# An Analysis of Summer Olympics Using Open Source Big Data Technologies

Shikhar Srivastava

*x18106960*

*Abstract*—Game events which occur on a global scale are important events in which players from all over the globe participate. Players participate in many different sports which are scheduled to happen over the course of many days. Olympics is one of the major sport events which occur after every four years in which many athletes participate from all over the world. This mega game event is further bifurcated in two Olympic events i.e. Summer Olympics and Winter Olympics with just a single difference that all the winter sports like skiing are scheduled to happen in Winter Olympics. It would be very interesting to analyze data recorded for this majestic game event and look at some insights. This project uses open source big data technologies like MapReduce, Sqoop, HBase, Hive, Pig and HDFS and tries to extract few useful and insightful information. This project also follows an approach that how these open source big data technologies can be efficiently used to store and process data which in our case is Olympic dataset.

*Index Terms*—MapReduce, MySQL, Sqoop, HBase, Hive, Pig, Hadoop, HDFS, Olympics data, R, Tableau

## I. INTRODUCTION

The world is witnessing an era of data explosion, an era where it is a known fact that data is a very important asset. As the volume of data is increasing exponentially, it became a need and necessity to be able to store and process this huge volume of data efficiently. With this explosion in data, new approaches emerged to store it, where people started to store huge volume data on distributed platforms rather than storing it in traditional relational databases. Many technologies also started to emerge which could manage data on a distributed platform. This project utilizes some of these popularly known open source technologies to analyze Olympic dataset and aims to showcase how these technologies can be leveraged to store and process huge datasets.

International sport events occur at a very large scale involving a large number of participants from many different countries. These participants compete in several events across many sports. Olympic games are one of the most popularly known sport events which takes place after every four years. Summer and Winter Olympics are the two types of Olympic events which are held where Winter Olympics are especially organized for winter sports. This project utilizes Olympic dataset which is downloaded from Kaggle [1] which contains player wise details like in which sport the player participated, in which Olympic, in which sport, nationality of the player, gender and few more details. Columns and their brief descrip-

tion is given in the appendix section of this document. The Olympic dataset is fed in different databases like MySQL, HBase, Hive, Pig and is then further analyzed by building queries. The objective of this project is to feed Olympic dataset to above mentioned databases and explore data to find some key insights, which are explained in the methodology section.

The rest of the document is structured in four different sections: Related Work, a concise literature review on relevant research papers; Methodology, explains all the steps followed for implementation of this project; Results, shows and explains results from the implementation.

## II. RELATED WORK

Olympics are ancient sport events which roughly started in around 5th century B.C. [2] when intellectuals known as Sophists travelled to Olympia to demonstrate combative oratory. This reflects that Olympics have rich history which dates back many hundreds of years. Major role in reviving the ancient Olympic event was played by Pierre de Coubertin [3] who provided a launch pad for modern Olympics. As the roots of Olympics come Greece, it faced negligence when the 1896 Olympic games were about to start. His name was no where mentioned and the greeks accused him of being a thief to steal their inheritance. Olympics are not just bound to happen in Greece, but now many other countries host Olympic events every four years. It greatly impacts the economic situation of host countries. A research has also been done to examine and evaluate methods and assumptions used by economic studies on Olympic host countries. Olympic games not just impact economy of host countries but they also have a social impact on which a research has been done for based on Sydney Olympics [4].

This project tries to analyze Olympic data to look at some insightful information. There are many open source big data technologies available to store, manage and process large scale data like MySQL, HBase, Hive, Pig etc. There have been many researches on these distributed platforms. One such research evaluated the performance of HBase on random reads and writes and further does a comparative study of HBase on HDFS and MySQL on HDFS. Similarly, Hive is also a famous platform for big data processing on top of HDFS. [5] provides a descriptive explanation of Hive architecture and why is it suitable to be used for large scale data on HDFS. Hive

also comes with almost SQL like querying language which is called HiveQL. Pig is another popular big data technology which is widely used to store and manage high volume of data. A new querying language was designed to operate on Pig by a group of researchers [6] and was called as Pig Latin which has the declarative style of SQL and procedural style of MapReduce. Storing, managing and processing of large data can be efficiently handled by above mentioned big data technologies, but what if there is a requirement to move data from RDBMS to hadoop ecosystem. Sqoop is very a robust framework which is widely used to transfer data from relational database like MySQL to HDFS. Sqoop works on JDBC drivers to transfer data into hadoop [7] [8]

## III. METHODOLOGY

This section provides a detailed explanation of data used in this project along with steps followed to prepare, store and analyze Olympic data.

