# Introduction

In the world of Data mining, the major task is to crawl, clean and retrieve meaningful information from the data. In order to retrieve meaningful information from the **WikiCFP** website, I have crawled it using **Java** and parse the data using **Jsoup: Java HTML Parser**. Then for cleaning **OpenRefine** is used. I have used **Hadoop** to gain an insight from the data. The processes followed throughout the accomplishment of each task is discussed here.

## 1. Data Crawling

In order to crawl the WikiCFP website for the conferences of data mining, machine learning, database and AI, I have used the provided Java code. In order to fetch the conferences of these 4 categories, I have used a string array named **searchCategory** and fetched 20 pages of each category. After fetching the content of each page, a function named **parseTableData()** is called. In this function, the content of one page is parsed. For data parsing, **Jsoup: Java HTML Parser** is used. While parsing the data of one page, it is observed that there are multiple tables on every page. Among them, table number 2 is the one in which the conference data is displayed. Therefore, table 2 is selected from the parsed data.

After selecting the correct table, each row is considered to fetch the conference information. In the table, each conference information is represented using 2 consecutive rows. Therefore, in order to fetch the information of one conference, two row is selected at a time. From Figure 1, the real scenario can be found, where each conference consists of two rows. The first row contains conference name and the down row contains time, location and deadline of the conference.



Figure 1: One conference entry

From each conference entry, the required information is collected as follows,
**1) Conference Acronym**: The conference acronym is found at table column number 0 of first row. So this column is selected for this information. There are some conference entries in which the conference acronym is **Expired CFPs**. These entries are ignored for further processing i.e. the conference name and location is not collected for these entries.

**2) Conference Name**: After getting the correct conference acronym, the conference name is fetched by selecting column number 1 of first row.

**3) Conference Location**: In order to fetch the conference location, the second row is selected. In second row, the conference location is found in column 1. Therefore, the conference location is selected from there. Though it is not required to select the time and deadline of the conference, these information are also selected to check if the required information is collected properly.

The three information is appended using tab between them in a string. Then all the conference entries of one page is returned to the calling routine where the returned string is written in the file named **wikicfp_crawl.txt**. The conference information of all categories i.e. Data mining, Artificial Intelligence, Databases and Machine Learning are stored in the same file.

In order to ensure the proper implementation of this portion, I have used netinstructions.com[1] to get an idea about how a java crawler works and Stackoverflow[2] to get an idea about how the table data are stored.

## 2. Data Cleaning

In the data cleaning portion, **OpenRefine** is used. Since the data were stored in tab separated format, when the file containing all conference data is uploaded in OpenRefine, all the data are automatically organized in a table format by OpenRefine as Figure 2.

