	0.014913	
10	,	推动	31	inf	0.013807	
11	。	我们	30	inf	0.013433	
12	。	加强	29	inf	0.013057	
13	我们	党	28	inf	0.012679	
14	体系	,	28	0.363636	0.012679	
15	制度	,	26	0.292135	0.011914	
16	,	加强	26	inf	0.011914	
17	,	必须	25	inf	0.011528	
18	,	在	24	inf	0.011138	
19	,	推进	24	inf	0.011138	
20	,	完善	22	inf	0.01035	

图 3.30 Bigrams.xls

◇ 文件统计信息输出文件 "E:\NLPIR-Parser\bin-win64\output\FileStat.xls"

文档名	总词频	总词数	用户词典总词频	用户词典总词数
十九大报告	13670	3137	1823	163

图 3.31 FileStat.xls

文件统计结果包括：文档名、总词频、总词数、用户词典总词频与用户词典总词数。

3.6 文本聚类

文本聚类能够从大规模数据中自动分析出热点事件，并提供事

件话题的关键特征描述。文本聚类适用于长文本和短信、微博等短文本的热点分析。

用户点击“文本聚类”，进入系统文本聚类功能模块。

Step1：“在语料源所在路径”选择语料源文件夹。

Step2：设置聚类参数（最大类数目与类中最大文档数），点击“聚类”，系统进行聚类。聚类结果自动保存在 output 目录下：

\\NLPIR-Parser\\output\\聚类结果，并于结果提示框呈现聚类结果。如下所示：

共有 24 个聚类类别，第一类别主题词汇：区块链 云平台 互联网平台 制造业增加值，第一类别文档总数 14 篇。



图 3.32 聚类

聚类结果文件有两种形式：网页和文件，都保存到本地文件夹：\\NLPIR-Parser\\output\\聚类结果，文件夹按照文件数量排序，名称包含

信息：文件数量与聚类特征词。

电脑 > 桌面 > NLPIR-Parser > output > 聚类结果

名称	修改日期	类型	大小
DocCount-2-沙雅赫梅托娃 哈萨克斯坦 ...	2018/7/20 10:36	文件夹	
DocCount-14-区块链 云平台 互联网平...	2018/7/20 10:38	文件夹	
ClusterResult.xml	2018/7/20 10:38	XML 文档	4 KB

图 3.33 聚类结果文件

用户可查看同属一个类别的多个文件。聚类详情文件名称包含：聚类特征词、媒体来源与新闻标题。

电脑 > 桌面 > NLPIR-Parser > output > 聚类结果 > DocCount-14-区块链 云平台 互联网平台 制造业增加值

名称	修改日期	类型	大小
区块链-CCTIME飞象网-报告称：未来五...	2018/7/20 10:38	文本文档	2 KB
区块链-白鲸出海-狙击IBM甲骨文亚马逊...	2018/7/20 10:38	文本文档	7 KB
区块链-北国网-开启区块链3.0金窝窝助...	2018/7/20 10:38	文本文档	3 KB
区块链-创业邦-谷歌创始人：我们在区块...	2018/7/20 10:38	文本文档	2 KB
区块链-东方资讯-谛听科技newifi云计算...	2018/7/20 10:38	文本文档	5 KB
区块链-环球网-汇桔网谢旭辉：用户的知...	2018/7/20 10:38	文本文档	6 KB
区块链-人民网-2018汽车后市场产业+区...	2018/7/20 10:38	文本文档	9 KB
区块链-人民网-人民日报-高质量发展，...	2018/7/20 10:38	文本文档	7 KB
区块链-人民网-人民日报海外版-探求金...	2018/7/20 10:38	文本文档	5 KB
区块链-人民网-人民日报-人民日报新论...	2018/7/20 10:38	文本文档	6 KB
区块链-腾讯财经-蚂蚁金服韩歆毅：ICO...	2018/7/20 10:38	文本文档	12 KB
区块链-腾讯大粤网-2018区块链场景应...	2018/7/20 10:38	文本文档	4 KB
区块链-证券时报-e公司-东旭蓝天：增资...	2018/7/20 10:38	文本文档	1 KB
一带一路-腾讯大秦网-首届“一带一路”科...	2018/7/20 10:38	文本文档	5 KB

图 3.34 聚类详情文件

3.7 文本分类

文本分类能够根据事先指定的规则和示例样本，自动从海量文档中识别并训练分类。NLPIR 深度文本分类，可以用于新闻分类、简历分类、邮件分类、办公文档分类、区域分类等诸多方面。此外还可以实现文本过滤，能够从大量文本中快速识别和过滤出符合特殊要求的信息，可应用于品牌报道监测、垃圾信息屏蔽、敏感信息审查等领域。

NLPIR 采用深度神经网络对分类体系进行了综合训练。演示平台目前训练的类别只是新闻的政治、经济、军事等。我们内置的算法支持类别自定义训练，该算法对常规文本的分类准确率较高，综合开放测试的 F 值接近 86%。

用户点击“文本分类”，进入系统文本分类功能模块。

文本分类有两种模式：专家规则分类与机器学习分类。

专家规则分类指的是根据事先人为制定的分类规则进行分类，比如“毒品”类别，我们可定义该类别的规则：“海洛因 快乐丸 罂粟可卡因 Pethidine 摇头丸 K 粉”，系统会根据文本中出现规（词）判定文本类别为：毒品。

机器学习分类是利用机器自动学习的能力，通过大量文本的训练，是系统具有分类的能力。比如我们准备军事、政治类别的大量语料，通过训练，机器自动学习类别特征，经过不断的语料训练，分类效果越来越精准。

➤ 专家规则分类

Step1:选择测试语料(分类测试语料为例)，点击“帮助”。

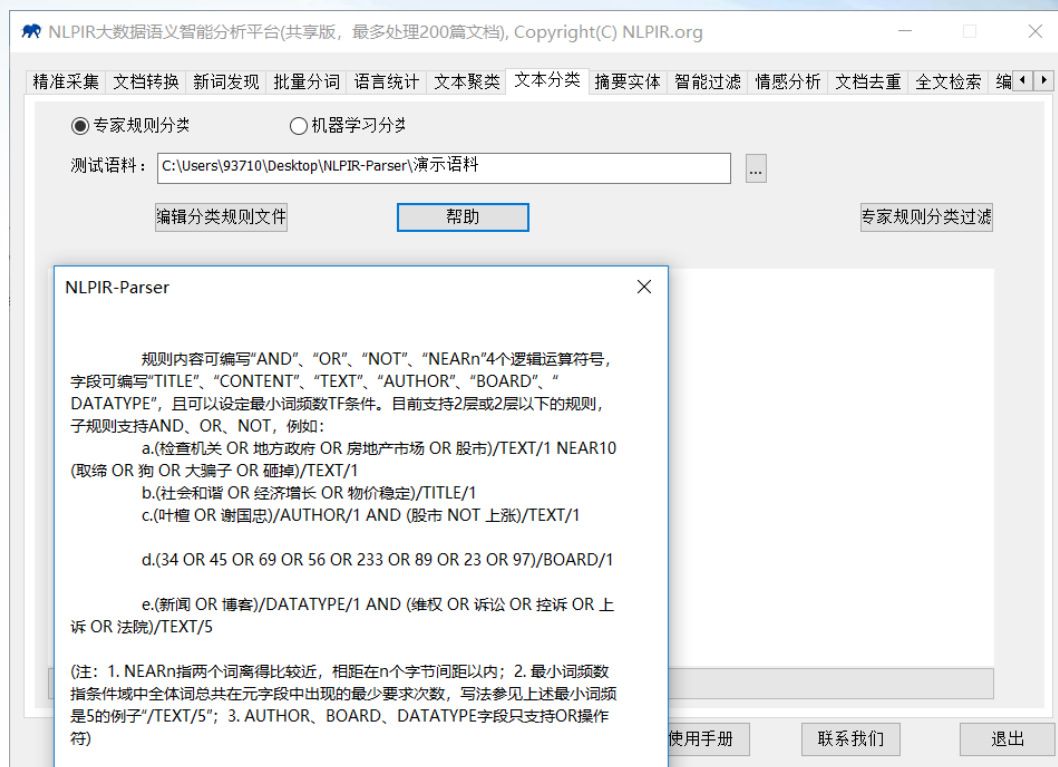


图 3.35 分类规则“帮助”

“帮助”里面详细介绍了分类规则的书写格式。例如：

(检察机关 OR 地方政府 OR 房地产市场 OR 股市) /TEXT/1
NEAR10(取缔 OR 狗 OR 大骗子 OR 砸掉) /TEXT/1

表示：在原文中，这两组词（组内任意一词与另一组内任意一词）距离在 10 个字节（5 个中文字）以内，共现频率至少为 1 次。

Step2: 点击“编辑分类规则文件”，系统弹出规则分类文件。用户可自定义分类规则。系统有默认的 rulelist，用户可在此基础上添加、修改、删除分类的规则。用户编辑完分类规则文件需要保存。

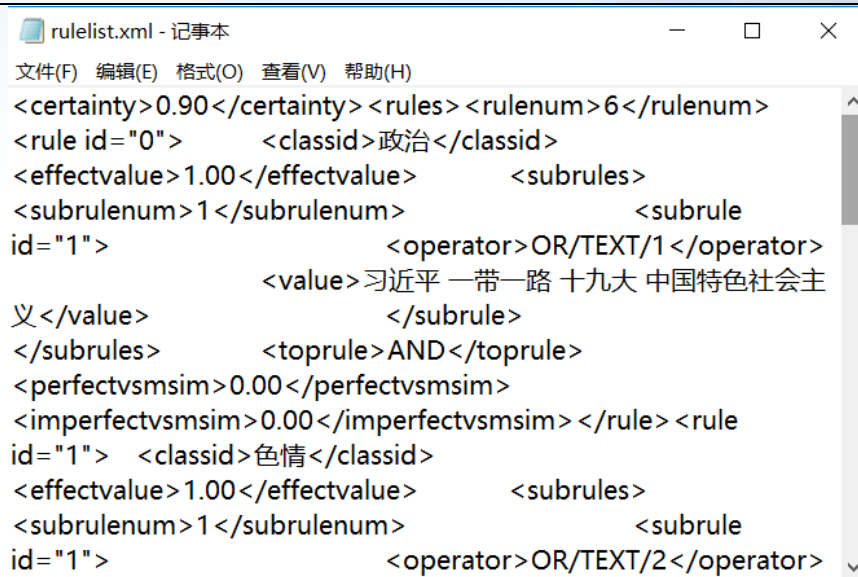


图 3.36 编辑分类规则文件

Step3: 点击“专家规则分类过滤”，系统进行分类分析。分类结果同样会呈现在结果提示框中，如下所示：



图 3.37 专家规则分类过滤

系统会将分类结果以网页和文件的同时自动保存至：\NLPIR-Parser\output\专家规则分类结果文件夹，并在分类结束后自动打开，

用户可直接查看与利用分类结果。

电脑 > 桌面 > NLPIR-Parser > output > 专家规则分类结果

名称	修改日期	类型	大小
政治	2018/7/19 14:59	文件夹	
RuleClassifier_Result.xml	2018/7/20 10:49	XML 文档	1 KB

