# Worcester Polytechnic Institute
# CS 585/DS 503 Big Data Management
# Project 2: Advanced big data operations and analytics

Group 9: Xiaoting Cui, Allison Rozet, Abhinav Singh

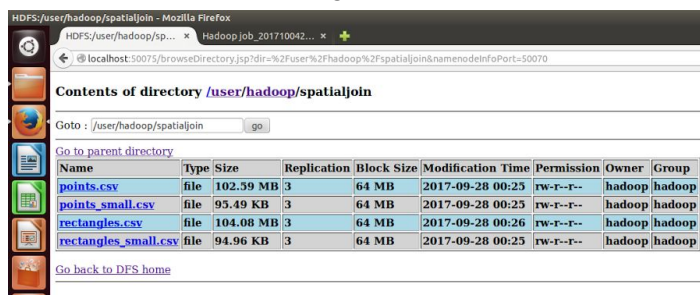October 6, 2017

## Problem 1: Spatial Join

We solved this problem three ways based on three different assumptions: (1) the rectangles data set was small and the points data set was large, (2) points was small and rectangles was large, and (3) both points and rectangles were large.

### Step 1 (Create the Two Datasets)

First, we do not assume that each point must be unique. This is because in the real, 3D world two points (or rectangles) might have the same x- and y-coordinates, but say, be on different floors of a building. We assume that points and edges can lie on the edge of the space, i.e. at 1 or 10000. Rectangles are generated so that both corners are within the space and are at most 20x20 in size. Finally, we generated integer coordinates. We discuss at the end of this section how to relax this final assumption.

To execute, run datacreation.jar group9.Main1. Use arguments such as: `-P points.csv -R rectangles.csv -np 11000000 -nr 4000000`

After creating the data sets, we loaded them into Hadoop as in Project 1. Our large points file has 11 million points, our large rectangles file has 4 million rectangles, our small points file has 10 thousand points, and our small rectangles file has 4 thousand rectangles. The screenshot below verifies that our large files exceed 100MB.
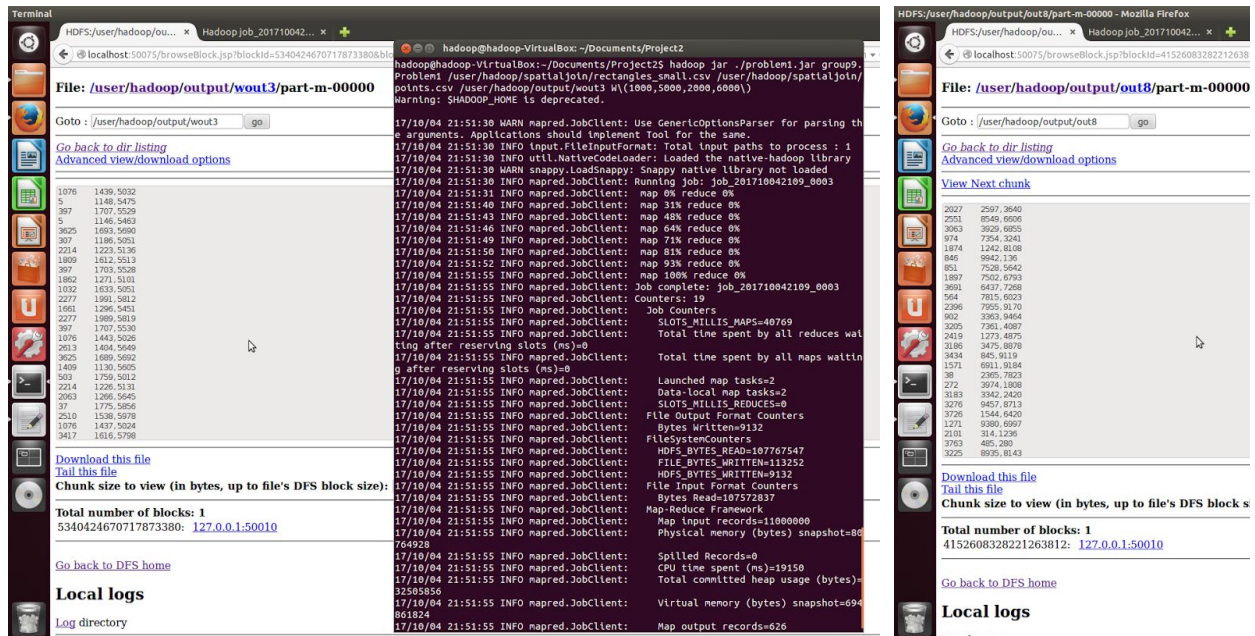
# Step 2 (MapReduce Job for Spatial Containment Join)