| All | | | Conference Acronym | Conference Name | Conference Location |
|---|---|---|---|---|---|
| ☆ | ⬒ | 1. | ICBCT 2017 | 2017 International Conference on Bioinformatics and Computing Technologies (ICBCT 2017) | Hong Kong, China |
| ☆ | ⬒ | 2. | SWM 2017 | 1st Workshop on Scholarly Web Mining | Cambridge, UK |
| ☆ | ⬒ | 3. | IFAC World Congress SIMCA 2017 | IFAC WORLD CONGRESS '2017 OPEN INVITED TRACK on SYSTEM IDENTIFICATION for MANUFACTURING CONTROL APPLICATIONS | Toulouse, France |
| ☆ | ⬒ | 4. | SPTM 2016 | Fourth International Conference of Security, Privacy and Trust Management | Chennai, India |
| ☆ | ⬒ | 5. | DKMP 2017 | Fifth International Conference on Data Mining & Knowledge Management Process | Dubai, UAE |
| ☆ | ⬒ | 6. | CCSEA 2017 | Seventh International Conference on Computer Science, Engineering and Applications | Dubai, UAE |
| ☆ | ⬒ | 7. | Big Data Analytics@IJCNN 2017 | Special Session: 'Large Datasets and Big Data Analytics: Theory, Methods, and Applications' at IJCNN 2017 | Anchorage, Alaska, USA |
| ☆ | ⬒ | 8. | ACSTY 2016 | Second International Conference on Advances in Computer Science and Information Technology | Chennai, India |
| ☆ | ⬒ | 9. | DMAP 2016 | Second International Conference on Data Mining and Applications | Vienna, Austria |
| ☆ | ⬒ | 10. | EAST - FLAIRS 2017 | lEArning from heterogeneouS data analyTics - FLAIRS | Marco Island, Florida, USA |
| ☆ | ⬒ | 11. | FLAIRS 2017 | FLAIRS-30: Special Track in Data Mining | Marco Island, Florida, USA |
| ☆ | ⬒ | 12. | Sideways 2017 | Social Media World Sensors 2017 | Madeira |
| ☆ | ⬒ | 13. | ICISDM 2017 | 2017 International Conference on Information System and Data Mining (ICISDM 2017)-SCOPUS, Ei Compendex | South Carolina, USA |
| ☆ | ⬒ | 14. | ITMLS 2017 | Intelligent Technologies and Methodologies of Learning Systems | Barcelona, Spain |
| ☆ | ⬒ | 15. | ITPF 2016 | The 5th international conference on Information Technology, Present and Future | Mashhad, Iran |
| ☆ | ⬒ | 16. | ICCSP 2017 | International Conference on Cryptography, Security and Privacy - Ei Compendex and Scopus | Wuhan, China |
| ☆ | ⬒ | 17. | ICKSE 2017 | 3rd International Conference on Knowledge and Software Engineering-Ei Compendex, Scopus & ISI CPCS | Paris, France |
| ☆ | ⬒ | 18. | WSDM Cup 2017 | WSDM Cup 2017: Knowledge Base Quality and Search | Cambridge, MA, USA |
| ☆ | ⬒ | 19. | KDD 2017 | Knowledge Discovery and Data Mining | Halifax, Nova Scotia, Canada |
| ☆ | ⬒ | 20. | ICOK 2017 | 2017 2nd International Conference on Knowledge (ICOK 2017) | Chengdu, Sichuan, China |
| ☆ | ⬒ | 21. | ICBDA 2017 | The 2017 IEEE International Conference on Big Data Analysis (ICBDA 2017) - Ei Compendex | Beijing, China |
| ☆ | ⬒ | 22. | IEA/AIE NAAD 2017 | Special Track on Novel Approaches to Anomaly Detection | Arras, France |

Figure 2: Data upload in OpenRefine

In order to clean the following steps are followed.

- **Trim white space**: At first, from all the columns the leading and trailing white spaces are removed. For this **Edit cells→Common transforms→Trim leading and trailing spaces**, steps are followed as shown in Figure 3.

- **Cleaning using Conference Location**: In order to clean the data, at first the the **Conference Location** column is selected. There were a number of inconsistencies in this column. In order to maintain consistency and make all the conference location to follow a specific format without loosing important information, **City name**, **Country** format is followed in general for conference location. While cleaning the data, the following scenarios are considered,

Figure 3: Trimming white spaces from columns

1. **N/A**: There are some conference entries in which the location is **N/A** or **N.A.** according to Figure 4. At first, these locations are considered as **Location Not Available**.



Figure 4: Removal of Location N/A to Location Not Available

2. **Online**: There are also some conference entries with location **Online** or **online** as Figure 5. Initially, we have also converted them to **Location Not Available**.

3. **Journal Name**: There are some conference entries which are actually call for journal article. In this case, the journal name is included as location. For example, in Figure 6, the location is shown as **Studies in Computational Intelligence**. These entries are also replaced by **Location Not Available**.

After observing the above 3 situations, there were total 110 entries with **Location Not**

Figure 5: Conferences with Location Online



Figure 6: Conferences with Journal name as location

**Available**. Since the location information is important, these entries are removed. The removal operation is shown in Figure 7.

After this operation, the clusters among the conference location data is determined by the Cluster option of OpenRefine. With proper observation, the entries were merged. The screen shots of the cluster options are given on Figure 8.

Then, there are some conference entries for which the location format was still inconsistent. For example, for some conference the location is **San Fransisco, CA, USA** and for some conference the location is **San Fransisco, California, USA**. Therefore, such entries are found out and converted to a single format. For example, for this case I have converted it to **San Fransisco, CA, USA**. In order to find such entries, from the **Conference Location column** →**Facet** →**Text Facet**. From the facet, the repetitions are found according to Figure 9.

