

# 今天来聊聊求职需要的 Python 技能

原创 某某白米饭 Python技术 1周前

每年的 3、4 月份都是金三银四跳槽季，企业一般也会选择在这个时期调整职工的薪资，小伙伴在这个时候也会心里痒痒，在招聘网站上看看是否有合适的机会，需要的 Python 技能是否符合年限等等情况。这里以招聘网站为例抓取魔都近一个月的招聘数据，生成柱状图与词云。

## 抓取招聘网站数据

首先将魔都近 1 个月的招聘职位都抓取出来，使用 requests 模块和 BeautifulSoup 模块

```
1  # -*- coding: utf-8 -*-
2  import requests
3  from bs4 import BeautifulSoup
4  import time
5  import random
6
7  urlFileName = 'urls.txt' # 存放招聘信息详情的URL文本
8  contentFileName = 'context.txt' # 存放抓取的内容
9  def getUrls2Txt(page_num):
10     p = page_num+1
11     for i in range(1, p):
12         urls = []
13         # 抓取魔都的
14         url = 'https://search.51job.com/list/020000,000000,0000,00,2,99,Python,2,'+str(i)+'.html?lan
15
16         html = requests.get(url)
17         soup = BeautifulSoup(html.content, "html.parser")
18         ps = soup.find_all('p', class_='t1')
19         for p in ps:
20             a = p.find('a')
21             urls.append(str(a['href']))
22         with open(urlFileName, 'a', encoding='utf-8') as f:
23             for url in urls:
24                 f.write(url+'\n')
25             s = random.randint(5, 30)
26             print(str(i)+'page done,'+str(s)+'s later')
27             time.sleep(s)
28
29  def getContent(url, headers):
30     record = ''
31     try:
32         html = requests.get(url, headers=headers)
```

```

33     soup = BeautifulSoup(ntml.content, 'html.parser')
34     positionTitle = str(soup.find('h1')['title']) # 标题
35     salary = soup.find_all('strong')[1].get_text() # 薪资
36     companyName = soup.find('p', class_='cname').get_text().strip().replace('\n', '').replace('查
37     positionInfo = soup.find(
38         'div', class_='bmsg job_msg inbox').get_text().strip().replace('\n', '').replace('分享',
39     record = positionTitle + '***' + salary + '***' + companyName + '***' + '***' + positionInfo
except Exception as e:
40     print('错误了')
41
42     return record
43
44
45 def main():
46     page_num = 93
47     getUrls2Txt(page_num)
48     user_Agent = 'Mozilla/5.0 (Macintosh; Intel Mac OS X 10_15_3) AppleWebKit/537.36 (KHTML, lik
49     headers = {'User-Agent': user_Agent}
50     with open(urlFileName, 'r', encoding='utf-8') as f:
51         urls = f.readlines()
52     i = 0
53     for url in urls:
54         url = url.strip()
55         if url != '':
56             record = getContent(url, headers)
57             with open(contentFileName, 'a', encoding='utf-8') as f:
58                 f.write(record + '\n')
59             i += 1
60             print(str(i)+'详情抓取完成')
61             time.sleep(1)
62
63     print('完成了')
64
65
66 if __name__ == '__main__':
67     main()

```

抓取网站内容结果图

## 分词

在这一步需要对招聘信息中的职位信息进行人工的初步删选，过滤掉常用字存入 `filterWords` 变量中，然后利用结巴分词(<https://github.com/fxsjy/jieba>)基于TF-IDF算法将职位信息进行分词，并统计技术词语出现的次数。

```

1  from jieba import analyse
2
3  fenCi = {}
4

