

# 司法大数据自动化标注与分析说明文档

---

## 司法大数据自动化标注与分析说明文档

- 1. 使用方法
  - 1.1 启动网页
  - 1.2 选择爬虫
  - 1.3 分词标注与保存
- 2. 具体实现
  - 2.1 项目结构
  - 2.2 前端实现
    - 2.2.1 时间选择器
    - 2.2.2 各种按钮
    - 2.2.3 单选框
    - 2.2.4 复选框
  - 2.3 爬虫实现
    - 2.3.1 反爬
    - 2.3.2 请求
    - 2.3.3 解析与下载
  - 2.4 分词实现
    - 2.4.1 方法分析
    - 2.4.2 分词过程
  - 2.5 前后连接
    - 2.5.1 初始
    - 2.5.2 爬虫
    - 2.5.3 分词
- 3. 不足
- 4. 分工

## 1. 使用方法

---

网页示意图：

司法大数据标注与分析

开始日期

结束日期

爬取数量

默认认为10

爬取案件

上传案例文件

案件文本：

中华人民共和国最高人民法院  
刑 事 判 决 书  
被告人辛志勇，化名张先奇、付仙阳，绰号“少林”“阿志”“智能”，男，汉族，1976年6月12日出生于江西省万载县，初中文化，无业，户籍地万载县××乡××村××组××号，暂住地广东省佛山市××区××街道××村××大街×××号出租屋。1996年10月22日因犯盗窃罪被判处有期徒刑十年，剥夺政治权利三年，2002年10月17日被假释，假释考验期至2004年2月23日止；2014年12月1日因犯职务侵占罪被判处有期徒刑一年六个月，2015年10月1日刑满释放。2017年5月19日因本案被逮捕。现在押。江西省宜春市中级人民法院审理宜春市人民检察院指控被告人辛志勇犯抢劫罪、盗窃罪一案，于2018年9月21日以（2018）赣09刑初5号刑事判决，认定被告人辛志勇犯抢劫罪，判处死刑，剥夺政治权利终身，并处没收个人全部财产；犯盗窃罪，判处有期徒刑十二年，并处罚金人民币十万元，决定执行死刑，剥夺政治权利终身，并处没收个人全部财产。宣判后，辛志勇提出上诉。江西省高级人民法院经依法开庭审理，于2019年3月1日以（2018）赣刑终293号刑事裁定，驳回上诉，维持原判，并依法报请本院核准。本院依法组成合议庭，对本案进行了复核，依法讯问了被告人。现已复核终结。  
经审理查明：  
一、抢劫事实  
2009年2月上旬，被告人辛志勇与同乡林啟粮、辛钧亮（均系另案被告人，已判刑）预谋前往江西省万载县经济条件较好的人家攫取财物。同月18日，辛志勇和辛钧亮从广东省回到万载县与林啟粮会合，后辛钧亮离开。同月23日上午，辛志勇与林啟粮携带菜刀、手套、电击防爆枪等作案工具，来到万载县××镇××村×组陈某1家屋后伺机作案。当日下午，辛志勇、林啟粮进入陈某1家，被陈的妻子李某1（被害人，殁年44岁）发现。辛志勇见状持陈家的木棍连续击打李某1的头部，林啟粮扼拍李某1的颈部，二人将李某1按倒在地，用围裙绳勒颈并用菜刀割颈，致李某1颈部动脉横断、大出血休克死亡。其间，林啟粮的手指被李某1咬伤。二人从李某1的身上和家中搜得现金750元、黄金耳环一副、黄金戒指、玉币和银链子等物后逃离陈家。随后，林啟粮打电话通知表兄卓秋明（另案被告人，已判刑）驾驶摩托车前往现场附近将二人接回。当晚，辛志勇、林啟粮、辛钧亮驾驶摩托车逃往广东省。辛志勇、林啟粮将所劫黄金首饰变卖并平分赃款，玉币由林啟粮留存。

开始分词

清空案例文本

清空所有标注

信息标注：

当事人	性别	民族	出生地	文化水平	案由	相关法院
不涉及地点名词						
<input type="checkbox"/> 郑某	<input type="checkbox"/> 梅花	<input type="checkbox"/> 孟某	<input checked="" type="checkbox"/> 辛志勇	<input type="checkbox"/> 卓秋明		
<input type="checkbox"/> 辛钧亮	<input type="checkbox"/> 阿志	<input type="checkbox"/> 林啟粮	<input type="checkbox"/> 刘某	<input type="checkbox"/> 初中文化		
<input type="checkbox"/> 男	<input type="checkbox"/> 汉族	<input type="checkbox"/> 桂秀泉	<input type="checkbox"/> 付仙阳	<input type="checkbox"/> 刑事判决		
<input type="checkbox"/> 陈均	<input type="checkbox"/> 张先奇	<input type="checkbox"/> 万载县	<input type="checkbox"/> 段某	<input type="checkbox"/> 何某		
<input type="checkbox"/> 李某	<input type="checkbox"/> 陈述	<input type="checkbox"/> 陈某	<input type="checkbox"/> 翟培江	<input type="checkbox"/> 方文军		
<input type="checkbox"/> 袁某	<input type="checkbox"/> 黄金首饰	<input type="checkbox"/> 一审判决				
涉及地点名词						
<input type="checkbox"/> 仙阳	<input type="checkbox"/> 西路	<input type="checkbox"/> 大街	<input type="checkbox"/> 江西省万载县	<input type="checkbox"/> 之日起		
<input type="checkbox"/> 北京	<input type="checkbox"/> 广东省佛山市	<input type="checkbox"/> 宜春市	<input type="checkbox"/> 瑞纳	<input type="checkbox"/> 最高人民法院		
<input type="checkbox"/> 宜春市中级人民法院	<input type="checkbox"/> 江西省宜春市中级人民法院	<input type="checkbox"/> 江西省宜春市	<input type="checkbox"/> 广东省广州市	<input type="checkbox"/> 仙海		
<input type="checkbox"/> 中华人民共和国最高人民法院	<input type="checkbox"/> 中华人民共和国	<input type="checkbox"/> 广州市	<input type="checkbox"/> 江西省高级人民法院	<input type="checkbox"/> 广东省		
<input type="checkbox"/> 佛山市	<input type="checkbox"/> 江西省					
案由						
<input type="checkbox"/> 抢劫	<input type="checkbox"/> 盗窃					
形容词						
<input type="checkbox"/> 具体	<input type="checkbox"/> 连续	<input type="checkbox"/> 准确	<input type="checkbox"/> 残忍	<input type="checkbox"/> 清		