After cleaning the data, using conference location column, I found that, there are some confer-
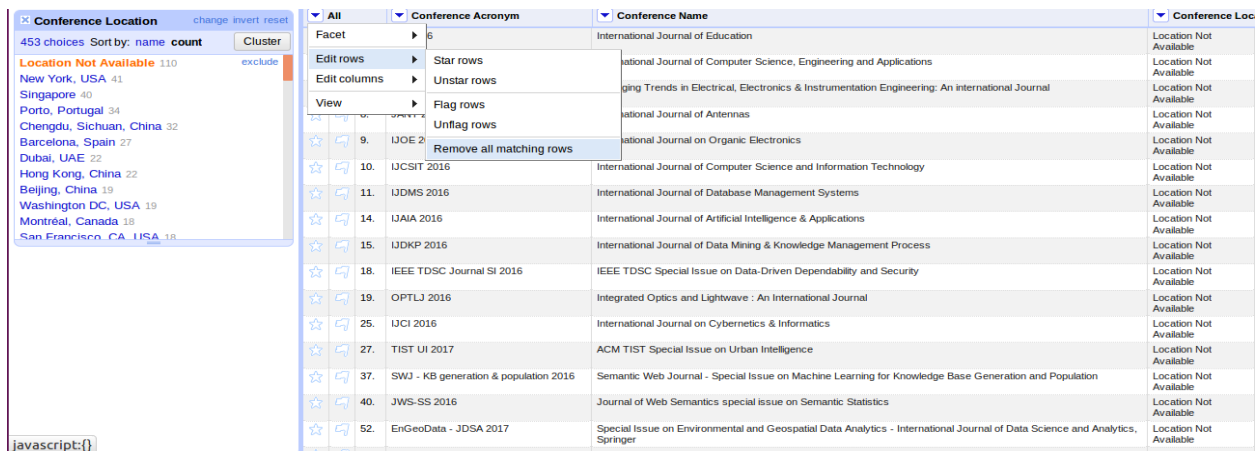
Figure 7: Removal of entries with **Location Not Available**



Figure 8: Cluster Options



Figure 9: Inconsistency in the name of locations

ence entries that are repeated multiple times. The reason behind this is, there are some conference which are included as Data mining and Artificial Intelligence categories or more. Therefore, in order to get a clear view of the data, I removed the repeated entries.

In order to accomplish this, at first I tried to merge duplicate rows according to same Conference Acronym. But in that case, some other conference entries might be removed. For example, Figure 10 is showing that conferences with Acronym **AI 2016** has 6 entries but only 2 different conferences. In that case, I removed the repeated 4 conferences using star. The operation is, **All →Facet →Facet by star**. From **Starred Rows** select **true**. Then only the starred rows will appear. After that **All →Edit rows →Remove all matching rows**.



Figure 10: Problem of removing repeated entries using conference acronym

After removing the repeated entries it becomes like Figure 11.



Figure 11: Scenario after removing repeated entries

The problem of conference acronym also occurs with conference name. Therefore, after observing the clusters manually, I clustered the same conferences together as shown in Figure 12.

Figure 12: Clustering according to conference name

There are some conference names in which the time and location of the conference are appended, I removed them as shown in Figure 13.



Figure 13: Correcting the conference name

Though, a thorough cleaning operation is performed on the Conference Location column, there are still some inconsistencies in the data. For example, same conference is entered once with English name. Then again with Spanish name as shown in Figure 14.

After removing repeated data and all inconsistencies now I have **1238** conference entries from initially collected **1600** conference entries. The number of location is reduced from **682** to **442**. In order to perform the data cleaning operation efficiently, I followed the videos of OpenRefine[3]. From this portion, the project is exported as **WikiCFP_CleanData.tsv** file.