### A. Data Acquisition and Preprocessing

Dataset used in this project has been downloaded from Kaggle website [1] which is a public dataset, so there are no ethical issues associated with dataset used. Although on above mentioned link there are two datasets in the form of CSV files which includes athelete_events.csv and noc_regions.csv, but for this project we are only using athelete_events.csv. CSV file of size 39.58 MB containing all player wise Olympic data was downloaded.

The downloaded dataset was in a pretty good condition and didn't need any cleaning as such. The only thing which was carried out as part of pre-processing was to replace NULL values in Medal column with 'No Medal' and remove all NAs from the dataset. This two step re-processing task was done in R by making use of tidyr library of R, (2_Step_Pre_Processing.R) is present in the compressed folder attached. The output file from this code was then copied and pasted in shared folder to use it in Ubuntu VM.

### B. Step by Step Process Followed

This section details all the steps followed to store and process Olympic data in different databases.

**Step 1:** File copied to shared folder was first again copied to hduser just to make sure that there are no conflicts with permissions. After this, MySQL was launched in which a new database was created named 'PDA_Project' in which a new table 'olympics' was created and data from the Olympic csv file was then loaded in this table. Screenshot in Fig. 1.

**Step 2:** 'olympics' table data stored in MySQL was then stored on HDFS to execute MapReduce shell script on this dataset.
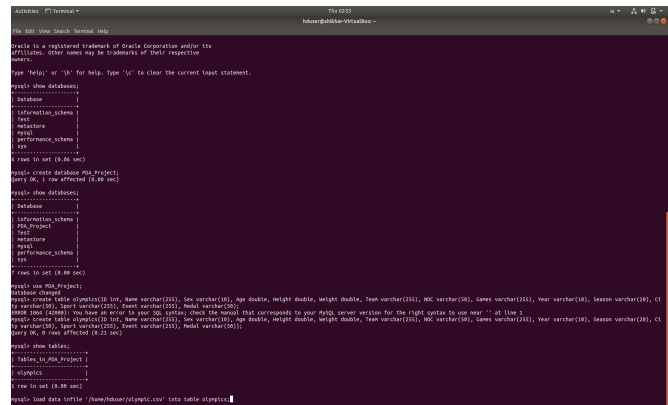


Fig. 1. Loading Data in MySQL



Fig. 2. Sqoop Output a

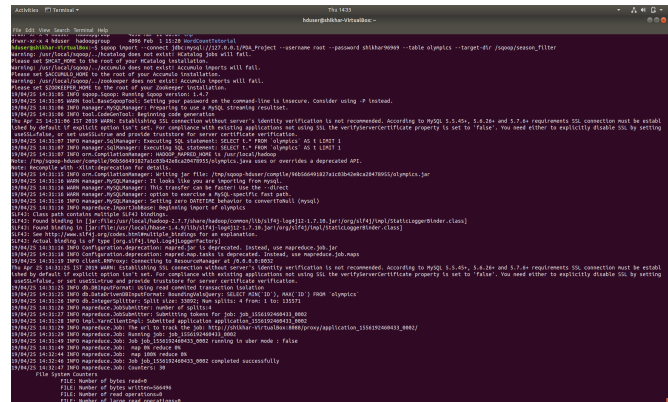MySQL table data was copied to HDFS with the help of Sqoop which stored MySQL table data in sqoop/season_filter directory in HDFS. Fig. 2 and Fig. 3 shows the successful execution of sqoop command.