保存案件与标注

## 1.1 启动网页

- 1. 在PyCharm中打开该项目，导入相关包
- 2. 在app.py文件中启动main函数
- 3. 在浏览器地址栏输入 <http://127.0.0.1:5000/> 即可进入网页

## 1.2 选择爬虫

在网页中输入开始日期、结束日期与爬取数量（默认为10），点击爬取案件按钮即可爬虫，此时将会跳

### 爬虫结果展示

爬取文书存放在项目的static/Downloads目录下，文书名称如下：

- 王某涉挪用资金罪刑事一审案件刑事判决书
- 陈学德涉非法经营罪刑事一审案件刑事判决书
- 陈春飞涉非法经营罪刑事一审案件刑事判决书
- 李振环涉帮助信息网络犯罪活动罪刑事一审案件刑事判决书
- 刘国扬故意伤害刑事一审刑事判决书
- 李刚危险驾驶罪一审刑事判决书
- 李某某故意伤害罪、故意伤害罪刑事一审刑事判决书
- 罗楚国帮助信息网络犯罪活动罪刑事一审刑事判决书
- 莫汉前非法捕捞水产品罪刑事一审刑事判决书
- 陈洪学交通肇事刑事一审刑事判决书

转到爬虫展示页，如图：

若未显示名称，则爬取失败

## 1.3 分词标注与保存

- 1. 手动输入文本或者上传案例文件（txt格式）
- 2. 点击开始分词按钮
- 3. 在分词结果中勾选标注
- 4. 点击保存案件与标注按钮即可保存案件文本与标注信息

PS：由于案例文本、标注的清空手动较为麻烦，故增加了一键清空文本与标注的功能按钮

## 2. 具体实现

以flask框架为基础，综合爬虫、HTML、CSS、JavaScript、jieba分词、BootStrap、Jinja2等技能实现该项目

### 2.1 项目结构

```
1 Program: .
2 |
3 |  app.py ..... 项目入口，连接前后端，启动网页
4 |  jiebaAnalysis.py ..... jieba分词
5 |  wenshu_spider.py ..... 爬虫源代码
6 |
7 |
8 |—static ..... 存放js, css, txt等文件
9 |   | dict.txt ..... jieba分词自定义词典
10 |   | 司法大数据自动化标注与分析.css ..... 网页除时间选择器的style
11 |   | 司法大数据自动化标注与分析.js ..... 网页除时间选择器的js源码
12 |   |
13 |   |—bootstrap ..... 用于构造时间选择器组件
14 |   |
15 |   |
16 |   |—Downloads ..... 文书爬取结果文件夹
17 |   |   刘国扬故意伤害刑事一审刑事判决书.txt
18 |   |   .....
19 |   |
20 |   |—jquery ..... 用于实现时间选择器功能
```

21			
22			
23		└js .....	用于实现时间选择器功能
24			
25			
26		└templates .....	网页HTML源文件
27		司法大数据自动化标注与分析.html .....	主体网页
28		文书爬取.html .....	爬虫结果展示网页

## 2.2 前端实现

前端页面主要包括以下几个重要组件：

时间选择器、各种按钮、文本域、单选框、复选框，其中文本域实现简单，此处不再赘述。

### 2.2.1 时间选择器

时间选择器的实现利用了Bootstrap框架，css和js代码直接copy了过来，放在了文件中，此处不再展示，下面为html代码

```
1 <div class="container">
2     <form action="/文书爬取" method="post" class="form-horizontal"
      role="form">
3         <div class="row">
4             <label for="dtp_input1" class="col-md-1 control-label">开始日期
      </label>
5             <div class="input-group date form_date col-md-3" id="begindate"
      data-date="" data-date-format="yyyy.MM.dd" data-link-field="dtp_input1"
      data-link-format="yyyy-mm-dd">
6                 <input class="form-control" name="beginDate" size="16"
      type="text" id="beginText" value="" readonly>
7                 <span class="input-group-addon"><span class="glyphicon
      glyphicon-remove"></span></span>
8                 <span class="input-group-addon"><span class="glyphicon
      glyphicon-calendar"></span></span>
9             </div>
10            <input type="hidden" id="dtp_input1" value="" />
11
12            <label for="dtp_input2" class="col-md-1 control-label">结束日期
      </label>
13            <div class="input-group date form_date col-md-3" id="enddate"
      data-date="" data-date-format="yyyy.MM.dd" data-link-field="dtp_input2"
      data-link-format="yyyy-mm-dd">
14                <input class="form-control" name="endDate" size="16"
      type="text" id="endText" value="" readonly>
15                <span class="input-group-addon"><span class="glyphicon
      glyphicon-remove"></span></span>
16                <span class="input-group-addon"><span class="glyphicon
      glyphicon-calendar"></span></span>
17            </div>
18            <input type="hidden" id="dtp_input2" value="" />
19
20            <label for="num" class="col-md-1 control-label">爬取数量</label>
21            <div class="input-group col-md-2">
22                <input type="text" id="num" name="num" class="text"
      placeholder="默认为10">
23            </div>
```

```

24         <div class="col-md-1 control-label">
25             <input type="submit" id="spiderButton" value="爬取案件" />
26         </div>
27     </div>
28 </form>
29 </div>