Figure 14: Repeated conference with the name written in different languages

# 3. Hadoop

In order to work with hadoop, at first I have followed all the steps of **Apache HADOOP**[4][5] tutorials. After running the **WordCount.java**, I have started working on the problems. For all the problems, in mapper portion, **StringTokenizer** is used to tokenize the string. Here as delimiter new line character is used, since in the **WikiCFP_CleanData.txt** file all the entires of conferences are in each new line. Then each new line is split using tab character as a delimiter. In each of the mapper function, I have checked, if current line's location is "Conference Location" or not. In order to prevent the column heading to be mistakenly considered by the program, I have deleted such entries.

- *Compute and plot the number of conferences per city. Which are the top 10 locations?*
  In this problem, the **Conference Location** is considered as **key**. For each location, **1** is stored as **value**. Here key is a Text element and value is a IntWritable element. In the reducer part, the values are calculated which is actually the frequency of locations and written with the location name. After the proper execution of the code, the output file is generated. The file is copied from the distributed file system to the local system. The data of this file is used to plot the graph through excel. The graph for the number of conferences per city is shown in Figure 15.

  The top 10 locations in which the highest number of conferences are taking place are shown in Figure 16. Figure 17 is showing the result of the hadoop program from terminal.

- *Output the list of conferences per city?*
  In this problem, the mapper has considered the **Conference Location** as **key** and the **Conference Name** as **value**. Therefore, both key and value are Text in this case. In the reducer, as **Iterable**, Text is considered which is the conference name. In each iteration, the conference names for a particular location are concatenated as a single string. Then the reducer writes the conference location and the conference names of that location (as a single string) in the output file. In order to solve this problem, I have used Text iterable and got

Figure 15: Plot showing number of conferences per city



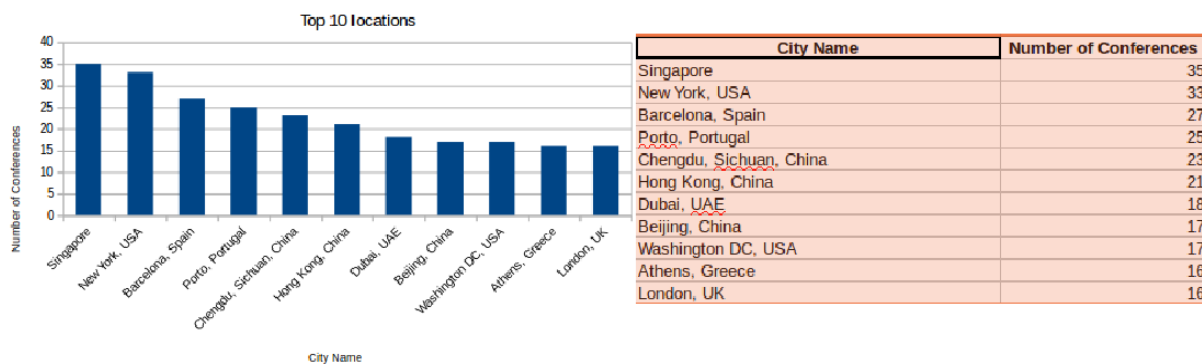| City Name | Number of Conferences |
|---|---|
| Singapore | 35 |
| New York, USA | 33 |
| Barcelona, Spain | 27 |
| Porto, Portugal | 25 |
| Chengdu, Sichuan, China | 23 |
| Hong Kong, China | 21 |
| Dubai, UAE | 18 |
| Beijing, China | 17 |
| Washington DC, USA | 17 |
| Athens, Greece | 16 |
| London, UK | 16 |

