

基于互联网挖掘的热点选股策略

——互联网大数据挖掘系列专题之（五）

报告摘要:

● 近年来，概念、热点炒作频繁

近年来，A 股市场上概念、热点的炒作频繁，例如最近的“沪港通”、“国企改革”、“互联网金融”、“一带一路”等，相关个股在某一概念、热点背景下，往往在中短期时间内有较好的市场表现。如何抓住这种热点、概念带来的投资机会，是当前量化研究的一个难点。本专题报告试图从大规模的互联网新闻的文本信息作为切入点，采用文本挖掘的方法研究热点、概念带来的投资机会。

● 基于互联网挖掘的热点选股策略构建

基于大规模的热门网站关于个股的新闻文本信息，本专题策略通过动态地构建热词词库以及关于热词的热度指标，采用文本挖掘方法对热词的变化进行了量化。基于量化后关于热词热度的变化指标，我们进行策略的构建。

互联网挖掘的热点选股策略构建：在历史回测期间内，将资金等分为若干份。动态地监控在每个交易日内是否有热点概念出现，如果有热点概念出现，则将其其中的一份没有持仓个股的资金配置当天出现的热点概念的个股，持有一段时间，不考虑个股在买进或者卖出时候出现涨停、跌停或者停牌的情况，个股在买进时以该个股当日的开盘价买进，卖出时以该个股在持有期末的收盘价卖出。

● 实证结果

历史回测结果显示，热点选股策略在取得了较优异的结果。相对沪深 300 指数，在整个回测期内，实现了约 19.00% 的年化超额收益率以及 112.08% 的年化绝对收益率；在胜率方面胜率，取得了 65.85% 的周胜率以及 81.82% 的月度胜率；从回撤角度看，在没有用沪深 300 进行对冲前，策略累计净值的回撤为 6.27%，在用沪深 300 进行对冲后，策略对冲后的净值的回撤为 3.59%。整体上而言，基于对大规模的热门网站的文本信息挖掘构建的热点选股策略，在历史回测期内取得相对较为优异的表现。

● 风险提示

本模型为采用纯量化方法，所推荐的个股未必具有实质性的利好，其股价表现还受到诸多因素影响，请结合基本面及自身判断进行恰当使用。

图 1 热点选股策略历史回测表现

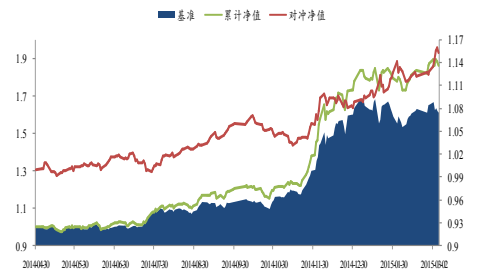


表 1 热点选股策略表现指标

指标	数值
年化绝对收益率	112.08%
年化相对收益率	19.00%
最大回撤	-6.27%
周胜率	65.85%
月胜率	81.82%

分析师：史庆盛 S0260513070004



020-87555888-8618



sqs@gf.com.cn

相关研究:

基于网络新闻热度的择时策

略—互联网大数据挖掘系列
专题之（一） 2014-06-25

公告披露背后隐藏的投资机

会—互联网大数据挖掘系列
专题之（二） 2014-06-26

倾听股吧之声—互联网大数据

挖掘系列专题之（三） 2014-06-27

那些年一起追过的财经小编

策略—互联网大数据挖掘系列
专题之（四） 2014-08-22

上市公司披露信息变更隐含

的投资机会—事件驱动策略
之（十四）>> 2014-12-26

目录索引

一、 前言.....	4
二、 文本挖掘框架流程.....	5
2.1 数据抓取平台组件框架.....	5
2.2 大规模互联网新闻获取介绍.....	6
2.3 基于互联网挖掘的热点选股策略框架.....	9
三、 基于互联网挖掘的热点选股策略实证.....	12
3.1 样本数据.....	12
3.2 策略构建原理.....	12
3.3 互联网挖掘的热点选股策略.....	12
四、 总结.....	13
4.1 总结.....	13
4.2 未来研究方向.....	14
4.3 工具推荐.....	14
风险提示.....	14

图表目录

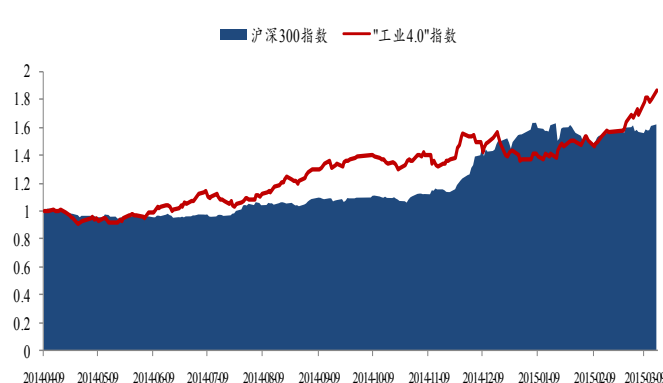
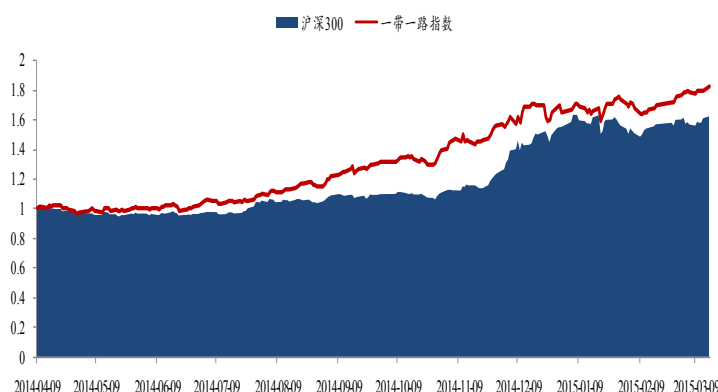
图 1 热点选股策略历史回测表现	1
图 2 "一带一路"概念指数走势	4
图 3 "工业 4.0"概念指数走势	4
图 4 "Duang"关键字提及的微博数变化趋势	5
图 5 "Duang"热词变化趋势图	5
图 6 数据抓取平台框架组件	6
图 7 新闻网个股资讯栏目	8
图 8 和讯网个股资讯栏目	8
图 9 个股资讯获取示例	8
图 10 新闻资讯提取一览	8
图 11 基于互联网挖掘的热点概念选股策略框架	9
图 12 txt 存储格式:个股代码_新闻发表时间_新闻标题	9
图 13 数据库存储字段:个股、时间、新闻标题、对应链接	9
图 14 热点概念传播与对应的股价关系	10
图 15 热词词库示例	10
图 16 热点概念热度指标构建步骤	11
图 17 "工业 4.0"热点概念热度变化	11
图 18 "云计算"热点概念热度变化	11
图 19 "雾霾"热点概念热度变化	11
图 20 "环保"热点概念热度变化	11
图 21 热点变化与对应标的股价变化示例 1	12
图 22 热点变化与对应标的股价变化示例 2	12
图 23 热点选股策略表现(相对沪深 300)	13
图 25 广发证券金融工程热点概念识别工具	14
表 1 热点选股策略表现指标	1
表 2 热点选股策略表现指标一览	13

一、前言

近年来，A股市场上概念、热点的炒作频繁，例如最近的“沪港通”、“国企改革”、“互联网金融”、“一带一路”等，相关个股在某一概念、热点背景下，往往在中短期时间内有较好的市场表现。热点、概念的形成往往是因为某一类信息，例如国家政策方针、龙头企业的相关产品信息在投资者之间互相传播，演进，达成共识最后带动相关标的的上涨。这种热点、概念的兴起带来的投资机会在往往在中短期内的市场表现比市场上的整体表现要优异得多，而如何研究这种热点、概念变化带来的投资机会是本专题策略研究的出发点。

图 2 “一带一路”概念指数走势

图 3 “工业 4.0”概念指数走势



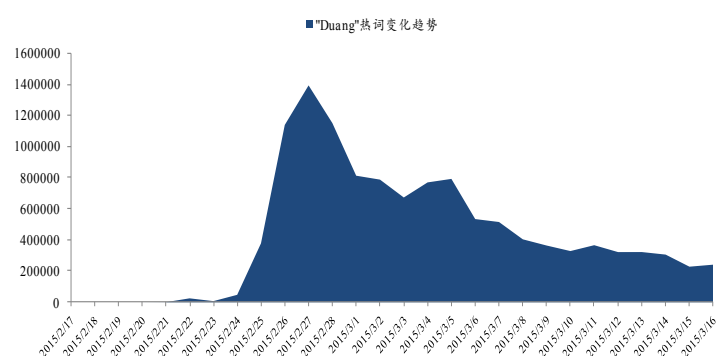
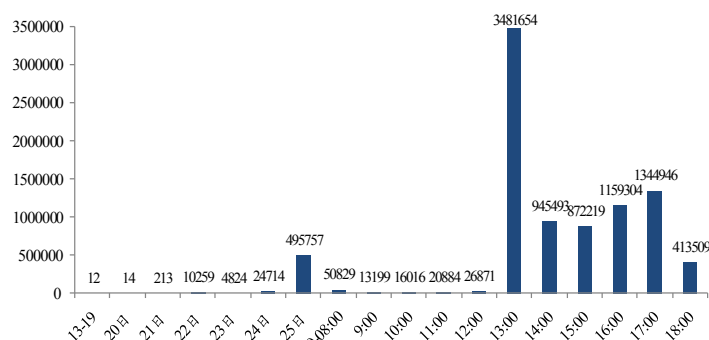
数据来源：广发证券发展研究中心、Wind

