

# CS209A Project Report

- 1. **architecture design:**
  - 1.1. **technology stack**
  - 1.2. **Design of database**
  - 1.3. **Front End**
    - 1.3.1. **Exhibition:**
  - 1.4. **Back End**
  - 1.5. **Data Crawling**
    - 1.5.1. **Result**
- 2. **Restful API**
- 3. **Insight**

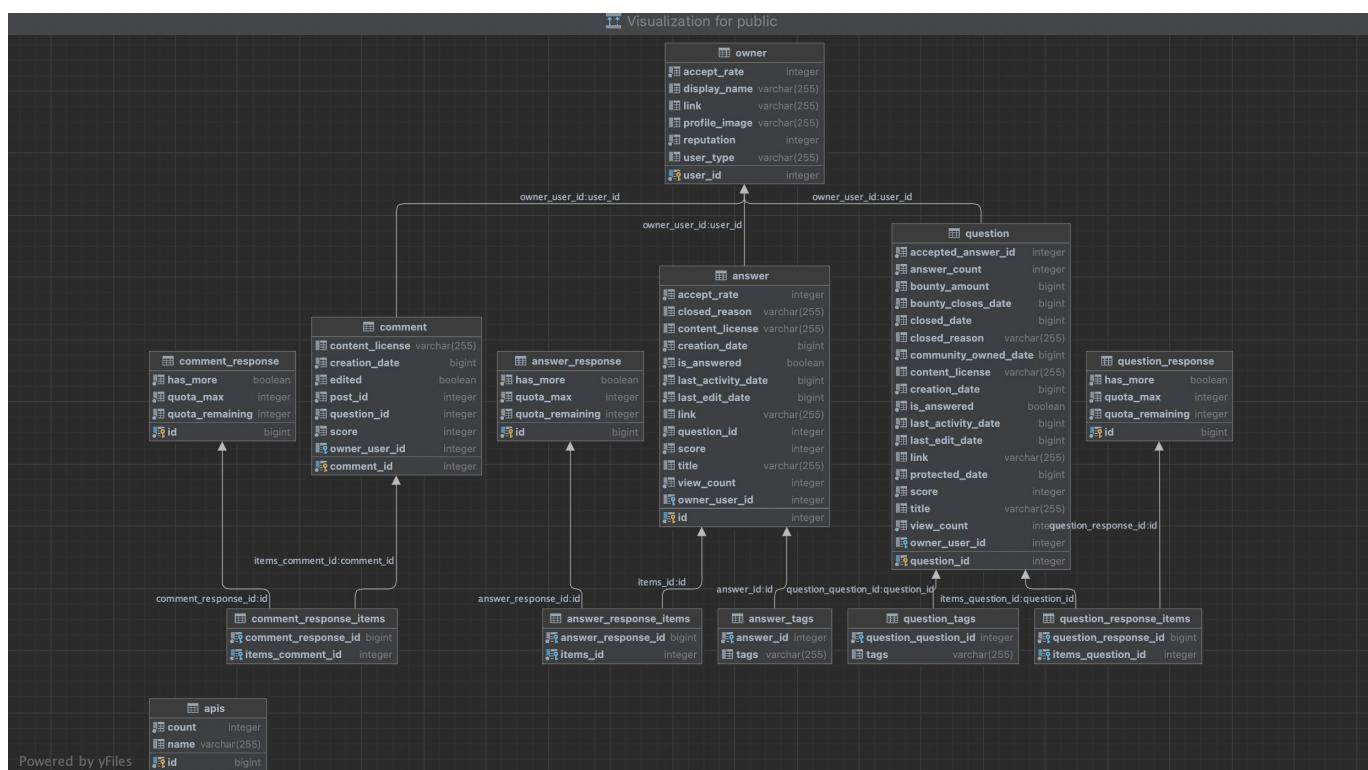
## 1. architecture design:

The architecture of this project mainly consists of two parts, with the front-end to accept data, show the visualization result and the back\_end to process the data and respond to the request of the front-end.

### 1.1. technology stack

- front end: React Axios echarts
- back end: SpringBoot JPA
- database: postgresSQL

### 1.2. Design of database



### 1.3. **Front End**

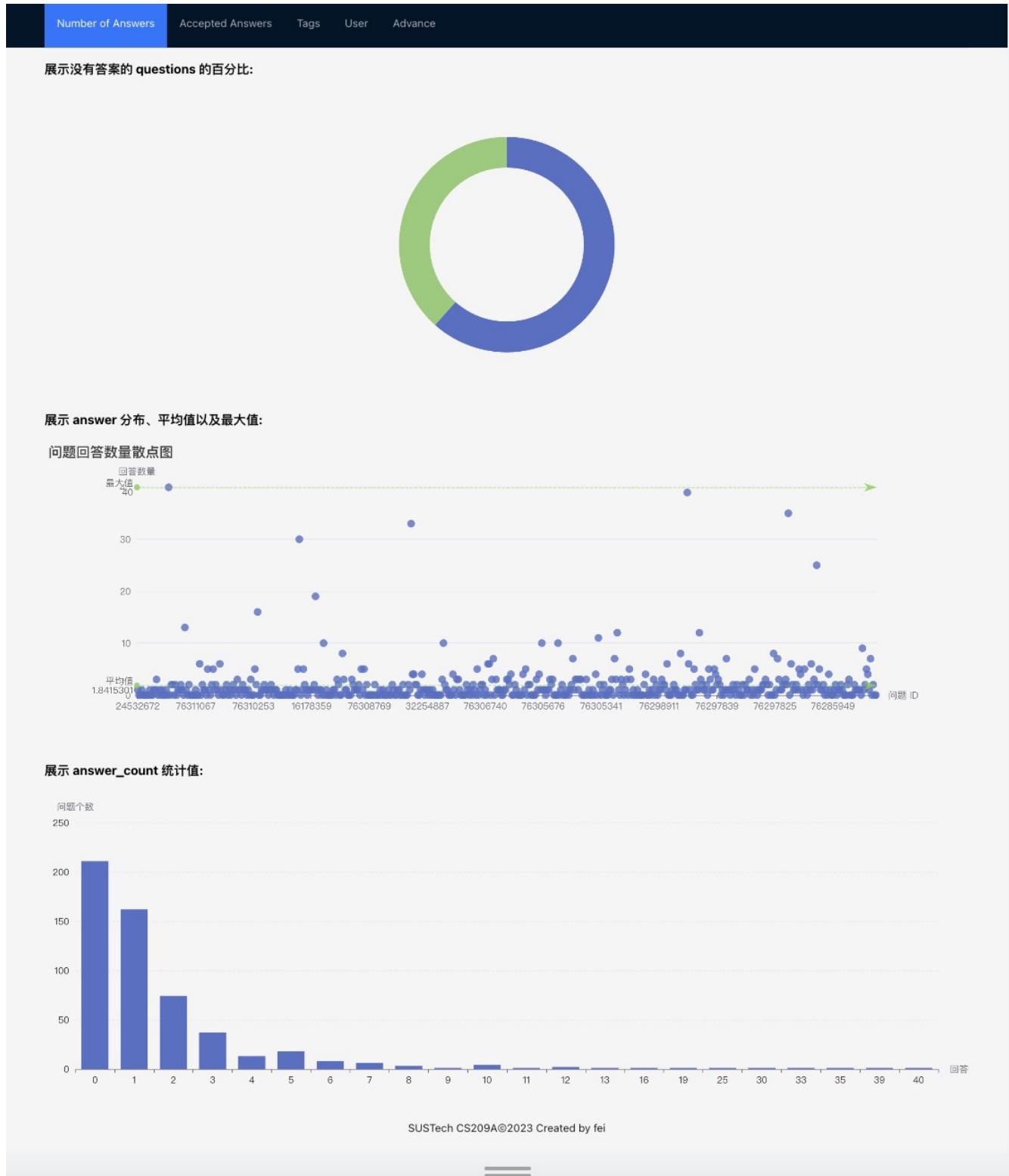
The front-end and back-end interaction was implemented using **axios**, a popular JavaScript library for making HTTP requests.

