

# **Problem statement:**

Analyzing the data and performing exploratory data analysis on the data set related to State University of New York at Buffalo. This data set cons0ists of data for related to the number of students registered, courses offered, time scheduling of classes in various hall rooms in UB campus and their maximum capacities.

### Goal:

Our main goal is to understand this data and come up with a list of relevant questions to realize important trends in this given data set. Also we came up with MR algorithms to extract "knowledge" from this data and provide answers to the set of questions designed.

### **Data Characteristics:**

We referred to bina\_classschedule.csv file provided by Professor Bina Ramamurthy. We worked on semi structured data in CSV format.

We designed a set of **9 complex questions each with 2 or more mappers** and came up with map reduce solutions to these problems in an attempt to give optimized solutions to these questions. This would help us comprehend the extensive data set in a constructive way which might reflect some important observations, thus helping to make better decisions in the future. This would surely help in the effective utilization of the classroom based on the courses offered in them, students registered etc.

Here is a brief explanation of each questions pertaining to its implementation details and conclusions derived and their usefulness.

**PROBLEM 1:** To find top 20 courses for which maximum number of students were registered.

# **Implementation details:**

### **Data Cleaning:**

We cleaned and filtered the field Students\_registered (field 7). We ignored those tuples where there were zero number of students registered or this field had a non -numeric or invalid entry. Also if the corresponding field of Hall Room has an invalid entry like "Unknown" and "Arr" (for halls).

### Map Reduce job: 3 MR Jobs

Job 1: Found the number of students registered for each course. Key: Name of course value: Number of students registered

#### Job 2: To sort the records on values

In Mapper 2 we swap the Key, value pairs obtained from stage 2 and in Reducer 2 display the swapped key value pair i.e. key: Number of students registered Value: Course Name. This stage will output the records in sorted order based on the initial values as Hadoop sorts the keys by default.

#### Job 3: To display top 20 records

In Mapper 3 we set the Key as "1" so that in Reducer we can put all the values in one Array. Accessing the first element of an array will give us the maximum value i.e. course for which the maximum number of students are registered. Thus we have accessed the first 20 records to display the top 20 courses for which maximum courses are registered.

### **Output:**

```
File Output Format Counters
 adoop@hadoop-VirtualBox:~/hadoop$ hdfs dfs -cat plo1/*
16/04/15 16:01:05 WARN util.NativeCodeLoader: Unable to load native-hadoop library for your platform... using builtin-java classes where applicable
154409 General Chemistry
139710 World Civilization 1
        World Civilization 2
103194 Writing 1
70018 Introductory Psychology
53553 College Physics
52191 Intro to Macroeconomics
        Intro to Microeconomics
 6938
         Special Topics
        Surv Calculus & Appl 1
        College Calculus 2
Psychological Statistics
         Evolutionary Biology
Intro to Financial Accounting
        Organic Chemistry I
Cell Biology
Human Nutrition
         Gen Chem for Engineers
        hadoop-VirtualBox:~/hadoop$
```

### **Conclusion:**

Maximum number of students have registered for General Chemistry up till now. Besides World Civilization 1 and World Civilization 2 are also in demand. These courses are more in demand at UB according the data provided and hence they should be allocated comparatively greater lecture rooms or lecture rooms having greater capacity so as to avoid mismanagement. These courses should be allocated a greater capacity as they seem to be more demand.

**PROBLEM 2**: Special Case: To find in which year semester had maximum students registered for General Chemistry.

#### **Implementation details:**

#### **Data Cleaning:**

We cleaned and filtered the field Students\_registered (field 7). We ignored those tuples where there were zero number of students registered or this field had a non -numeric or invalid entry. Also if the corresponding field of Hall\_Room has an invalid entry like "Unknown" and "Arr" (for halls). Also in this case we filtered out the data on the basis of the field course\_name and considered only those fields that hade "General Chemistry" in their course\_name field as we are considering it as a special case.

### Map Reduce job: 2 MR Jobs

Job 1: Found the number of students registered for General chemistry for each year semester respectively.

Key: Year semester field value: Number of students registered

Job 2: To sort the records on values

In Mapper 2 we swap the Key, value pairs obtained from stage 2 and in Reducer 2 display the swapped key value pair i.e. key: Number of students registered Value: Year semester. This stage will output the records in sorted order based on the initial values as Hadoop sorts the keys by default.