数据来源：广发证券发展研究中心、Wind

在当前互联网大数据时代下，信息的传播渠道已经越来越丰富，微博、微信、热门财经网站等已经成为了信息传播的主要途径。举个例子而言，最近兴起的一个热点词汇--“Duang”。“Duang”这个词好像一夜之间，其实就是一夜之间，“Duang”这个词火了，火得一塌糊涂，这个热词的形成毫无疑问是因为互联网的传播。简要讲述一下这个热词的形成过程，“Duang”这个词最初是来自一个网友在Bilibili上传的一段关于成龙代言的一段广告视频，视频比较搞笑，同时涉及成龙，又是在春节期间，因此就慢慢地就在网友之间传播起来了，其中在2月26日早上9:00至19:00这段时间内在新浪微博中被提及826万次，也即是说每秒平均被提及229次。从中可以看到互联网的传播对于热点、概念的形成起到了至关重要的作用。因此在针对研究热点、概念的产生、形成变化过程，可以从互联网上大规模的新闻文本信息作为切入点，而如何获取到大规模的关于个股相关资讯的新闻文本，并从这些文本信息中挖掘出当前可能的一些热点、概念，并识别出对应的标的带来的投资机会则是本专题策略研究的难点。

图 4 "Duang"关键字提及的微博数变化趋势

图 5 "Duang"热词变化趋势图



数据来源：广发证券发展研究中心、数据化管理

数据来源：广发证券发展研究中心、新浪

对大规模的互联网新闻文本信息的挖掘与研究，我们广发金工在该领域进行了比较深入的研究，在专题策略上也取得了一系列的研究成果。例如之前采用文本挖掘的方法对上市公司公告披露背后的投资机会进行了统计分析以及实证，得到了较好的就结果，具体可见《公告披露背后隐藏的投资机会—互联网大数据挖掘系列专题之（二）》专题报告；从股吧、个股的新闻热度、上市公司信息变更、财经频道的荐股信息等角度对文本信息进行挖掘，具体可见《基于网络新闻热度的择时策略—互联网大数据挖掘系列专题(一)》、《倾听股吧之声，洞察大盘趋势—互联网大数据挖掘系列专题(三)》、《那些年一起追过的财经小编选股策略—互联网财经频道文本挖掘策略》、《上市公司披露信息变更隐含的投资机会—事件驱动策略之(十四)》等相关的专题策略报告。

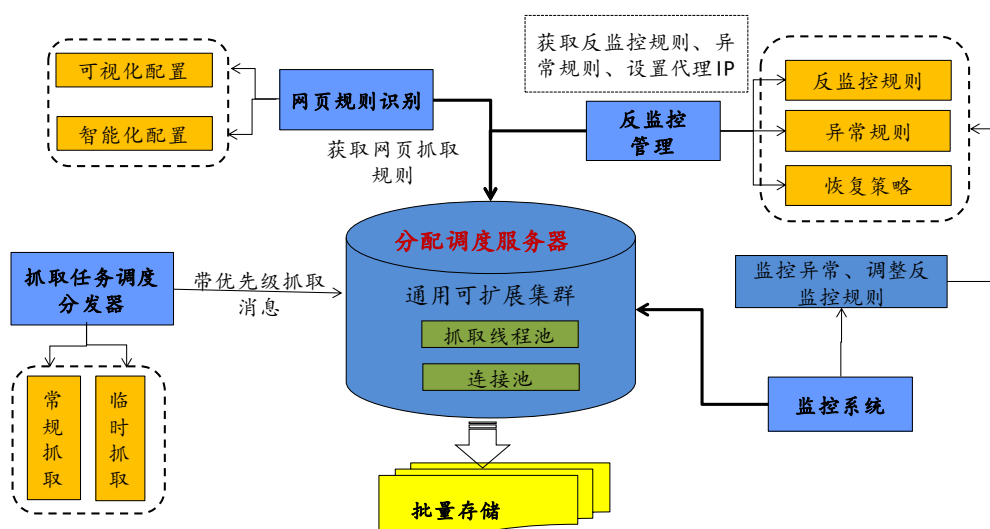
二、 文本挖掘框架流程

2.1 数据抓取平台组件框架

在之前的互联网大数据挖掘的专题系列报告中，我们广发金工搭建了一个完善的数据抓取平台框架。这里我们首先对整个数据抓取的平台框架进行详细地介绍，然后再介绍基于此平台上搭建的此次专题策略的框架流程。

首先，介绍一下我们金融工程小组搭建的完善的数据抓取平台框架。

图 6 数据抓取平台框架组件



数据来源：广发证券发展研究中心

上图刻画了我们金融工程小组搭建的完善的数据抓取平台组件的整个框架流程。整个抓取平台的平台核心部分是中间的分配调度服务器，辅助搭建的模块有四个，分别为抓取任务调度器、网页规则识别、反监控系统以及监控系统。

分配调度服务器功能为负责将所需数据从互联网上抓取下来，然后以指定形式的格式批量存储下来，这里可以将记录以特定的标题格式以 txt 的形式存储于本地或者以数据库的形式存储相关的文本信息。如果是新闻文本的标题作为存储格式的一部分的话，对于在 Windows 系统下注意 txt 文件名中的一些非法字符的处理以及由于网页本身代码的特殊性，导致的一些非法字符的处理。抓取任务调度分发器负责分配抓取的网站的调度，包括一些常规的抓取以及一些临时的抓取(链接失效、断网重新抓取等任务)调度、网页规则识别负责根据抓取任务调度器分配到的网站，调用指定网站的网页内容获取规则。监控系统负责监控网页的异常（例如链接失效、断网、网页加载过慢等情况）、反监控系统负责监控网站的防盗链等问题(例如 IP 频繁访问造成的 IP 被禁等问题)。

2.2 大规模互联网新闻获取介绍

这里简要介绍一下，本专题策略对大规模互联网信息文本信息的抓取过程。因为本专题策略目的是研究热点、概念变化带来的投资机会，因此在数据获取上应该尽可能地获取到市场上关于个股的所有新闻文本信息。然而，“理想是美好的，现实是残酷的”，获取互联网上所有关于个股的新闻文本信息，无论是在硬件条件以及效率上来说都是不现实的。由于热点、概念的传播、形成往往是通过一些比较热门的财经网站渠道，因为这些热门财经网站上每天都有大量的浏览用户，这些信息在这些浏览用户之间传播后，逐渐地形成共识，进而产生热点、概念。因此，在数据源上可以选择一些用户浏览量大的热门财经网站的，利用这些热门网站的新闻热点作为数据源切入点。本专题策略的数据源分别选择了新浪网、和讯网、腾讯网、东方财富网这四个热门网站上的关于个股的新闻资讯，这几个网站上据我们广发金

融工程小组对相关资讯栏目的统计,每一天关于全市场的个股新闻平均总共有 15000 条左右的新闻量。

以下以新浪网为例,简要介绍一下个股新闻资讯爬取的主要过程。

- 1) 首先获取到关于个股资讯的网址:

http://vip.stock.finance.sina.com.cn/corp/view/vCB_AllNewsStock.php?symbol=sz000001&Page=1;

- 2) 从新浪的个股资讯网站上可以查看到每天每个个股的相关的新闻的数据,而且从这个个股资讯网站上可以查看到该个股历史上所有的关于该个股的新闻的文本信息数据。

- 3) 通过2)中的解析后有两种方法可以获取到个股指定日期的新闻资讯,一种为通过每次查询后的结果获取到的结果URL,然后比较不同URL之间的固定参数以及可变参数的关系,找到URL之间的规律,另一种为通过模拟鼠标操作,例如翻页的操作等来获取到个股的相关资讯;

