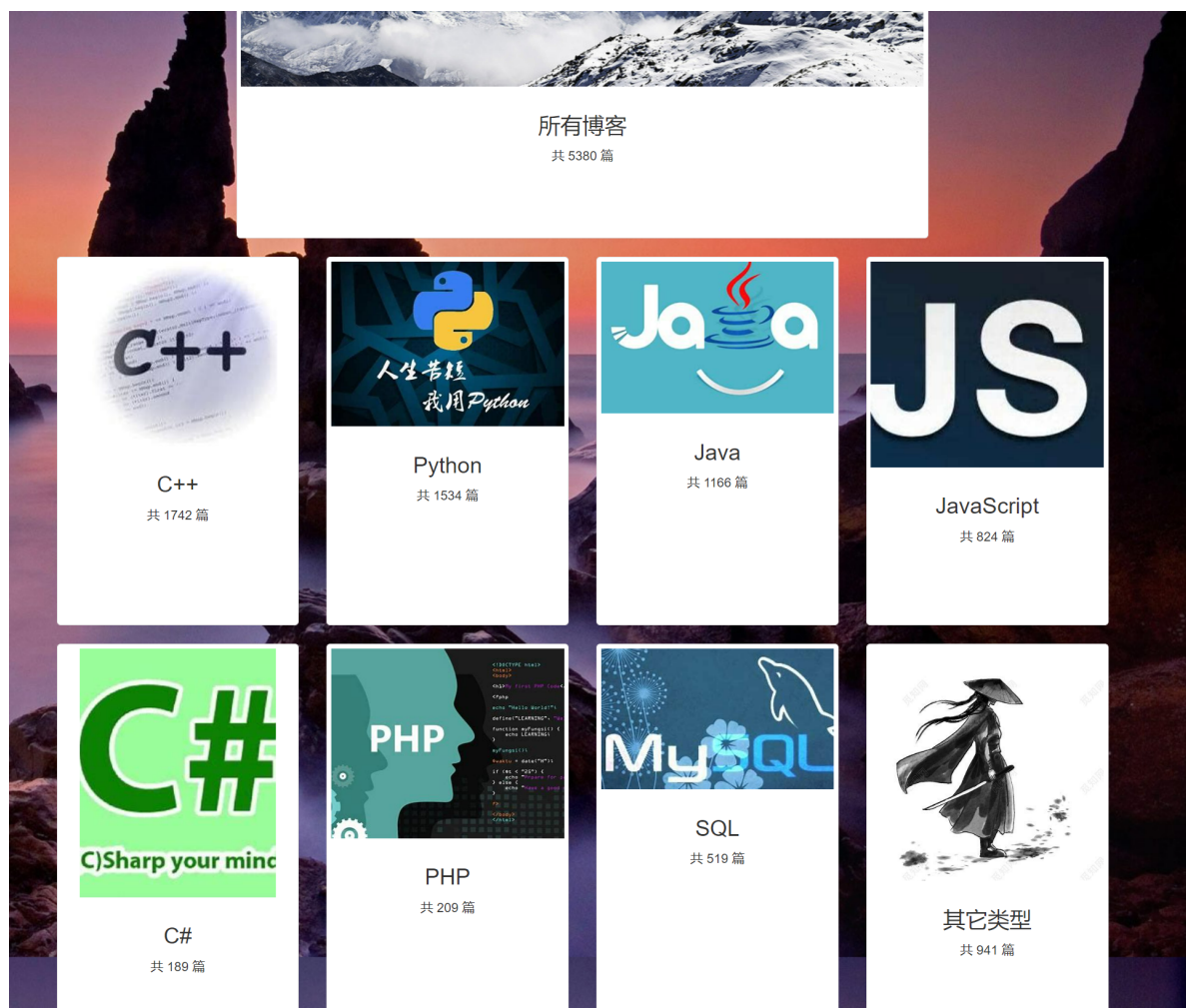


# CSDN爬取数据分析

## 样本情况

- 从CSDN网站的语言分类平台上爬取共5385条博客，博客内容包括c++，c#等7种编程语言，其它数据归类到其它类型中。
- 主要爬取的数据包括：文章题目、正文、发表日期时间、点赞收藏量、作者头像、点赞量等