Figure 16: Plot showing top 10 cities with highest number of conferences
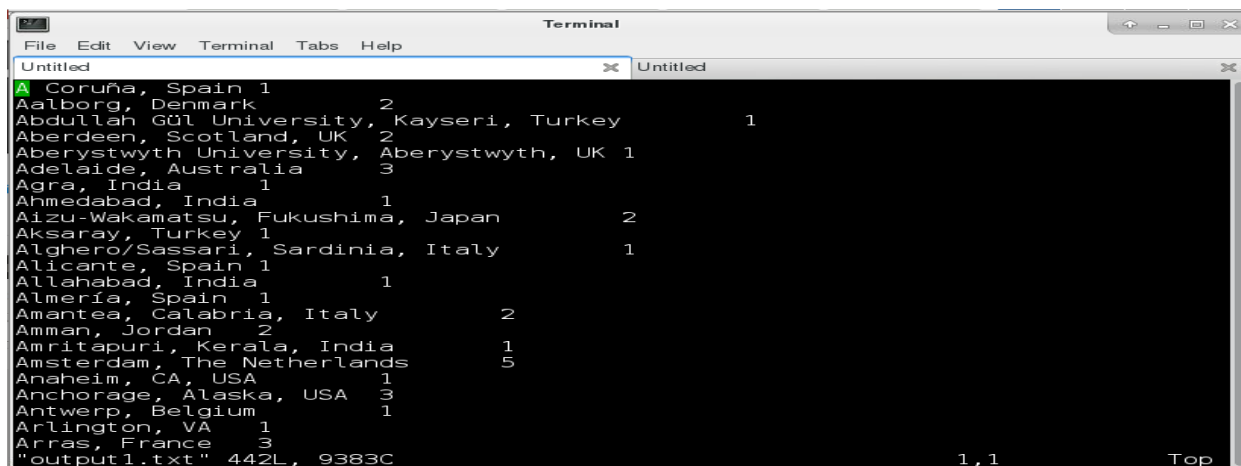


Figure 17: Result of Problem 1

the initial idea about it from CONVERGE Blog[6]. Figure 18 is showing the result generated by the hadoop program in terminal.



Figure 18: Result of Problem 2

- **For each conference regardless of the year (e.g., KDD), output the list of cities.** In this problem, at first the conference acronym is separated from the year. Then only **Conference Acronym** is emitted as **key** and the **Conference Location** is emitted as **value** by the mapper. In the reducer part, the conference location is again considered as text iterable. For a particular conference/conference acronym, the locations are concatenated in one string which is then emitted with the conference acronym as output. Figure 19 is showing the result generated by the hadoop program in terminal.



Figure 19: Result of Problem 3

- **For each city compute and plot a time series of number of conferences per year.** In order to solve this part of the problem, in the mapper part, the Year information is collected from the conference acronym. While separating the year information from the acronym, I considered the last 4 characters of it. However, I find that there are some discrepancies.

For example, for conference acronym like "ACTION15", considering the last 4 digits will give "ON15". In order to remove this problem, I have checked the first 2 letters. If they are not numbers, then I replaced them by "20", assuming that the record must be showing information for conference that is going to take place by year $2000 - 2099$. Then this Year information is concatenated with conference location information. This merged information of **Conference location and Year** is emitted as **key**. Then for each such key a value of **1** is emitted as **value** from the mapper. Then in the reducer part, for each conference location and year key, the frequency is calculated and emitted as output. Here as iterable the frequency is considered as **IntWritable**. Figure 20 is showing the result generated by the hadoop program in terminal.
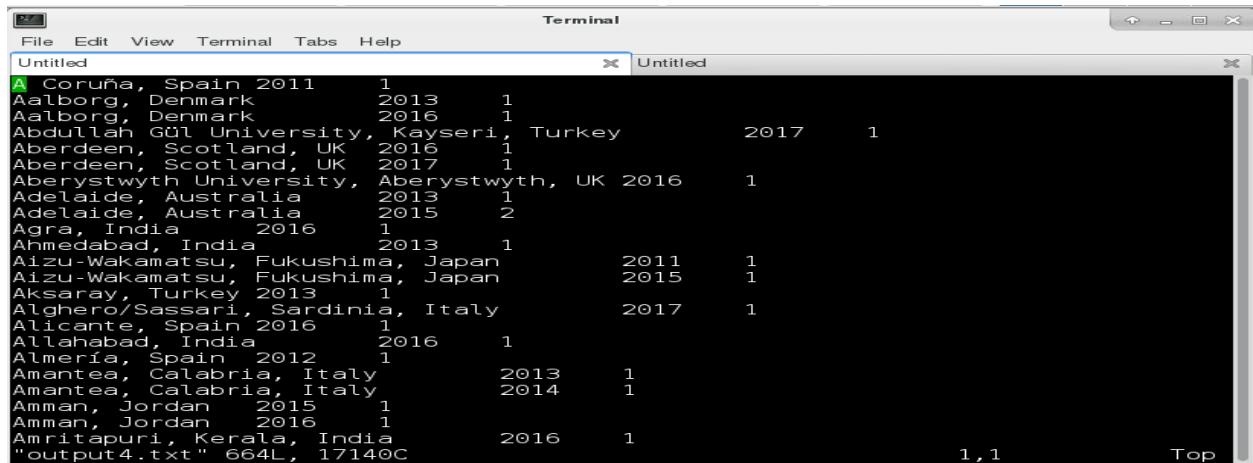