图 3.38 分类过滤结果文件

➤ 机器学习分类

Step1:选择训练分类，点击“训练”按钮，系统进行类别特征的自学习；系统目前已有 7 大类分类训练语料，用户仍可自定义在此基础上进行语料的更新或类别的定义。

> NLPIR-Parser > 训练分类用文本

名称	修改日期	类型
交通	2017/12/11 18:09	文件夹
教育	2017/12/11 18:09	文件夹
经济	2017/12/11 18:09	文件夹
军事	2017/12/11 18:09	文件夹
体育	2017/12/11 18:09	文件夹
艺术	2017/12/11 18:09	文件夹
政治	2017/12/11 18:09	文件夹

图 3.39 训练分类用文本



图 3.40 训练

如上所示，系统将训练结果以网页的形式呈现在提示框中，总计频率为 186964，共有 1000 个特征词，第一个特征词为“会谈”，在 9 篇文档中出现共 22 次，权重值为 11。

Step2: 选择测试语料（以十九大报告为例），点击“机器学习分类过滤”，系统进行分类分析。分类结果如下：

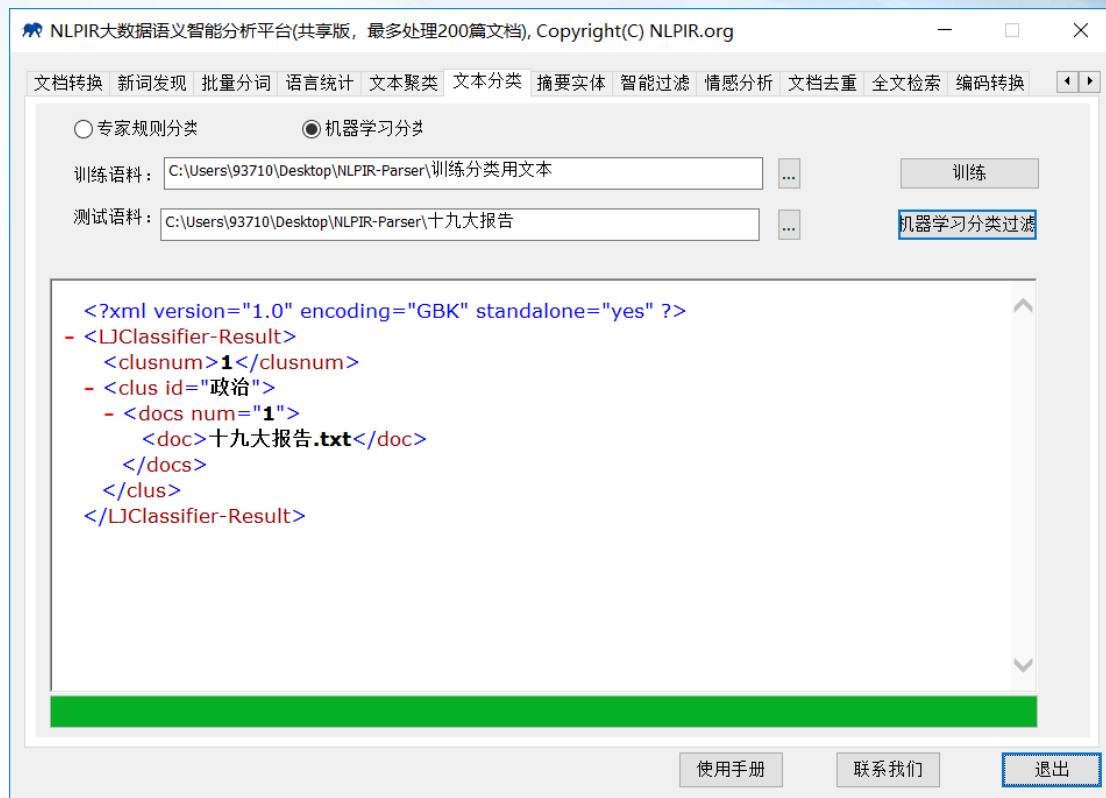


图 3.41 分类过滤

系统会将分类结果以网页和文件的同时自动保存至：\NLPIR-Parser\output\机器学习分类结果，并自动打开文件夹。



图 3.42 分类结果文件

3.8 摘要实体

自动摘要能够对单篇或多篇文章，自动提炼出内容的精华，方便用户快速浏览文本内容。实体提取能够对单篇或多篇文章，自动提炼

出内容摘要，抽取人名、地名、机构名、时间及主题关键词；方便用户快速浏览文本内容。

用户首先点击“摘要实体”，进入系统摘要实体功能模块。

Step1: 选择语料源目录（以十九大报告为例）；用户可自定义摘要长度（默认最大为 250），摘要最大压缩率和关键词数量；

Step2: 点击“摘要与实体抽取”，系统进行提取与分析，并显示摘要和关键词的结果。点击“上一篇”、“下一篇”按钮，可实现结果的快速浏览。



图 3.43 摘要与实体抽取

摘要实体结果包括：自动摘要和实体抽取（关键词、人、时间、地点、国家与机构）。

十九大报告分析结果：

摘要（摘要长度定义为 300 的结果）：要长期坚持、不断发展我国社会主义民主政治，积极稳妥推进政治体制改革，推进社会主义民主政治制度化、规范化、法治化、程序化，保证人民依法通过各种途径和形式管理国家事务，管理经济文化事业，管理社会事务，巩固和发展生动活泼、安定团结的政治局面。成立中央全面依法治国领导小组，加强对法治中国建设的统一领导。

实体抽取：

关键词（关键词数量定义为 10 的分析结果）：发展#建设#人民#中国#国家#政治#社会#文化#经济#创新#

时间：2017 年 10 月 18 日#现在#当前#近代#一九二一年#一九四九年#今天#未来#本世纪中叶#千年#二〇二〇年#二〇三五年#现代#当今#冬#当代#清明#

国：中国#

人物：习近平#金山银#向发力#德治相#言代法#安邦定#强国强#来海#晏河清#高强#

地点：中国#台湾#北京#京津冀#中华人民共和国#惠民#澳门#香港#长江#澳门特别行政区#亚洲#杭州#香港特别行政区#厦门#南海#古田#亚丁#安新#

机构：中国共产党#党中央#联合国#中共中央#

3.9 智能过滤

智能过滤能够对文本内容进行语义智能过滤审查，内置国内最全

词库，智能识别多种变种：形变、音变、繁简等多种变形，且实现语义精准排歧。

用户首先点击“智能过滤”，进入系统智能过滤功能模块。

（1）导入关键词

系统已内置约 10 类近 4 万关键词，用户仍可根据需求添加自己的关键词。

Step1: 选择关键词文件，在“关键词列表文件”中选择文件，点击“编辑”，系统弹出关键词文件，用户可编辑关键词列表，编辑完成后保存并关闭文件。

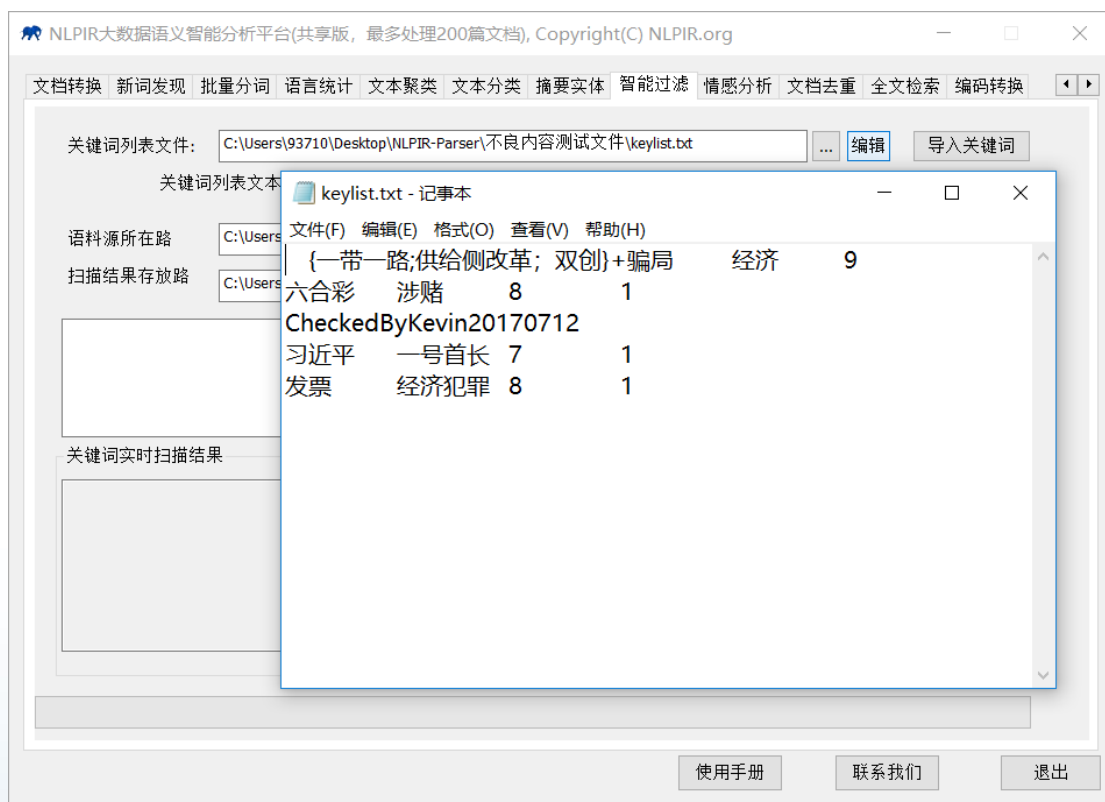


图 3.44 编辑关键词

Step2: 点击“导入关键词”，系统显示导入关键词成功。

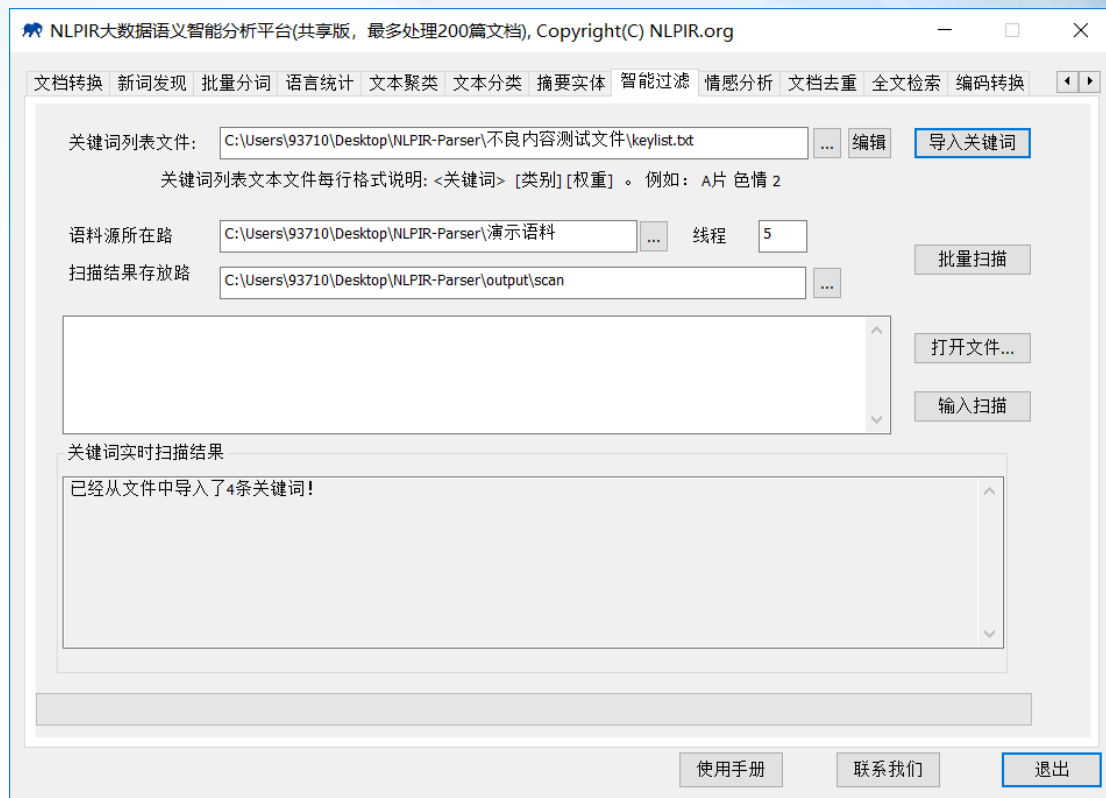


图 3.45 导入关键词成功

（2）批量扫描

Step1: 选择语料源: \NLPIR-Parser\不良内容测试文件(系统默认, 用户可定义自己的过滤语料), 系统会指定默认的“扫描结果存放路径”为: \NLPIR-Parser\output\scan。用户也可以指定其它输出路径。