Axios provided a convenient and efficient way to send requests from the front-end to the back-end API endpoints and handle the responses.

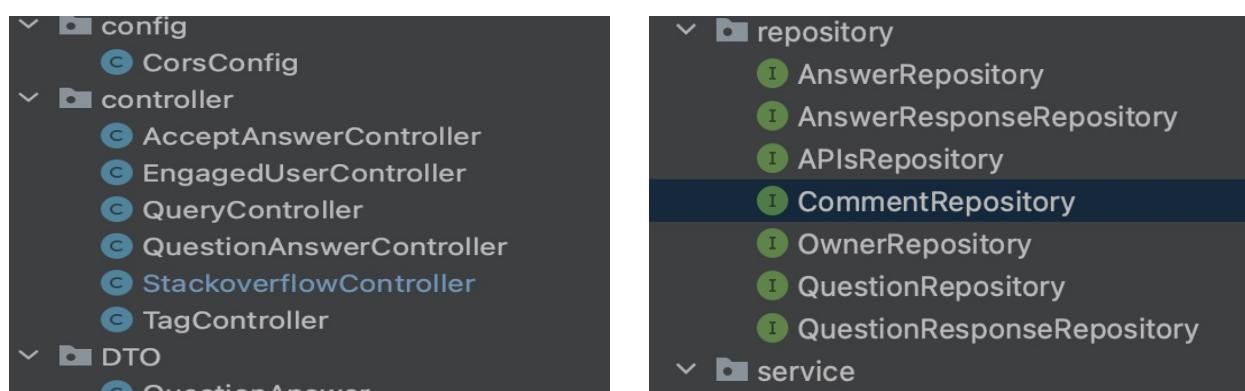
After the front-end sends a request to the endpoint, it gets the data and draws a suitable chart using the **echarts** library.

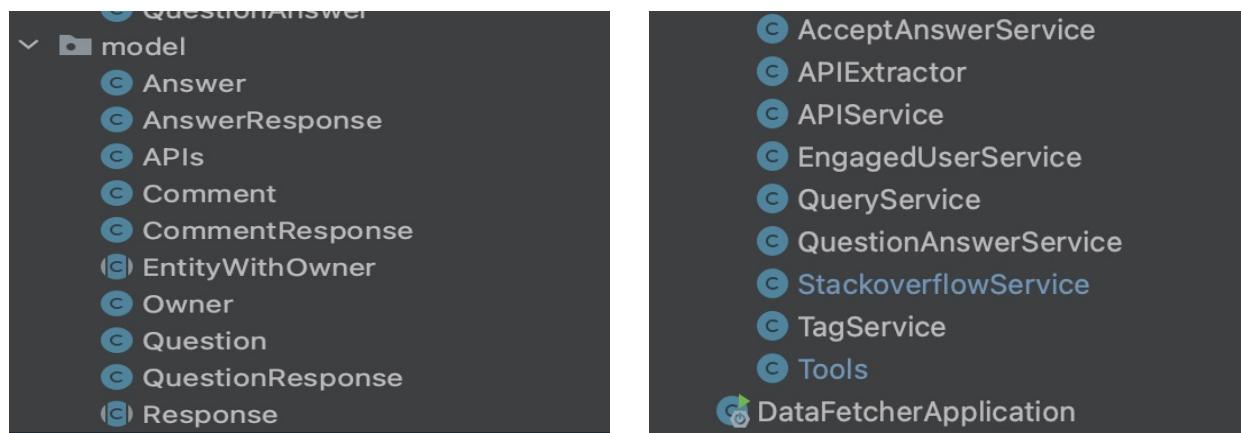
#### 1.3.1. **Exhibition:**

The top navigation bar allows you to select different categories of statistics to view.



## 1.4. Back End





- **CorsConfig:** Solve the cross-domain problem of front and back-end communication
- **StackoverflowService:** Data Crawling
- **AcceptAnswerController:** Provide API related to **Accepted Answers**
- **QuestionAnswerController:** Provide API related to **Number of Answers**
- **TagController:** Provide API related to **Tags**
- **EngagedUserController:** Provide API related to **Users**

## 1.5. Data Crawling

**StackoverflowService** is a service class for getting data from the Stack Overflow API. It contains the following main methods:

**Apache HttpClient:** The CloseableHttpClient and HttpClients classes from the Apache HttpClient library are used to create an HTTP client for making requests to the Stack Overflow API.

### Data Source:

[official Stack Overflow REST API documentation](#)

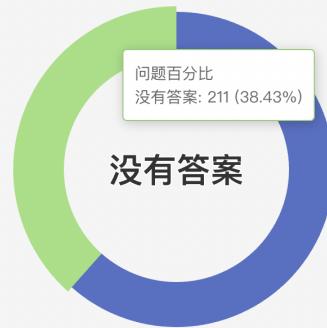
### The basic idea of the crawler is as follows:

A StackoverflowService class is defined to encapsulate the methods and logic for crawling data. A RestTemplate is used to initiate HTTP requests and the Stack Overflow API is called to get the data. The GET request is mainly used and the corresponding parameters are passed to specify the type of data to be obtained, sorting method, tags and other information. Based on the return result of the API, the JSON data is parsed using ObjectMapper and converted to Java objects. The parsed data is saved to the corresponding database table through the methods provided by the JpaRepository interface. For multi-page data, the data is paged using a loop traversal and saved to the database. Page crawling is achieved by controlling the parameters of page number and amount of data per page. Process the returned data if needed, such as extracting the APIs mentioned in the question and counting them. Use the Apache HttpClient library to create and manage HTTP clients for HTTP requests.