### **Output:**

```
3903 :Fall 2003_General Chemistry
3420 :Fall 2004_General Chemistry
3505 :Fall 1995_General Chemistry
3505 :Fall 1995_General Chemistry
3506 :Spring 2004_General Chemistry
3709 :Fall 2006_General Chemistry
3709 :Fall 2006_General Chemistry
3700 :Fall 2005_General Chemistry
3700 :Spring 2005_General Chemistry
4023 :Spring 2006_General Chemistry
4024 :Spring 2007_General Chemistry
40404 :Spring 2007_General Chemistry
407 :Spring 2007_General Chemistry
4090 :Fall 2010_General Chemistry
4111 :Fall 2007_General Chemistry
4127 :Fall 2008_General Chemistry
4134 :Spring 2015_General Chemistry
4238 :Spring 2016_General Chemistry
4238 :Spring 2016_General Chemistry
4314 :Spring 2016_General Chemistry
4314 :Spring 2016_General Chemistry
4316 :Spring 2016_General Chemistry
4317 :Fall 2009_General Chemistry
4318 :Spring 2016_General Chemistry
4319 :Spring 2016_General Chemistry
4310 :Spring 2016_General Chemistry
4311 :Spring 2016_General Chemistry
4312 :Spring 2016_General Chemistry
4313 :Spring 2016_General Chemistry
4314 :Spring 2016_General Chemistry
4315 :Spring 2016_General Chemistry
4316 :Spring 2016_General Chemistry
4317 :Fall 2016_General Chemistry
4318 :Spring 2016_General Chemistry
4319 :Spring 2016_General Chemistry
4310 :Spring 2016_General Chemistry
4311 :Spring 2016_General Chemistry
4312 :Spring 2016_General Chemistry
4313 :Spring 2016_General Chemistry
4314 :Spring 2016_General Chemistry
4315 :Spring 2016_General Chemistry
4316 :Spring 2016_General Chemistry
4317 :Fall 2016_General Chemistry
4318 :Spring 2016_General Chemistry
4319 :Spring 2016_General Chemistry
4319 :Spring 2016_General Chemistry
4310 :Spring 2016_General Chemistry
4311 :Spring 2016_General Chemistry
4311 :Spring 2016_General Chemistry
4312 :Spring 2016_General Chemistry
4313 :Spring 2016_General Chemistry
4314 :Spring 2016_General Chemistry
4315 :Spring 2016_General Chemistry
4316 :Spring 2016_General Chemistry
4317 :Spring 2016_General Chemistry
4318 :Spring 2016_General Chemistry
4318 :Spring 2016_General Chemistry
4318 :Spring 2016_Ge
```

#### **Conclusion:**

Maximum students had registered for General Chemistry last year in Fall semester i.e. Fall 2015. Also as we can observe from the output greater number of students have registered for this subject in Fall and Spring semesters. Hence this subject should be allocated larger or greater number of classrooms in Fall and Spring semesters. The course capacity should be probably increased for these semesters so that students can get the benefits of the course.

**PROBLEM 3:** Which hall utilized the maximum time since the year 2000.

# **Implementation details:**

# **Data Cleaning:**

We cleaned and filtered the field Students\_registered (field 7). We ignored those tuples where there were zero number of students registered or this field had a non -numeric or invalid entry. Also if the corresponding field of Hall\_Room has an invalid entry like "Unknown" and "Arr" (for halls). We have also ignored those tuples that have their time field as "Unknown" or "Before 8: 00 AM". We have considered only those years since 2000.

### Map Reduce job: 3 MR Jobs

Job 1: Found the time utilized by each hall.

Key: Name of hall value: Time utilized (set to 1) as all the classes have a duration of one hour.

Job 2: To sort the records on values

In Mapper 2 we swap the Key, value pairs obtained from stage 2 and in Reducer 2 display the swapped key value pair i.e. key: Name of hall Value: time. This stage will output the records in sorted order based on the initial values as Hadoop sorts the keys by default.

#### Job 3: To display topmost record

In Mapper 3 we set the Key as "1" so that in Reducer we can put all the values in one Array. Accessing the first element of an array will give us the maximum value i.e. the hall which utilized maximum time.