```

其中增添了结束日期不能早于开始日期的设定，并且将语言设置为中文，实现如下：

```

1  $('form_date').datetimepicker({
2      language: 'zh-CN', //设置语言为中文
3      weekStart: 1,
4      todayBtn: 1,
5      autoclose: 1,
6      todayHighlight: 1,
7      startView: 2,
8      minView: 2,
9      forceParse: 0
10 }).on('changeDate', function (e) {
11     var BeginTime = $("#beginText").val();
12     $("#enddate").datetimepicker("setStartDate", BeginTime); //设置结束
    时间只能从开始时间选择起
13 });

```

## 2.2.2 各种按钮

按钮的重点在于功能的实现，以下介绍了实现一些按钮功能的js代码，另外一些按钮需要后端支持

```

1  /**
2   * 上传文件并在textarea中显示
3   */
4  window.onload = function() {
5      /**
6       * 上传函数
7       * @param fileInput DOM对象
8       * @param callback 回调函数
9       */
10     var getFileContent = function (fileInput, callback) {
11         if (fileInput.files && fileInput.files.length > 0 &&
12             fileInput.files[0].size > 0) {
13             //下面这一句相当于jQuery的: var file = $("#upload").prop('files')
14             [0];
15             var file = fileInput.files[0];
16             if (window.FileReader) {
17                 var reader = new FileReader();
18                 reader.onloadend = function (evt) {
19                     if (evt.target.readyState === FileReader.DONE) {
20                         callback(evt.target.result);
21                     }
22                 };
23             }
24             // 包含中文内容用utf-8编码
25             reader.readAsText(file, 'utf-8');
26         }
27     };
28 }

```

```

27     /**
28      * upload内容变化时载入内容
29      */
30     document.getElementById('uploadButton').onchange = function () {
31         var textArea = document.getElementById('textArea');
32         getFileContent(this, function (str) {
33             textArea.value = str;
34         });
35     };
36 };
37
38 /**
39  * 清空案例文本
40  */
41 function clearContent(){
42     var textArea = document.getElementById('textArea'); // 文本域文档元素：案
    件文本
43     textArea.value = "";
44 }
45
46 /**
47  * 清空标注
48  */
49 function clearMark(){
50     var tab__contentList =
    document.getElementsByClassName("tab__content");//长度为7
51     for(var i = 0; i < tab__contentList.length; i++){
52         var checkboxes = tab__contentList[i].getElementsByTagName("input");
53         for(var j = 0; j < checkboxes.length; j++){
54             checkboxes[j].checked = false;
55         }
56     }
57 }
58
59 /**
60  * 保存案件文本和标注
61  */
62 function saveFile(){
63     saveTxt();
64     saveJSON();
65 }
66 // 保存案件文本.txt
67 function saveTxt(){
68     var textArea = document.getElementById('textArea'); // 文本域文档元素：案
    件文本
69     var textAreaContent = textArea.value; // 具体内容
70     if(textAreaContent){ // 有内容
71         var blob = new Blob([textAreaContent], {type:
    "text/plain;charset=utf-8"});
72         saveAs(blob, "案件文本.txt");
73         //textArea.value="";// 清空输入框
74     }else{ // 无内容
75         window.alert("你尚未输入案件文本，请重新输入");
76     }
77 }
78 // 保存标注.json
79 function saveJSON(){

```

```

80     var tab__contentList =
document.getElementsByClassName("tab__content");//长度为7
81     var markList = new Array(tab__contentList.length);
82     for(var i = 0; i < tab__contentList.length; i++){
83         var checkboxes = tab__contentList[i].getElementsByTagName("input");
84         markList[i] = "";
85         for(var j = 0; j < checkboxes.length; j++){
86             if(checkboxes[j].checked){
87                 var label = checkboxes[j].parentNode;
88                 if(markList[i].length != 0) {
89                     markList[i] += ("、" + label.innerText);
90                 }else {
91                     markList[i] = label.innerText;
92                 }
93             }
94         }
95     }
96
97     var nameList = ["当事人", "性别", "民族", "出生地", "文化水平", "案由", "相关
法院"];
98     var textContents = [];
99     for (var i = 0; i < 7; i++){
100         var textJSONString = '"' + nameList[i] + ':' + '"' + markList[i] +
''';
101         textContents.push(textJSONString);
102     }
103     var textListJSONString = "{" + textContents.join(",") + "}"; // JSON字
字符串
104     var blob = new Blob([textListJSONString], {type:
"text/plain;charset=utf-8"});
105     saveAs(blob, "标注.json");
106 }