Step2: 点击“批量扫描”, 系统开始进行不良信息过滤。

智能过滤扫描结果以 txt 格式文件存放, 文件名与源语料中的文件名一致。扫描统计结果 KeyScanStatResult.xls 放入 NLPIR-Parser\output 目录下并自动打开。扫描详情结果存放路径: \NLPIR-Parser\output\scan, 扫描完成时自动为用户打开该目录。

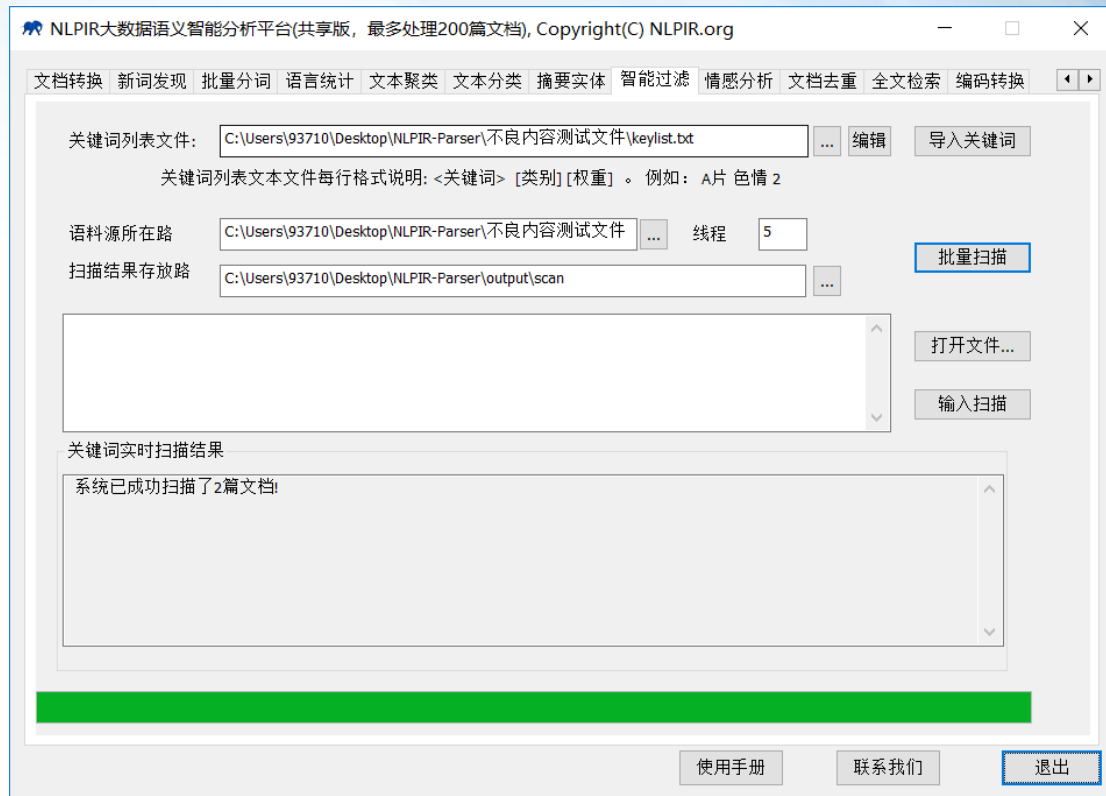


图 3.46 批量扫描

KeyScanStatResult.xls 包括：关键词、类别、权重与命中次数。
扫描详情文件会在原文中标出扫描结果。

测试时间:	Fri Jul 20 11:06:17 2018			
扫描的记录:	0	条记录		
扫描所花:	1971.02	秒		
处理速度:	0	条/秒		
命中的规则:	10	命中的记录:	0	疑似敏感率:
规则编号	关键词	类别	权重	命中次数
10277	犯罪	敏感词	1	1
16636	六合彩	涉赌	8	1
19314	骗局	complex	-1	1
21446	涉赌	涉赌	2	1
22095	首长	涉领导人	7	1
30905	习近平	complex	-1	1
45769	发票	经济犯罪	8	1
45770	一带一路	complex	-1	1
45771	供给侧改	complex	-1	1
45772	双创	complex	-1	1

图 3.47 扫描过滤结果统计

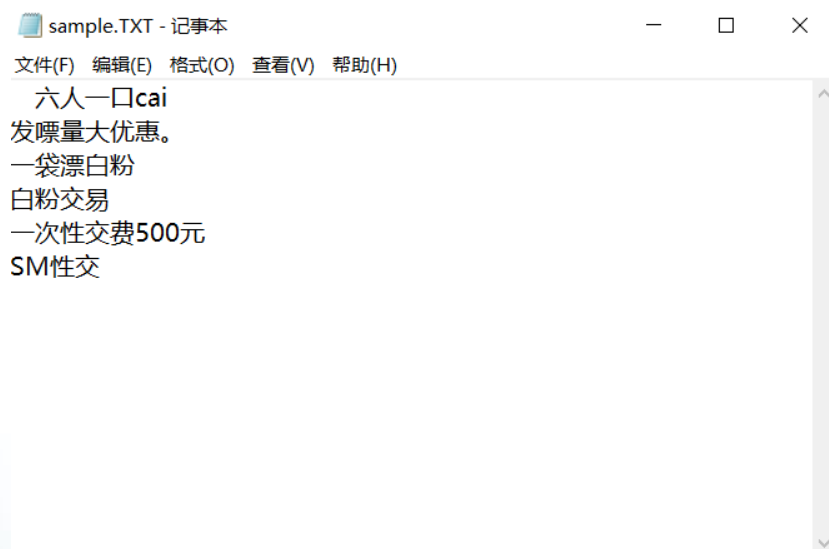


图 3.48 原文扫描结果

(3) 输入扫描

Step1: 点击“打开文件”或者直接将扫描文本粘贴至文本框中；**Step2:** 点击“输入扫描”，结果如下：

输入文本：六人一口彩（六合彩的形变）和法轮功

扫描结果：不良得分、命中不良内容、不良类别与命中文字

六人一口彩：

不良得分：16，命中不良内容：[形变]六合彩→六人一口彩，

不良类别：涉赌，命中文字：六人一口彩。

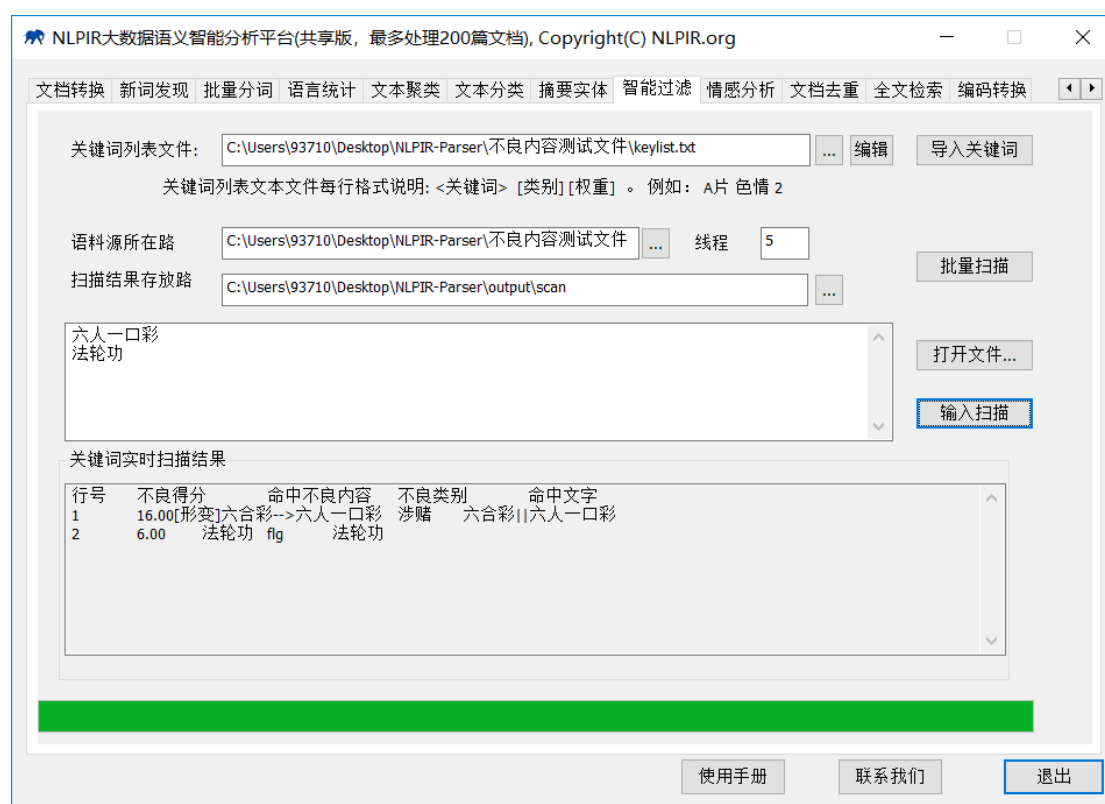


图 3.49 输入扫描

3.10 情感分析

情感分析，针对事先指定的分析对象，系统自动分析海量文档的情感倾向：情感极性及其情感值测量，并在原文中给出正负面的得分和句子样例。NLPIR 情感分析的情感分类丰富，不仅包括正、负两面，还包括好、乐、惊、怒、恶、哀和惧的具体情感属性。NLPIR 还提供关于特定人物的情感分析，并能计算正负面的具体得分。