以第一种方法为例,解析查询后获取的网址的规律。

通过对网页编码的解析,我们可以获取到每次查询后的网址的实际网址信息,例如想获取到平安银行(000001.SZ)的个股资讯,可以看到实际的网址就是:

http://vip.stock.finance.sina.com.cn/corp/view/vCB_AllNewsStock.php?symbol=sz000001&Page=1,而如果在翻页操作后,例如翻到第二页,可以看到网址为:

http://vip.stock.finance.sina.com.cn/corp/view/vCB_AllNewsStock.php?symbol=sz000001&Page=2,因此可以看到,对于单个个股资讯的获取,可以通过page参数的变化来获取到历史上每一天的个股资讯新闻文本信息,而对于不同的个股,例如万科A(000002.SZ),个股资讯的网址为:

http://vip.stock.finance.sina.com.cn/corp/view/vCB_AllNewsStock.php?symbol=sz000002&Page=1,对不同个股资讯的查询可以通过symbol参数的变化来获取到不同个股历史上所有的新闻文本信息。通过以上的解析,可以知道这个URL的规律如下:

固定参数部分:

http://vip.stock.finance.sina.com.cn/corp/view/vCB_AllNewsStock.php

可变参数部分: symbol、Page

其中Page表示如果查询结果返回的是多页情况,可以通过Page的变化获取到不同页面的结果,symbol指的是不同的个股的代码,格式为:"上市交易所(sz or sh)+个股上市代码"

- 4) 通过3)中对查询网址的规律的查询,我们可以根据URL的规律批量获取到市场上所有的个股在历史上的所有的个股资讯信息,并将将个股资讯信息提取出来。

基于上述对新浪网个股资讯获取规则的解析,我们可以得到不同个股在历史上所有的个股资讯信息,通过对需要提取的信息的网页编码格式进行解析后,可以获取到不同的个股新闻文本信息。

图 7 新闻网个股资讯栏目

资讯与公告	个股资讯	行业资讯	公司公告	年度报告	中期报告
个股相关资讯:					
2015-03-10 18:31	招商地产:长轴善舞 协力同行				
2015-03-10 18:15	万科北京区域CEO辞职创业 地产行业吸引力消退?				
2015-03-10 15:00	王若郁离毛大庆 万科三巨头绝配还是富斗				
2015-03-10 14:35	房地产行业:中国土地市场开局不利 路到尽头?				
2015-03-10 11:12	房地产行业:成交延续回暖势头 后续政策存在进一步宽松可能				
2015-03-10 09:15	离职潮袭房企 万科人才战略博弈				
2015-03-10 05:59	楼市低迷倒逼开发商谋变 标杆房企现高管离职潮				
2015-03-10 05:59	毛大庆离职创业就万科力推郁离押宝“75”后新人吗未来				
2015-03-10 04:04	万科五个副总裁中有三个75后 要拿年轻人“赌明天”				
2015-03-10 02:30	毛大庆告别万科创业要做啥 将创建中国版WeWork				
2015-03-10 01:49	万科8个副总裁中有三个75后 要拿年轻人“赌明天”				
2015-03-10 01:49	毛大庆离职背后:非住宅业务成万科发展新方向				
2015-03-10 01:40	毛大庆版WeWork的头号难题				
2015-03-10 01:11	北京万科告别“毛大庆时代” 坊间曾传其因言获罪				
2015-03-09 16:53	郁离:把万科未来赌在年轻人身上				
2015-03-09 14:17	万科郁离:没有回款的销售都是耍流氓				
2015-03-09 11:32	万科毛大庆离职创业 德东家慷慨投资创新工场				
2015-03-09 10:37	摩根大通减持万科企业314万股 套现5343万港币				
2015-03-09 09:55	[互动]万科A:度假物业是公司正探索的方向				

数据来源: 广发证券发展研究中心、新浪网

图 8 和讯网个股资讯栏目

平安银行资讯聚合	百度一下	谷歌一下
查看: 全部 个股新闻 行业新闻 研报 博客 聚吧 专题 视频		
[个股新闻] 合约陆续有来 中纺国际借势攀升 和讯网	2015-03-17	
[个股新闻] 谢娜: 接班人的隐与立 接力	2015-03-17	
[个股新闻] 恒指逼近24,000点 留意034915、沽193 阿思达克财经网	2015-03-17	
[个股新闻] 3.15银行满意度调查: 超三成投资者遭理财产品夸大 和讯银行	2015-03-17	
[个股新闻] 平保《2015.版》获准认购平安银行非公开发行 阿思达克财经网	2015-03-17	
[个股新闻] 中国平安认购平安银行45%-50%定增股份 证券时报网	2015-03-17	
[个股新闻] 近半银行未披露理财产品运作信息 北京商报	2015-03-17	
[个股新闻] 平安银行邵平: 二季度成商业银行不良率转折点 21世纪经济报道	2015-03-17	
[个股新闻] 股份制商业银行理财能力更强 上海金融报	2015-03-17	
[个股新闻] 平安银行 金融时报	2015-03-17	

数据来源: 广发证券发展研究中心、和讯网

图 9 个股资讯获取示例

万科合生创展身陷质量门: 墙皮脱落 装修用纸
2015年03月12日 09:59 经济参考报 教育资讯 0人参与 0条评论
墙皮脱落 装修用“纸”
万科、合生创展: 身陷“质量门”
除墙皮外, 近年因精装修引发维权的企业不在少数。北京的合生世界花园项目日前继续曝出墙皮脱落、房间漏水严重。
一位业主表示, 房间内的墙皮在触碰后很容易脱落, 而脱落处露出了里面建筑材料之间存在的缝隙。甚至缝隙处还可以伸一条完整的木工板条来封堵进去。与此同时, 也出现了大量入门门开裂的情况。
而更让业主们难以接受的情况是, 这个商住两用的项目在规划之初曾经设计除大堂以外各层均设有两个公用卫生间, 并计入公摊。当业主们按照合同缴纳公摊费用之后, 却发现原本规划中的公用卫生间已经经开发商更改为住宅出售出去。业主们计算了相关损失认为, “消失的卫生间”大约涉及近9000万元左右的公摊费用。此外, 业主们还向当地住建部门投诉该项目在临水临电以及供暖收费等方面存在问题。
房企龙头也在2014年曝出“质量门”。报道显示, 万科被投诉的楼盘是位于广东佛山南海区的万科金域蓝湾, 该项目最早于2010年开盘。根据业主向媒体爆料称, 其于2011年购买了该项目二期一套精装修单元房, 收房时发现精装修房的装修都是用“纸”做的。