**Step 3:** This step involves filtering Olympic data for just summer Olympics in Season column. For this a MapReduce java program was written in eclipse and jar file was exported to
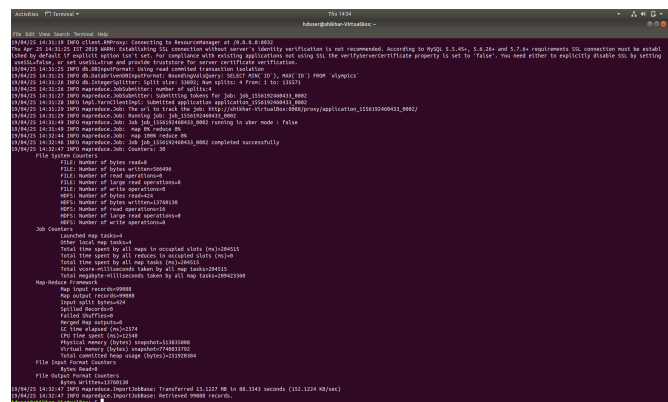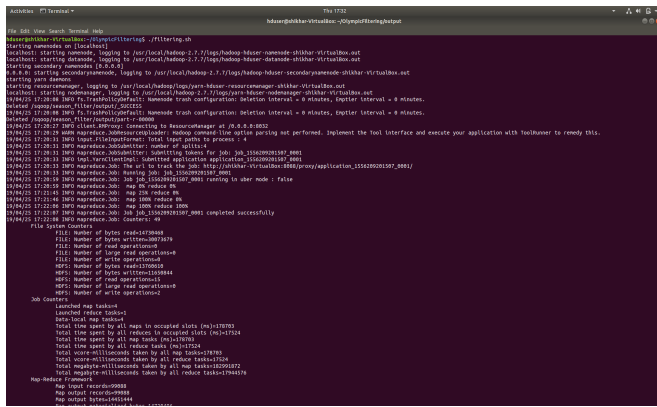


Fig. 3. Sqoop Output b
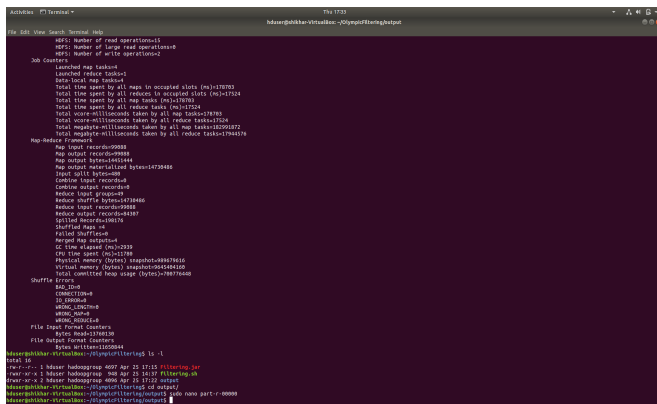
Fig. 4. MapReduce Shell Script a



Fig. 6. MapReduce Output



Fig. 5. MapReduce Shell Script b



Fig. 7. Processed MapReduce Output

hduser in 'OlympicFiltering' directory. A shell script was then written to execute MapReduce jar file on Olympic data which was sqooped from MySQL to HDFS. Jar file and the shell script written are present in the compressed folder uploaded. Fig. 4 and Fig. 5 show the execution of MapReduce task on HDFS Olympic data.

**Step 4:** Output from MapReduce task is shown in Fig. 6. The output data from MapReduce task was not looking good because of the space before each ID and the headers were also not present. The space before IDs were because a blank field was included in the Reducer class while writing the output after filtering. Due to this, output file from MapReduce task was then processed to make it look like a neat output. Fig. 7 shows the neat file after processing output from MapReduce.

**Step 5:** This Step involves putting organized MapReduce output data in HDFS to further load it in HBase. Script used to load MapReduce output data in HDFS is shown in Fig. 8. A separate directory was created named 'HBASE' to store MapReduce output data in HDFS.