## 分析结论

### 话题一：主流语言

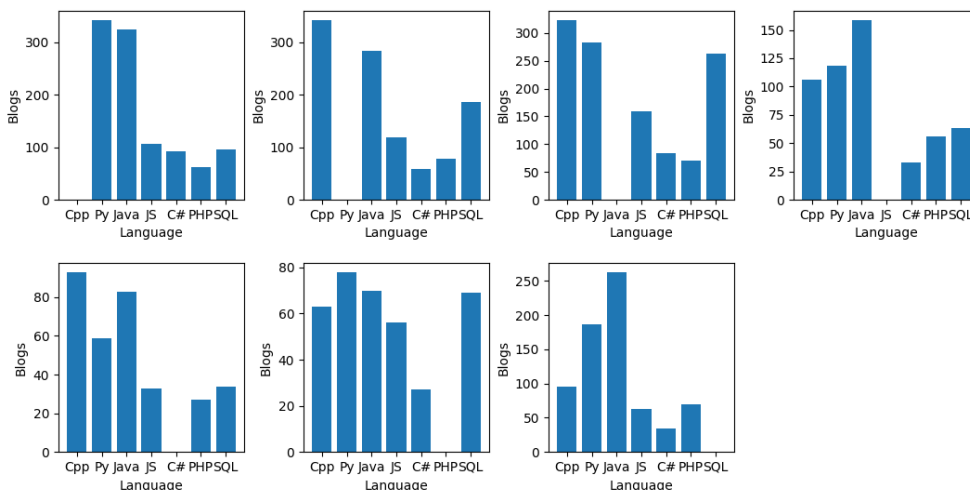
- 下面这段代码统计了在带有所有带有某种语言标签的数据种带有其它语言标签的数量，并据此绘制了柱状统计图。以图一为例，在所有带有 Cpp 标签的数据中，带有 Python 和 Java 标签的博客都达到了300条以上，而其它语言标签数量均在100条以下

```
1 import json
2 from scipy.integrate import quad
3 import numpy as np
4 import matplotlib.pyplot as plt
5 import dexplore as dxp
6 import pandas as pd
7 with open('all_data.json','r',encoding='utf-8') as JsonFile:
8     data = json.load(JsonFile)
```

```

9 all_dict = {'Cpp': {'Cpp': 0, 'Python': 0, 'Java': 0, 'JavaScript': 0,
  'Csharp': 0, 'PHP': 0, 'SQL': 0}, 'Python': {'Cpp': 0, 'Python': 0, 'Java':
  0, 'JavaScript': 0, 'Csharp': 0, 'PHP': 0, 'SQL': 0}, 'Java': {'Cpp': 0,
  'Python': 0, 'Java': 0, 'JavaScript': 0, 'Csharp': 0, 'PHP': 0, 'SQL': 0},
  'JavaScript': {'Cpp': 0, 'Python': 0, 'Java': 0, 'JavaScript': 0, 'Csharp':
  0, 'PHP': 0, 'SQL': 0}, 'Csharp': {'Cpp': 0, 'Python': 0, 'Java': 0,
  'JavaScript': 0, 'Csharp': 0, 'PHP': 0, 'SQL': 0}, 'PHP': {'Cpp': 0,
  'Python': 0, 'Java': 0, 'JavaScript': 0, 'Csharp': 0, 'PHP': 0, 'SQL': 0},
  'SQL': {'Cpp': 0, 'Python': 0, 'Java': 0, 'JavaScript': 0, 'Csharp': 0,
  'PHP': 0, 'SQL': 0}}
10 for each in data:
11     for lan in each['language']:
12         if lan == 'Others':
13             continue
14         for lan_to in each['language']:
15             if lan_to == 'Others' or lan == lan_to:
16                 continue
17             all_dict[lan][lan_to] += 1
18 num=1
19 for each in all_dict:
20     plt.subplot(2, 4, num)
21     num+=1
22     lst = list(all_dict[each].values())
23     plt.bar(range(7), lst)
24     plt.xticks(range(7), ['Cpp', 'Py', 'Java', 'JS', 'C#', 'PHP', 'SQL'])
25     plt.xlabel('Language')
26     plt.ylabel('Blogs')
27 plt.show()