```

### 2.2.3 单选框

单选框实现不难，但要通过多个单选框实现选项卡有些困难，重点在于CSS代码

```

1  <div class="tab-wrap" id="jiebas">
2
3      <!-- active tab on page load gets checked attribute -->
4      <input type="radio" id="party" name="Mark" class="tab" checked>
5      <label for="party">当事人</label>
6
7      <input type="radio" id="sex" name="Mark" class="tab">
8      <label for="sex">性别</label>
9
10     <input type="radio" id="nationality" name="Mark" class="tab">
11     <label for="nationality">民族</label>
12
13     <input type="radio" id="birthPlace" name="Mark" class="tab">
14     <label for="birthPlace">出生地</label>
15
16     <input type="radio" id="education" name="Mark" class="tab">
17     <label for="education">文化水平</label>
18
19     <input type="radio" id="cause" name="Mark" class="tab">
20     <label for="cause">案由</label>

```

```
21
22         <input type="radio" id="court" name="Mark" class="tab">
23         <label for="court">相关法院</label>
24     </div>
```

CSS代码:

```
1  .tab-wrap {
2      transition: 0.3s box-shadow ease;
3      border-radius: 6px;
4      max-width: 100%;
5      display: flex;
6      flex-wrap: wrap;
7      position: relative;
8      list-style: none;
9      background-color: #fff;
10     margin: 40px 0;
11     box-shadow: 0 1px 3px rgba(0, 0, 0, 0.12), 0 1px 2px rgba(0, 0, 0,
12     0.24);
13 }
14 .tab-wrap:hover {
15     box-shadow: 0 12px 23px rgba(0, 0, 0, 0.23), 0 10px 10px rgba(0, 0, 0,
16     0.19);
17 }
18 .tab__content {
19     padding: 10px 25px;
20     background-color: transparent;
21     position: absolute;
22     width: 100%;
23     z-index: -1;
24     opacity: 0;
25     left: 0;
26     transform: translateY(-3px);
27     border-radius: 6px;
28 }
29 .tab {
30     display: none;
31 }
32 .tab:checked:nth-of-type(1) ~ .tab__content:nth-of-type(1) {
33     opacity: 1;
34     transition: 0.5s opacity ease-in, 0.8s transform ease;
35     position: relative;
36     top: 0;
37     z-index: 100;
38     transform: translateY(0px);
39     text-shadow: 0 0 0;
40 }
41 .tab:checked:nth-of-type(2) ~ .tab__content:nth-of-type(2) {
42     opacity: 1;
43     transition: 0.5s opacity ease-in, 0.8s transform ease;
44     position: relative;
45     top: 0;
46     z-index: 100;
47     transform: translateY(0px);
48     text-shadow: 0 0 0;
49 }
```



```
49 .tab:checked:nth-of-type(3) ~ .tab__content:nth-of-type(3) {
50     opacity: 1;
51     transition: 0.5s opacity ease-in, 0.8s transform ease;
52     position: relative;
53     top: 0;
54     z-index: 100;
55     transform: translateY(0px);
56     text-shadow: 0 0 0;
57 }
58 .tab:checked:nth-of-type(4) ~ .tab__content:nth-of-type(4) {
59     opacity: 1;
60     transition: 0.5s opacity ease-in, 0.8s transform ease;
61     position: relative;
62     top: 0;
63     z-index: 100;
64     transform: translateY(0px);
65     text-shadow: 0 0 0;
66 }
67 .tab:checked:nth-of-type(5) ~ .tab__content:nth-of-type(5) {
68     opacity: 1;
69     transition: 0.5s opacity ease-in, 0.8s transform ease;
70     position: relative;
71     top: 0;
72     z-index: 100;
73     transform: translateY(0px);
74     text-shadow: 0 0 0;
75 }
76 .tab:checked:nth-of-type(6) ~ .tab__content:nth-of-type(6) {
77     opacity: 1;
78     transition: 0.5s opacity ease-in, 0.8s transform ease;
79     position: relative;
80     top: 0;
81     z-index: 100;
82     transform: translateY(0px);
83     text-shadow: 0 0 0;
84 }
85 .tab:checked:nth-of-type(7) ~ .tab__content:nth-of-type(7) {
86     opacity: 1;
87     transition: 0.5s opacity ease-in, 0.8s transform ease;
88     position: relative;
89     top: 0;
90     z-index: 100;
91     transform: translateY(0px);
92     text-shadow: 0 0 0;
93 }
94 .tab:first-of-type:not(:last-of-type) + label {
95     border-top-right-radius: 0;
96     border-bottom-right-radius: 0;
97 }
98 .tab:not(:first-of-type):not(:last-of-type) + label {
99     border-radius: 0;
100 }
101 .tab:last-of-type:not(:first-of-type) + label {
102     border-top-left-radius: 0;
103     border-bottom-left-radius: 0;
104 }
105 .tab:checked + label {
106     background-color: #fff;
```

```

107     box-shadow: 0 -1px 0 #fff inset;
108     cursor: default;
109 }
110 .tab:checked + label:hover {
111     box-shadow: 0 -1px 0 #fff inset;
112     background-color: #fff;
113 }
114 .tab + label {
115     box-shadow: 0 -1px 0 #eee inset;
116     border-radius: 6px 6px 0 0;
117     cursor: pointer;
118     display: block;
119     text-decoration: none;
120     color: #333;
121     flex-grow: 3;
122     text-align: center;
123     background-color: #f2f2f2;
124     user-select: none;
125     text-align: center;
126     transition: 0.3s background-color ease, 0.3s box-shadow ease;
127     height: 50px;
128     box-sizing: border-box;
129     padding: 15px;
130 }
131 .tab + label:hover {
132     background-color: #f9f9f9;
133     box-shadow: 0 1px 0 #f4f4f4 inset;
134 }
135