### **Output:**

```
hadoop@hadoop-VirtualBox:~/hadoop$ hdfs dfs -cat plo3/*
16/04/15 17:15:40 WARN util.NativeCodeLoader: Unable to load native-hadoop library for your platform... using builtin-java classes where applicable
19934 Baldy
hadoop@hadoop-VirtualBox:~/hadoop$
hadoop@hadoop-VirtualBox:~/hadoop$
```

#### **Conclusion:**

Baldy hall seems to be more occupied as compared to other halls. The schedule for this hall should be decided well in advance as the hall seems to be occupied.

**PROBLEM 4:** To find the trends of time utilized by Baldy hall for consecutive years since 2000.

# **Implementation details:**

### **Data Cleaning:**

We cleaned and filtered the field Students\_registered (field 7). We ignored those tuples where there were zero number of students registered or this field had a non -numeric or invalid entry. Also if the corresponding field of Hall\_Room has an invalid entry like "Unknown" and "Arr" (for halls). We have also ignored those tuples that have their time field as "Unknown" or "Before 8: 00 AM". We have considered only those years since 2000.

### Map Reduce job: 2 MR Jobs

JOB 1: To find the time utilized by Baldy hall in each year.

Key: Year Value: Time utilized (set to 1 as all the classes take one hour)

JOB2: To find the change in time utilization for consecutive years:

In Mapper2 we have set the key twice and the reducer 2 displays the change in time utilized for those consecutive years for which the key was set.

# **Output:**

```
nadoop@hadoop-VirtualBox:~/hadoop$ hdfs dfs -cat plo4/*
16/04/15 18:24:58 WARN util.NativeCodeLoader: Unable to load native-hadoop library for your platform... using builtin-java classes where applicable
000-2001
 902-2003
003-2004
004-2005
                   -61
-60
005-2006
 906-2007
007-2008
                   22
-66
111
 908-2009
 009-2010
2010-2011
011-2012
012-2013
 913-2014
015-2016
  doop@hadoop-VirtualBox:~/hadoop$
```

#### **Conclusion:**

There was a remarkable increase in utilization of Baldy hall from year 2016 to 2017 by 1762 hours. Roughly the time utilized by Baldy has increased over the years. Even if there was a decrease in the time utilization that was comparatively less.

**PROBLEM 5:** What was the peak time for which most number of classes were held for each year.

#### **Implementation details:**

#### **Data Cleaning:**

We cleaned and filtered the field Students\_registered (field 7). We ignored those tuples where there were zero number of students registered or this field had a non -numeric or invalid entry. Also if the corresponding field of Hall\_Room has an invalid entry like "Unknown" and "Arr" (for halls). We have also ignored those tuples that

have their time field as "Unknown" or "Before 8: 00 AM". We have considered only those years since 2000. Besides we have also filtered out the day's field (field 3), we have ignored those tuples that have "Uknwn" or "Arr" in the days field.

### Map Reduce job: 2 MR Jobs

JOB 1: To get how many classes are held respective to each time period.

Key: Year + time period e.g. 2000\_3:00AM-3:59AM Value: Number of classes (We have calculated the number of classes) based on the length of the day's field i.e. if the entry is M-F we have assigned the length 5, if MTW the length assigned is 3.

JOB2: To get the time for which maximum classes are held in each year respectively.

Mapper 2 passes the entire value to the reducer where the Reducer2 calculates the time for which maximum number of classes held keeping the year unique.

# **Output:**

#### **Conclusion:**

As we can see from the above output the peak time for almost all years is 9:00AM - 9:59 AM. This slot of the day seems to be very busy in all years and there seem to be many classes going on during this time period. Over all the classrooms in UB campus will be filled with students during this period as this is the peak lecture period. Probably more class rooms should be allocated to this time period to avoid mismanagement.

**PROBLEM 6:** To find the top 5 halls having the maximum capacity for the recent year.

### **Implementation details:**

### **Data Cleaning:**

We cleaned and filtered the field Students\_registered (field 7). We ignored those tuples where there were zero number of students registered or this field had a non -numeric or invalid entry. Also if the corresponding field of Hall\_Room has an invalid entry like "Unknown" and "Arr" (for halls). Also we have checked that if the number of students registered for the corresponding tuple is greater than zero.

### Map Reduce job: 3 MR Jobs

Job 1: To find the capacity for each hall.

Key: Name of hall value: Maximum capacity.