用户首先点击“情感分析”，进入系统情感分析功能模块。

➤ 单个分析：对单个对象做情感分析

Step1：选择语料源；输入分析对象：区块链；

Step2：点击“单个分析”，系统开始以“区块链”为分析对象进行情感分析。

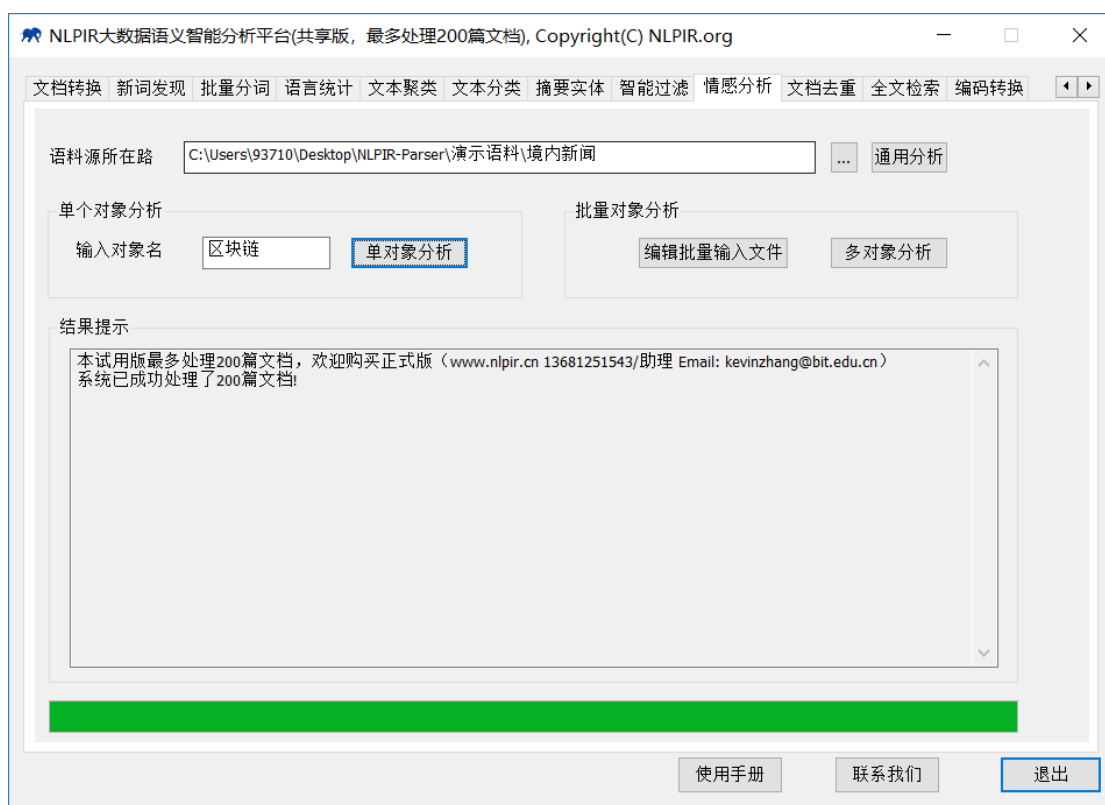


图 3.50 情感分析

情感分析结果默认存放路径：NLPIR-Parser\output，情感分析有两个分析结果，sentiment-rank.xls(系统分析完毕后自动打开)和sentiment-detail.txt，前者是统计结果，后者是分析详情结果。

情感分析统计结果包括：文档总数、正面数量及占比，每一篇文档的正负面得分与排序。情感分析详情结果会在原文本中显示情感分析的详情：对象、得分、原文等。

情感分析-区块链											
文档标题		200	负面总数	2	负面占比	1.00%	正面总数	12	正面占比	6.00%	
标题	出处	发表时间	情感得分	正面得分	负面得分	原始链接	本地文件名				
银保监会 百度新闻搜索-网易		2018/7/8 15:37	-2	1	-3	http://mc	区块链-财联社-银保监会国际部负责人范文仲：不要神化区块链.txt				
区块链板 百度新闻搜索-网易		2018/7/10 13:49	-1	0	-1	http://mc	区块链-上海证券报·中国证券网-区块链板块拉升宣亚国际涨停.txt				
相关传销 百度新闻搜索-网易		2018/7/10 10:14	1	77	-76	http://ne	区块链-人民网-相关传销平台已超3000家“区块链”竟成“上亿大坑”.txt				
蚂蚁金服 百度新闻搜索-新浪		2018/7/8 15:03	6	10	-4	http://fj	区块链-新浪财经-蚂蚁金服韩敏毅：ICO不是真正的区块链金融服务.txt				
【三言】 百度新闻搜索-新浪		2018/7/6 8:31	11	15	-4	http://fj	区块链-三言财经-【三言】中国区区块链50城之长沙：将设30个产业基金.txt				
世界热点 百度新闻搜索-新浪		2018/7/7 2:04	17	57	-40	http://fj	区块链-证券日报-世界杯点燃区块链预测“引线”宣称“只赢不输”难避.txt				
快讯：区 百度新闻搜索-新浪		2018/7/10 13:29	18	19	-1	http://fj	区块链-新浪财经-快讯：区块链概念股集体走强宣亚国际涨停.txt				
区块链投 百度新闻搜索-网易		2018/7/5 8:32	20	38	-18	http://mc	区块链-新京报-区块链投资切勿盲目跟风“网红”的诱惑.txt				
合茂数字 百度新闻搜索-腾讯		2018/7/2 14:59	28	31	-3	http://bj	一带一路-北国网-合茂数字发起成立“一带一路数字经济联盟”.txt				
梅地卡尔 百度新闻搜索-网易		2018/6/29 19:50	38	52	-14	http://mc	区块链-证券日报-梅地卡尔利用区块链技术发力医疗数据交换系统.txt				
区块链+人 百度新闻搜索-网易		2018/7/6 9:21	39	55	-16	http://mc	区块链-钛媒体-区块链+人工智能如何重新定义世界黑科技?.txt				
区块链跨 百度新闻搜索-新浪		2018/7/5 9:18	41	66	-25	http://fj	区块链-每日经济新闻-区块链跨国“国域”：中韩两国都希望寻求外部空.txt				
迅雷发布 百度新闻搜索-新浪		2018/7/9 18:18	54	86	-32	http://fj	区块链-21世纪经济报道-迅雷发布区块链文件协议TCFS陈磊：技术会发生.txt				
央行数字货币 百度新闻搜索-网易		2018/7/1 8:13	92	113	-21	http://mc	区块链-21世纪经济报道-央行数字货币研究所蒋国庆：法定数字货币与发				

图 3.51 sentiment-rank

```
sentiment_SingleObject_Detail.txt - 记事本
文件(F) 编辑(E) 格式(O) 查看(V) 帮助(H)
区块链-证券日报-梅地卡尔利用区块链技术发力医疗数据交换系统.txt 的分析结果:
<?xml version="1.0" encoding="utf-8" standalone="yes"?>
<LJSentiment-Result>
  <result>
    <object>区块链</object>
    <polarity>38.00</polarity>
    <positivepoint>52.00</positivepoint>
    <negativepoint>-14.00</negativepoint>
    <sentenceclue>
      <titlesentenceclue><![CDATA[
<object>区块链</object>-证券日报-梅地卡尔<neg value="-1">利用</neg><object>区块链</object><pos value="1">技术</pos>发力医疗数据交换<pos value="1">系统</pos>]]></titlesentenceclue>
      <contentsentenceclue><![CDATA[
<title>梅地卡尔<neg value="-1">利用</neg><object>区块链</object><pos value="1">技术</pos>发力医疗数据交换<pos value="1">系统</pos><!--<hr/>
</wtitle><href>http:money.html</whref><publish_time>2018-06-29 19:50:00</wpublish_time><source>证券日报</wsource><file_name><object>区块链</object>-证券日报-梅地卡尔<neg value="-1">利用</neg><object>区块链</object><pos value="1">技术</pos>发力医疗数据交换
```

图 3.52 sentiment-detail

对象：区块链，情感得分：38，正面得分：52，负面得分：-14

► 批量对象分析

Step1: 选择语料源; 点击“编辑批量输入文件”, 用户可自定义多个分析对象与分析条件。

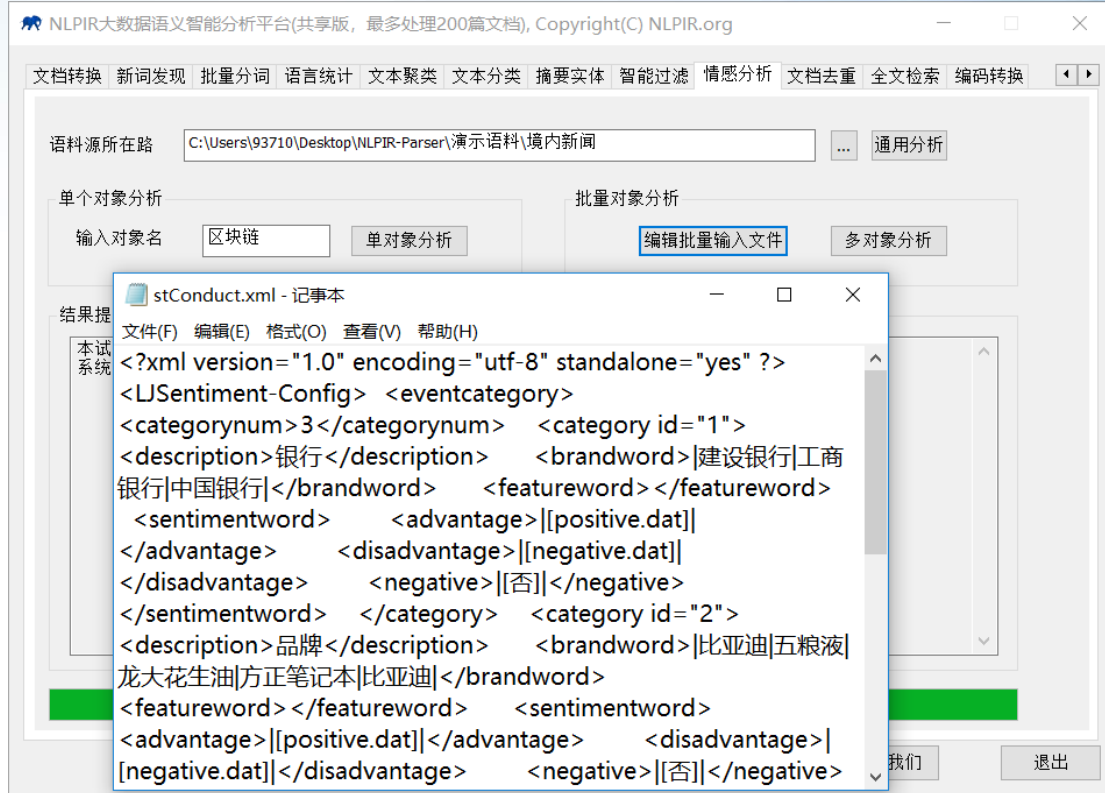


图 3.53 定义批量对象

Step2: 点击“多对象分析”，系统开始对多个对象进行情感分析。

批量分析同样有两个结果文件(\\NLPIR-Parser\\output), sentiment-rank.xls(系统分析完毕后自动打开)和 sentiment-detail.txt, 前者是统计结果, 后者是分析详情结果。

[illegible]

图 3.54 情感批量分析结果

3.11 文档去重

文档去重能够快速准确地判断文件集合或数据库中是否存在相同或相似内容的记录，同时找出所有的重复记录。

用户首先点击“文档去重”，进入系统文档去重功能模块。

Step1：选择语料源；选择结果文件存放路径。

Step2：点击“开始查重”，系统即刻开始查重处理，并输出查重结果文件 RepeatFile（NLPIR-Parser\bin-win64\output\RepeatFile.txt）
查重结果会显示在结果提示框中，如下图所示：