**Step 6:** After putting MapReduce output data in HDFS, HBase was then started and launched. A new table was created in HBase to store the output from MapReduce. Fig. 9 shows

screenshot for creating a new table in HBase and Fig. 10 then shows data loading in HBase. The HBase table data is also shown in Fig. 11

**Step 7:** Now, with an aim to pull MapReduce output data in Hive from MySQL using Sqoop for further analysis, a new MySQL table was created named 'summer_olympics'. Fig. 12 shows the creation of a new table and some sample records from that table. This table holds data for only summer Olympics as data put into this table is coming from MapRe-



Fig. 8. MapReduce Output to HDFS for HBASE



Fig. 9. Table Created in HBase

Fig. 10.  Importing MapReduce Output Data in HBASE



Fig. 11.  HBase Table Data

duce output.

**Step 8:** This involved importing data from newly created summer Olympics table from MySQL to Hive. But, the command was giving unusual errors when the script was executed. Fig. 13 show the first command which threw error when trying to sqoop data from MySQL to Hive. This command was supposed to create a new table in Hive and then put MySQL
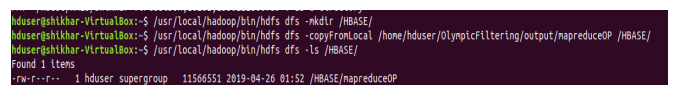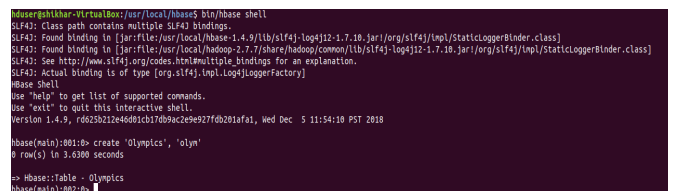


Fig. 12.  Summer Olympics Table Created in MySQL



Fig. 13.  First Sqoop Error Command



Fig. 14.  First Sqoop Command Error

data in that table.

It gave an error in executeScript while importing data from MySQL as shown in Fig. 14. With an assumption that error is caused due to primary key which was present in MySQL table being imported, a new table was created in Hive manually and then Sqooping was tried again, but this time create hive table was removed from the script. This also didn't work out and this command gave new error of HDFS directory already exists as shown in Fig. 15.

The HDFS directory was then manually deleted as shown in Fig. 16 and command was executed again. Unfortunately, this also didn't work out and gave same error of executeScript failure.

**Step 9:** Due to the errors faced while sqooping data from MySQL to Hive, data was manually inserted in Hive table (shown in Fig. 17) and queries were then executed on this table.

**Step 10:** In this step four queries were written to analyze data loaded in Hive. These queries and its output will be discussed in detail in Results section.
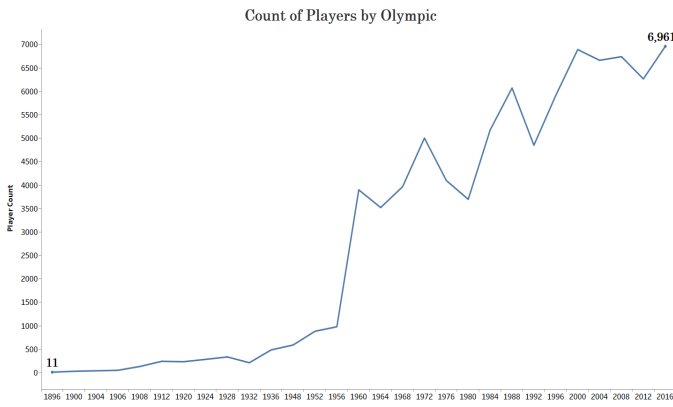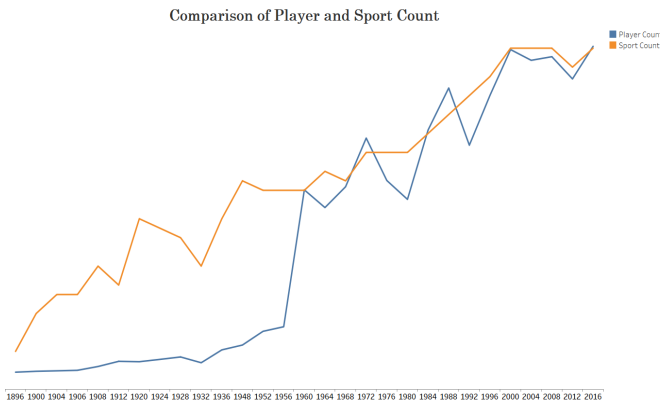