数据来源: 广发证券发展研究中心、新浪网

图 10 新闻资讯提取一览

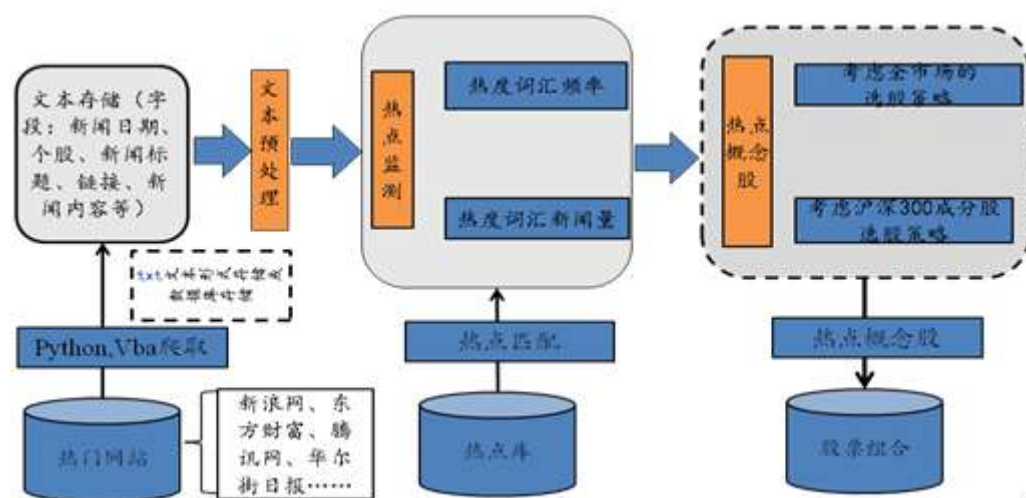
万科合生创展身陷质量门: 墙皮脱落 装修用纸
<pre> 1 <!-- 标题 --> 2 <!-- 副标题 --> 3 <!-- 正文 --> 4 <!-- 正文 --> 5 <!-- 正文 --> 6 <!-- 正文 --> 7 <!-- 正文 --> 8 <!-- 正文 --> 9 <!-- 正文 --> 10 <!-- 正文 --> 11 <!-- 正文 --> 12 <!-- 正文 --> 13 <!-- 正文 --> 14 <!-- 正文 --> 15 <!-- 正文 --> 16 <!-- 正文 --> 17 <!-- 正文 --> 18 <!-- 正文 --> 19 <!-- 正文 --> 20 <!-- 正文 --> 21 <!-- 正文 --> 22 <!-- 正文 --> 23 <!-- 正文 --> 24 <!-- 正文 --> 25 <!-- 正文 --> 26 <!-- 正文 --> 27 <!-- 正文 --> 28 <!-- 正文 --> 29 <!-- 正文 --> 30 <!-- 正文 --> 31 <!-- 正文 --> 32 <!-- 正文 --> 33 <!-- 正文 --> 34 <!-- 正文 --> 35 <!-- 正文 --> 36 <!-- 正文 --> 37 <!-- 正文 --> 38 <!-- 正文 --> 39 <!-- 正文 --> 40 <!-- 正文 --> 41 <!-- 正文 --> 42 <!-- 正文 --> 43 <!-- 正文 --> 44 <!-- 正文 --> 45 <!-- 正文 --> 46 <!-- 正文 --> 47 <!-- 正文 --> 48 <!-- 正文 --> 49 <!-- 正文 --> 50 <!-- 正文 --> 51 <!-- 正文 --> 52 <!-- 正文 --> 53 <!-- 正文 --> 54 <!-- 正文 --> 55 <!-- 正文 --> 56 <!-- 正文 --> 57 <!-- 正文 --> 58 <!-- 正文 --> 59 <!-- 正文 --> 60 <!-- 正文 --> 61 <!-- 正文 --> 62 <!-- 正文 --> 63 <!-- 正文 --> 64 <!-- 正文 --> 65 <!-- 正文 --> 66 <!-- 正文 --> 67 <!-- 正文 --> 68 <!-- 正文 --> 69 <!-- 正文 --> 70 <!-- 正文 --> 71 <!-- 正文 --> 72 <!-- 正文 --> 73 <!-- 正文 --> 74 <!-- 正文 --> 75 <!-- 正文 --> 76 <!-- 正文 --> 77 <!-- 正文 --> 78 <!-- 正文 --> 79 <!-- 正文 --> 80 <!-- 正文 --> 81 <!-- 正文 --> 82 <!-- 正文 --> 83 <!-- 正文 --> 84 <!-- 正文 --> 85 <!-- 正文 --> 86 <!-- 正文 --> 87 <!-- 正文 --> 88 <!-- 正文 --> 89 <!-- 正文 --> 90 <!-- 正文 --> 91 <!-- 正文 --> 92 <!-- 正文 --> 93 <!-- 正文 --> 94 <!-- 正文 --> 95 <!-- 正文 --> 96 <!-- 正文 --> 97 <!-- 正文 --> 98 <!-- 正文 --> 99 <!-- 正文 --> 100 <!-- 正文 --> 101 <!-- 正文 --> 102 <!-- 正文 --> 103 <!-- 正文 --> 104 <!-- 正文 --> 105 <!-- 正文 --> 106 <!-- 正文 --> 107 <!-- 正文 --> 108 <!-- 正文 --> 109 <!-- 正文 --> 110 <!-- 正文 --> 111 <!-- 正文 --> 112 <!-- 正文 --> 113 <!-- 正文 --> 114 <!-- 正文 --> 115 <!-- 正文 --> 116 <!-- 正文 --> 117 <!-- 正文 --> 118 <!-- 正文 --> 119 <!-- 正文 --> 120 <!-- 正文 --> 121 <!-- 正文 --> 122 <!-- 正文 --> 123 <!-- 正文 --> 124 <!-- 正文 --> 125 <!-- 正文 --> 126 <!-- 正文 --> 127 <!-- 正文 --> 128 <!-- 正文 --> 129 <!-- 正文 --> 130 <!-- 正文 --> 131 <!-- 正文 --> 132 <!-- 正文 --> 133 <!-- 正文 --> 134 <!-- 正文 --> 135 <!-- 正文 --> 136 <!-- 正文 --> 137 <!-- 正文 --> 138 <!-- 正文 --> 139 <!-- 正文 --> 140 <!-- 正文 --> 141 <!-- 正文 --> 142 <!-- 正文 --> 143 <!-- 正文 --> 144 <!-- 正文 --> 145 <!-- 正文 --> 146 <!-- 正文 --> 147 <!-- 正文 --> 148 <!-- 正文 --> 149 <!-- 正文 --> 150 <!-- 正文 --> 151 <!-- 正文 --> 152 <!-- 正文 --> 153 <!-- 正文 --> 154 <!-- 正文 --> 155 <!-- 正文 --> 156 <!-- 正文 --> 157 <!-- 正文 --> 158 <!-- 正文 --> 159 <!-- 正文 --> 160 <!-- 正文 --> 161 <!-- 正文 --> 162 <!-- 正文 --> 163 <!-- 正文 --> 164 <!-- 正文 --> 165 <!-- 正文 --> 166 <!-- 正文 --> 167 <!-- 正文 --> 168 <!-- 正文 --> 169 <!-- 正文 --> 170 <!-- 正文 --> 171 <!-- 正文 --> 172 <!-- 正文 --> 173 <!-- 正文 --> 174 <!-- 正文 --> 175 <!-- 正文 --> 176 <!-- 正文 --> 177 <!-- 正文 --> 178 <!-- 正文 --> 179 <!-- 正文 --> 180 <!-- 正文 --> 181 <!-- 正文 --> 182 <!-- 正文 --> 183 <!-- 正文 --> 184 <!-- 正文 --> 185 <!-- 正文 --> 186 <!-- 正文 --> 187 <!-- 正文 --> 188 <!-- 正文 --> 189 <!-- 正文 --> 190 <!-- 正文 --> 191 <!-- 正文 --> 192 <!-- 正文 --> 193 <!-- 正文 --> 194 <!-- 正文 --> 195 <!-- 正文 --> 196 <!-- 正文 --> 197 <!-- 正文 --> 198 <!-- 正文 --> 199 <!-- 正文 --> 200 <!-- 正文 --> 201 <!-- 正文 --> 202 <!-- 正文 --> 203 <!-- 正文 --> 204 <!-- 正文 --> 205 <!-- 正文 --> 206 <!-- 正文 --> 207 <!-- 正文 --> 208 <!-- 正文 --> 209 <!-- 正文 --> 210 <!-- 正文 --> 211 <!-- 正文 --> 212 <!-- 正文 --> 213 <!-- 正文 --> 214 <!-- 正文 --> 215 <!-- 正文 --> 216 <!-- 正文 --> 217 <!-- 正文 --> 218 <!-- 正文 --> 219 <!-- 正文 --> 220 <!-- 正文 --> 221 <!-- 正文 --> 222 <!-- 正文 --> 223 <!-- 正文 --> 224 <!-- 正文 --> 225 <!-- 正文 --> 226 <!-- 正文 --> 227 <!-- 正文 --> 228 <!-- 正文 --> 229 <!-- 正文 --> 230 <!-- 正文 --> 231 <!-- 正文 --> 232 <!-- 正文 --> 233 <!-- 正文 --> 234 <!-- 正文 --> 235 <!-- 正文 --> 236 <!-- 正文 --> 237 <!-- 正文 --> 238 <!-- 正文 --> 239 <!-- 正文 --> 240 <!-- 正文 --> 241 <!-- 正文 --> 242 <!-- 正文 --> 243 <!-- 正文 --> 244 <!-- 正文 --> 245 <!-- 正文 --> 246 <!-- 正文 --> 247 <!-- 正文 --> 248 <!-- 正文 --> 249 <!-- 正文 --> 250 <!-- 正文 --> 251 <!-- 正文 --> 252 <!-- 正文 --> 253 <!-- 正文 --> 254 <!-- 正文 --> 255 <!-- 正文 --> 256 <!-- 正文 --> 257 <!-- 正文 --> 258 <!-- 正文 --> 259 <!-- 正文 --> 260 <!-- 正文 --> 261 <!-- 正文 --> 262 <!-- 正文 --> 263 <!-- 正文 --> 264 <!-- 正文 --> 265 <!-- 正文 --> 266 <!-- 正文 --> 267 <!-- 正文 --> 268 <!-- 正文 --> 269 <!-- 正文 --> 270 <!-- 正文 --> 271 <!-- 正文 --> 272 <!-- 正文 --> 273 <!-- 正文 --> 274 <!-- 正文 --> 275 <!-- 正文 --> 276 <!-- 正文 --> 277 <!-- 正文 --> 278 <!-- 正文 --> 279 <!-- 正文 --> 280 <!-- 正文 --> 281 <!-- 正文 --> 282 <!-- 正文 --> 283 <!-- 正文 --> 284 <!-- 正文 --> 285 <!-- 正文 --> 286 <!-- 正文 --> 287 <!-- 正文 --> 288 <!-- 正文 --> 289 <!-- 正文 --> 290 <!-- 正文 --> 291 <!-- 正文 --> 292 <!-- 正文 --> 293 <!-- 正文 --> 294 <!-- 正文 --> 295 <!-- 正文 --> 296 <!-- 正文 --> 297 <!-- 正文 --> 298 <!-- 正文 --> 299 <!-- 正文 --> 300 <!-- 正文 --> 301 <!-- 正文 --> 302 <!-- 正文 --> 303 <!-- 正文 --> 304 <!-- 正文 --> 305 <!-- 正文 --> 306 <!-- 正文 --> 307 <!-- 正文 --> 308 <!-- 正文 --> 309 <!-- 正文 --> 310 <!-- 正文 --> 311 <!-- 正文 --> 312 <!-- 正文 --> 313 <!-- 正文 --> 314 <!-- 正文 --> 315 <!-- 正文 --> 316 <!-- 正文 --> 317 <!-- 正文 --> 318 <!-- 正文 --> 319 <!-- 正文 --> 320 <!-- 正文 --> 321 <!-- 正文 --> 322 <!-- 正文 --> 323 <!-- 正文 --> 324 <!-- 正文 --> 325 <!-- 正文 --> 326 <!-- 正文 --> 327 <!