# CS209A Project Report

## Stack Overflow Q&A Data w.r.t. Java Programming

### Student Name and Student ID

丁健乐　　12111517

林俊杰　　12111511

# 1　　Contributions

| Member | Contributions | Ratio |
|---|---|---|
| 丁健乐 | 1.　Data getting and storage | 50% |
| 林俊杰 | 1.　Data analysis and visualization | 50% |

# 2　　Problem Analysis

Number of answers (15 points)

1. Percentage of unanswered questions

2. For each question, the mean and maximum values of the answers

3. Distribution of the number of answers for each question

Acceptable answer (15 points)

1. Percentage of questions with acceptable answers

2. Distribution of problem resolution times (problem release times -- release times with acceptable answers)

3. The percentage of questions in which an answer that was not accepted received more likes than an answer that was acceptable

Label (15 points)

1. Tags that often appear with java tags

2. The hashtag or hashtag combination that gets the most likes

3. The TAB or combination of tabs that get the most views

Users (15 points)

1. The distribution of the total number of different users Posting questions, answering questions, or commenting on answers, per question

2. The most active users who regularly participate in discussions

# 3  The Description of the Collected Data

First of all, we obtained 500 basic information from https://api.stackexchange.com/2.3/questions, Including owner_id, question_id, is_answered answer_numbers, question_posting_time, tags, upvote, views, user_count, comment. At the same time, we also store every user who asks a question. Then we added from https://api.stackexchange.com/2.3/questions/ {ids} / answers to get accepted_answer_id, answer_posting_time,most_upvote_answer_id. We also store every user who answers the question

Finally, we go through https://api.stackexchange.com/2.3/questions/ {ids} / comments added to every question the user comments. So that's the combination of the data that we got.

The SQL statements below are the tables we built to store the answer.

```
 1  CREATE TABLE question
 2  (
 3      owner_id              VARCHAR,
 4      question_id           INTEGER PRIMARY KEY,
 5      is_answered           BOOLEAN,
 6      answer_numbers        INTEGER,
 7      accepted_answer_id    INTEGER,
 8      question_posting_time INTEGER,
 9      answer_posting_time   INTEGER,
10      most_upvote_answer_id INTEGER,
11      tags                  VARCHAR,
12      upvote                INTEGER,
13      views                 INTEGER,
14      user_count            INTEGER,
15      comment               INTEGER
16  );
17
18  CREATE TABLE owner
19  (
20      owner_id         VARCHAR PRIMARY KEY,
21      question_numbers INTEGER,
22      answer_numbers   INTEGER,
```

```
23     comment_numbers   INTEGER
24 );
```

# 4      The Architecture Design of Project

In terms of getting the data, I use Java's HttpURLConnection to send the HTTP request and process the JSON response data returned from the request. A BufferedReader in is then created to read the connected input stream. Determine whether to use GZIP decompression based on the content encoding in the response header. The code uses a BufferedReader to read the data line by line, appending each line to the response until the data read is empty. Next, the code converts the JSON data in response to the JSONObject object 'json'. Finally, we loop through each element of the items array to get the corresponding JSONObject object 'item' for each element.

```java
1 URL url = new URL(apiUrl + "?page=" + page + "&pagesize=" + pageSize +
   apiParams);
2 HttpURLConnection connection1 = (HttpURLConnection) url.openConnection();
3 connection1.setRequestMethod("GET");
4 BufferedReader in;
5 if ("gzip".equals(connection1.getContentEncoding())) {
6     in = new BufferedReader(
7         new InputStreamReader(new
   GZIPInputStream(connection1.getInputStream())));
8 } else {
9     in = new BufferedReader(new
   InputStreamReader(connection1.getInputStream()));
10 }
11 StringBuilder response = new StringBuilder();
12 String line;
13 while ((line = in.readLine()) != null) {
14     response.append(line);
15 }
16 in.close();
17 JSONObject json = new JSONObject(response.toString());
18 JSONArray items = json.getJSONArray("items");
```

In terms of storing data, after I get the data from the website, I select the data I need and save it in the postgresql database. I choose data insertion as an example:

```java
1 private static void insertData(Connection connection, String ownerId, int
   questionId,
```

```
2       boolean isAnswered, int answerCount,
3       int acceptedAnswerId, int questionPostingTime, int answerPostingTime,
4       int most_upvote_answer_id, String tags,
5       int upvote, int views, int userCount, int comments) throws SQLException,
   SQLException {
6       String insertQuery = "INSERT INTO question
   (owner_id,question_id,is_answered," +
7
   "answer_numbers,accepted_answer_id,question_posting_time,answer_posting_time,"
   +
8           "most_upvote_answer_id,tags,upvote,views,user_count,comment) VALUES
   (?,?,?,?,?,?,?,?,?,?,?,?,?)";
9       PreparedStatement statement = connection.prepareStatement(insertQuery);
10      statement.setString(1, ownerId);
11      statement.setInt(2, questionId);
12      statement.setBoolean(3, isAnswered);
13      statement.setInt(4, answerCount);
14      statement.setInt(5, acceptedAnswerId);
15      statement.setInt(6, questionPostingTime);
16      statement.setInt(7, answerPostingTime);
17      statement.setInt(8, most_upvote_answer_id);
18      statement.setString(9, tags);
19      statement.setInt(10, upvote);
20      statement.setInt(11, views);
21      statement.setInt(12, userCount);
22      statement.setInt(13, comments);
23      statement.executeUpdate();
24      statement.close();
25 }
```

I use PreparedStatement to store the insertQuery, and add the required part into the statement and execute the statement.