```



据此可以得出下面的结论：

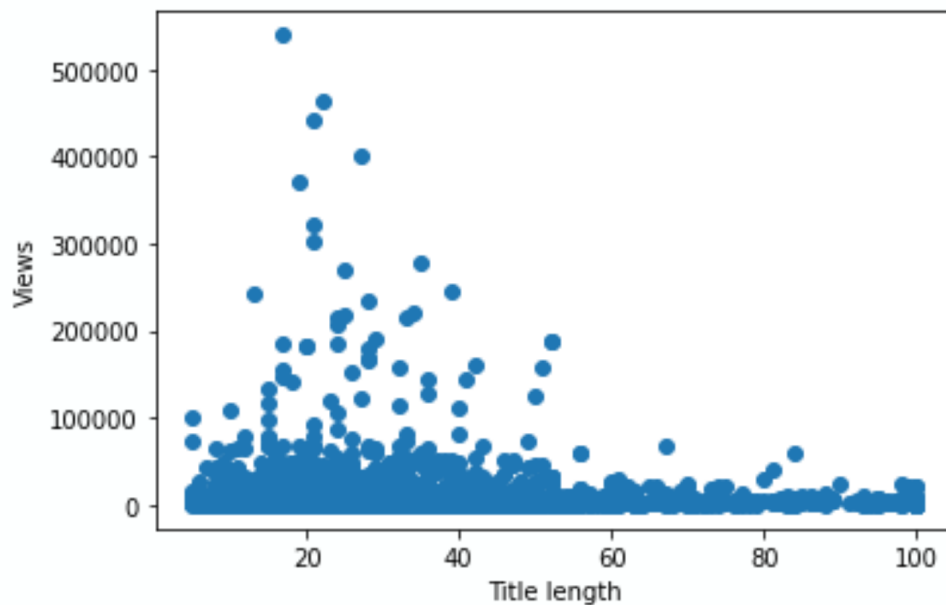
- Cpp、Java、JavaScript 在七种语言中热度较高，覆盖面较广，有两点数据支撑：
  1. 从标签总数来看，带有这三种语言标签的数量均在1000条以上，远大于其它语言
  2. 从语言标签的包含关系来看，无论是在哪种语言标签下，这三种语言均占据了数量前三的位置
- Java 与 JavaScript 和 SQL 常常被放到一起讨论，说明这三种语言的联系相对更紧密

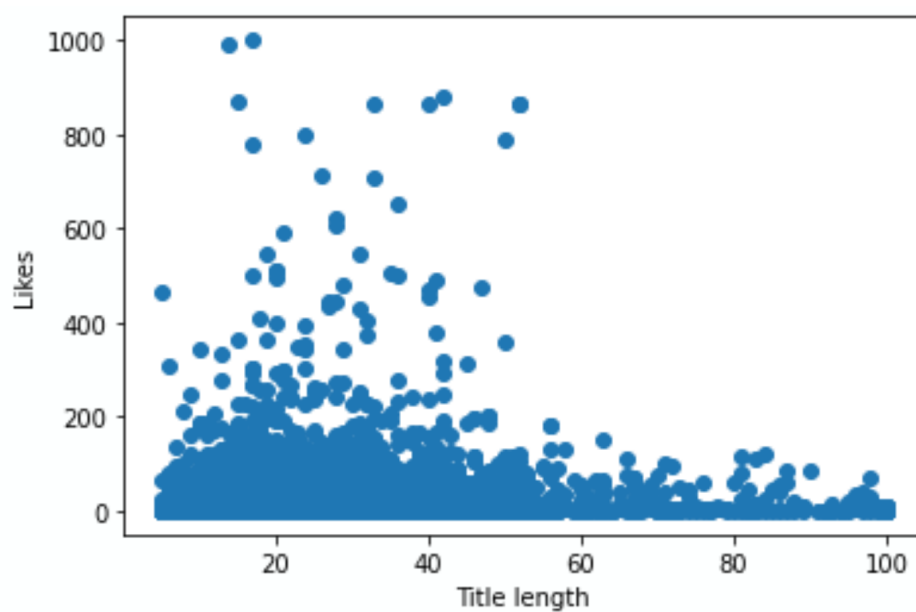
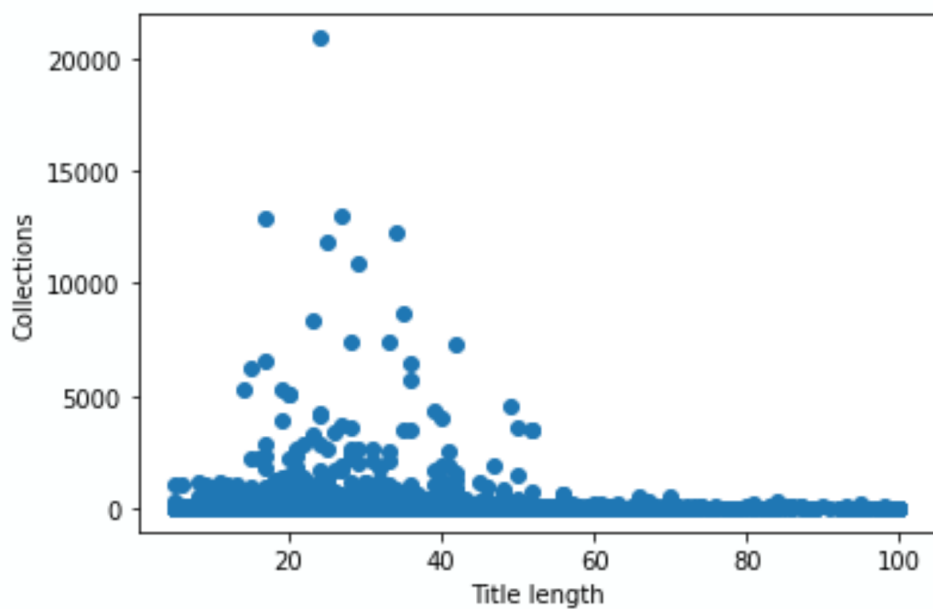
## 话题二：文章受欢迎程度的影响因素