-- 正文 --> 328 <!-- 正文 --> 329 <!-- 正文 --> 330 <!-- 正文 --> 331 <!-- 正文 --> 332 <!-- 正文 --> 333 <!-- 正文 --> 334 <!-- 正文 --> 335 <!-- 正文 --> 336 <!-- 正文 --> 337 <!-- 正文 --> 338 <!-- 正文 --> 339 <!-- 正文 --> 340 <!-- 正文 --> 341 <!-- 正文 --> 342 <!-- 正文 --> 343 <!-- 正文 --> 344 <!-- 正文 --> 345 <!-- 正文 --> 346 <!-- 正文 --> 347 <!-- 正文 --> 348 <!-- 正文 --> 349 <!-- 正文 --> 350 <!-- 正文 --> 351 <!-- 正文 --> 352 <!-- 正文 --> 353 <!-- 正文 --> 354 <!-- 正文 --> 355 <!-- 正文 --> 356 <!-- 正文 --> 357 <!-- 正文 --> 358 <!-- 正文 --> 359 <!-- 正文 --> 360 <!-- 正文 --> 361 <!-- 正文 --> 362 <!-- 正文 --> 363 <!-- 正文 --> 364 <!-- 正文 --> 365 <!-- 正文 --> 366 <!-- 正文 --> 367 <!-- 正文 --> 368 <!-- 正文 --> 369 <!-- 正文 --> 370 <!-- 正文 --> 371 <!-- 正文 --> 372 <!-- 正文 --> 373 <!-- 正文 --> 374 <!-- 正文 --> 375 <!-- 正文 --> 376 <!-- 正文 --> 377 <!-- 正文 --> 378 <!-- 正文 --> 379 <!-- 正文 --> 380 <!-- 正文 --> 381 <!-- 正文 --> 382 <!-- 正文 --> 383 <!-- 正文 --> 384 <!-- 正文 --> 385 <!-- 正文 --> 386 <!-- 正文 --> 387 <!-- 正文 --> 388 <!-- 正文 --> 389 <!-- 正文 --> 390 <!-- 正文 --> 391 <!-- 正文 --> 392 <!-- 正文 --> 393 <!-- 正文 --> 394 <!-- 正文 --> 395 <!-- 正文 --> 396 <!-- 正文 --> 397 <!-- 正文 --> 398 <!-- 正文 --> 399 <!-- 正文 --> 400 <!-- 正文 --> 401 <!-- 正文 --> 402 <!-- 正文 --> 403 <!-- 正文 --> 404 <!-- 正文 --> 405 <!-- 正文 --> 406 <!-- 正文 --> 407 <!-- 正文 --> 408 <!-- 正文 --> 409 <!-- 正文 --> 410 <!-- 正文 --> 411 <!-- 正文 --> 412 <!-- 正文 --> 413 <!-- 正文 --> 414 <!-- 正文 --> 415 <!-- 正文 --> 416 <!-- 正文 --> 417 <!-- 正文 --> 418 <!-- 正文 --> 419 <!-- 正文 --> 420 <!-- 正文 --> 421 <!-- 正文 --> 422 <!-- 正文 --> 423 <!-- 正文 --> 424 <!-- 正文 --> 425 <!-- 正文 --> 426 <!-- 正文 --> 427 <!-- 正文 --> 428 <!-- 正文 --> 429 <!-- 正文 --> 430 <!-- 正文 --> 431 <!-- 正文 --> 432 <!-- 正文 --> 433 <!-- 正文 --> 434 <!-- 正文 --> 435 <!-- 正文 --> 436 <!-- 正文 --> 437 <!-- 正文 --> 438 <!-- 正文 --> 439 <!-- 正文 --> 440 <!-- 正文 --> 441 <!-- 正文 --> 442 <!-- 正文 --> 443 <!-- 正文 --> 444 <!-- 正文 --> 445 <!-- 正文 --> 446 <!-- 正文 --> 447 <!-- 正文 --> 448 <!-- 正文 --> 449 <!-- 正文 --> 450 <!-- 正文 --> 451 <!-- 正文 --> 452 <!-- 正文 --> 453 <!-- 正文 --> 454 <!-- 正文 --> 455 <!-- 正文 --> 456 <!-- 正文 --> 457 <!-- 正文 --> 458 <!-- 正文 --> 459 <!-- 正文 --> 460 <!-- 正文 --> 461 <!-- 正文 --> 462 <!-- 正文 --> 463 <!-- 正文 --> 464 <!-- 正文 --> 465 <!-- 正文 --> 466 <!-- 正文 --> 467 <!-- 正文 --> 468 <!-- 正文 --> 469 <!-- 正文 --> 470 <!-- 正文 --> 471 <!-- 正文 --> 472 <!-- 正文 --> 473 <!-- 正文 --> 474 <!-- 正文 --> 475 <!-- 正文 --> 476 <!-- 正文 --> 477 <!-- 正文 --> 478 <!-- 正文 --> 479 <!-- 正文 --> 480 <!-- 正文 --> 481 <!-- 正文 --> 482 <!-- 正文 --> 483 <!-- 正文 --> 484 <!-- 正文 --> 485 <!-- 正文 --> 486 <!-- 正文 --> 487 <!-- 正文 --> 488 <!-- 正文 --> 489 <!-- 正文 --> 490 <!-- 正文 --> 491 <!-- 正文 --> 492 <!-- 正文 --> 493 <!-- 正文 --> 494 <!-- 正文 --> 495 <!-- 正文 --> 496 <!-- 正文 --> 497 <!-- 正文 --> 498 <!-- 正文 --> 499 <!-- 正文 --> 500 <!-- 正文 --> 501 <!-- 正文 --> 502 <!-- 正文 --> 503 <!-- 正文 --> 504 <!-- 正文 --> 505 <!-- 正文 --> 506 <!-- 正文 --> 507 <!-- 正文 --> 508 <!-- 正文 --> 509 <!-- 正文 --> 510 <!-- 正文 --> 511 <!-- 正文 --> 512 <!-- 正文 --> 513 <!-- 正文 --> 514 <!-- 正文 --> 515 <!-- 正文 --> 516 <!-- 正文 --> 517 <!-- 正文 --> 518 <!-- 正文 --> 519 <!-- 正文 --> 520 <!-- 正文 --> 521 <!-- 正文 --> 522 <!-- 正文 --> 523 <!-- 正文 --> 524 <!-- 正文 --> 525 <!-- 正文 --> 526 <!-- 正文 --> 527 <!-- 正文 --> 528 <!-- 正文 --> 529 <!-- 正文 --> 530 <!-- 正文 --> 531 <!-- 正文 --> 532 <!-- 正文 --> 533 <!-- 正文 --> 534 <!-- 正文 --> 535 <!-- 正文 --> 536 <!-- 正文 --> 537 <!-- 正文 --> 538 <!-- 正文 --> 539 <!-- 正文 --> 540 <!-- 正文 --> 541 <!-- 正文 --> 542 <!-- 正文 --> 543 <!-- 正文 --> 544 <!-- 正文 --> 545 <!-- 正文 --> 546 <!-- 正文 --> 547 <!-- 正文 --> 548 <!-- 正文 --> 549 <!-- 正文 --> 550 <!-- 正文 --> 551 <!-- 正文 --> 552 <!-- 正文 --> 553 <!-- 正文 --> 554 <!-- 正文 --> 555 <!-- 正文 --> 556 <!-- 正文 --> 557 <!-- 正文 --> 558 <!-- 正文 --> 559 <!-- 正文 --> 560 <!-- 正文 --> 561 <!-- 正文 --> 562 <!-- 正文 --> 563 <!-- 正文 --> 564 <!-- 正文 --> 565 <!-- 正文 --> 566 <!-- 正文 --> 567 <!-- 正文 --> 568 <!-- 正文 --> 569 <!-- 正文 --> 570 <!-- 正文 --> 571 <!-- 正文 --> 572 <!-- 正文 --> 573 <!-- 正文 --> 574 <!-- 正文 --> 575 <!-- 正文 --> 576 <!-- 正文 --> 577 <!-- 正文 --> 578 <!-- 正文 --> 579 <!-- 正文 --> 580 <!-- 正文 --> 581 <!-- 正文 --> 582 <!-- 正文 --> 583 <!-- 正文 --> 584 <!-- 正文 --> 585 <!-- 正文 --> 586 <!-- 正文 --> 587 <!-- 正文 --> 588 <!-- 正文 --> 589 <!-- 正文 --> 590 <!-- 正文 --> 591 <!-- 正文 --> 592 <!-- 正文 --> 593 <!-- 正文 --> 594 <!-- 正文 --> 595 <!-- 正文 --> 596 <!-- 正文 --> 597 <!-- 正文 --> 598 <!-- 正文 --> 599 <!-- 正文 --> 600 <!-- 正文 --> 601 <!-- 正文 --> 602 <!-- 正文 --> 603 <!-- 正文 --> 604 <!-- 正文 --> 605 <!-- 正文 --> 606 <!-- 正文 --> 607 <!-- 正文 --> 608 <!-- 正文 --> 609 <!-- 正文 --> 610 <!-- 正文 --> 611 <!-- 正文 --> 612 <!-- 正文 --> 613 <!-- 正文 --> 614 <!-- 正文 --> 615 <!-- 正文 --> 616 <!-- 正文 --> 617 <!-- 正文 --> 618 <!-- 正文 --> 619 <!-- 正文 --> 620 <!-- 正文 --> 621 <!-- 正文 --> 622 <!-- 正文 --> 623 <!-- 正文 --> 624 <!-- 正文 --> 625 <!-- 正文 --> 626 <!-- 正文 --> 627 <!-- 正文 --> 628 <!-- 正文 --> 629 <!-- 正文 --> 630 <!-- 正文 --> 631 <!-- 正文 --> 632 <!-- 正文 --> 633 <!-- 正文 --> 634 <!-- 正文 --> 635 <!-- 正文 --> 636 <!-- 正文 --> 637 <!-- 正文 --> 638 <!-- 正文 --> 639 <!-- 正文 --> 640 <!-- 正文 --> 641 <!-- 正文 --> 642 <!-- 正文 --> 643 <!-- 正文 --> 644 <!-- 正文 --> 645 <!-- 正文 --> 646 <!-- 正文 --> 647 <!-- 正文 --> 648 <!-- 正文 --> 649 <!-- 正文 --> 650 <!-- 正文 --> 651 <!-- 正文 --> 652 <!-- 正文 --> 653 <!-- 正文 --></pre>