Job 2: To sort the records on values

In Mapper 2 we swap the Key, value pairs obtained from stage 2 and in Reducer 2 display the swapped key value pair i.e. key: Maximum capacity Value: Hall Name. This stage will output the records in sorted order based on the initial values as Hadoop sorts the keys by default.

### Job 3: To display top 5 records

In Mapper 3 we set the Key as "1" so that in Reducer we can put all the values in one Array. Accessing the first element of an array will give us the maximum value i.e. hall which has maximum capacity. Thus we have accessed the first 5 records to display the top 5 halls having the maximum capacity.

# **Output:**

```
hadoop@hadoop-VirtualBox:~/hadoop$
```

#### **Conclusion:**

Nsc Hall has the maximum capacity in the UB campus. Hence comparatively greater classes must be being held in this hall or this hall must be utilized for conducting events or seminars attended by huge people. Besides Baldy,knox and Dfn Hall also have greater capacity, as we all know Baldy hall has a Gym which required comparatively more space. As we have seen from Problem 3 Baldy hall has utilized the maximum time overall, also it has a comparatively greater capacity which might mostly be the reason for it being used for a greater time period.

**PROBLEM 7:** To find the trends and compare the number of students registered for online course v/s those on campus for a particular course for each year.

#### **Implementation details:**

#### Clean up:

We cleaned and filtered the field Students\_registered (field 7). We ignored those tuples where there were zero number of students registered or this field had a non -numeric or invalid entry. Also if the corresponding field of Hall\_Room has an invalid entry like "Unknown" and "Arr" (for halls). Also we have checked that if the number of students registered for the corresponding tuple is greater than zero. We filtered the data to consider the courses offered only after 2009 as online courses were offered only after 2009.

### Map Reduce job: 2 MR Jobs

JOB 1: To find the number of students registered for online and on campus course (for each course respectively)

Key: Course\_name+Year Value: Online or campus depending on the case

JOB2: Filtering those courses which offer both online and on campus course for that year.

Mapper 2 passes the entire value to the reducer where the Reducer2 calculates and makes the comparison between the number of students registered online and on campus for the same course in each year.

# **Output:**

```
ddv Statistical Technique_2016
Adv Theory Const & Devel_2015
                                   Online:1
                                                     OnCampus:7
Advanced Design Theory_2012
Advanced Design Theory_2013
                                   Online:4
                                                      OnCampus:9
Advanced Design Theory_2014
                                   Online:1
                                                     OnCampus:12
Advanced Design Theory 2015
                                   Online:5
                                                     OnCampus:2
Advanced Design Theory_2016
Advanced Standing Seminar_2011
                                   Online:2
                                                     OnCampus:0
                                   Online:10
                                                     OnCampus:50
Advanced Standing Seminar_2015 Online:17
Advanced Standing Seminar_2016 Online:10
                                                     OnCampus:50
                                                     OnCampus:14
Aging Population & Family_2016 Online:23
                                                     OnCampus:0
Analog Circuits_2014 Online:7
                                            OnCampus:61
Analysis/Quantity Rsrch 2_2014 Online:13
                                                     OnCampus:17
Appl Child Devl&Learning_2016
                                   Online:0
                                                     OnCampus:40
Assess MH Counseling_2013
                                                     OnCampus:16
Assess MH Counseling 2014
                                   Online:17
Assess MH Counseling 2015
                                                     OnCampus:13
Assess MH Counseling 2016
                                   Online:19
                                                     OnCampus:13
Assess Sec Lang Profency_2013 Online:4
Bankruptcy Law_2009 Online:35 C
                                                     OnCampus:32
Bankruptcy Law_2009
                                            OnCampus:29
```

#### **Conclusion:**

This gives us a clear idea about which course had more demand online or on campus for each year. As we can see the output above for the subject Assess MH Counseling the demand for online courses is greater as compared to off campus and greater number of students are registering online. Thus allocating larger space classrooms such courses will probably not be wise decision. However, some courses like Analog Circuits have greater demand on campus hence they should positively be considered while allocating classes. This surely will give us a clear idea about the courses whose online demand is increasing, and probably increasing the capacity would help as people are considering taking up the online course option. Also a point worth noting is that online courses were offered in UB since the year 2009 but not before that.