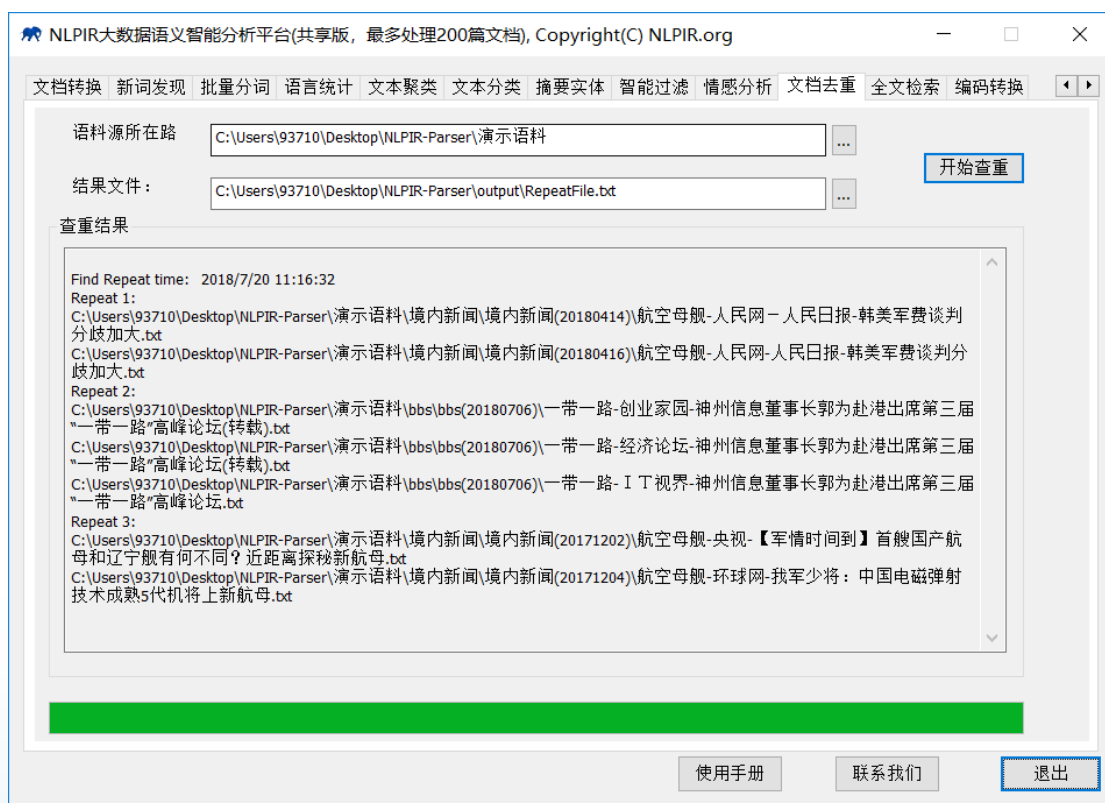


图 3.55 文档去重

RepeatFile 文档去重分析结果包括：重复文档数量统计(共有 5 片文档重复)，重复文档标题与重复文档路径。

```
Find Repeat time: 2018/7/20 11:16:32
Repeat 1:
C:\Users\93710\Desktop\NLPIR-Parser\演示语料\境内新闻\境内新闻(20180414)\航空母舰-人民网-人民日报-韩美军费谈判分歧加大.txt
C:\Users\93710\Desktop\NLPIR-Parser\演示语料\境内新闻\境内新闻(20180416)\航空母舰-人民网-人民日报-韩美军费谈判分歧加大.txt
Repeat 2:
C:\Users\93710\Desktop\NLPIR-Parser\演示语料\bbs\bbs(20180706)\一带一路-创业家园-神州信息董事长郭为赴港出席第三届“一带一路”高峰论坛(转载).txt
C:\Users\93710\Desktop\NLPIR-Parser\演示语料\bbs\bbs(20180706)\一带一路-经济论坛-神州信息董事长郭为赴港出席第三届“一带一路”高峰论坛(转载).txt
C:\Users\93710\Desktop\NLPIR-Parser\演示语料\bbs\bbs(20180706)\一带一路-IT视界-神州信息董事长郭为赴港出席第三届“一带一路”高峰论坛.txt
Repeat 3:
C:\Users\93710\Desktop\NLPIR-Parser\演示语料\境内新闻\境内新闻(20171202)\航空母舰-央视-【军情时间到】首艘国产航母和辽宁舰有何不同? 近距离探秘新航母.txt
C:\Users\93710\Desktop\NLPIR-Parser\演示语料\境内新闻\境内新闻(20171204)\航空母舰-环球网-我军少将: 中国电磁弹射技术成熟5代机将上新航母.txt
```

图 3.56 RepeatFile

3.12 全文检索

全文检索支持文本、数字、日期、字符串等各种数据类型，多字段的高效搜索，支持 AND/OR/NOT 以及 NEAR 邻近等查询语法，支持维语、藏语、蒙语、阿拉伯、韩语等多种少数民族语言的检索。可以无缝地与现有文本处理系统与数据库系统融合。

支持的典型查询语法包括：

Sample1: [FIELD] title [AND] 解放军

Sample3: [FIELD] content [AND] 甲型 H1N1 流感

Sample4: [FIELD] content [NEAR] 张雁灵 解放军

Sample5: [FIELD] content [OR] 解放军 甲流

Sample6: [FIELD] title [AND] 解放军 [FIELD] content [NOT]

甲流

用户首先点击“全文检索”，进入系统全文检索功能模块。

➤ 建立索引

Step1: 选择语料文件夹（以十九大报告为例）；

Step2: 选择是否“增量”，增量是指在历史索引的基础上需要对新增部分文件的内容建立索引。系统在历史索引基础上新增索引，不选择增量，系统将以预料源为基础重新建立索引。点击“建立索引”，

系统对语料快速建立压缩索引。

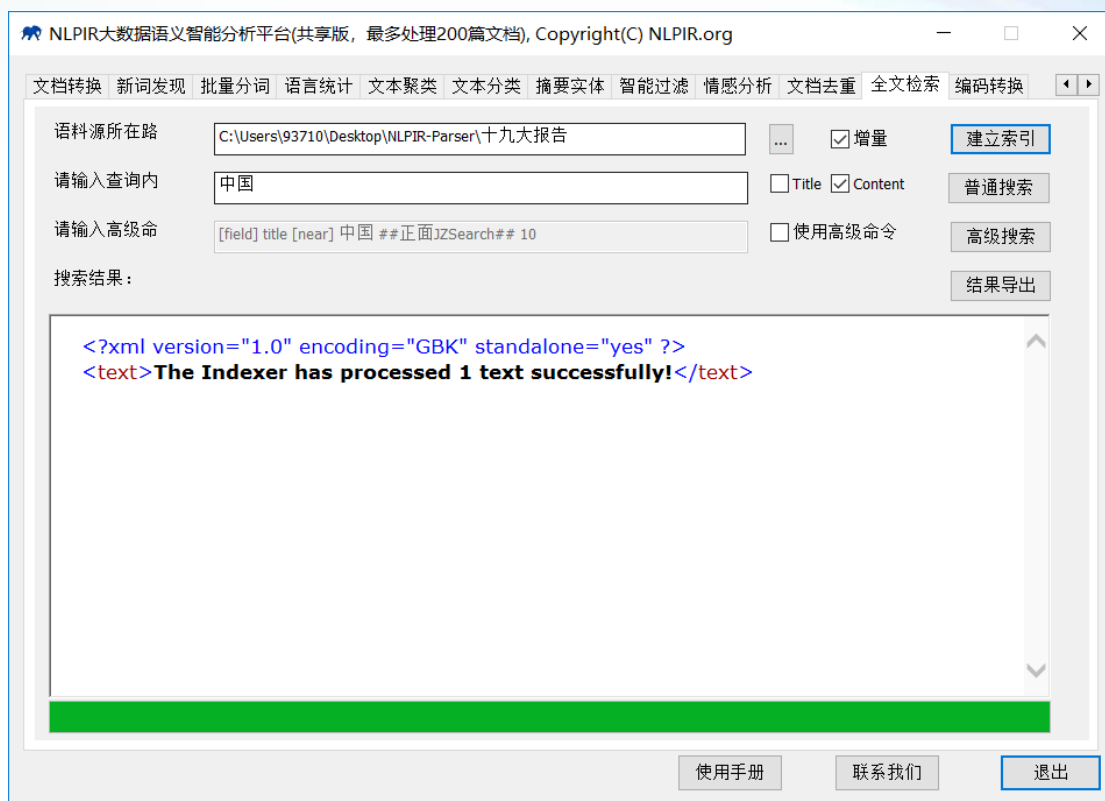


图 3.57 建立索引

➤ 普通检索

Step1: 输入查询关键词（中国），选择“Title”（标题查询）与“content”（内容查询），两者可同时选择。

Step2: 点击“普通检索”。搜索结果框会呈现查询结果，并配以相似得分。检索结果文件（\NLPIR-Parser\output\搜索结果JZSearch-result）以网页形式保存。

检索结果包括：文档总量统计、标题、内容与相似得分。



图 3.58 普通检索

Step3: 点击“结果导出”，系统将检索目标文档导出\NLPIR-Parser\output\搜索结果\中国，并自动打开文件目录。



图 3.59 结果导出

电脑 > 桌面 > NLPIR-Parser > output > 搜索结果 > 中国

名称	修改日期	类型	大小
航空母舰-KT天鹰营销-厉害了--康婷, 未...	2018/7/20 11:21	文本文档	5 KB
航空母舰-参考消息-港媒: 中国首艘国产...	2018/7/20 11:21	文本文档	3 KB
航空母舰-参考消息-美媒分析: 中国会否...	2018/7/20 11:21	文本文档	4 KB
航空母舰-参考消息网-美军官自称其航母...	2018/7/20 11:21	文本文档	4 KB
航空母舰-参考消息网-美专家预测: 中国...	2018/7/20 11:21	文本文档	6 KB
航空母舰-车买买网-看完中国航母下水, ...	2018/7/20 11:21	文本文档	9 KB
航空母舰-传媒-如何讲好我们的故事.txt	2018/7/20 11:21	文本文档	15 KB
航空母舰-第一电动-电动出游并不难 比...	2018/7/20 11:21	文本文档	10 KB
航空母舰-封面-中国航母史 两艘航母为...	2018/7/20 11:21	文本文档	10 KB
航空母舰-观察者网-美军罗斯福号航母今...	2018/7/20 11:21	文本文档	5 KB
航空母舰-光明日报-永不褪色的精神礼赞...	2018/7/20 11:21	文本文档	30 KB
航空母舰-国防部网站-中国正自主开展设...	2018/7/20 11:21	文本文档	3 KB
航空母舰-环球网-我军少将: 中国电磁弹...	2018/7/20 11:21	文本文档	16 KB
航空母舰-环球网-乌克兰媒体称中国正在...	2018/7/20 11:21	文本文档	4 KB

图 3.60 搜索结果

➤ 高级检索

Step1: 点击“使用高级命令”，输入高级命令。例: [field] content
[AND] 中国 人民 表示: 搜索内容字段中同时包含“中国”和“人民”
的文档，采用该语法信息过滤将更有针对性;

Step2: 点击“高级检索”，系统将进行高级检索。搜索结果框会呈现查询结果，并配以相似得分。检索结果文件(\\NLPIR-Parser\\output\\搜索结果 JZSearch-result) 以网页形式保存。