本质上，网站返回到本地电脑的是一个代码文件，代码文件经过浏览器的重新解析识别后，才会以各种非常规则漂亮的形式显示在我们的视野面前。网站返回的网页源代码其实是一种树状结构的，有根节点、子节点、元素、属性等。而HTMLDOM能够将这种树状结构的网页源代码存储到本地，从而能够方便快捷地利用HTMLDOM的相关语句以及语法对代码进行解析，从而获取到需要的文本信息。这对于大规模获取热门网站上的新闻的文本信息，效率上是有非常大的提高的。

2.3 基于互联网挖掘的热点选股策略框架

基于上述的完善的数据抓取平台，我们构建了此次专题策略的整体框架。

图 11 基于互联网挖掘的热点概念选股策略框架



数据来源：广发证券发展研究中心

整个信息变更的事件驱动策略流程分为三大模块，分别为个股新闻数据的抓取存储、热词词库的构建以及热度指标的构建、策略构建。

个股新闻抓取与存储模块主要是利用VBA+Python编程从指定的热门网站上获取到所有的关于个股新闻的文本信息，并将获取到的文本信息以指定的格式存储到本地磁盘，供后续建模调用处理。

图 12 txt 存储格式:个股代码_新闻发表时间_新闻标题

股票代码	新闻发表时间	新闻标题	新闻内容
600187	2015-02-04	中海发展(600187) 2015年第一次临时股东大会决议公告	中海发展(600187) 2015年第一次临时股东大会决议公告
600187	2015-02-04	中海发展(600187) 2015年第一次临时股东大会决议公告	中海发展(600187) 2015年第一次临时股东大会决议公告
600187	2015-02-04	中海发展(600187) 2015年第一次临时股东大会决议公告	中海发展(600187) 2015年第一次临时股东大会决议公告
600187	2015-02-04	中海发展(600187) 2015年第一次临时股东大会决议公告	中海发展(600187) 2015年第一次临时股东大会决议公告
600187	2015-02-04	中海发展(600187) 2015年第一次临时股东大会决议公告	中海发展(600187) 2015年第一次临时股东大会决议公告
600187	2015-02-04	中海发展(600187) 2015年第一次临时股东大会决议公告	中海发展(600187) 2015年第一次临时股东大会决议公告
600187	2015-02-04	中海发展(600187) 2015年第一次临时股东大会决议公告	中海发展(600187) 2015年第一次临时股东大会决议公告
600187	2015-02-04	中海发展(600187) 2015年第一次临时股东大会决议公告	中海发展(600187) 2015年第一次临时股东大会决议公告
600187	2015-02-04	中海发展(600187) 2015年第一次临时股东大会决议公告	中海发展(600187) 2015年第一次临时股东大会决议公告
600187	2015-02-04	中海发展(600187) 2015年第一次临时股东大会决议公告	中海发展(600187) 2015年第一次临时股东大会决议公告

图 13 数据库存储字段:个股、时间、新闻标题、对应链接

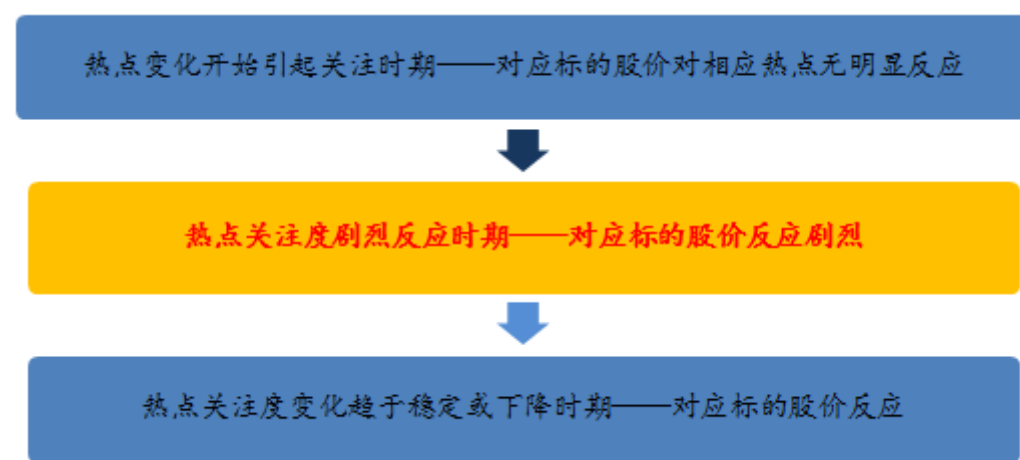
Stock	Title	publ_date	Link
600187	中海发展(600187) 2015年第一次临时股东大会决议公告	2015-02-04	http://finance.sina.com.cn/stock/2015-02-04/019C190086.shtml
600187	中海发展(600187) 2015年第一次临时股东大会决议公告	2015-02-04	http://finance.sina.com.cn/stock/2015-02-04/019C190086.shtml
600187	中海发展(600187) 2015年第一次临时股东大会决议公告	2015-02-04	http://finance.sina.com.cn/stock/2015-02-04/019C190086.shtml
600187	中海发展(600187) 2015年第一次临时股东大会决议公告	2015-02-04	http://finance.sina.com.cn/stock/2015-02-04/019C190086.shtml
600187	中海发展(600187) 2015年第一次临时股东大会决议公告	2015-02-04	http://finance.sina.com.cn/stock/2015-02-04/019C190086.shtml
600187	中海发展(600187) 2015年第一次临时股东大会决议公告	2015-02-04	http://finance.sina.com.cn/stock/2015-02-04/019C190086.shtml
600187	中海发展(600187) 2015年第一次临时股东大会决议公告	2015-02-04	http://finance.sina.com.cn/stock/2015-02-04/019C190086.shtml
600187	中海发展(600187) 2015年第一次临时股东大会决议公告	2015-02-04	http://finance.sina.com.cn/stock/2015-02-04/019C190086.shtml
600187	中海发展(600187) 2015年第一次临时股东大会决议公告	2015-02-04	http://finance.sina.com.cn/stock/2015-02-04/019C190086.shtml
600187	中海发展(600187) 2015年第一次临时股东大会决议公告	2015-02-04	http://finance.sina.com.cn/stock/2015-02-04/019C190086.shtml