```

## 2.2.4 复选框

由于复选框的内容需要分词结果，故需要通过jinja2来返回渲染后的模板来实现动态加载复选框

html代码：

```

1  <div class="tab__content">
2      <div>
3          <h3>不涉及地点名词</h3>
4          {% for jieba1 in jieba1s %}
5              <label class="checkbox-inline">
6                  <input type="checkbox">{{ jieba1 }}
7              </label>
8          {% endfor %}
9      </div>
10
11     <div>
12         <h3>涉及地点名词</h3>
13         {% for jieba2 in jieba2s %}
14             <label class="checkbox-inline">
15                 <input type="checkbox">{{ jieba2 }}
16             </label>
17         {% endfor %}
18     </div>
19

```

```

20     <div>
21         <h3>案由</h3>
22         {% for jieba3 in jieba3s %}
23         <label class="checkbox-inline">
24             <input type="checkbox">{{ jieba3 }}
25         </label>
26         {% endfor %}
27     </div>
28
29     <div>
30         <h3>形容词</h3>
31         {% for jieba4 in jieba4s %}
32         <label class="checkbox-inline">
33             <input type="checkbox">{{ jieba4 }}
34         </label>
35         {% endfor %}
36     </div>
37 </div>
38

```

## 2.3 爬虫实现

由于裁判文书网爬取困难，最终我们选择爬取威科先行网站，通过分析发现对url='<https://law.wkinfo.com.cn/csi/search>'，发送携带检索参数的POST请求即可获取带有文书内容的Response

### 2.3.1 反爬

利用谷歌开发者工具分析发现，此网站反爬很低，只要发送一次请求即可获取文书信息。可以说，只要你买了网站的数据库就可以爬取很多文本。由于我们未购买数据库，只是试用，每天可访问的文书数量有限，同时为了反爬，我们随机使用小组成员账号来爬取文书。

Cookie:

该网站的Cookie有下面特点：

1. 对于同一用户，只有一项Cookie即Hm\_lpvt\_fecce484974a74c6d10f421b6d3bd395不稳定，是不断变化的，研究发现其值为当前时间戳，通过time.time()即可获得
2. 对于同一用户，其余Cookie相对稳定，长期不变，故直接登录组员的几个账号来获取Cookie，放入列表中，使用时通过random.choice()随机选取

User-Agent:

由于只需要一次请求，故一个User-Agent足以

Sleep:

由于只需要一次请求，故不需要sleep

### 2.3.2 请求

发送POST请求，重点在于headers的Cookie和表单数据的构建，Cookie上面已经介绍。表单数据即检索条件，我们只需要设置爬取时间和数量即可，实现如下

```

1  def spider(beginDate, endDate, num):
2      """
3      :param beginDate: "2020.01.20"
4      :param endDate: "2021.12.25"

```

```

5         :param num:
6         :return:
7         ""
8         dateList = [] #表单中的时间参数
9         date = "judgmentDate:['+beginDate+' TO '+endDate+']"
10        dateList.append(date)
11
12        #小组成员Cookie
13        Cookies = ['autologin=true; username=2034885420@qq.com; userConfig=
{"moduleList":
[], "userStaffType":0, "isIpUser":0, "parentId":0, "sourceSiteUrl":null, "sources
iteName":null, "clientSource": "自主注
册", "trial":true, "expire":false, "caseView":false, "advanced1SearchExpires": "2
", "proximateLatestExpires": "48", "proximateSearch":false}; check=valid;
TY_SESSION_ID=827ba7ea-6794-4eaf-abf1-8f66eb227b74; connect.sid=s:ZXBCNi-
URL3-e59bsiqzeRJNCckhEXEX.zuDa4Mc1AW4qwi4yAyCr88dNkXS19xsGOxBGpudIEYk;
userInfo=
{"id": "1000290499", "username": "2034885420@qq.com", "password": "0x02000000511e
c6b5b749cd46de1cf27e985d76cd0b5f608d2d7f8172158330542a1718a781d215aaee94b1de
4b5377df41109118", "userType": "normal", "email": "2034885420@qq.com", "userLang"
: "cn", "userPageSize": 25, "isSend": true, "sendLang": "cn", "recieveEmails":
[], "groupName": "law", "libraryCode": "law,taa,hr,HKBo1d", "licences": 1, "telepho
ne": "18355442634", "conf": ""}; loginin=true;
loginId=e6029baaa3b9484eb96ecf84287938f6;
acw_tc=78e2b62616426938197051079e9f34ff3064942aaf94eb0cf973020c11;
Hm_lpvtfecce484974a74c6d10f421b6d3bd395=',
14        'check=valid; TY_SESSION_ID=346b2eee-b5d7-4e77-8e10-
fd5498f42ede; autologin=true; username=632101734@qq.com;
acw_tc=78cebd9816427533357544001e2808a8cd92fc1f6388c9ef189f60038f;
connect.sid=s:Z3qGUVaiBnAsqnk40KZY5ELvtrB3GHCV.L4+emsCG+duBmRk5BwwAHOM2VyyE5
yjj+/7BXr5f1Hg; userInfo=
{"id": "1000290723", "username": "632101734@qq.com", "password": "0x0200000030923
68a749aa28c92e2afa813943ebb1d1d136f1baa5e9a3533b889a410991b300166b3365384146
50dec260dcf5280", "userType": "normal", "email": "632101734@qq.com", "userLang": "
cn", "userPageSize": 25, "isSend": true, "sendLang": "cn", "recieveEmails":
[], "groupName": "law", "libraryCode": "law,taa,hr,HKBo1d", "licences": 1, "telepho
ne": "18169538061", "conf": ""}; loginin=true;
loginId=298901ca248749449df3198c9dba1df4; userConfig={"moduleList":
[], "userStaffType":0, "isIpUser":0, "parentId":0, "sourceSiteUrl":null, "sources
iteName":null, "clientSource": "自主注
册", "trial":true, "expire":false, "caseView":false, "advanced1SearchExpires": "2
", "proximateLatestExpires": "48", "proximateSearch":false};
Hm_lvt_fecce484974a74c6d10f421b6d3bd395=1642745554,1642747324,1642753336,164
2753375; Hm_lpvtfecce484974a74c6d10f421b6d3bd395='
15        ]
16
17        cookie = random.choice(Cookies)
18
19        headers = {
20            'Accept': 'application/json, text/plain, */*',
21            'Accept-Encoding': 'gzip, deflate, br',
22            'Accept-Language': 'zh-CN,zh;q=0.9',
23            'Appversion': '2021.03.17.1',
24            'Cache-Control': 'no-cache',
25            'Connection': 'keep-alive',
26            'Content-Type': 'application/json;charset=utf-8',
27            'Cookie': (cookie+str(int(time.time()))).encode('utf-8'),
28            'Host': 'law.wkinfo.com.cn',