***Small rectangles file; large points file***.  This is a map-only job and essentially a broadcast-join. We cache the small rectangles file so that we can broadcast it to every mapper, along with the optional window. In the setup() method, we retrieve the window corner coordinates and read in the rectangles file. In the map() method, if a point falls within or on the window, we iterate through each rectangle and write out the (rectangle id, point) if the point falls within or on the rectangle.

Please use code similar to the following to execute the jar file.
```
hadoop jar ./problem1.jar group9.Problem1 <small rectangles input>
<large points input> <output> <optional W\(xbl,ybl,xtr,ytr\)>
```



The screenshots above are sample outputs with a window (left) and without a window (right).



The screenshot on the left verifies that our program with a window was successful. The screenshot on the next page verifies it without a window.

***Small points file; large rectangles file***. This is also a map-only job and analogous to the previous solution. The main difference is that we perform the check for the points that fall within the window in the setup() method rather than the map() method so that it only has to be done once rather than every time a new rectangle is read.

Please use code similar to the following to execute the jar file.

```
hadoop jar ./problem1v2.jar group9.Problem1v2 <small points input>
<large rectangles input> <output> <optional W\(xbl,ybl,xtr,ytr\)>
```

Again, two screenshots to verify successes.



**Both files are large**. This is a one map-reduce job with a separate Mapper class for points and rectangles. We employ the "striping" method similar to the lecture on matrix multiplication. Our first attempt was to "stripe" by x-value, but Hadoop failed with the two large datasets. We further divided our "stripe" into "top" (y>5000) and "bottom" (y<=5000).

We join the points and rectangles with the key <x-value:top/bottom>. In addition, rectangles were marked with their id for the secondary sort (see next paragraph). The map output value for a point was simply the y-value. In the point mapper, we also filtered only the points that fell within or on the window. The map output value for a rectangle was the y-values for the bottom-left and top-right corners. Although not required for the correct result (since the point mapper filtered points in/on the window), we chose to also filter rectangles that overlapped the window to speed up the sort/shuffle and read/write to disk in the reducer.

After performing a few experiments, we determined that each key receives about twice as many rectangles than points in the reducer. Thus, we want to ensure that the reducer reads all points first and stores them in main memory. Next, as each rectangle arrives, it checks each point, similar to our previous solution (small points file, large rectangles file). We had to overwrite the partitioner, grouping comparator, and sort comparator in order to implement the secondary sort.

Please use code similar to the following to execute the jar file.

```
hadoop jar ./problem1v3.jar group9.Problem1v3 <small points input>
<large rectangles input> <output> <optional W\(xbl,ybl,xtr,ytr\)>
```

Final screenshots are above.  The table below summarizes the performance of each job.

| Job | W(1000,5000,2000,6000) | No window input |
|---|---|---|
| Small rectangles; large points | 24 sec | 23 min, 36 sec |
| Small points; large rectangles | 28 sec | 48 min, 40 sec |
| Large points; large rectangles | 42 sec | 21 min, 54 sec |

We conclude that the "striping" method for large points and rectangles is best because it scales well, even though it had the slowest performance for a window. Indeed, if the data was event larger, we could continue striping the y-coordinates into smaller sections. We might generalize that as the window increases, the performance of the map-only broadcast join methods deteriorates much more rapidly.

Finally, we discuss how our code for this problem would be adjusted for real number coordinates rather than integers. First, the data generation requires the Random instance to generate nextFloat() rather than nextInt() with some simple arithmetic to ensure that the domain and range is still [1, 10000]. In the Hadoop programs, we would parseFloat() rather than parseInt() and change the types of all applicable data structures.