图 3.61 高级检索

Step3: 点击“结果导出”，系统将检索目标文档导出，并自动打开文件目录。

电脑 > 桌面 > NLPIR-Parser > output > 搜索结果 > [field] content [AND] 中国人民

名称	修改日期	类型	大小
航空母舰-参考消息网-美专家预测：中国...	2018/7/20 11:22	文本文档	6 KB
航空母舰-传媒-如何讲好我们的故事.txt	2018/7/20 11:22	文本文档	15 KB
航空母舰-光明日报-永不褪色的精神礼赞...	2018/7/20 11:22	文本文档	30 KB
航空母舰-环球网-我军少将：中国电磁弹...	2018/7/20 11:22	文本文档	16 KB
航空母舰-科技日报-辽宁舰已具备初始作...	2018/7/20 11:22	文本文档	9 KB
航空母舰-科技日报-张召忠：中国双航母...	2018/7/20 11:22	文本文档	5 KB
航空母舰-每日经济新闻-昨天航母刚刚下...	2018/7/20 11:22	文本文档	8 KB
航空母舰-人民网-第十三届中国大学生年...	2018/7/20 11:22	文本文档	7 KB
航空母舰-人民网-军事频道-国防部回应...	2018/7/20 11:22	文本文档	2 KB
航空母舰-人民网-军事频道-专家：首艘...	2018/7/20 11:22	文本文档	4 KB
航空母舰-人民网-人民日报海外版-海外...	2018/7/20 11:22	文本文档	4 KB
航空母舰-人民网-人民日报海外版-海外...	2018/7/20 11:22	文本文档	7 KB
航空母舰-人民网-人民日报海外版-海洋...	2018/7/20 11:22	文本文档	11 KB
航空母舰-人民网-人民日报海外版-人民...	2018/7/20 11:22	文本文档	5 KB

图 3.62 结果导出

3.13 编码转换

编码转换功能，自动识别内容的编码，并把编码统一转换为 GBK 编码。目前支持 Unicode/BIG5/UTF-8 等编码自动转换为简体的 GBK，同时将繁体 BIG5 和繁体 GBK 进行繁简转化。

用户首先点击“编码转换”，进入系统编码转换功能模块。

➤ 转换为 GBK 编码

Step1：选择语料源：\NLPIR-Parser\编码转换测试文本，系统指定输出路径：\NLPIR-Parser\bin-win64\output\GBK。

Step2：点击“转换为 GBK 编码”。系统自动识别给定的 BIG5 文件，GBK 以及 UTF-8,Unicode 文件，最终转化为简体 GBK 编码的文件。转换结果提示框将显示转换结果，并将编码转换结果文件夹自动打开，用户可直接查看与使用转换后的文件。

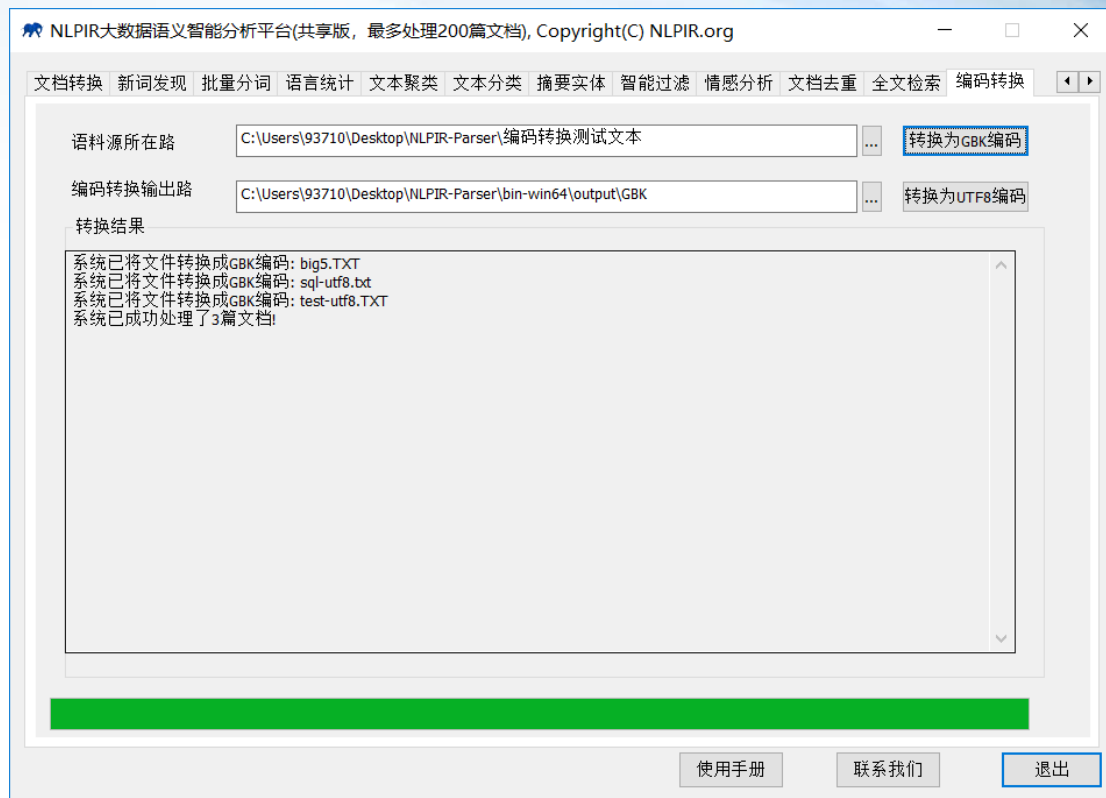


图 3.63 转换为 GBK 编码

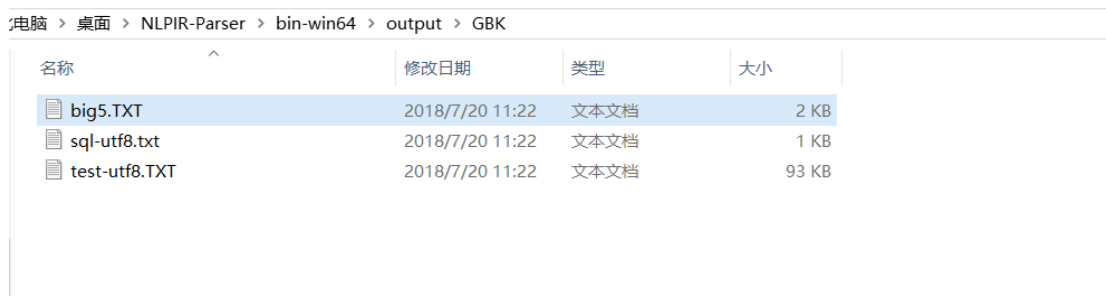


图 3.64 转换为 GBK 编码

➤ 转换为 UTF8 编码

Step1: 选择语料源: \NLPIR-Parser\编码转换测试文本, 系统指定输出路径: \NLPIR-Parser\bin-win64\output\UTF8。

Step2: 点击“转换为 UTF8 编码”。系统自动识别给定的 BIG5 文件, GBK 以及 UTF-8,Unicode 文件, 最终转化为简体 UTF8 编码的文件。转换结果提示框将显示转换结果, 并将编码转换结果文件夹自动打开, 用户可直接查看与使用转换后的文件。

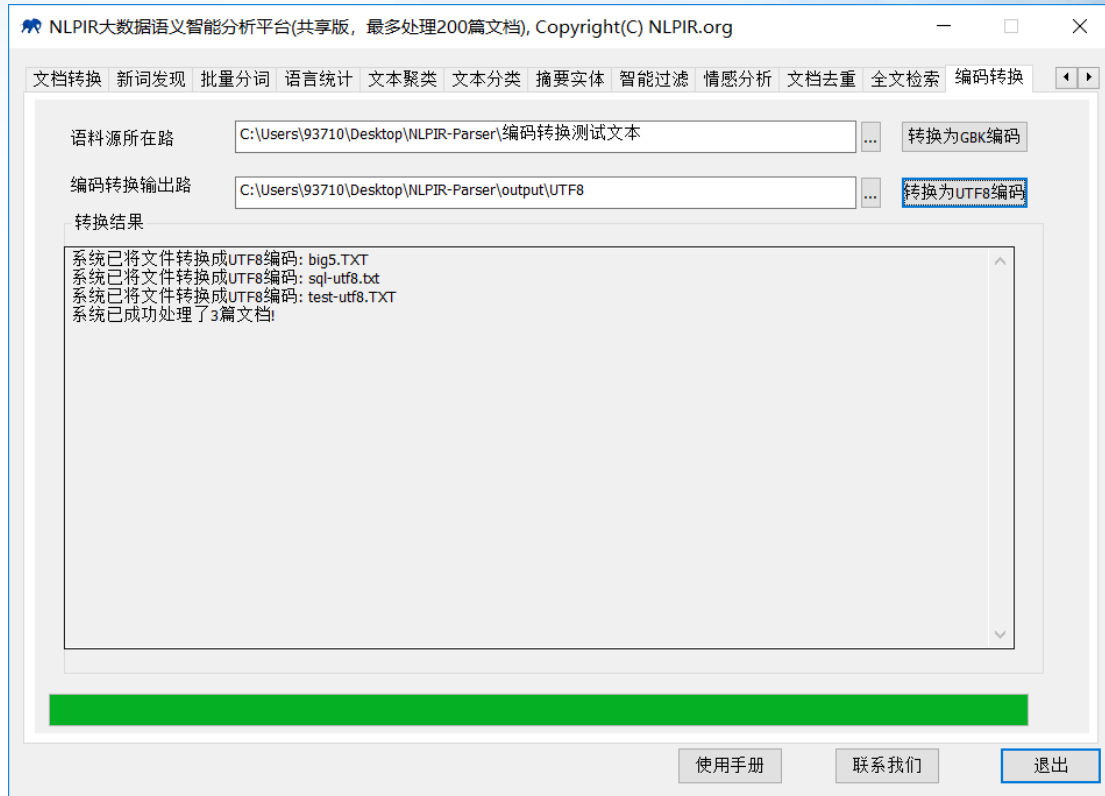


图 3. 65 转换为 UTF8 编码

电脑 > 桌面 > NLPIR-Parser > output > UTF8

名称	修改日期	类型	大小
big5.TXT	2018/7/20 11:23	文本文档	3 KB
sql-utf8.txt	2018/7/20 11:23	文本文档	1 KB
test-utf8.TXT	2018/7/20 11:23	文本文档	127 KB

图 3. 66 转换为 UTF8 编码

四、应用示范案例

4.1 十九大报告语义智能分析

2017 年 10 月 18 日，中国共产党第十九次全国代表大会在北京隆重召开，习近平代表第十八届中央委员会向大会作报告。这份沉甸甸的报告总结了自十八大以来我国的发展进程，党的引领脚步，人民

的生活改变.....以及未来如何开启新时代、谱写新篇章。

如何精准解读这份报告，我们采用自然语言处理工具 nlpir-paser，通过挖掘十九大报告的关键词、概念新词、内容图谱等语义智能处理技术，带你深度感受十九大精神。

➤ 关键词提取

十九大报告全文 3 万余字，本文使用 NLPIR 对十九大报告进行关键词提取，以期揭示十九大报告的核心要点。关键词 top100 结果展示如下：



图 4.1 关键词 top100

由于篇幅所限，本文只展示了部分关键词提取的结果，关键词词云图分析结果比较充分地展示了十九大报告的核心概念。

➤ 词频统计

分析结果显示，词频统计 top10 的关键词分别为：“中国特色社会主义”、“中华民族伟大复兴”、“依法治国”、“全面建成小康社会”、“中国梦”、“人民当家作主”、“美好生活”、“现代化