### 1.5.1. Result

[What percentage of questions don't have any answer?](#)

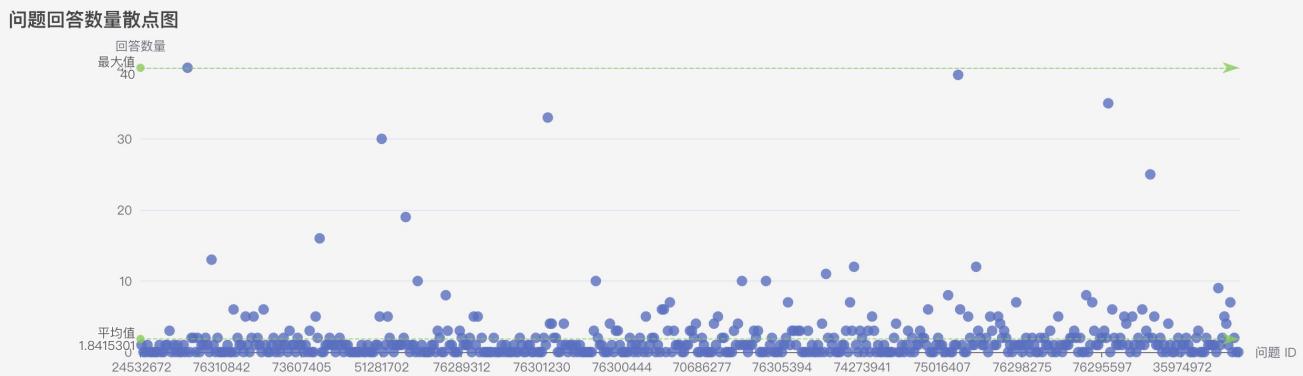
展示没有答案的 questions 的百分比:



- 61.57% questions have answer.

#### What is the average and maximum number of answers?

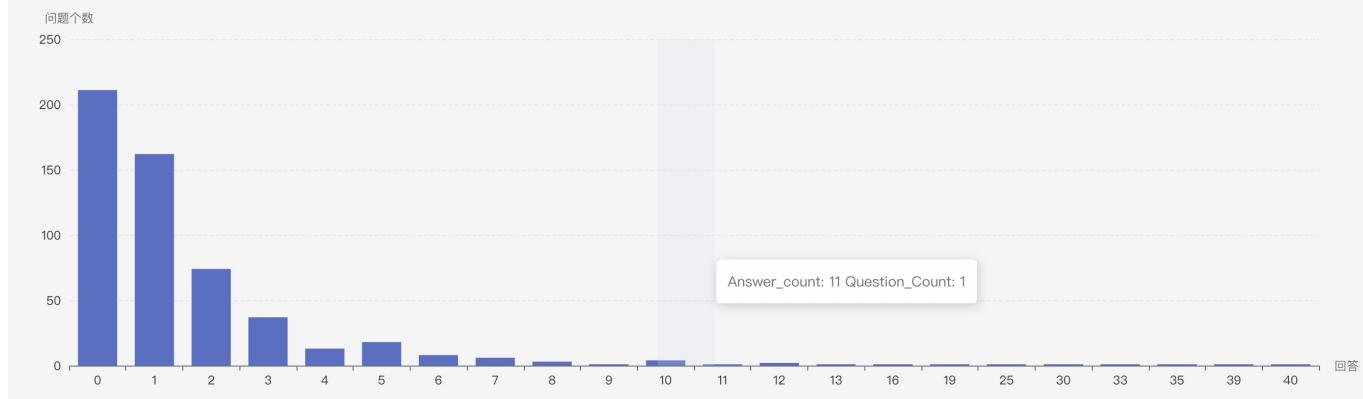
展示 answer 分布、平均值以及最大值:



- average:1.8415301
- maximum:40

#### What is the distribution of the number of answers?

展示 answer\_count 统计值:



- Most of the data had only 5 and less answers. Overall, the number of questions decreases as the number of answers increases.

**What percentage of questions have accepted answers (one question could only have one accepted answer)?**

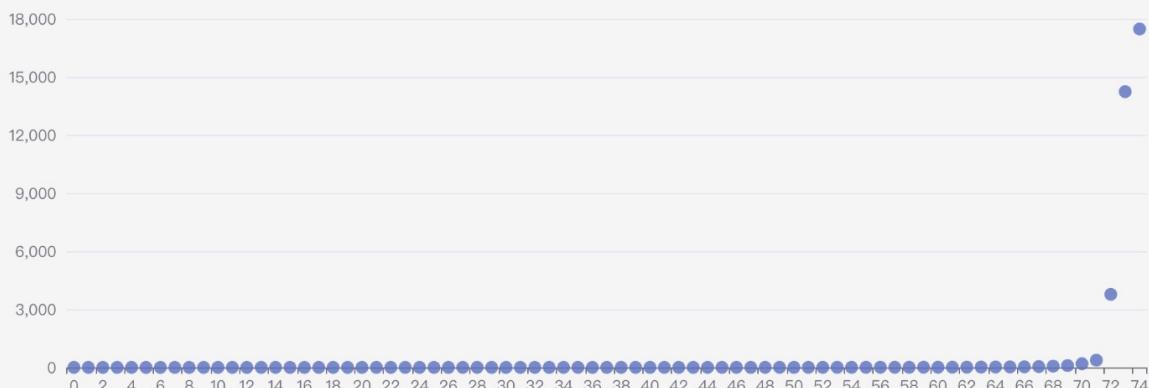
展示有 accepted answer 的问题的百分比:



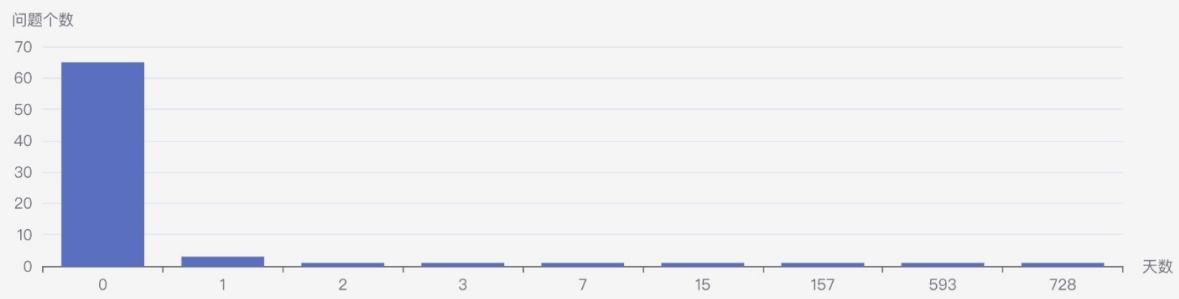
- 20.77% questions have answer.

**What is the distribution of question resolution time (i.e., the duration between the question posting time and the posting time of the accepted answer)?**

展示问题从提出到解决 (answer accepted time – question post time) 的时间间隔分布(单位:h):



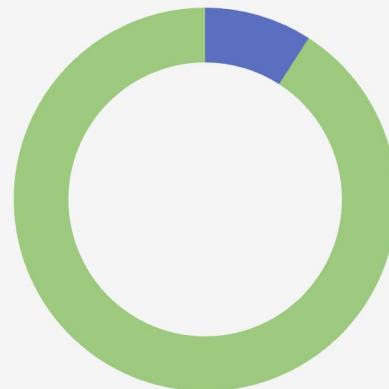
展示问题从提出到解决 (answer accepted time – question post time) 的时间间隔分布(单位:天):



- most question can be solved within one day.