```

```

5 def main():
6
7     # 负责过滤的词语，这里只列出了几个
8     filterWords = ['熟悉', '熟练', '经验', '优先', '应用开发', '相关', '工作', '开发', '能力', '负责', '技
9
10    # 结巴分词基于 TF-IDF 算法的关键词
11    tfidf = analyse.extract_tags
12
13    for zpInfo in open('context.txt', 'r', encoding='utf-8'):
14
15        if zpInfo.strip() == '':
16            continue
17        # 详情数据是用&&&分割的
18        infos = zpInfo.split("&&&")
19        words = tfidf(infos[-1])
20
21        words = [x.upper() for x in words if x.upper() not in filterWords]
22
23        for word in words:
24            num = fenCi.get(word, 0) + 1
25            fenCi[word] = num
26
27    print(sorted(fenCi.items(), key=lambda kv: (kv[1], kv[0]), reverse=True))
28    print('分出了' + str(len(fenCi)) + '了词语')
29
30
31 if __name__ == '__main__':
32     main()

```

分词结果图

## 技能图表

在分词中，分出了 12663 个词这些词大多都是常用字，需要进一步筛选出多个高频的 Python 技能利用 matplotlib 模块画出柱状图。

```

1 import matplotlib.pyplot as plt
2 import numpy as np
3
4 plt.rcParams['font.sans-serif'] = ['Arial Unicode MS']
5 params = {
6     'axes.labelsize': '14',
7     'xtick.labelsize': '14',
8     'ytick.labelsize': '13',
9     'lines.linewidth': '2',
10    'legend.fontsize': '20',

```

```
11     'figure.figsize': '26, 24'
12 }
13 plt.style.use("ggplot")
14 plt.rcParams.update(params)
15
16 # 筛选分词中高频的
17 barDir = {
18     'PYTHON': 2283,
19     'LINUX': 981,
20     '算法': 658,
21     '运维': 530,
22     '数据库(MySql,Sql,Redis等)': 1021,
23     'SHELL': 996,
24     '数据分析/挖掘': 695,
25     'WEB': 454,
26     '测试用例': 515,
27     'MATLAB': 221,
28     'PERL': 209,
29     'HIVE': 122,
30     'HADOOP': 176,
31     'SPARK': 146,
32     'TENSORFLOW': 136,
33     '多线程': 127,
34     'AI': 106,
35     'SAS': 104,
36     '视觉/图像处理': 180,
37     '人工智能': 170,
38     'HTTP': 90,
39     'DOCKER': 82,
40     'DJANGO': 82,
41 }
42
43 fig, ax = plt.subplots(figsize=(20, 10), dpi=100)
44
45 # 添加刻度标签
46 labels = np.array(list(barDir.keys()))
47 ax.barh(range(len(barDir.values())), barDir.values(), tick_label=labels, alpha=1)
48
49 ax.set_xlabel('Python技术词的次数', color='k')
50 ax.set_title('Python工作高频技术词')
51
52
53 # 为每个条形图添加数值标签
54 for x, y in enumerate(barDir.values()):
55     ax.text(y + 0.5, x, y, va='center', fontsize=14)
56
57 # 显示图形
58 plt.show()
59
```

## 词云

最后将分词数据生成一个词云，将 Python 图标作为底图使用。

```
1 def getWorldCloud():
2     # 底层图片路径
3     path_img = "python.jpg"
4     background_image = np.array(Image.open(path_img))
5
6     wordcloud = WordCloud(
7         # 字体路径
8         font_path="/System/Library/Fonts/STHeiti Light.ttc",
9         background_color="white",
10        mask=background_image).generate(" ".join(list(fenCi.keys()))))
11    image_colors = ImageColorGenerator(background_image)
12    plt.imshow(wordcloud.recolor(color_func=image_colors), interpolation="bilinear")
13    plt.axis("off")
14    plt.show()
```

最后生成的词云图

## 总结

本文主要是从招聘网站抓取 Python 工作职责并生成柱状图和词云，展示企业需要哪些 Python 技能，从而在面试前学会并运用这些技能。在生成最后结果的过程中存在 2 点不完美的情况，一点是存在人工筛选另一个是在分词中没有完全过滤掉通用字。随着小编的 Python 技能树的增长，有理由相信在不久这 2 种情况将完全避免。

示例代码：求职需要的 Python 技能

**PS：** 公号内回复「Python」即可进入 Python 新手学习交流群，一起 100天计划！

-END-

**Python 技术**  
**关于 Python 都在这里**