数据来源：广发证券发展研究中心

数据来源：广发证券发展研究中心

热词词库的构建以及热度指标的构建模块整体上来说主要分类为两大部分，一部分是热词词库的构建，在热词词库构建方面，根据我们广发金工在文本挖掘领域的深入的研究，已经积累了大量的历史上相关的热点、概念词库，同时也基于每天的市场行情的变化，动态地往热词词库中添加新的可能的热点词汇或者删除部分已经失效的热点词汇，从而基于热词词库来动态地研究每一天可能出现的热点、概念，并通过建模来研究这种热点、概念对应标的的表现情况。热词、概念的热度指标的构建，可以通过扫描热词词库中每一个热词在指定的日期内的新闻量，如果出现新闻量或者是热词的频率跟过去的一段时间内相比，有一个继续上升的趋势，并且上升的速率超过给定的阈值，则说明该热词可能成为当前的一个热点，相应的标的则可能存在的上升机会。

图 14 热点概念传播与对应的股价关系



数据来源：广发证券发展研究中心

图 15 热词词库示例

热词词库示例				
一带一路	人脑工程	海洋执法	人工角膜	安防服务
沪港通概念	基金三方销售	沪港通行软件	易信	智能表
奶牛养殖	云停车场	土豆涨价	病毒检测芯片	上海金改
江北新区	碳海绵	GEP	航空租赁	百度金融
辉钼	新疆旅游	翡翠涨价	新疆电源项目	上海国资委重组
重庆水资源	五一出境游	聚四氟乙烯	浦东前滩	硅铁
海藻炼油	参股申万	殡葬	海上丝路	低碳经济
沪新通	铂矿	沪港通结算银行	京津冀	黄金概念
大麻概念	临近空间	四人行	粤港澳	燃料电池
帕拉米韦	烟花爆竹	长三角	次新股	水利建设
疫苗存储	老八股	珠三角	含H股	电子支付
农业综改	参股外资金融	武汉规划	含可转债	海水淡化
参股银河证	上海崇明岛	丝绸之路	军工航天	抗流感
维生素	智能交通	信息安全	免疫医疗	智能机器
文化振兴	空气治理	电商	特岗	央企改革
节能	污水治理	网贷	钛金属	全息概念
金融改革	固废治理	油气改革	海外工程	绿色照明
农村金融	装饰园林	苹果	聚氨基酯	赛马
宽带提速	体育	生物智能		菠菜

数据来源：广发证券发展研究中心

策略构建模块主要是根据通过动态地监控热词词库中每个热词当前时刻的热度值，然后根据当前时刻的热度值跟该热词在过去一段时间的热度值进行比较，如果当前的热词相比过去一段时间的热度值的变化率超过一定的阈值则将该热词选定为

当前时点的热点，并对该热词对应的标的进行组合的配置，持有一段时间，在持有期末卖出的个股。

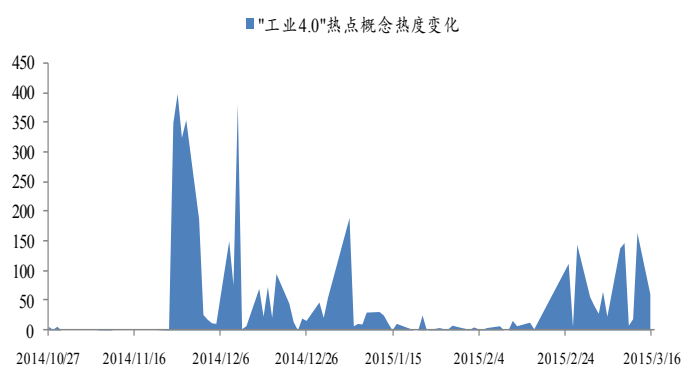
图 16 热点概念热度指标构建步骤



数据来源：广发证券发展研究中心

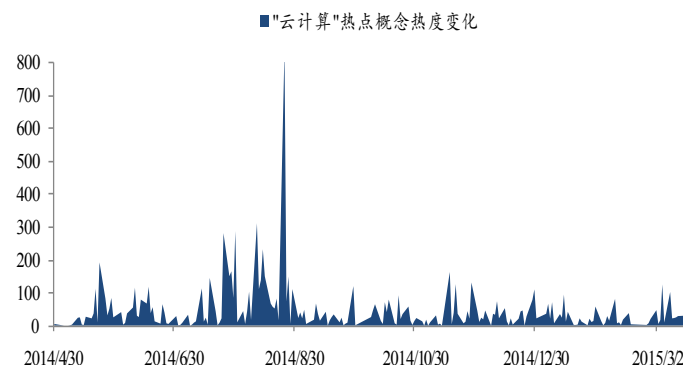
热点概念热度计算过程中，首先对抓取下来的新闻文本信息进行预处理：每个新闻对应的网站进行标注，对在之前抓取下来内容为空的新闻对应的链接进行重新抓取；其次对热词词库中的每一个热词计算热度指标，并将其存储到对应的热度词库矩阵中；最后根据热度词库中的每个热度的变化识别出每一个交易日对应的热点概念标的，然后配置对应的标的。

图 17 "工业 4.0"热点概念热度变化



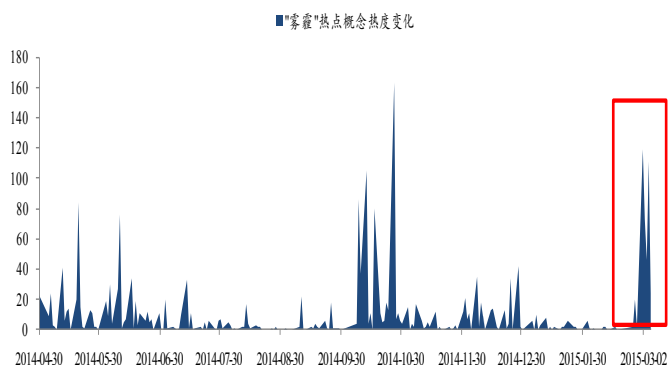
数据来源：广发证券发展研究中心

图 18 "云计算"热点概念热度变化



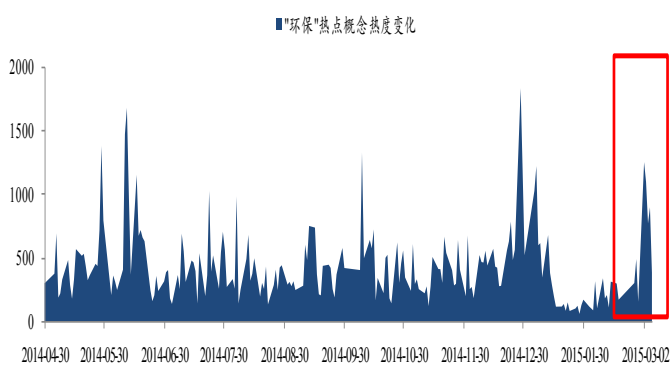
数据来源：广发证券发展研究中心

图 19 "雾霾"热点概念热度变化



数据来源：广发证券发展研究中心

图 20 "环保"热点概念热度变化

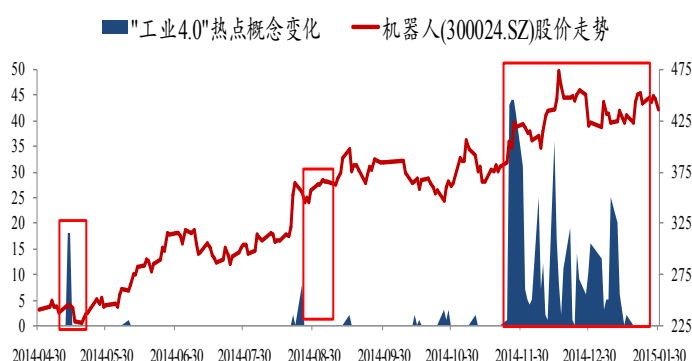


数据来源：广发证券发展研究中心

从图 19、图 20 可以看出，近期由于柴静的纪录片“穹顶之下”中关于雾霾的报道，使得人们对“雾霾”、“环保”等的关注度在短期内就呈现一个剧烈上升的趋势。

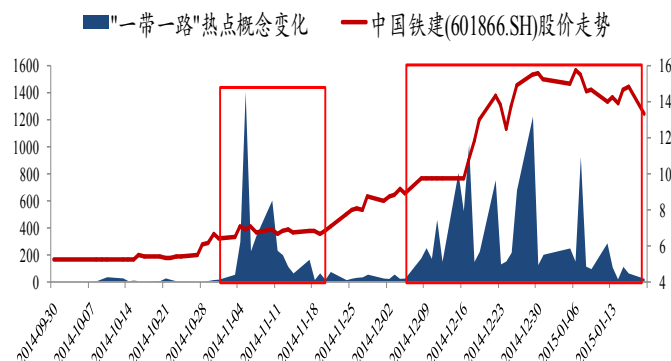
势，而相应的个股在短期内也可能比同期的大盘表现要更加地优异。

图 21 热点变化与对应标的股价变化示例 1



数据来源：广发证券发展研究中心

图 22 热点变化与对应标的股价变化示例 2



数据来源：广发证券发展研究中心

从图 21、22 可以看出，每当对应的热点改变呈现比较剧烈的一个上升趋势时，相对应的个股也往往会表现出一个上涨的趋势。

三、 基于互联网挖掘的热点选股策略实证

3.1 样本数据

由上述的策略框架可得，本专题策略的数据主要来源于新浪网、和讯网、腾讯网、东方财富网这四个热门网站。所需的数据包括热门网站上关于个股新闻资讯的历史文本信息以及 A 股市场上所有个股的历史交易数据，包括开盘价、收盘价、成交量、交易状态等数据。本专题策略选择实证分析的样本区间为 2014 年 4 月 30 日至今。策略的历史收益比较基准为沪深 300 指数。