Figure 20: Result of Problem 4

For this problem, I have plotted the top 11 cities. The plots for the top 11 cities are shown in Figure 21 - 23.
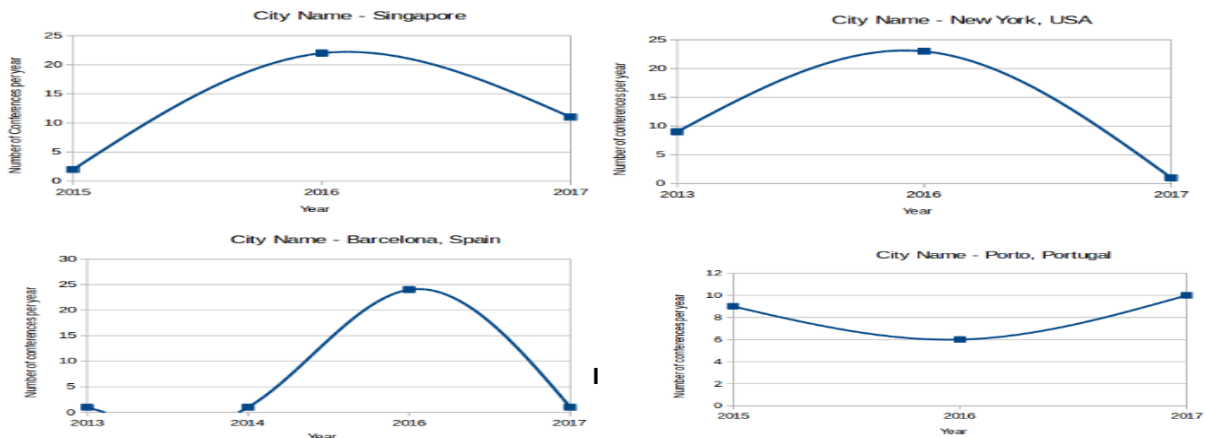