What percentage of questions have non-accepted answers (i.e., answers that are not marked as accepted) that have received more upvotes than the accepted answers?

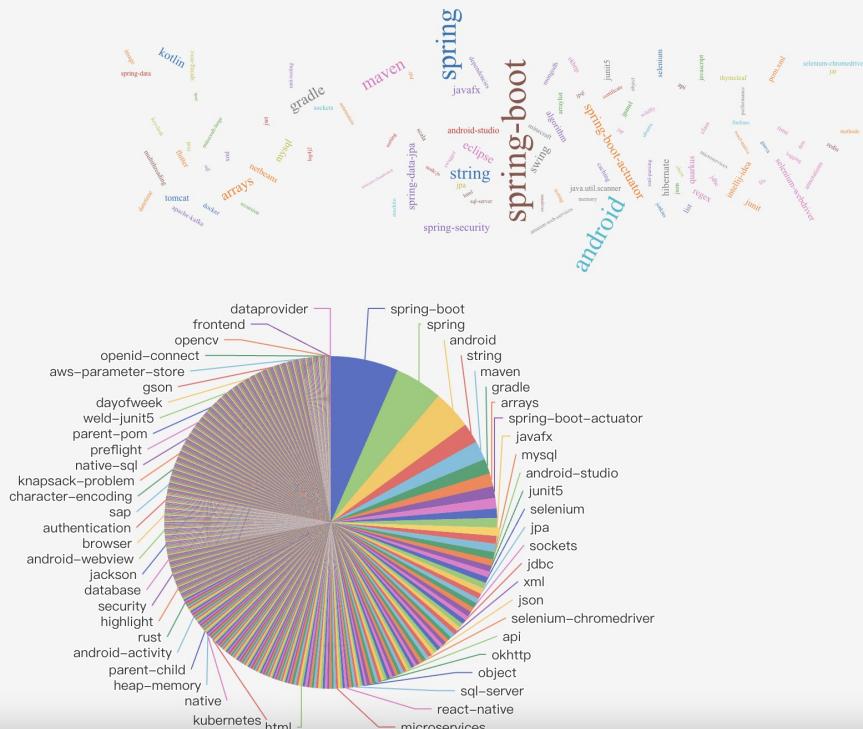
展示含有 non-accepted answer 的 upvote 数高于 accepted answer 的问题的百分比:



- 9.09%
  - Consistent with intuitive expectations, answers with more votes are more likely to be adopted

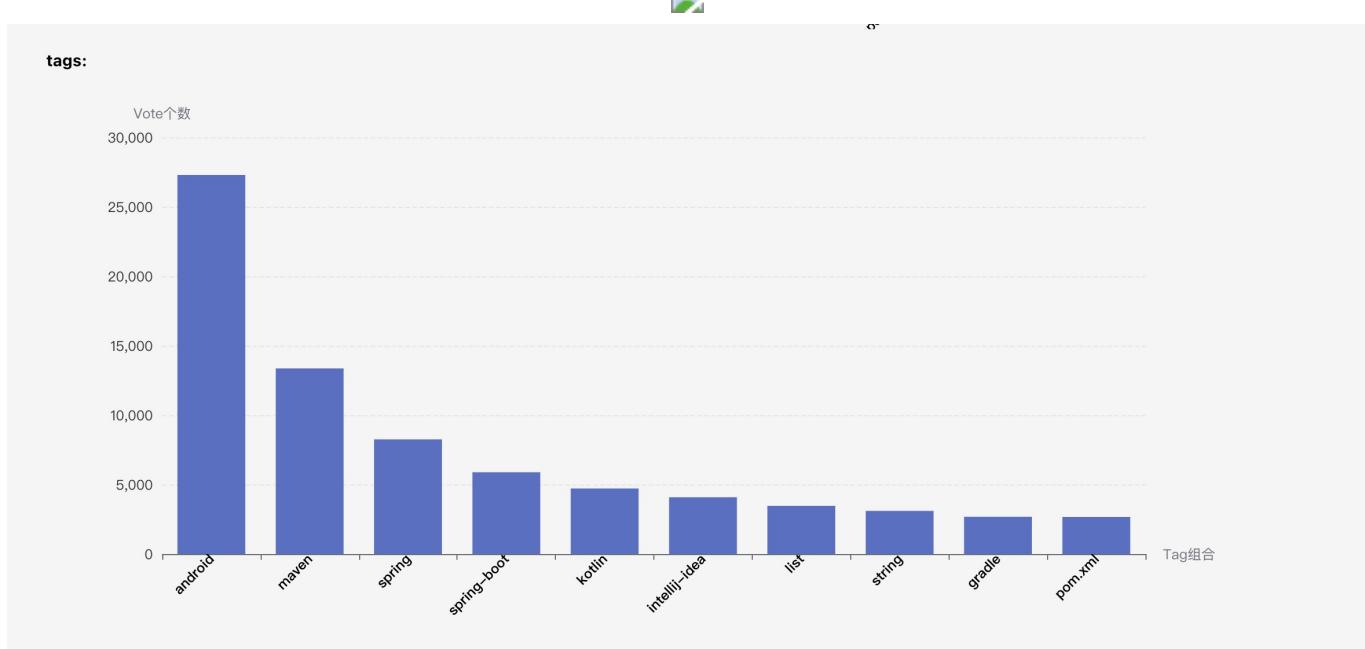
**Which tags frequently appear together with the java tag?**

展示哪些 tags 经常和 Java tag 一起出现:



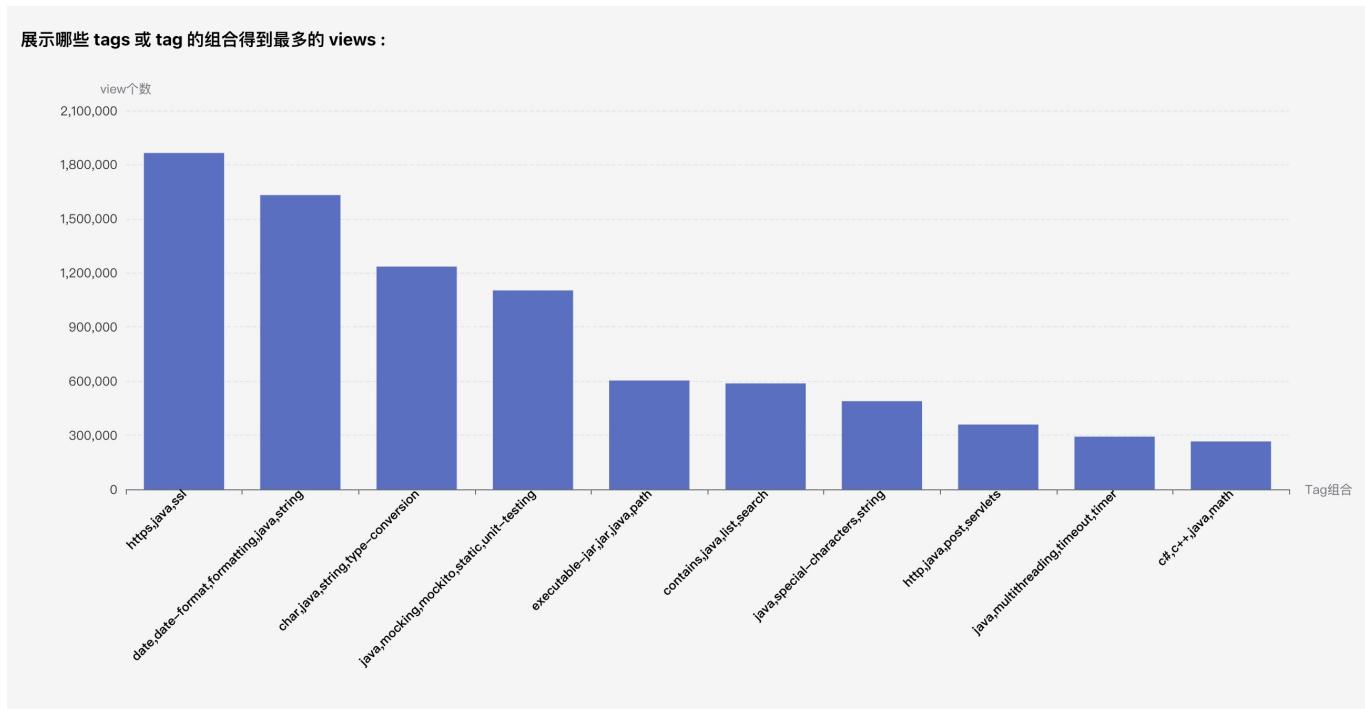
- Spring Boot
  - spring
  - android
  - It looks like the Springboot framework is very popular at the moment

## Which tags or tag combinations receive the most upvotes?



- tag combinations : executable-jar,jar,java,path

### Which tags or tag combinations receive the most views?



- tag combinations : arrylist,collections,duplicates,java,list
- Probably because the list is really used very often

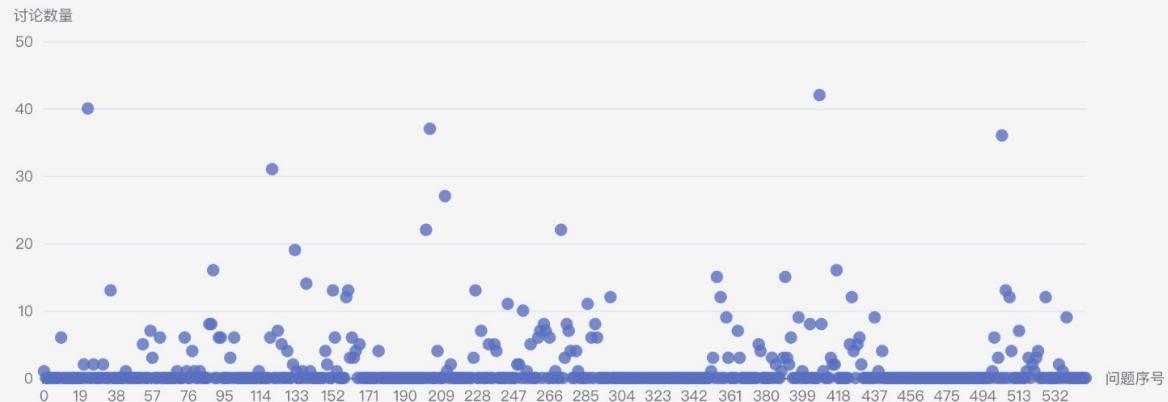
### Many users could participate in a thread discussion. What is the distribution of such participation (i.e., the number of distinct users who post the question, answers, or comments in a thread)?

merge:

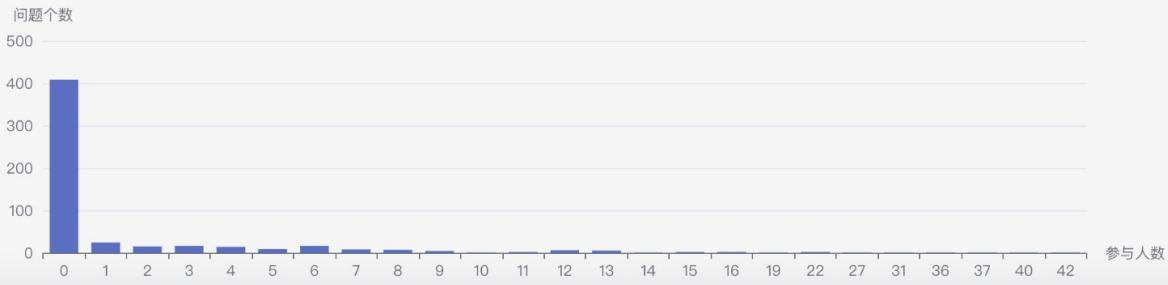
展示参与 Thread 讨论的用户数量的分布:

总体分布 :

散点图(不同问题的参与者者数量):



直方图(不同参与者数量的频数):

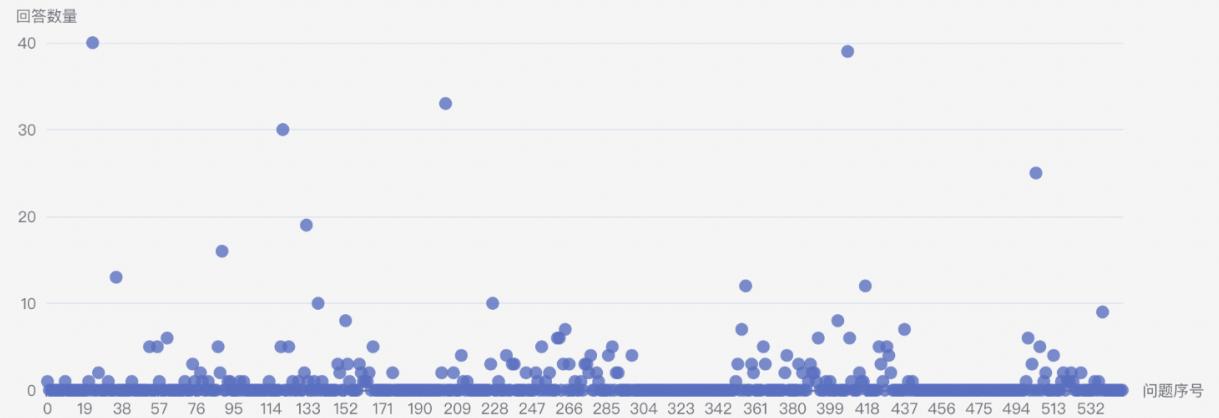


- There are many issues that no one is involved in discussing. Most of the questions with people participating in the discussion had less than 10 participants.