```

```

29         'Identification': '_2a66579079e311ecbc3ba7fea9a4723d',
30         'Referer': 'https://law.wkinfo.com.cn/judgment-documents/list?
mode=advanced',
31         'sec-ch-ua-mobile': '?0',
32         'Sec-Fetch-Dest': 'empty',
33         'Sec-Fetch-Mode': 'cors',
34         'Sec-Fetch-Site': 'same-origin',
35         'User-Agent': 'Mozilla/5.0 (Windows NT 10.0; Win64; x64)
AppleWebKit/537.36 (KHTML, like Gecko) Chrome/97.0.4692.71 Safari/537.36
Edg/97.0.1072.62',
36         'X-Tingyun-Id': 'tn6win9ZeY4;r=96286256',
37     }
38
39     params = {
40         "indexId": "law.case",
41         "query": {
42             "queryString": "typeOfCase:刑事 AND typeOfDecision:((001))",
43             "filterQueries": [],
44             "filterDates": dateList
45         },
46         "searchScope": {"treeNodeIds": []},
47         "relatedIndexQueries": [],
48         "sortOrderList": [{"sortKey": "judgmentDate", "sortDirection":
"DESC"}],
49         "pageInfo": {"limit": int(num), "offset": 0},
50         "otherOptions": {
51             "requireLanguage": "cn",
52             "relatedIndexEnabled": False,
53             "groupEnabled": False,
54             "smartEnabled": True,
55             "buy": False,
56             "summaryLengthLimit": 100,
57             "advanced": True,
58             "synonymEnabled": True,
59             "isHideBigLib": 0,
60             "relatedIndexFetchRows": 5,
61             "proximateCourtID": "",
62             "module": ""
63         },
64         "chargingInfo": {"useBalance": True}}
65
66     url = 'https://law.wkinfo.com.cn/csi/search'
67     result = requests.session().post(url, headers=headers, json=params)
68
69     if result.status_code == 200:
70         content = result.content
71         return htmlParser(content) #返回标题列表，用于前端展示爬虫结果
72     else:
73         print(result.content.decode('utf-8'))
74         return []

```

### 2.3.3 解析与下载

对请求的结果通过json.loads()转化为字典，获取文书列表documentList，遍历文书列表，对每个document进行正则分析等来提取文书内容，并保存至static/Downloads目录下

```

1 def htmlParser(content):

```

```

2     titleList = []
3     contentDict = json.loads(content,strict=False)
4     documentList = contentDict['documentList']
5     for i in range(len(documentList)):
6         titleList.append(wenshuDownload(documentList[i]))
7     return titleList #返回标题列表
8
9 def wenshuDownload(document):
10     title = document['title'] #文书标题
11     additionalFields = document['additionalFields'] #文书内容
12     res = title + '\r' #待保存的结果
13     for key,value in additionalFields.items(): #正则分析
14         if key == "judgmentDate":
15             continue
16         if key == "instancecode":
17             continue
18         if key == "referenceLevel":
19             continue
20         if key == "updateContentStatus":
21             continue
22         if key == "product":
23             continue
24         if value == None:
25             continue
26         if value=='':
27             continue
28         r = re.compile(r'\s+')
29         value = r.sub("\r", value) #正则表达式替换
30         res+=value
31     path = r'static/Downloads/'+title+'.txt' #保存路径
32     with open(path,'w',encoding='utf-8') as file:
33         file.write(res)
34     return title #返回标题

```

## 2.4 分词实现

通过分析多份法律文书，最终决定分为4类词性：{"不涉及地点名词": nr, "涉及地点名词": ns, "案由": r, "形容词": adj}，同时添加自定义词典以增强分词的通用性。

### 2.4.1 方法分析

#### 2.4.1.1 out方法

获取前端界面得到文书的内容，然后仅需调用该方法即可得到分词结果。

传出处理后的分词结果，形式为字典，其中value值为一个列表，其中的元素为某一类型的词语的集合，每个value值对应的key值为该类型词语的名称

```

1     """传出结果，形式为字典，key为value中词语类型，value为列表"""
2     def out(content):
3         jieba.load_userdict("static/dict.txt")
4         # 补充由于结构较为复杂难以分出的案由
5         content = content.replace(" ", "").replace("\n", "")
6         words = pseg.lcut(content, use_paddle=True)
7         wordsList = []
8         # 初步分词结果
9         flagList = []