根据收集的数据，文章受欢迎程度主要可以由文章的点赞量、收藏量和阅读量来衡量

下段代码绘制了点赞量、收藏量、阅读量和标题长度的关系

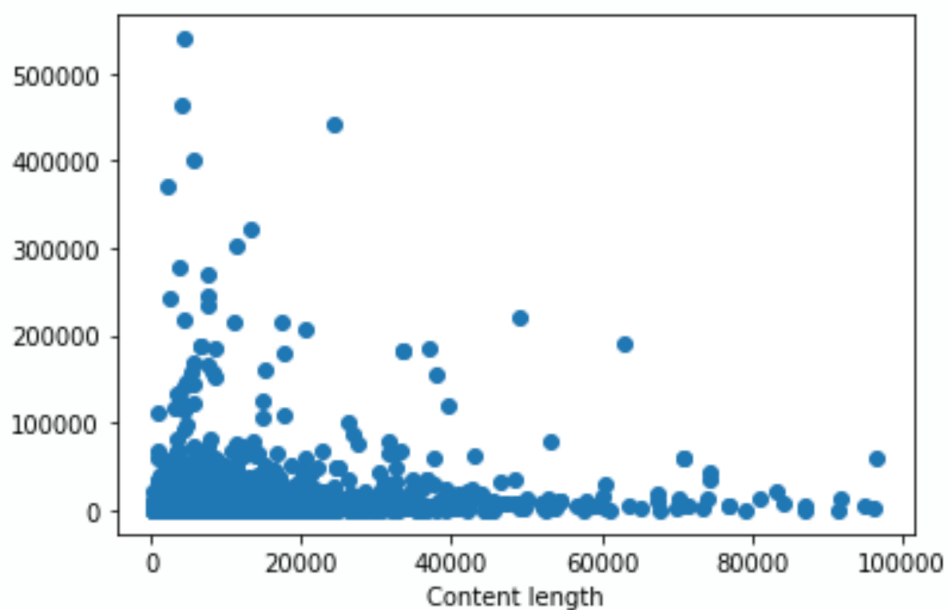
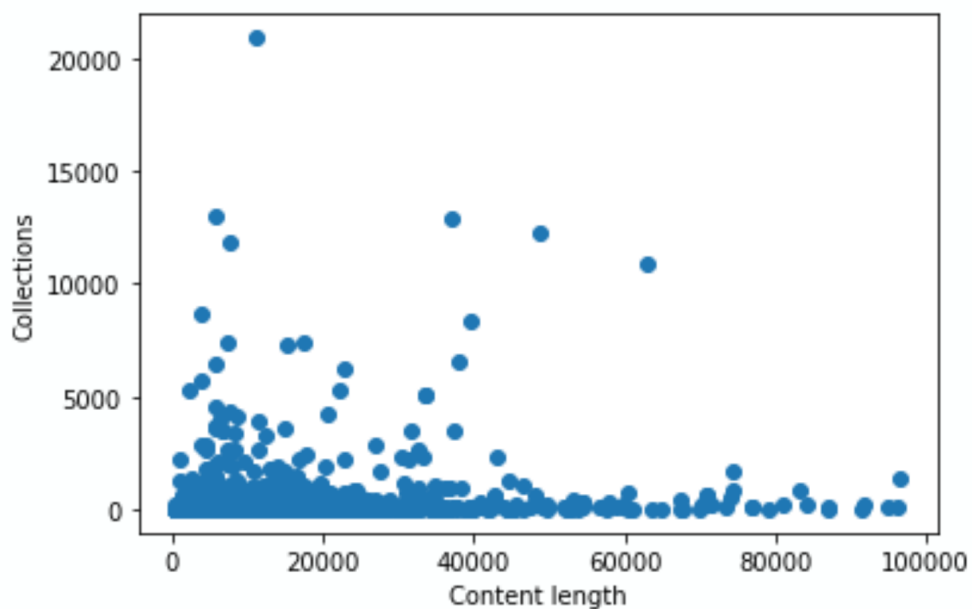
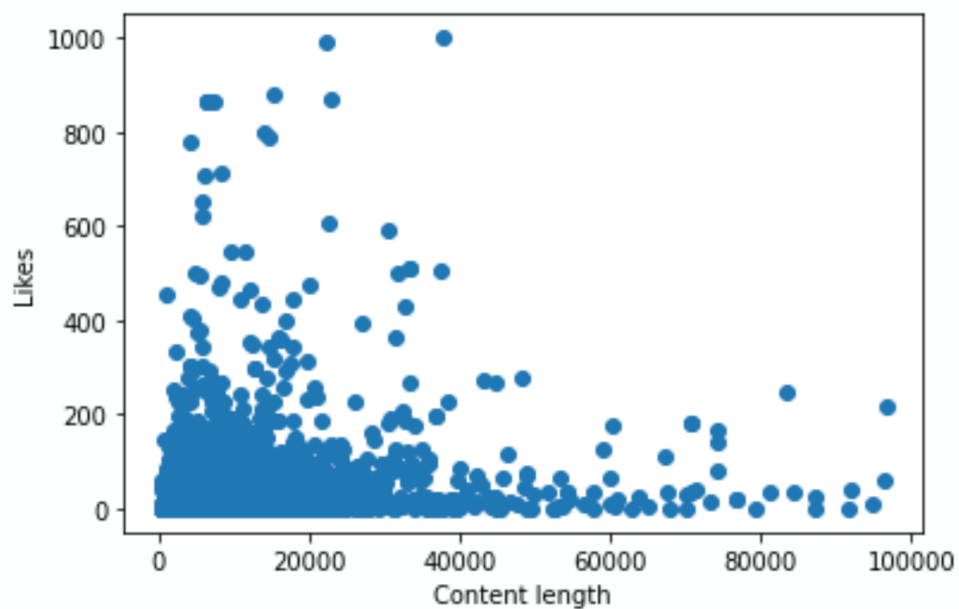
```
1 import matplotlib.pyplot as plt
2 import json
3 import random
4 with open('all_data.json','r',encoding='utf-8') as JsonFile:
5     data = json.load(JsonFile)
6 x=[]
7 y=[]
8 rec_index = []
9 num = 0
10 for each in data:
11     x.append(len(each['Title']))
12     y.append(int(each['Views'])) #还可以是Likes和Collections
13
14 plt.scatter(x,y)
15 plt.show()
```