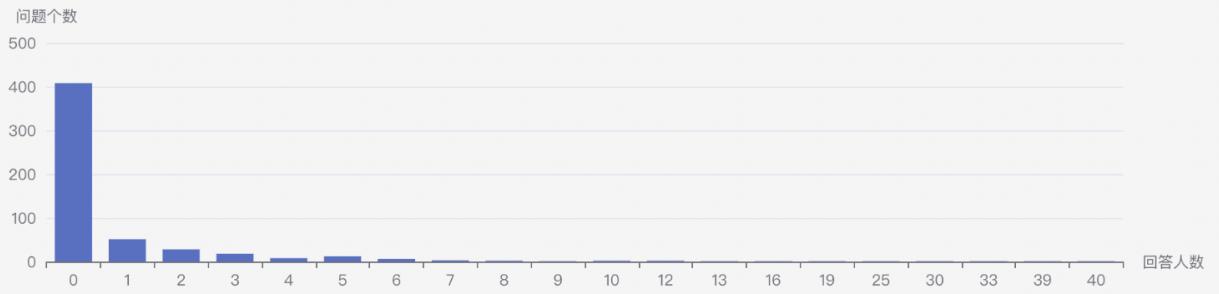
answer:

回答者分布：

散点图(不同问题的回答者数量)：



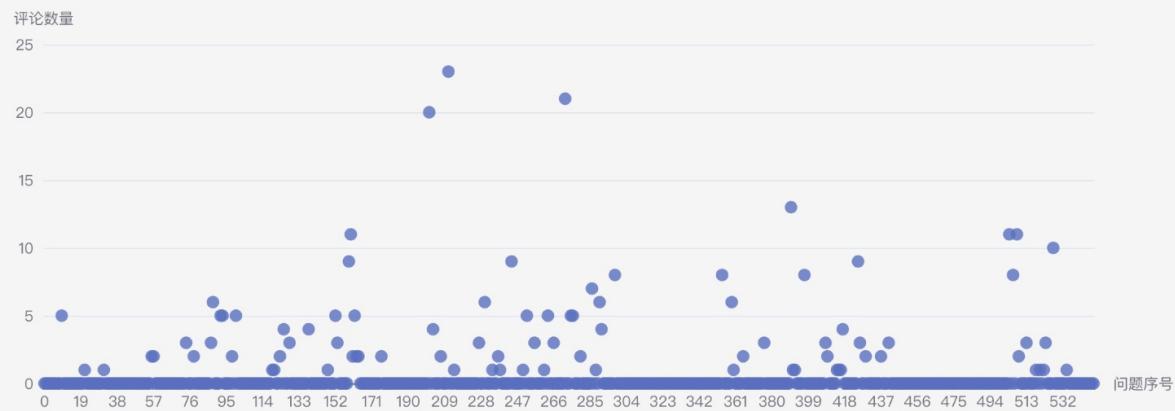
直方图(不同回答者数量的频数)：



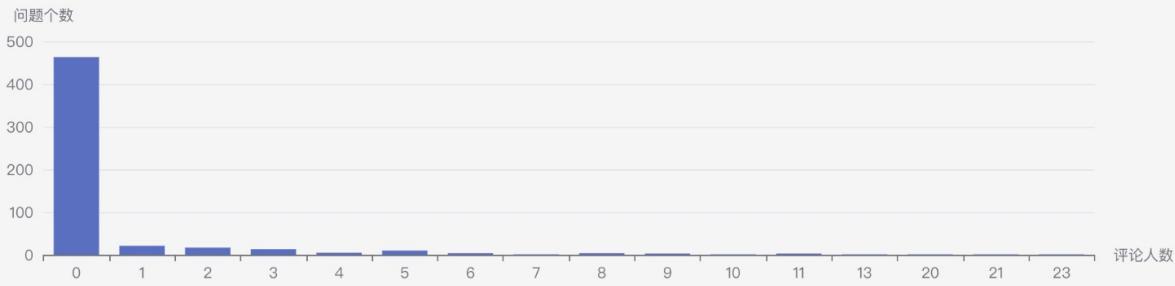
comment:

评论者分布：

散点图(不同问题的评论者数量)：

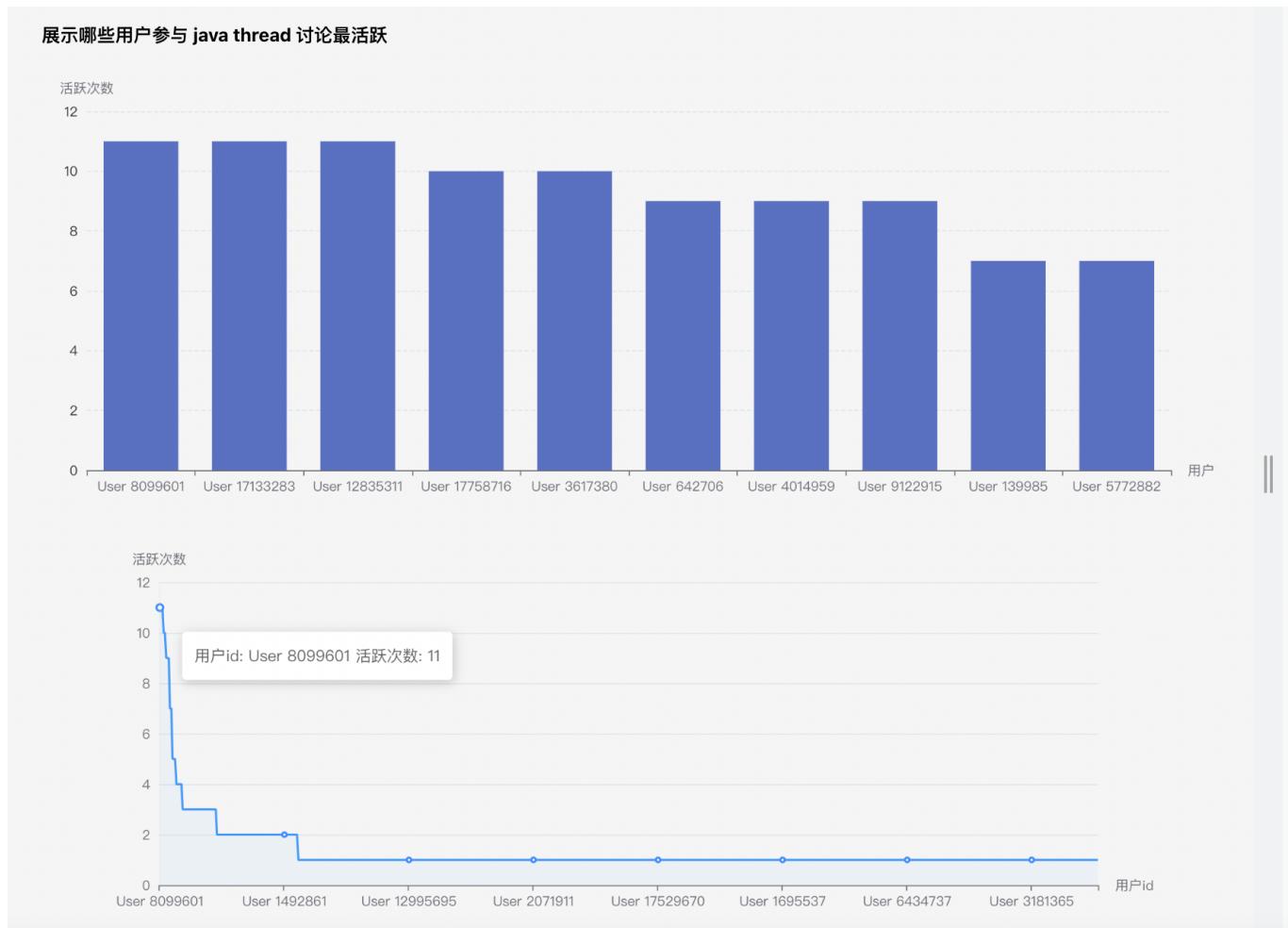


直方图(不同评论者数量的频数)：



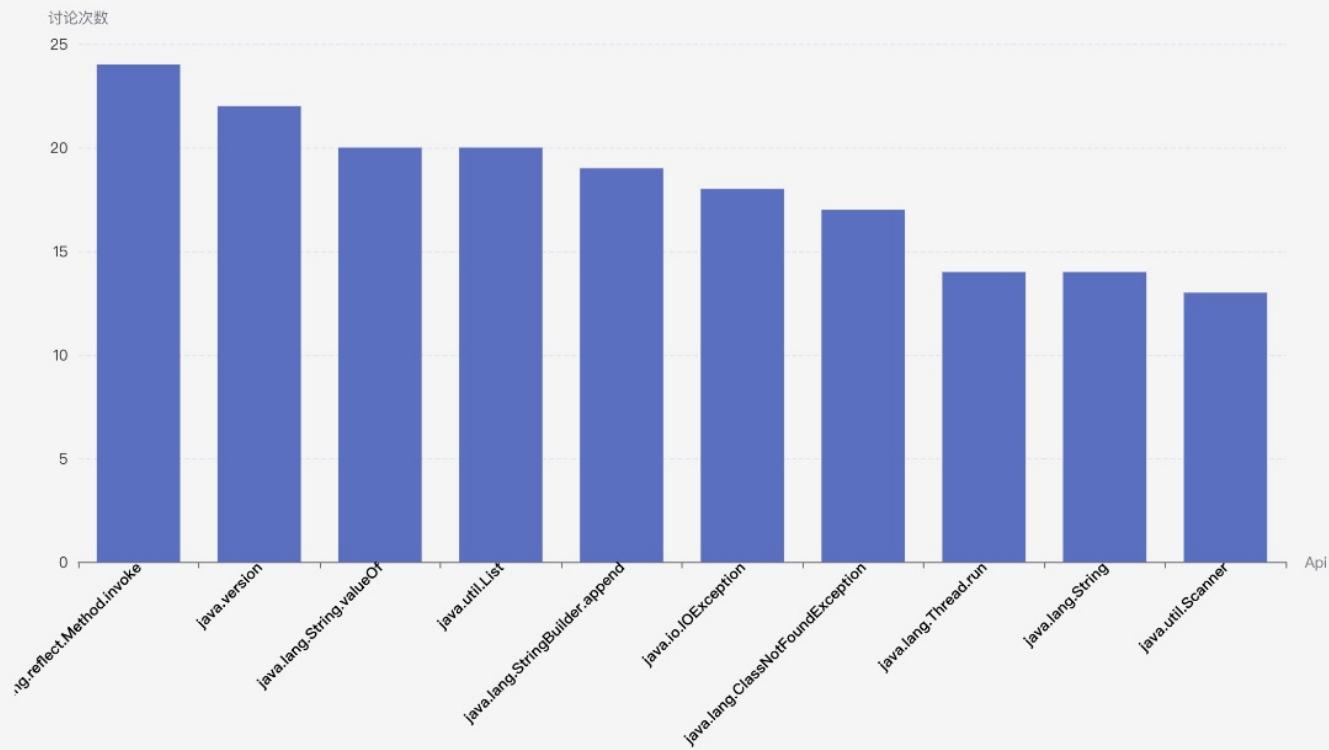
- Questions with a higher number of responses tend to have a higher number of comments. The reverse is also true.

**Which are the most active users who frequently participate in thread discussions?**



- Among the users with active records, most of them have only 1 active count.

### Advance: Frequently discussed Java APIs

**展示哪些 Java APIs 在 Stack Overflow 上最频繁被讨论**

- `java.lang.reflect.invoke`

## 2. Restful API

Header Query Body 认证 预执行脚本 后执行脚本 一键压测 NEW

导入参数 导出参数

	参数名	参数值	必填	类型	参数描述			
<input checked="" type="checkbox"/>	question	47363490	<input checked="" type="checkbox"/>	String	参数描述,用于生成文档			
<input type="checkbox"/>	参数名	参数值,支持mock字段变	<input checked="" type="checkbox"/>	String	参数描述,用于生成文档			
<input type="checkbox"/>	参数名	参数值 支持mock字段变	<input checked="" type="checkbox"/>	String	参数描述 用于生成文档			

实时响应 请求头(10) 响应头(6) Cookies 成功示例 响应示例 200 20:12:58 42.00ms 2.02kb

美化 原生 预览 断言与校验 可视化 utf8 绑定响应结果到变量?

```

1 [
2   {
3     "tags": [],
4     "owner": {
5       "reputation": 1,
6       "link": "https://stackoverflow.com/users/21938958/dicluu",
7       "user_id": 21938958,
8       "user_type": "registered",
9       "profile_image": "https://www.gravatar.com/avatar/bbaeb9f81c1787c14b10f4c36971ced5?s=256&d=identicon&r=PG&f=y&so-version=2",
10      "display_name": "Dicluu",
11      "accept_rate": 0
12    },
13    ...
  
```

- find the answer according to the question id

Header Query Body 认证 预执行脚本 后执行脚本 一键压测 NEW

导入参数 导出参数

	参数名	参数值	必填	类型	参数描述			
<input checked="" type="checkbox"/>	tags	c	<input checked="" type="checkbox"/>	String	参数描述,用于生成文档			
<input checked="" type="checkbox"/>	参数名	参数值,支持mock字段变	<input checked="" type="checkbox"/>	String	参数描述,用于生成文档			

实时响应 请求头(10) 响应头(6) Cookies 成功示例 响应示例 200 20:15:59 267.00ms 0.80kb

美化 原生 预览 断言与校验 可视化 utf8 绑定响应结果到变量?

```

1 [
2   {
3     "tags": [
4       "java",
5       "python",
6       "c++",
7       "c",
8       "hash"
9     ],
10    "owner": {
11      "reputation": 547,
12      "link": "https://stackoverflow.com/users/2419921/devesh-saini",
13      "user_id": 2419921,
14      ...
  