About data analysis, we use SQL statements to search for the required data from the database, and store it into JavaScript files and Sync to visualization side. The code above is an example of the whole process:

```
1 public static void ParticipationDistribution(String str) throws SQLException {
2      String filePath = "src/main/resources/static/js/Users/ParticipationDistribut
3
4      Statement stmt = con.createStatement();
5      ResultSet rs0 = stmt.executeQuery("select count(*)\n"
6          + "from (select (answer_numbers + comment + 1) as sum\n"
7          + "       from question) l\n"
8          + "where sum >= 1\n"
```

```java
 9             + "  and sum <= 10;");
10      rs0.next();
11      int first = rs0.getInt(1);
12
13      ResultSet rs1 = stmt.executeQuery("select count(*)\n"
14          + "from (select (answer_numbers + comment + 1) as sum\n"
15          + "      from question) l\n"
16          + "where sum > 10\n"
17          + "  and sum <= 20;");
18      rs1.next();
19      int second = Math.round(rs1.getFloat(1));
20
21      ResultSet rs2 = stmt.executeQuery("select count(*)\n"
22          + "from (select (answer_numbers + comment + 1) as sum\n"
23          + "      from question) l\n"
24          + "where sum > 20\n"
25          + "  and sum <= 30;");
26      rs2.next();
27      int third = Math.round(rs2.getFloat(1));
28
29      ResultSet rs3 = stmt.executeQuery("select count(*)\n"
30          + "from (select (answer_numbers + comment + 1) as sum\n"
31          + "      from question) l\n"
32          + "where sum > 30\n"
33          + "  and sum <= 40;");
34      rs3.next();
35      int fourth = Math.round(rs3.getFloat(1));
36
37      ResultSet rs4 = stmt.executeQuery("select count(*)\n"
38          + "from (select (answer_numbers + comment + 1) as sum\n"
39          + "      from question) l\n"
40          + "where sum > 40;");
41      rs4.next();
42      int fifth = Math.round(rs4.getFloat(1));
43
44      try (BufferedReader reader = new BufferedReader(new FileReader(filePath))) {
45          StringBuilder content = new StringBuilder();
46          String line;
47          while ((line = reader.readLine()) != null) {
48              if (line.contains("let " + str + " =")) {
49                  Map<String, Integer> values = extractValues(line);
50                  if (values.containsKey(">= 1 & <= 10")) {
51                      values.put(">= 1 & <= 10", first);
52                  }
53                  if (values.containsKey("> 10 & <= 20")) {
54                      values.put("> 10 & <= 20", second);
55                  }
```

```
56                  if (values.containsKey("> 20 & <= 30")) {
57                      values.put("> 20 & <= 30", third);
58                  }
59                  if (values.containsKey("> 30 & <= 40")) {
60                      values.put("> 30 & <= 40", fourth);
61                  }
62                  if (values.containsKey("> 40")) {
63                      values.put("> 40", fifth);
64                  }
65                  String modifiedLine = generateModifiedLine(str, values);
66                  line = line.replaceFirst("let " + str + " =.*", modifiedLine);
67              }
68              content.append(line).append("\n");
69          }
70          try (BufferedWriter writer = new BufferedWriter(new FileWriter(filePath)
71              writer.write(content.toString());
72          }
73          System.out.println("JavaScript file modified successfully.");
74      } catch (IOException e) {
75          System.out.println(
76              "An error occurred while modifying the JavaScript file: " + e.getMes
77      }
78 }
```

# 5    The insights from the Data Analysis Results

After conducting the data analysis on Java programming, several insights have emerged. Here are some of the key findings:

1. What percentage of questions don't have any answer?

    Almost two-thirds of the questions asked were not answered.

    "No answer": 332, "Has answer": 168

2. What is the average and maximum number of answers?

    The maximum number of answers to answered questions was 40 and the mean was 1.832.

    "Maximum": 40, "Average": 2

3. What is the distribution of the number of answers?

    ">= 4 & <= 6": 27, ">= 10": 19, ">= 1 & <= 3": 238, "= 0": 212, ">= 7 & <= 9": 4

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

    "Has accepted answer": 76, "No accepted answer": 424

5. 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)?

   ">= 0, <= 1000": 15, "> 100000": 17, "> 1000, <= 10000": 24, "> 10000, <= 100000": 20

   The unit of time is seconds.

6. 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?

   "More Upvotes Non-accepted Answer": 26, "Less Upvotes Non-accepted Answer": 50

7. Which tags frequently appear together with the java tag?

   Most Frequently Asked Topics: The analysis revealed that certain topics of Java programming are asked more frequently than others. The top five most asked topics are:

   A. Spring-boot  B. Spring  C. Android  D. Spring-boot-actuator  E. Maven

8. Which tags or tag combinations receive the most upvotes?

   A. List  B. Jar  C. Arraylist  D. Java-stream  E. Java8

9. Which tags or tag combinations receive the most views?

   A. Spring-boot  B. Spring  C. Android  D. Spring-boot-actuator  E. Maven

10. 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)?

    "> 40": 1, "> 20 & <= 30": 5, "> 10 & <= 20": 33, ">= 1 & <= 10": 457, "> 30 & <= 40": 4

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

    "6140624", "28659413", "7079914", "19118721", "28657233"...

# 6    Code

https://github.com/isLinHehea/CS209A_Project