```

```

10     # 初步分词结果对应的词性
11     length = 0
12     for word, flag in words:
13         length += 1
14         wordsList.append(str(word))
15         flagList.append(str(flag))
16     nr = get_about_nr(wordsList, flagList, length)
17     ns = get_about_ns(wordsList, flagList, length)
18     # v = get_v(wordsList, flagList, length)
19     r = get_r(wordsList, flagList, length)
20     adj = get_adj(wordsList, flagList, length)
21     result = {"不涉及地点名词": nr, "涉及地点名词": ns, "案由": r, "形容词": adj}
22     return result

```

#### 2.4.1.2 get\_about\_nr方法

得到不包含地点的名词（考虑到名词数量过多）如图所示

```

1  #在引入自定义词典后，wordsList为分词的初始结果，flagList为分词得到的词性的初始结果，
   length为列表长度
2  def get_about_nr(wordsList, flagList, length):

```

后续方法与该方法传入参数相同

#### 2.4.1.3 get\_about\_ns方法

得到包含地点的名词

```

1  def get_about_ns(wordsList, flagList, length):

```

#### 2.4.1.4 get\_r方法

得到案由

```

1  def get_r(wordsList, flagList, length):

```

#### 2.4.1.5 get\_adj方法

得到形词

```

1  def get_adj(wordsList, flagList, length):

```

#### 2.4.1.6 getStrFromTxt & showRes方法

由于本地测试时，传入为txt格式的文本文件，该方法将从txt文件中获得文本，返回形式为字符串，主要目的是为处理由于编码不同产生的部分问题，该方法仅供本地测试使用  
将本地分词结果进行输出，便于本地测试

### 2.4.2 分词过程

#### 2.1 直接使用jieba分词结果的部分

人名部分

#### 2.2 结合jieba分词结果及对部分字符串的解析部分

文化程度部分；民族部分；性别部分

## 2.3 结合jieba分词时构建的自定义词典和对字符串的解析

地名；相关法院名称；案由

地名采用向后匹配策略，由于部分地名无法识别（例如西藏自治区那曲市），故添加至自定义词典

相关法院采用添加前继位名词后匹配地点名词的方式

案由构成的分析，首先进行字符串解析，由于形式较为复杂，且不存在普适性，故全部添加至自定义词典并保存为re词性（dict.txt文件）

## 2.5 前后连接

利用flask框架实现前后段连接

### 2.5.1 初始

项目启动后，直接在浏览器输入<http://127.0.0.1:5000/>来发出GET请求即可进入初始界面，路由实现如下

```
1 @app.route('/', methods = ['GET', 'POST'])
2 def jiebas():
3     if request.method == 'GET':
4         return render_template('司法大数据自动化标注与分析.html', text='')
```

### 2.5.2 爬虫

通过初始页面的爬虫按钮来提交爬取日期和数量，发出POST请求，app.py调用wenshu\_spider.py来爬取文书，并返回Jinja2渲染文书标题后的爬虫展示页面

```
1 @app.route('/文书爬取', methods = ['GET', 'POST'])
2 def crawl():
3     if request.method == 'GET':
4         return render_template('司法大数据自动化标注与分析.html', text='') # 返回模板
5     else:
6         #获取参数值
7         beginDate = request.form.get('beginDate')
8         endDate = request.form.get('endDate')
9         num = request.form.get('num')
10        #设置参数默认值
11        if beginDate == "":
12            beginDate = "2002.12.25"
13        if endDate == "":
14            endDate = "2021.12.25"
15        if num == "":
16            num = "10"
17        titleList = wenshu_spider.spider(beginDate, endDate, num)
18        return render_template('文书爬取.html', titles=titleList) #Jinja2渲染
```

### 2.5.3 分词

通过初始页面的分词按钮来提交案件文本，发出POST请求，app.py调用jiebaAnalysis.py来进行分词，并返回Jinja2渲染分词结果后的页面



```

1  @app.route('/', methods = ['GET', 'POST'])
2  def jiebas():
3      if request.method == 'GET':
4          return render_template('司法大数据自动化标注与分析.html', text='') # 返回模板
5      else:
6          text = request.form.get('textArea') #获取案件文本
7          result = jiebaAnalysis.out(text)
8          jieba1s = result.get('不涉及地点名词')
9          jieba2s = result.get('涉及地点名词')
10         jieba3s = result.get('案由')
11         jieba4s = result.get('形容词')
12         return render_template('司法大数据自动化标注与分析.html', text=text,
13                                jieba1s=jieba1s, jieba2s=jieba2s, jieba3s=jieba3s, jieba4s=jieba4s) #Jinja2渲染

```

```

1  def spider(beginDate, endDate, num):
2      """
3      :param beginDate: "2020.01.20"
4      :param endDate: "2021.12.25"
5      :param num:
6      :return:
7      """
8      dateList = [] #表单中的时间参数
9      date = "judgmentDate:[" + beginDate + ' TO ' + endDate + ']'
10     dateList.append(date)
11
12     #小组成员Cookie
13     Cookies = ['autologin=true; username=2034885420@qq.com; userConfig=
{"moduleList":
[], "userStaffType":0, "isIpUser":0, "parentId":0, "sourceSiteUrl":null, "sourceSiteName":null, "clientSource": "自主注册", "trial":true, "expire":false, "caseView":false, "advanced1SearchExpires": "2", "proximateLatestExpires": "48", "proximateSearch":false}; check=valid; TY_SESSION_ID=827ba7ea-6794-4eaf-abf1-8f66eb227b74; connect.sid=s:ZXBCNi-URL3-e59bsiqzeRJNCckhEXEX.zuDa4Mc1AW4qwi4yAyCr88dNkXS19xsGOxBGpudIEYk; userInfo=
{"id": "1000290499", "username": "2034885420@qq.com", "password": "0x02000000511ec6b5b749cd46de1cf27e985d76cd0b5f608d2d7f8172158330542a1718a781d215aaee94b1de4b5377df41109118", "userType": "normal", "email": "2034885420@qq.com", "userLang": "cn", "userPageSize": 25, "isSend": true, "sendLang": "cn", "recieveEmails": [], "groupName": "law", "libraryCode": "law,taa,hr,HKBo1d", "licences": 1, "telephone": "18355442634", "conf": ""}; loginin=true; loginId=e6029baaa3b9484eb96ecf84287938f6; acw_tc=78e2b62616426938197051079e9f34ff3064942aaf94eb0cf973020c11; Hm_lpvtfecce484974a74c6d10f421b6d3bd395=',