3.2 策略构建原理

初始资金：1；

策略原理：在历史回测期间内，将资金等分为若干份。动态地监控在每个交易日内是否有热点概念出现，如果有热点概念出现，则将其中的一份没有持仓个股的资金配置当天出现的热点概念的个股，持有一段时间，不考虑个股在买进或者卖出时候出现涨停、跌停或者停牌的情况，个股在买进时以该个股当日的开盘价买进，卖出时以该个股在持有期末的收盘价卖出。

交易费用：千分之二费用，在卖出时候收取。

资金投资权重：等权投资于个股；

3.3 互联网挖掘的热点选股策略

基于以上的交易策略原理，将策略应用于沪深300指数中，具体的交易策略原理为：初始资金为1，在历史回测期内，将资金等分为10份，通过扫描每一日热词词库

中每个热词的热度，将该日热词的热度与过去一段交易日内该热词的热度进行比较，选择相对增长率最大的前两个热点作为当期的热点概念，并用没有配置个股的一份资金配置该热点所对应的个股，持有期为10个交易日，资金等权投资于该个股。个股在买进时以买进日的开盘价买进，在持有期末以收盘价卖出。

图 23 热点选股策略表现(相对沪深 300)



数据来源：广发证券发展研究中心

表 2 热点选股策略表现指标一览

指标	数据
年化绝对收益率	112.08%
年化相对收益率	19.00%
最大回撤	-6.27%
周胜率	65.85%
月胜率	81.82%

数据来源：广发证券发展研究中心

从上述图与表可以看出，热点选股策略在历史回测期内取得了较优异的结果。相对沪深 300 指数，在整个回测期内，实现了约 19.00%的年化超额收益率以及 112.08%的年化绝对收益率；在胜率方面胜率，取得了 65.85%的周胜率以及 81.82%的月度胜率；从回撤角度看，在没有用沪深 300 进行对冲前，策略累计净值的回撤为 6.27%，在用沪深 300 进行对冲后，策略对冲后的净值的回撤为 3.59%。整体上而言，基于对大规模的热门网站的文本信息挖掘构建的热点选股策略，在历史回测期内取得相对较为优异的表现。

四、 总结

4.1 总结

基于大规模的互联网财经文本信息，本专题策略通过采用文本挖掘的方法研究

这些文本信息中所包含的热点、概念变化带来的投资机会。通过构建热词的热度变化指标，建立量化投资策略，考虑了热点、概念变化带来的投资机会。

通过所构建的量化策略，对热点、概念变化所引起的相应的个股的价格的变化进行了历史回测，实证结果表明，所构建的量化策略在回测期间内，表现优异，主要结论有：

- 本专题对互联网大规模的新闻数据的文本信息进行了量化处理，建立了基于热点概念识别的文本挖掘选股策略；
- 基于所建立的量化策略，可以动态地识别出每日可能的热点、概念并筛选出对应个股；

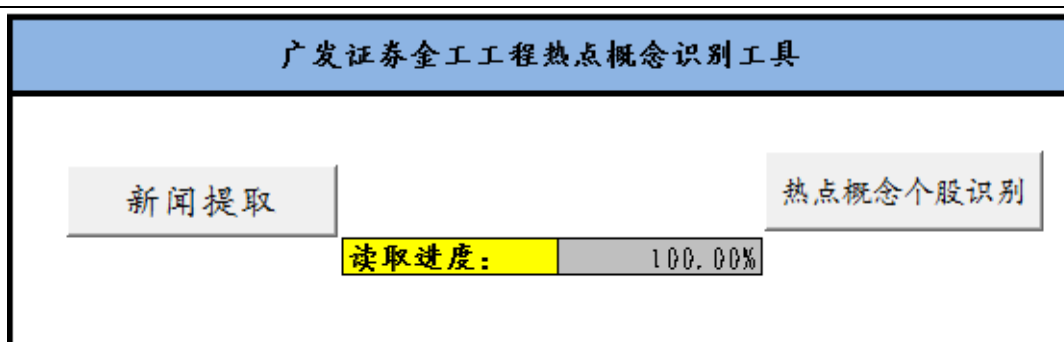
4.2 未来研究方向

- 动态更新并维护热点概念词库；
- 研究基于全市场的热点概念文本挖掘选股策略；
- 基于互联网大规模的新闻数据，每日定时推送相应的宏观热点(国内、国外)、热点概念以及对应的个股；
- 个性化需求定制；

4.3 工具推荐

基于对大规模的互联网热门网站上的个股信息相关网页的网页格式的研究，我们开发出了基于互联网热门网站的热点概念识别工具。这个工具能根据用户对不同网站的选择或者是基于所跟踪的所有的热门网站，根据本专题策略关于热点概念构建的算法，自动地识别出当前市场上可能的热点概念，并输入对应热点概念的标的个股，同时该工具能够显示用户自定义的热点概念在历史的走势情况，以图表的形式呈现。目前，该工具已经开发接近完善，即将外发给客户在本地使用，欢迎发邮件与我们索取或交流。

图 24 广发证券金融工程热点概念识别工具



数据来源：广发证券发展研究中心

风险提示

本模型为采用纯量化方法，所推荐的个股未必具有实质性的利好，其股

价表现还受到诸多因素影响，请结合基本面及自身判断进行恰当使用。

广发证券—行业投资评级说明

- 买入： 预期未来 12 个月内，股价表现强于大盘 10%以上。
- 持有： 预期未来 12 个月内，股价相对大盘的变动幅度介于-10%~+10%。
- 卖出： 预期未来 12 个月内，股价表现弱于大盘 10%以上。

广发证券—公司投资评级说明

- 买入： 预期未来 12 个月内，股价表现强于大盘 15%以上。
- 谨慎增持： 预期未来 12 个月内，股价表现强于大盘 5%-15%。
- 持有： 预期未来 12 个月内，股价相对大盘的变动幅度介于-5%~+5%。
- 卖出： 预期未来 12 个月内，股价表现弱于大盘 5%以上。

联系我们

	广州市	深圳市	北京市	上海市
地址	广州市天河北路183号大都会广场5楼	深圳市福田区金田路4018号安联大厦15楼A座03-04	北京市西城区月坛北街2号月坛大厦18层	上海市浦东新区富城路99号震旦大厦18楼
邮政编码	510075	518026	100045	200120
客服邮箱	gfyf@gf.com.cn			
服务热线	020-87555888-8612			

免责声明

广发证券股份有限公司具备证券投资咨询业务资格。本报告只发送给广发证券重点客户，不对外公开发布。

本报告所载资料的来源及观点的出处皆被广发证券股份有限公司认为可靠，但广发证券不对其准确性或完整性做出任何保证。报告内容仅供参考，报告中的信息或所表达观点不构成所涉证券买卖的出价或询价。广发证券不对因使用本报告的内容而引致的损失承担任何责任，除非法律法规有明确规定。客户不应以本报告取代其独立判断或仅根据本报告做出决策。

广发证券可发出其它与本报告所载信息不一致及有不同结论的报告。本报告反映研究人员的不同观点、见解及分析方法，并不代表广发证券或其附属机构的立场。报告所载资料、意见及推测仅反映研究人员于发出本报告当日的判断，可随时更改且不予通告。

本报告旨在发送给广发证券的特定客户及其它专业人士。未经广发证券事先书面许可，任何机构或个人不得以任何形式翻版、复制、刊登、转载和引用，否则由此造成的一切不良后果及法律责任由私自翻版、复制、刊登、转载和引用者承担。