经济体系”、“人民军队”、“小康社会”。这些高频词汇基本概括了十九大报告中的基础概念。



图 4.2 词频统计

► 新词发现

“人类命运共同体”，“新征程”，“现代化经济体系”，“社会主要矛盾转化”，“历史性变革”……

十九大报告中出现的不少新的“关键词”，这些新词展示了新理念、新观点，给予了重大时代课题明确的回答，在实践上作出了新部署。



图 4.3 十九大新词

4.2 文章风格对比：方文山 VS 汪峰

不同人的文章风格不同，汪峰的摇滚歌词给人奔放、热烈的情感激荡，而方文山中国风歌词则会给我们造成委婉、缠绵悱恻的心湖涟漪。这类文章风格主观感受的差别能否经得起科学实验的验证或证明呢？再者，文学、艺术等多个领域都存在文章作品对比与评价的争议，造成了很多不良的影响。通过技术能否为此提供一个评估的新维度或方法呢？我们通过 nlpir-paser 进行语言统计与分析、情感分析与词曲语言广度分析（信息熵）来进行文章风格的对比分析。

➤ 词频广度分析

通过歌词数目对比，通过工具可以得出以下方文山与汪峰对比：

（比率=方文山/汪峰，平均用词=总词数/歌曲数）

表 4.1 方文山和汪峰用词分析

	总词数	歌曲数	平均用词
方文山	8195	200	40.975
汪峰	2270	127	17.874
比率	3.610	1.574	2.292

可以很明显的看出方文山所用词汇数量远远多于汪峰。通过平均用词可以发现方文山比汪峰用词广度大。每首歌曲方文山是汪峰的用词量的二倍。

➤ 情感对比分析

将方文山和汪峰的形容词作为情感分析的主要词汇。



图 4.4 方文山（左）和汪峰（右）的情感词汇词云图

从形容词上统计方文山和汪峰，可以看出汪峰是一种激进的用词，负向很明显“孤独”“破碎”，正向“美丽”“坚强”，这些对生命的感悟的词汇。汪峰多写生命的感悟，同时把摇滚歌手那种想表达的孤单，力量感，表达出来。而方文山的形容词性则以比较温柔的情感词为主“温柔”“美丽”“简单”。这里也能说明两个作词人风格不同，方文山多写爱情和亲情。通过比对能很明显的发现两个作词人词风不同。

➤ 信息熵分析

信息熵公式： $H(X) = -\sum_{i=1}^n P(X) \log P(X)$ 。信息熵用来表示作词人用词的广度。用词数量越小，信息熵越小。通过用词信息熵进行加和来比较方文山和汪峰的用词广度。

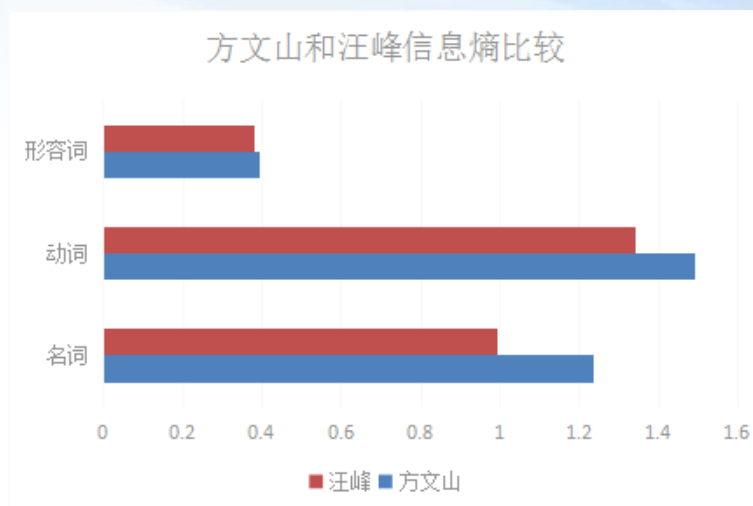


图 4.5 信息熵对比

可以看出汪峰词作在三组词性上的信息熵均小于方文山。同时验证了汪峰的词作中用词信息量较少。可以推理出汪峰词作多重复性词汇，方文山用词量大，广泛。

4.3 《红楼梦》作者前后同一性识别

《红楼梦》前八十回和后四十回到底是不是同一个作者？我们都知道《红楼梦》的作者有两个：曹雪芹写了前八十回，高鹗续写了后四十回。然而，红学上关于《红楼梦》的作者争议一直很大，存在着很多种版本。我们将利用大数据语义智能分析工具 `nlpir-paser`，通过语言统计、概率计算与文本相似度分析来进行《红楼梦》前后作者同一性判别。

➤ 虚词统计

每个人的写作都有些小习惯，虽然文章前后说的内容会有差别。但是每个人使用虚词的顺序与数量可能存在着差异。

将《红楼梦》120 回按顺序均分为 3 组，使用 NLPIR 统计出文言

虚词的词频，再对不同组数据之间进行 KL 距离计算。第一组将 120 回按顺序均分为三等份即第 1 回-第 40 回、第 41 回-第 80 回、第 81-第 120 回。这 3 组数据中部分虚词以及该词的概率如表所示：

表 4.2 三组虚词统计分析

词	第1回-第40回		第41回-第80回		第81回-第120回	
	词频	概率	词频	概率	词频	概率
了	5981	0.199712836	7740	0.213299529	6710	0.206786033
的	3854	0.128689729	5156	0.142089454	5269	0.162377885
不	3063	0.102277281	3805	0.104858489	3510	0.108169743
是	2293	0.076566048	2975	0.081985284	3039	0.093654658
一	2202	0.073527448	2750	0.075784716	1953	0.060186755
着	1607	0.053659677	1855	0.051120236	2112	0.065086752
便	1075	0.035895552	1272	0.035053876	1295	0.03990878
在	1026	0.034259383	1089	0.030010748	1253	0.038614441
就	935	0.031220783	1101	0.030341445	817	0.025177972
儿	899	0.030018699	1108	0.030534351	1143	0.035224506
好	786	0.026245492	956	0.026345523	939	0.028937718
之	747	0.024943235	658	0.018133216	243	0.007488675
呢	601	0.020068118	515	0.014192411	719	0.022157848
因	571	0.019066382	724	0.019952049	363	0.011186785
再	395	0.013189529	456	0.012566484	262	0.008074209
可	385	0.012855616	362	0.009976024	254	0.007827668
罢	328	0.010952317	354	0.009755556	407	0.012542759
把	324	0.010818753	364	0.010031141	420	0.012943388
方	266	0.008882062	284	0.007826494	59	0.001818238
往	253	0.008447976	243	0.006696613	140	0.004314463
别	250	0.008347803	314	0.008653237	165	0.005084902
向	212	0.007078937	203	0.00559429	119	0.003667293
亦	171	0.005709897	144	0.003968363	28	8.63E-04
比	160	0.005342594	211	0.005814755	110	0.003389935

➤ KL 距离

KL 距离（相对熵）可以衡量两个随机分布之间的距离，当两个随机分布相同时，它们的相对熵为零，当两个随机分布的差别增大时，它们的相对熵也会增大。所以相对熵（KL 散度）可以用于比较文本的相似度。

从下表中可以观察到第一行中 1-40 与 81-120 的 KL 值是 1-40 与 41-80 的 KL 值的十倍。由于当两个随机分布的差别增大时，它们的相对熵也会增大。所以 1-40 与 81-120 的相似性比 1-40 与 41-80 低。

表 4.3 三组 KL 距离分析

回数 \ KL 值	回数	KL 值	回数
回数	KL 值	回数	KL 值
1-40	0	0.008	0.082
41-80	0.007	0	0.06
81-120	0.051	0.049	0

可以看出前八十回的各组数据的 KL 值与后四十回的数据的 KL 值有不同程度的差距。后四十回之间的 KL 值比其他组得 KL 值要小，说明后四十回的相似度较高。可以大胆猜测后四十回是出自于另外一个人。

五、联系我们

需要购买 NLPIR 大数据语义智能分析平台正式版本，或者需要使用 NLPIR 各类二次开发包，可以通过以下方式联系到我们：

大数据搜索与挖掘实验室（北京市海量语言信息处理与云计算应用工程技术研究中心）

地址：北京海淀区中关村南大街 5 号 100081

电话：13681251543(商务助手电话)

Email: kevinzhang@bit.edu.cn

MSN: pipy_zhang@msn.com;

网站: <http://www.nlpir.org> (自然语言处理与信息检索共享平台)

<http://www.bigdataBBS.com> (大数据论坛)

微博:<http://www.weibo.com/drkevinzhang/>

微信公众号: 大数据千人会

Beijing Engineering Research Center of Massive Language Information Processing and Cloud Computing Application

Beijing Institute of Technology

Add: No.5, South St.,Zhongguancun,Haidian District,Beijing,P.R.C PC:100081

Tel: 13681251543(Assistant)

Email: kevinzhang@bit.edu.cn

MSN: pipy_zhang@msn.com;

Website: <http://www.nlpir.org> (Natural Language Processing and Information Retrieval Sharing Platform)

<http://www.bigdataBBS.com> (Big Data Forum)

Twitter:<http://www.weibo.com/drkevinzhang/>

Subscriptions: Thousands of Big Data Experts

六、附录

6.1 下载途径

NLPIR-Parser 系统的多种下载途径:

1、GitHub: <https://github.com/NLPIR->

[team/NLPIR/tree/master/NLPIR-Parser](https://github.com/NLPIR/tree/master/NLPIR-Parser)

【有可能国内访问国外网址受限】

2、官方网站下载：

链接：<http://www.nlpir.org/NLPIR-Parser.zip>

打开浏览器，复制下载链接，即可启动下载。

3、百度网盘：

链接：<https://pan.baidu.com/s/1Khxt0nEQxI7FfaVrfXOOMw> 密

码：4nyr 【有可能开大会期间会被误封】

4、也可以百度各软件下载平台，下载 NLPIR-Parser。

访问 NLPIR-Parser 目录即可。

注：用户在 github 上下载 NLPIR-Parser 文件时需要专门的下载工具，建议使用 svn 工具下载文件。百度网盘下载量大时，需要安装百度网盘客户端。