```

```

14         'check=valid; TY_SESSION_ID=346b2eee-b5d7-4e77-8e10-
fd5498f42ede; autoLogin=true; username=632101734@qq.com;
acw_tc=78cebd9816427533357544001e2808a8cd92fc1f6388c9ef189f60038f;
connect.sid=s:Z3qGUVaiBnAsqnk40KZY5ELvtrB3GHCv.L4+emsCG+duBmRk5BwwAH0M2VyyE5
yjj+/7BXR5f1Hg; userInfo=
{"id":"1000290723","username":"632101734@qq.com","password":"0x0200000030923
68a749aa28c92e2afa813943ebb1d1d136f1baa5e9a3533b889a410991b300166b3365384146
50dec260dcf5280","userType":"normal","email":"632101734@qq.com","userLang":"
cn","userPageSize":25,"isSend":true,"sendLang":"cn","receiveEmails":
[],"groupName":"law","libraryCode":"law,taa,hr,HKBoId","licences":1,"telepho
ne":"18169538061","conf":""}; loginin=true;
loginId=298901ca248749449df3198c9dba1df4; userConfig={"moduleList":
[],"userStaffType":0,"isIpUser":0,"parentId":0,"sourceSiteUrl":null,"sources
iteName":null,"clientSource":"自主注
册","trial":true,"expire":false,"caseview":false,"advanced1SearchExpires":"2
","proximateLatestExpires":"48","proximateSearch":false};
Hm_lvt_fecce484974a74c6d10f421b6d3bd395=1642745554,1642747324,1642753336,164
2753375; Hm_lpv_t_fecce484974a74c6d10f421b6d3bd395=
15     ]
16
17     cookie = random.choice(Cookies)
18
19     headers = {
20         'Accept': 'application/json, text/plain, */*',
21         'Accept-Encoding': 'gzip, deflate, br',
22         'Accept-Language': 'zh-CN,zh;q=0.9',
23         'Appversion': '2021.03.17.1',
24         'Cache-Control': 'no-cache',
25         'Connection': 'keep-alive',
26         'Content-Type': 'application/json;charset=utf-8',
27         'Cookie': (cookie+str(int(time.time()))).encode('utf-8'),
28         'Host': 'law.wkinfo.com.cn',
29         'Identification': '_2a66579079e311ecbc3ba7fea9a4723d',
30         'Referer': 'https://law.wkinfo.com.cn/judgment-documents/list?
mode=advanced',
31         'sec-ch-ua-mobile': '?0',
32         'Sec-Fetch-Dest': 'empty',
33         'Sec-Fetch-Mode': 'cors',
34         'Sec-Fetch-Site': 'same-origin',
35         'User-Agent': 'Mozilla/5.0 (Windows NT 10.0; Win64; x64)
AppleWebKit/537.36 (KHTML, like Gecko) Chrome/97.0.4692.71 Safari/537.36
Edg/97.0.1072.62',
36         'X-Tingyun-Id': 'tN6win9ZeY4;r=96286256',
37     }
38
39     params = {
40         "indexId": "law.case",
41         "query": {
42             "queryString": "typeOfCase:刑事 AND typeOfDecision:((001))",
43             "filterQueries": [],
44             "filterDates": dateList
45         },
46         "searchScope": {"treeNodeIds": []},
47         "relatedIndexQueries": [],
48         "sortOrderList": [{"sortKey": "judgmentDate", "sortDirection":
"DESC"}],
49         "pageInfo": {"limit": int(num), "offset": 0},
50         "otherOptions": {

```

```

51         "requireLanguage": "cn",
52         "relatedIndexEnabled": False,
53         "groupEnabled": False,
54         "smartEnabled": True,
55         "buy": False,
56         "summaryLengthLimit": 100,
57         "advanced": True,
58         "synonymEnabled": True,
59         "isHideBigLib": 0,
60         "relatedIndexFetchRows": 5,
61         "proximateCourtID": "",
62         "module": ""
63     },
64     "chargingInfo": {"useBalance": True}}
65
66 url = 'https://law.wkinfo.com.cn/csi/search'
67 result = requests.session().post(url, headers=headers, json=params)
68
69 if result.status_code == 200:
70     content = result.content
71     return htmlParser(content) #返回标题列表，用于前端展示爬虫结果
72 else:
73     print(result.content.decode('utf-8'))
74     return []

```

### 3. 不足

1. 由于我们未购买威科先行数据库，故爬取数量有限，每天只允许爬取几百份案例
2. 对于单文本多犯人的案例，最好多次分词标注，一次性标注的话json文件较乱
3. 虽然已经扩展了分词词典，但可能仍存在难以识别的地名、人名

### 4. 分工

组长：

于欣博 201250165

前后端结合

组员：

马子睿 201250163

分词

刘云辉 201250166

爬虫、前端、前后端结合