```

- find the problem according to the tag

GET <http://localhost:8080/query/accepted-answer?question=5127687> 发送

Header Query Body 认证 预执行脚本 后执行脚本 一键压测 NEW

导入参数  导出参数

	参数名	参数值	必填	类型	参数描述			
<input checked="" type="checkbox"/>	question	5127687	<input type="checkbox"/>	String	参数描述,用于生成文档	<input type="checkbox"/>	<input checked="" type="checkbox"/>	<input type="checkbox"/>
<input checked="" type="checkbox"/>	参数名	参数值,支持mock字段变	<input checked="" type="checkbox"/>	String	参数描述,用于生成文档	<input type="checkbox"/>	<input checked="" type="checkbox"/>	<input type="checkbox"/>

实时响应 请求头(10) 响应头(6) Cookies 成功示例 响应示例 200 ⌚ 20:17:22 78.00ms 0.51kb 🌐 ⌚

美化 原生 预览 断言与校验 可视化 utf8 utf8 ⌚ ⌚ 绑定响应结果到变量?

```

1 {
2   "tags": [],
3   "owner": {
4     "reputation": 42856,
5     "link": "https://stackoverflow.com/users/238884/michael-lorton",
6     "user_id": 238884,
7     "user_type": "registered",
8     "profile_image": "https://www.gravatar.com/avatar/febf321631229c7408834722691c06e9?s=256&d=identicon&r=PG",
9     "display_name": "Michael Lorton",
10    "accept_rate": 60
11  },
12  "is_answered": false,
13  ...

```

- find the accepted answer according to the question id

GET <http://localhost:8080/query/comment/question?id=5127687> 发送

Header Query Body 认证 预执行脚本 后执行脚本 一键压测 NEW

导入参数  导出参数

	参数名	参数值	必填	类型	参数描述			
<input checked="" type="checkbox"/>	id	5127687	<input type="checkbox"/>	String	参数描述,用于生成文档	<input type="checkbox"/>	<input checked="" type="checkbox"/>	<input type="checkbox"/>
<input checked="" type="checkbox"/>	参数名	参数值,支持mock字段变	<input checked="" type="checkbox"/>	String	参数描述,用于生成文档	<input type="checkbox"/>	<input checked="" type="checkbox"/>	<input type="checkbox"/>

实时响应 请求头(10) 响应头(6) Cookies 成功示例 响应示例 200 ⌚ 21:03:59 42.00ms 0.43kb 🌐

美化 原生 预览 断言与校验 可视化 utf8 utf8 ⌚ ⌚ 绑定响应结果到变量?

```

1 {
2   "owner": {
3     "reputation": 42856,
4     "link": "https://stackoverflow.com/users/238884/michael-lorton",
5     "user_id": 238884,
6     "user_type": "registered",
7     "profile_image": "https://www.gravatar.com/avatar/febf321631229c7408834722691c06e9?s=256&d=identicon&r=PG",
8     "display_name": "Michael Lorton",
9     "accept_rate": 60
10  },
11  "edited": false,
12  "score": 1

```

- find the comment according to the question id

Header: Query Body 认证 预执行脚本 后执行脚本 一键压测 NEW

参数名: id, 参数值: 120445855; 必填: No, 类型: String, 描述: 参数描述, 用于生成文档

参数名: , 参数值: 参数值, 支持 mock 字段变, 必填: Yes, 类型: String, 描述: 参数描述, 用于生成文档

实时响应: 请求头(10) 响应头(6) Cookies 成功示例 响应示例 200 21:11:41 85.00ms 0.42kb

美化 原生 预览 断言与校验 可视化 utf8 编辑 搜索 绑定响应结果到变量?

```

1 {
2   "owner": {
3     "reputation": 2861,
4     "link": "https://stackoverflow.com/users/2481153/hungnm2",
5     "user_id": 2481153,
6     "user_type": "registered",
7     "profile_image": "https://www.gravatar.com/avatar/c185ac8babb726f8c78fdfb061a3a360?s=256&d=identicon&r=PG",
8     "display_name": "HungNM2",
9     "accept_rate": 0
10    },
11    "edited": false,
12    "score": 0,
13    "comment": {
14      "id": 120445855,
15      "content": "This is a test comment.",
16      "owner": {
17        "reputation": 2861,
18        "link": "https://stackoverflow.com/users/2481153/hungnm2",
19        "user_id": 2481153,
20        "user_type": "registered",
21        "profile_image": "https://www.gravatar.com/avatar/c185ac8babb726f8c78fdfb061a3a360?s=256&d=identicon&r=PG",
22        "display_name": "HungNM2",
23        "accept_rate": 0
24      },
25      "parent_id": null,
26      "created_at": "2023-05-23T12:00:00Z",
27      "updated_at": "2023-05-23T12:00:00Z"
28    }
29  }
30}

```

- find comment by comment id

Header: Query Body 认证 预执行脚本 后执行脚本 一键压测 NEW

参数名: id, 参数值: 18729480, 必填: Yes, 类型: String, 描述: 参数描述, 用于生成文档

参数名: , 参数值: 参数值, 支持 mock 字段变, 必填: Yes, 类型: String, 描述: 参数描述, 用于生成文档

实时响应: 请求头(10) 响应头(6) Cookies 成功示例 响应示例 200 21:02:35 124.00ms 0.31kb

美化 原生 预览 断言与校验 可视化 utf8 编辑 搜索 绑定响应结果到变量?

```

1 [
2   {
3     "reputation": 37,
4     "link": "https://stackoverflow.com/users/18729480/pluggarose",
5     "user_id": 18729480,
6     "user_type": "registered",
7     "profile_image": "https://lh3.googleusercontent.com/-fRiFiolJ6U/AAAAAAAIAAI/AAAAAAA/AMZuucn0kcE1UYFtEwSZHVzow3XpZuHI3Q/s96-c/photo.jpg?sz=256",
8     "display_name": "PluggaRose",
9     "accept_rate": 0
10   }
11 ]

```

- find user by userid

### 3. **Insight**

The user activity frequency on stackoverflow follows a \*\*\*pyramid structure\*\*\*, which can provide us with the following insights:

1. User Distribution: The pyramid structure indicates an imbalance in user groups. Only a small portion of users are highly active on the platform, while the majority of users are less engaged or passive.
2. Key Users: The top of the pyramid represents highly active users who are often opinion leaders, content creators, or frequent participants in discussions and interactions on social media platforms. These users are crucial for platform growth and community development.
3. Long Tail Effect: The bottom of the pyramid consists of a large number of users with relatively low activity levels, but together, they contribute a significant portion of the overall activity. This reflects the long tail effect, where the cumulative contribution of many low-activity users can surpass that of a few high-activity users.
4. User Engagement: The pyramid structure of user activity suggests that platforms may need to take measures to enhance user engagement and activity levels. By providing engaging content, facilitating user interaction, and implementing reward mechanisms, platforms can stimulate user interest and participation to mitigate a decline in activity.
5. User Segmentation and Personalization: Understanding the distribution of user activity can help platforms segment users and provide personalized recommendations. By analyzing user behavior and interests, tailored content and experiences can be offered to users with different activity levels, thereby increasing their engagement.

The user question answerd frequency and most view tag ,upvote tag and api on stackoverflow follows a **pyramid structure**, which can provide us with the following insights:

1. These features of Java are probably the most commonly used.eg.list,springboot,java.lang.reflect.invoke.
2. Some api may not be perfect yet or easy to user, so was asked a lot of times.
3. springboot is probably the most popular framework now.