Fig. 15.  Second Sqoop Trial

Fig. 16. Manually Deleting Target Directory



Fig. 18. Count of Sports by Year and City



Fig. 17. Inserting Data in Hive Table

**Step 11:** In this last step, three .pig files were created in which pig latin queries were written and executed on MapReduce output dataset which was pulled in these .pig files. The output from all .pig files was stored in three different folders. All three .pig files and output files are present in the compressed folder uploaded.

## IV. EVALUATION AND RESULTS

In this section, output from all the tasks performed are discussed in detail.

### A. R Processing - MySQL Data Loading - MapReduce

After downloading data from Kagggle, a short code was written in R to clean it. The R Code is used to perform pre-processing task is present in the compressed folder. After a little bit of pre-processing, the data is then put in shared folder to use it in Ubuntu VM. In Ubuntu Vm, olympic data file is further moved from the shared directory to hduser and then moved to HDFS using '-copyFromLocal' command. MapReduce task was built in eclipse for filtering Olympic dataset for just summer Olympics. This MapReduce task was

built to filter the Season column in dataset for just summer (file present in compressed folder) and this dataset will further be fed to other databases for analyzing. After filtering data was left with 84307 records in 15 columns holding details of summer Olympics from year 1896 to 2016. This data was further organized a little bit as the MapReduce output was having leading spaces and header was also missing.

### B. HBase

HBase in this project is only used to store the MapReduce output in it. To store MapReduce output in HBase, data was first moved to HDFS and then pulled in HBase. No other task was performed in HBase.

### C. Hive

There were few errors while sqooping data from MySQL to Hive (as explained in Methodology section), due to which sqoop command was not working and a table was manually created in Hive in which data was manually inserted. Four queries were written in Hive to analyze Olympic data.

*1) Query 1:* First query which was written in Hive gives an output of total number of sports for each summer Olympic held. The output from this query contained fields like year, city and count of sports which was grouped by city and year. Screenshot has already been attached in Methodology section. Out files is present in compressed folder. Fig. 18 represents the output from this query where we can clearly see that the number of sports increased through the year and in the last summer Olympic held in year 2016, events related to 34 sports were organized. The color in charts reflect the number of sports and higher the number of sports, darker is the color. Similar color logic has been followed for top chart as well.

*2) Query 2:* In second Hive query, count of players who participated in Olympoic games from year 1896 to 2016 are shown. In the Fig. 19 we can clearly get an idea that the

Fig. 19. Count of Players by Year



Fig. 20. Sport Count vs. Player Count



Fig. 21. Last 4 Summer Olympics



Fig. 22. Number of Events by Olympic

number of players participating in summer Olympic games has increased since its inception. Also, we can see that there had been a sudden increase in the number of players participating from year 1956 to 1960. The chart also shows that first Olympic game started with just 11 players.

*3) Query 3:* The third query shows a comparison of count of players and count of sports. This output was achieved by joining output from first and second query in Hive. In Fig. 20 we can see that there is a steady increase in both numbers through the years till 2016. However, the interesting thing to notice here is that although there is a sudden increase in number of players from year 1956 to 1960, the number of sports organized in these years remains the same.