分析数据发现，标题长度在20到40个字符之间的文章更受欢迎，无论是在点赞量、收藏量还是阅读量上都明显更高，文章点赞量相比其它数据受标题长度影响更小一些。

采用相同的方法、进一步统计了点赞量、收藏量以及阅读量和正文长度的关系：



发现受欢迎的文章字符数大多在40000以内。

据此可以得出下面的结论：

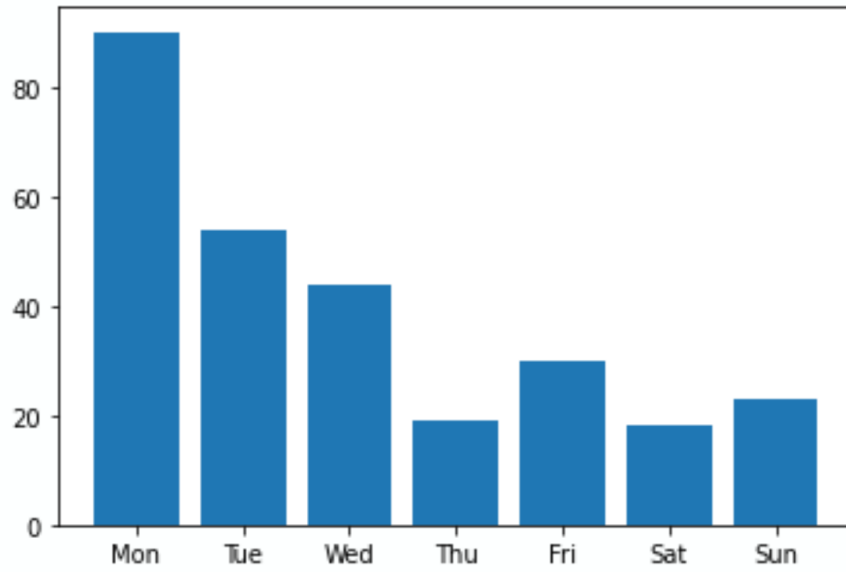
- 受欢迎的文章、标题和正文长度均适中（20-40字、40000以内）。总体来看，字数较少的文章更容易受到欢迎，点赞量相比收藏和观看受文章与标题长度影响比较小

## 话题三：基于作者发布习惯的分析

通过将样本按照作者归类、可以观察到不同作者发表文章时间的特点，下面代码统计了每个作者发表文章数量最多的一天。

```
1 import matplotlib.pyplot as plt
2 import json
3 import random
4 import re
5 import datetime
6 with open('all_data.json','r',encoding='utf-8') as JsonFile:
7     data = json.load(JsonFile)
8     author_dict={}
9     for each in data:
10         if each['Author'] in author_dict:
11             author_dict[each['Author']] += 1
12         else:
13             author_dict.update({each['Author']:1})
14     author_list=sorted(author_dict.items(), key = lambda d:d[1],reverse=True)
15     hot_author_list=[]
16     for each in author_list:
17         if each[1]>=3:
18             hot_author_list.append(each)
19     hot_author_dict = {}
20     for aut in hot_author_list:
21         aut_time_list=[]
22         for each in data:
23             if each['Author']==aut[0]:
24                 aut_time_list.append(each['Date'])
25             aut_time_list.sort(reverse=True)
26             hot_author_dict.update({aut[0]:aut_time_list})
27     author_week_dict={}
28     for aut in hot_author_dict:
29         week_dict = [0,0,0,0,0,0,0]
30         for each in hot_author_dict[aut]:
31             day = datetime.date(int(re.match(r'(\d+)-(\d+)-(\d+)',
each).group(1)),int(re.match(r'(\d+)-(\d+)-(\d+)',
each).group(2)),int(re.match(r'(\d+)-(\d+)-(\d+)',
each).group(3))).weekday()
32             week_dict[day] += 1
33             author_week_dict.update({aut:week_dict})
34     author_day = [0,0,0,0,0,0,0]
35     for each in author_week_dict:
36         max = 0
37         da = 0
38         for x in author_week_dict[each]:
39             if x>max:
40                 max=x
41                 da = author_week_dict[each].index(x)
42     author_day[da] += 1
```

```
43 plt.xticks(range(7), ['Mon', 'Tue', 'Wed', 'Thu', 'Fri', 'Sat', 'Sun',])
44 plt.bar(range(7), author_day)
45 plt.show()
46 print(author_day)
```



统计的结果略微有些违背直觉，CSDN的博主们普遍更喜欢在一周的上半段时间发布文章，而非通常认为的会在周末发文章，周一是大部分作者选择发布文章的高峰期。