6.2 Github 下载

首先，打开 github 上 NLPIR-Parser 文件下载地址，复制该地址：

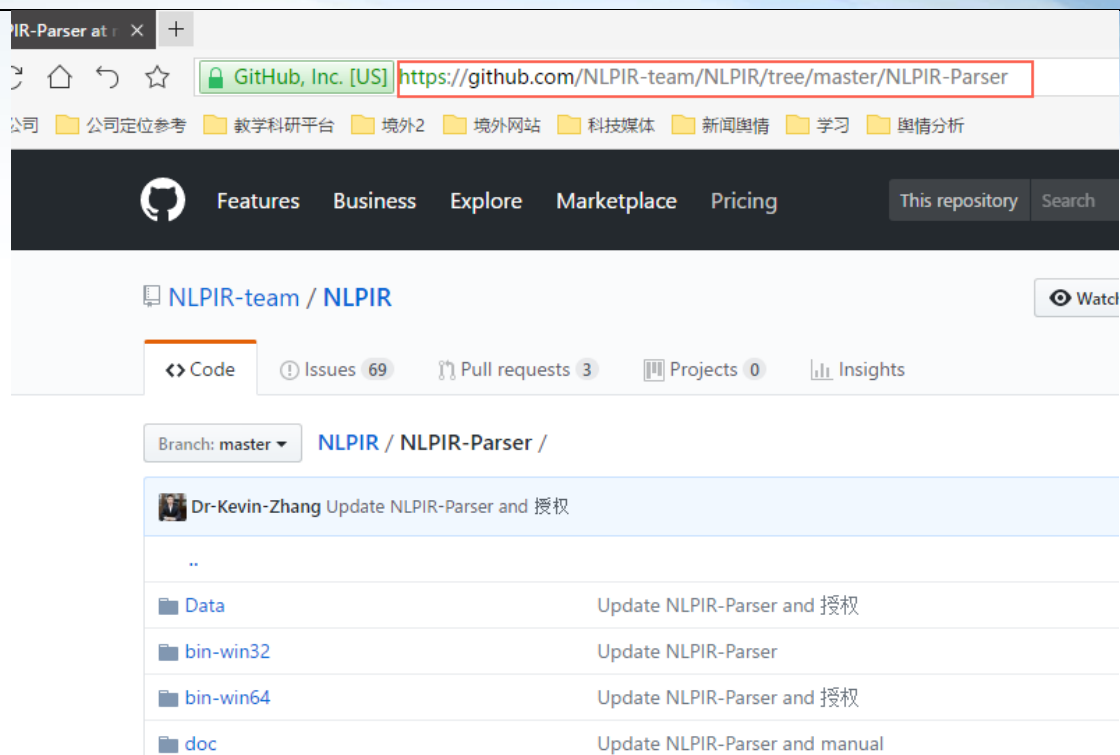


图 6.1 github 网址

然后，点击鼠标右键 SVN Checkout，弹出以下窗口，文件下载地址已经自动复制。

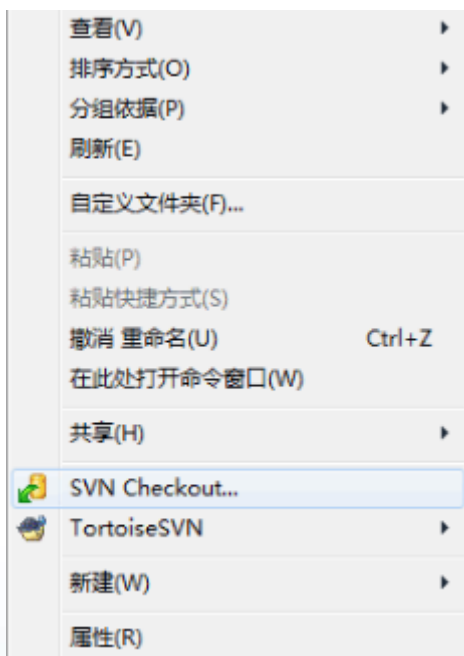


图 6.2 右键 “svn checkout”

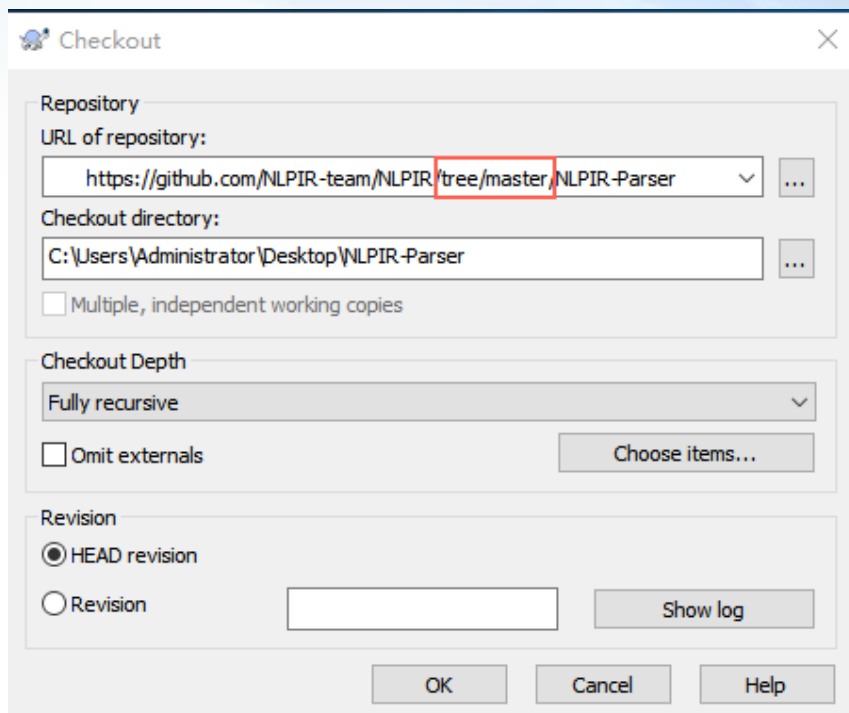


图 6.3 自动复制网址

最后，将地址中的“/tree/master/”修改为“/trunk/”，选择文件存放地址（桌面 desktop 或其他地址），点击“ok”，文件下载启动，下载完毕后点击“ok”，文件下载完毕。

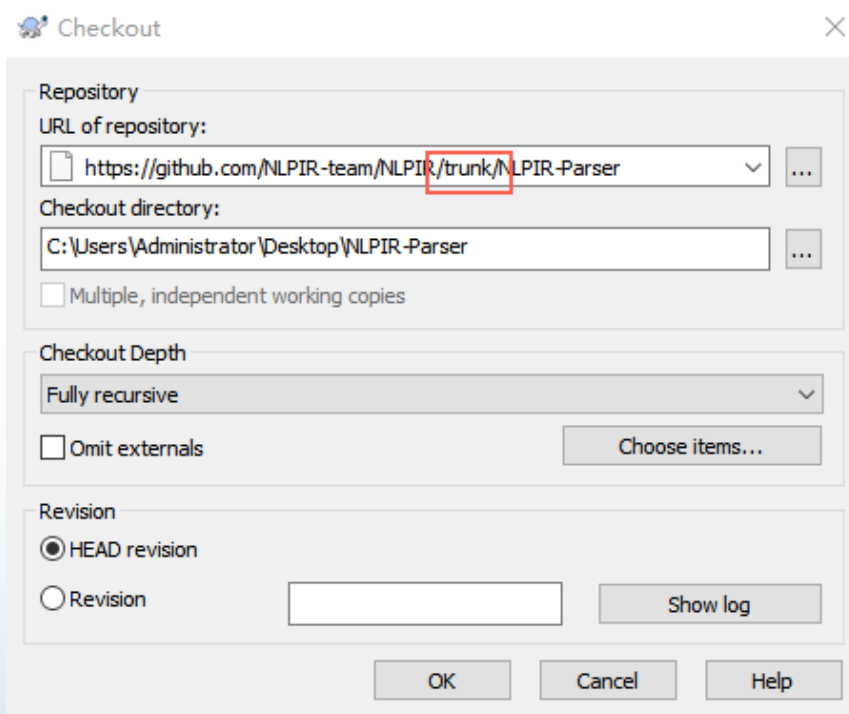
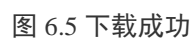


图 6.4 修改地址



首先，在浏览器打开 **NLPIR-Parser** 文件链接。输入密码。

图 6.6 打开连接

然后，打开 NLPIR 大数据语义智能分析文件夹，找到 NLPIR-Parser 文件目录。



图 6.7 文件目录

接下来，将 NLPIR-Parser 文件保存在自己的百度网盘账户中。



图 6.8-1 保存文件



图 6.8-2 保存文件

下一步，打开百度网盘客户端（下载量大推荐）或在线网盘，登录自己的账号，找到上一步保存的文件。

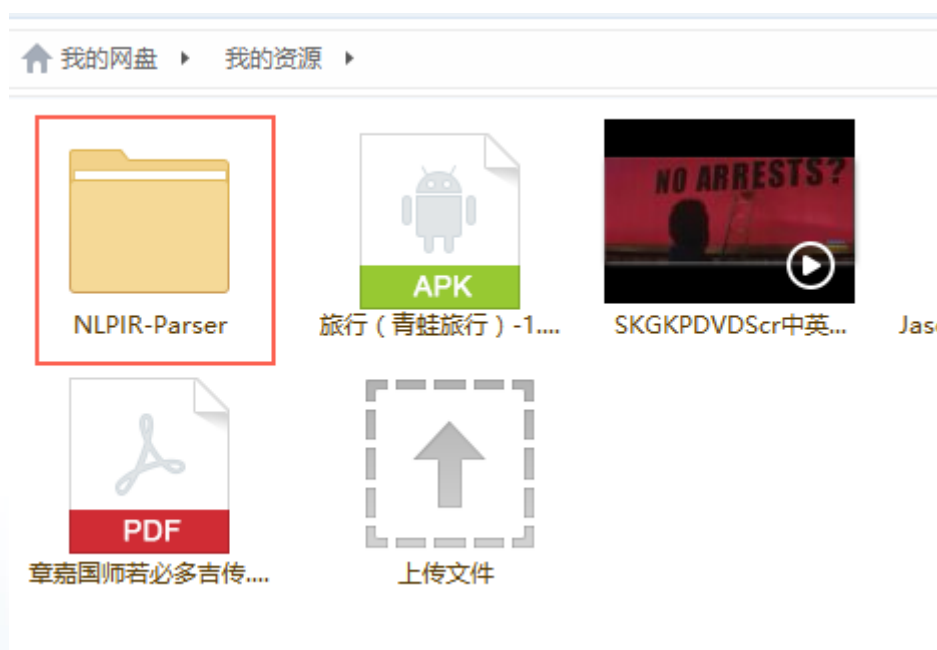


图 6.9 寻找文件

最后，右击文件，选择下拉列表的“下载”，定义文件下载地址，文件下载即可启动。

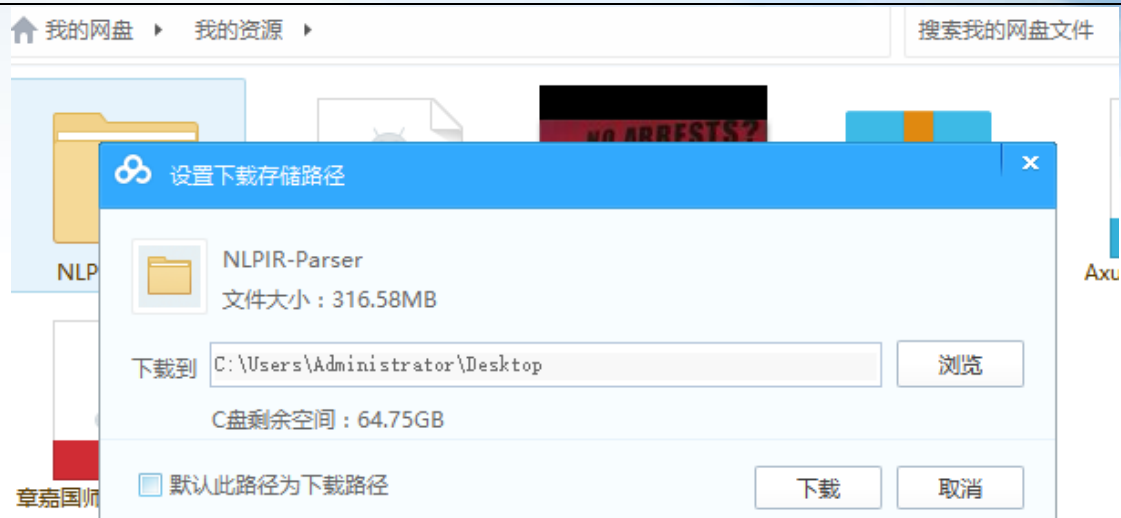


图 6.10 下载地址



图 6.11 下载文件