*4) Query 4:* Fourth query written in Hive filters data just for last four summer Olympics i.e. summer Olympics held in the year 2016, 2012, 2008 and 2004. A gender wise and country wise analysis was carried out on top of this filtered dataset. Fig. 21 shows an exploratory analysis based on gender and country for last four summer Olympics. Major observations after bulding charts on this dataset are: The number of medals won by females are less than number of medals won by males in each of the four years and it is observed that USA from past four years tops the chart bagging maximum number of
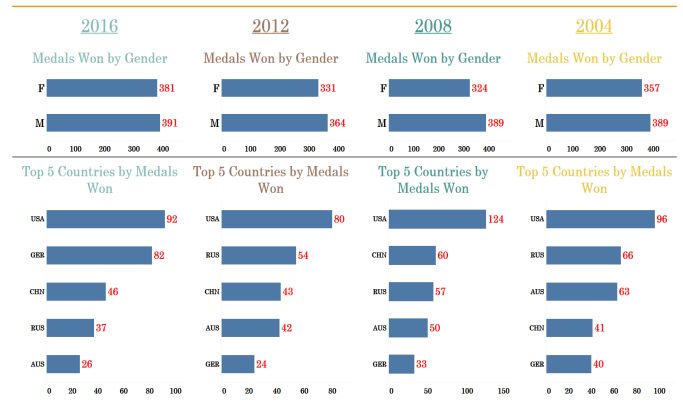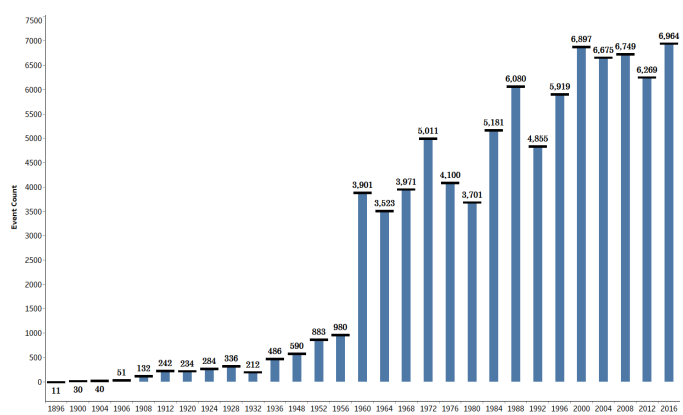
medals followed by Russia, China, Germany and Australia.

*D. Pig*

At first Pig was launched in distributed mode data was pulled into Pig from HDFS, but then it was noticed that when while pulling data from HDFS, Pig was trying to connect to server and every time the script was executed it iterated 10 times to try to connect to server which in this case is not possible. So, rather than launching pig in distributed mode, it was launched in local mode and MapReduce output file was then pulled form local system. A total of 3 queries were written in 3 different .pig files and similarly output was stored in 3 different folders.

*1) Query 1:* First query written in Pig was to get an idea about the number of total events which tool place in each Olympic across multiple sports. Here events for an example refer to different types of races like 100m, 400m etc. under sport running. Like in one of the Hive queries, here also we can see that after year 1956, there had been a considerable rise in number of events with last Olympic in year 2016 witnessing the highest ever number of events. Refer Fig. 22
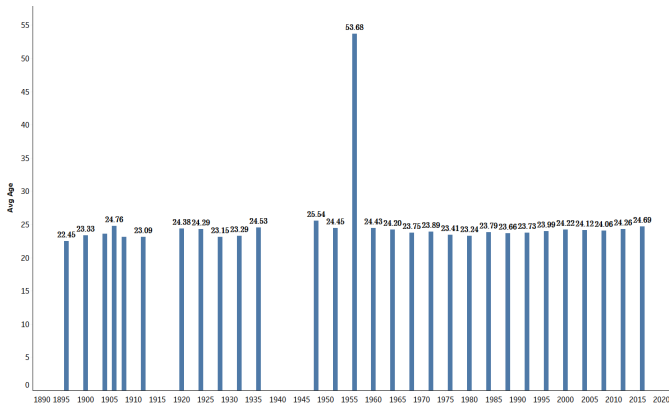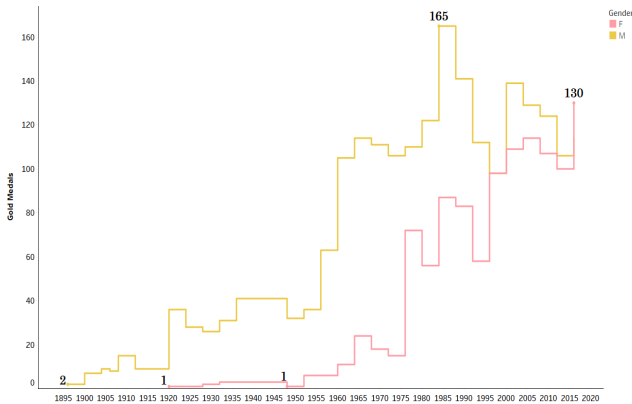
Fig. 23. Avg. Age by Olympic