**PROBLEM 8:** To find which hall was most efficiently utilized for each semester in each year since the year 2000.

### **Implementation details:**

# Clean up:

We cleaned and filtered the field Students\_registered (field 7). We ignored those tuples where there were zero number of students registered or this field had a non -numeric or invalid entry. Also if the corresponding field of Hall\_Room has an invalid entry like "Unknown" and "Arr" (for halls). Also we have checked that if the number of students registered for the corresponding tuple is greater than zero. Besides we have considered years since only 2000.

JOB 1: To find how efficiently utilized each class was used for each semester year. Here efficiency is a measure being calculated on the basis of difference between max capacity-registered students. Lesser the difference more efficiently has the class been used.

Key: Year Semester+Hall Name Value: Efficiency (Max capacity-Registered students).

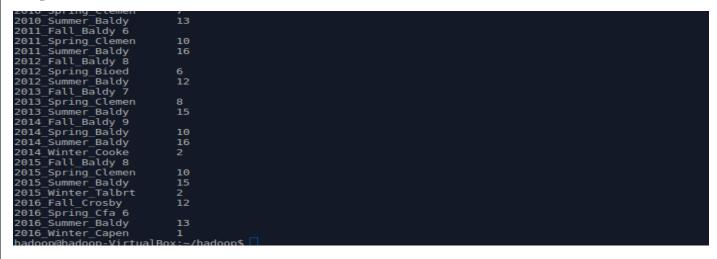
JOB2: To find out how many classes were there in each hall which were efficiently utilized. To filter out those entries whose difference between Max capacity- registered students is less than or equal to 25.

Mapper 2 filters out those classes that have a difference of 25 or less for Max capacity-Registered students and passes on to the second Reducer

JOB3: To find a hall for each year with the maximum number of most efficiently utilized rooms.

Mapper 3 passes the entire value to the reducer where the Reducer2 calculates the hall which had the maximum classes efficiently utilized.

# **Output:**



#### **Conclusion:**

On observing the output we can see that the most efficiently utilized hall was Capen in Winter 2016 over the years as there is difference of only one. This gives a clear idea about the how efficiency of the hall varied each year and in which year which hall was the most efficiently utilized Also we can see from the above output that Baldy is most frequently most efficiently utilized hall for each semester. From problem3 we know that Baldy also was utilized for maximum time besides it has comparatively greater capacity, from this we can probably conclude that Baldy seems to be always occupied thus making it efficiently utilized majority of the times. Also we can see that the halls are more efficiently utilized in the Fall semester comparatively, from problem 9 we will see that maximum students register during this semester thus being the causes for halls being mostly efficiently utilized. increasing its efficiency

**PROBLEM 9:** To find the semester for which generally greater students registered.

### **Implementation details:**

### Clean up:

We cleaned and filtered the field Students\_registered (field 7). We ignored those tuples where there were zero number of students registered or this field had a non -numeric or invalid entry. Also if the corresponding field of Hall\_Room has an invalid entry like "Unknown" and "Arr" (for halls). Also we have checked that if the number of students registered for the corresponding tuple is greater than zero.

Map Reduce job: 3 MR Jobs

Job 1: Found the capacity for each hall.

Key: semester l value: Maximum capacity.

Job 2: To sort the records on values

In Mapper 2 we swap the Key, value pairs obtained from stage 2 and in Reducer 2 display the swapped key value pair i.e. key: Maximum registered Value: Semester. This stage will output the records in sorted order based on the initial values as Hadoop sorts the keys by default.

### Job 3: To display topmost record

In Mapper 3 we set the Key as "1" so that in Reducer we can put all the values in one Array. Accessing the first element of an array will give us the maximum value i.e. the semester for which maximum students have registered.

# **Output:**

Bytes Written=13
hadoop@hadoop-VirtualBox:~/hadoop\$ hdfs dfs -cat plott1/\*
16/04/15 23:05:05 WARN util.NativeCodeLoader: Unable to load native-hadoop library for your platform... using builtin-java classes where applicable
3745058 Fall
hadoop@hadoop-VirtualBox:~/hadoop\$
hadoop@hadoop-VirtualBox:~/hadoop\$

# **Conclusion:**

Greater number of students have registered for Fall semester. The Fall semester must be probably offering greater number of courses or having greater capacity. Hence many students must be registering for classes in Fall as compared to Spring, Winter and Summer.