Figure 21: Plot for Singapore with rank 1, New York, USA with rank 2, Barcelona, Spain with rank 3 and Porto, Portugal with rank 4
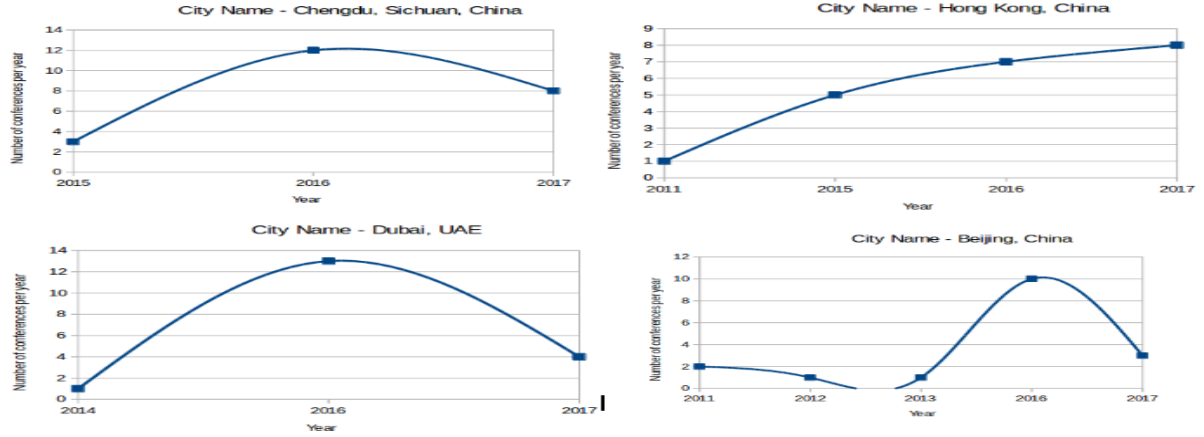
Figure 22: Plot for Chengdu, Sichuan, China with rank 5, Hong Kong, China with rank 6, Dubai, UAE with rank 7 and Beijing, China with rank 8
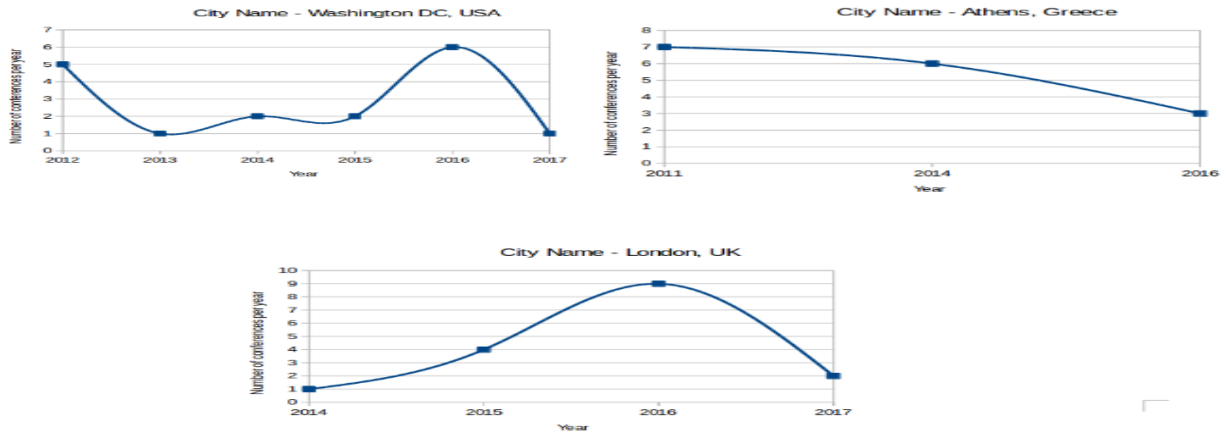


Figure 23: Plot for Washington DC, USA with rank 9, Athens, Greece with rank 10 and London, UK with rank 11

## Commands

The solution code for each problem is compiled and run using the commands in Table 1.

Table 1: Commands for the Hadoop program

| Task | Command |
| --- | --- |
| Copy .java file | /extra/netid/hadoop/hadoop-2.7.3/bin |
| Code compilation | hadoop com.sun.tools.javac.Main Code_File_Name.java |
| Jar creation | jar cf Code_File_Name.jar Code_File_Name*.class |
| Run application | hadoop jar Code_File_Name.jar Code_File_Name input output |
| Copy output file from HDFS | hadoop fs -copyToLocal output/part-r-00000 /extra/netid |
| Copy output file to local system | scp netid@machineid:/extra/netid/part-r-00000 /home/csgrads/netid |

# References

[1] StackOverflow, "How to parse HTML table using jsoup?" Last accessed: November 16, 2016. [Online]. Available: http://stackoverflow.com/questions/24772828/how-to-parse-html-table-using-jsoup

[2] Net Instructions, "How to make a simple web crawler in Java," Last accessed: November 16, 2016. [Online]. Available: http://www.netinstructions.com/how-to-make-a-simple-web-crawler-in-java/

[3] Google, "OpenRefine," Last accessed: November 16, 2016. [Online]. Available: http://openrefine.org/

[4] Apache Hadoop - The Apache Software Foundation, "Hadoop: Setting up a Single Node Cluster," Last accessed: November 16, 2016. [Online]. Available: https://hadoop.apache.org/docs/r2.7.0/hadoop-project-dist/hadoop-common/SingleCluster.html

[5] Apache Hadoop, "MapReduce Tutorial," Last accessed: November 16, 2016. [Online]. Available: https://hadoop.apache.org/docs/r2.7.0/hadoop-mapreduce-client/hadoop-mapreduce-client-core/MapReduceTutorial.html

[6] CONVERGE Blog, "How to Write a MapReduce Program," Last accessed: November 16, 2016. [Online]. Available: https://www.mapr.com/blog/how-write-mapreduce-program