Fig. 24. Avg. Age by Olympic

*2) Query 2:* Second query checks the average age for each Olympic. Surprisingly, average age for year 1956 Olympic is very high as compared to other Olympics (Fig. 23). For other Olympics, average age is almost the same.

*3) Query 3:* The third query digs into data to find out number of gold medals by gender. In Fig. 24 we can clearly see two patterns, one showing that there is an increase in medals won by males in females both over the period of different Olympics starting from year 1896 to 2016. Second pattern which can be observed is that the medals won by males in almost every Olympic is more than that of females. In year 1960, there was a considerable increase in number of medals won by makes whereas, year 1976 witnessed the same trend for females.

## V. Conclusion and Future Work

As the world is moving towards an era where data will continue to grow exponentially, it has become mandatory to be able to and be prepared to store, manage and process this huge volume of data. Storing data on distributed platforms rather on traditional relational databases is the new approach which many organizations are currently following. This project leverages the power of some popularly known Big Data open source technologies and systematically tries to understand their use whilst getting to know some common errors which could be faced while implementing these technologies for handling data.

Olympic dataset which was downloaded from Kaggle was used for this project which was first pre-processed a little but before using it for data analysis in Ubuntu. This dataset was first loaded into MySQL and then sqooped to HDFS after which a MapReduce filtering job was applied on it. Output from MapReduce was organized and was stored in a directory. This output dataset was further pulled in HBase where it was just stored. Later, MapReduce output data was tried to be pulled in Hive using Sqoop but faced some errors due to which it was then manually pulled in Hive after creating a table, following which four Hive queries were built on this data. Lastly, MapReduce output data was then pulled in Pig and three queries were built on it.

Major observations while building charts on query output from Hive and Pig are: All charts are following an increasing trend which particularly means that every Olympic starting from year 1896 saw an increase in number of participants, number of events, number of sports and number of medals won; for most of the charts year 1960 saw great jump in above mentioned numbers somehow; as far as last four Olympics are concerned, USA tops the chart in bagging most number of medals followed by Russia, China, Germany and Australia. Future work proposed would be to analyse data using HBase, building Pig queries for data present in HDFS rather than using Pig in local mode and using sqoop in all cases for importing data from MySQL into Hive or Pig.

## References

[1] 120 years of Olympic history: athletes and results. [Online]. Available: https://kaggle.com/heesoo37/120-years-of-olympic-history-athletes-and-results. [Accessed: 28-Apr-2019].

[2] N. Spivey and N. J. Spivey, The Ancient Olympics. Oxford University Press, 2005.

[3] A. Guttmann, The Olympics: A History of the Modern Games. University of Illinois Press, 2002.

[4] G. Waitt, Social impacts of the Sydney Olympics, Annals of Tourism Research, vol. 30, no. 1, pp. 194215, Jan. 2003.

[5] A. Thusoo et al., Hive: A Warehousing Solution over a Map-reduce Framework, Proc. VLDB Endow., vol. 2, no. 2, pp. 16261629, Aug. 2009.

[6] C. Olston, B. Reed, U. Srivastava, R. Kumar, and A. Tomkins, Pig Latin: A Not-so-foreign Language for Data Processing, in Proceedings of the 2008 ACM SIGMOD International Conference on Management of Data, New York, NY, USA, 2008, pp. 10991110.

[7] D. Vohra, Using Apache Sqoop, in Pro Docker, D. Vohra, Ed. Berkeley, CA: Apress, 2016, pp. 151183.

[8] M. S. S. Aravinth and A. H. Begam, An Efficient HADOOP Frameworks SQOOP and Ambari for Big Data Processing, vol. 1, no. 10, p. 4.