# Problem 2: Custom Input Format

## Step 1 (Upload the Dataset)

We uploaded the dataset provided on Canvas. First, we created a folder on the Hadoop file system:

```
hadoop fs -mkdir /user/hadoop/project2/input
```



For loading dataset into hadoop, type:

```
hadoop fs -put project2-airfield-dataset.txt
/user/hadoop/project2/input
```



To verify if the file is inside the correct path, type:

```
hadoop fs -ls /user/hadoop/project/input
```



## Step 2 (MapReduce Job with a Custom Input Format)

Our JsonInputFormat class extends the FileInputFormat class.  (To view the code, please open JsonInputFormat.java.)  We have overridden the RecordReader class. The overridden method nextKeyValue() has our core logic for reading JSON input file. We read in each line in the text file from "{" symbol to "}" symbol and output the whole record as one Text line to the mapper.

We have one mapper and one reducer for this problem.  The input for our mapper is <Text,Text>, and the output key-value pair is <IntWritable,IntWritable>.  For the reducer (and combiner), the input is <IntWritable,IntWritable> and the output is <IntWritable,IntWritable>.

Making the key an IntWritable sorts our results numerically by elevation, which will be helpful to the user reading the output. The code is straightforward and analogous to a simple word count program.  Our mapper, reducer, and `main()` code can be found in JsonOverider.java.

To execute the jar file, we created an output folder in Hadoop.

```
hadoop fs -mkdir /user/hadoop/project2/output
```



Please use code similar to the following to execute the jar file.

```
hadoop jar problem2.jar project2.JsonOverider
/user/hadoop/project2/input /user/hadoop/project2/output
```

We executed the program both with and without a combiner to compare the performance. In the case of executing the job with a combiner, it takes about **13 seconds** to finish the job. In case of executing the job without the combiner, it takes about **20 seconds**. The combiner improves performance by decreasing shuffle/sort and I/O costs.

The following screenshot is the log on the terminal after executing the job *with* a combiner.

After executing the job *without* a combiner:



Finally, we have included a screenshot of sample output after executing the jar file. The number on the right (key) is the elevation, and the number on the left (value) is the count. For example, there were 36 records that listed an elevation of 3.

# Problem 3: K-Means Clustering

## Step 1 (Creation of Dataset) :

We created a dataset of 2D points based on some randomly chosen centers. This way, the data is more likely to be "clustered" than the uniform distribution of points we generated in Problem 1. Our data set contains 11 million points and is over 100MB (see the screenshot below). We also tested our code on a smaller set of 300 points.

To execute, run datacreation.jar group9.Main3. Use arguments such as: `-P kmeans.csv -k 110 -m 100000`



| Name | Type | Size | Replication | Block Size | Modification Time | Permission | Owner | Group |
|------|------|------|-------------|------------|-------------------|------------|-------|-------|
| data.csv | file | 101.82 MB | 3 | 64 MB | 2017-09-28 15:26 | rw-r--r-- | hadoop | hadoop |
| test.csv | file | 2.81 KB | 3 | 64 MB | 2017-09-28 15:26 | rw-r--r-- | hadoop | hadoop |

Go back to DFS home

To put data into HDFS:
```
hadoop fs -put /home/yifan/Desktop/data.txt /home/yifan/hw2/
```

We created another file that stores k random initial centers, put it into HDFS, and updated it during each iteration. This file is created each time we run a Map-Reduce Java class (when iteration is 0).

The following screenshot are initial k centers. Each line is a 2D point with real number coordinates.
```
7220.096,7346.627
1949.7604,7158.033
6671.5957,275.3204
7784.408,3030.697
6186.0757,7020.345
6230.969,2487.7893
```

## Step 2 (Clustering the Data)
1.single-iteration: We set the number of centers to 6. First, we created a file that contains the initial k centers. In the mapper, the points in data.csv are read in line by line. We use the k center files to find the nearest center for each point, returning Key (center)-Pair (point that belongs to this center). The reducer takes the output of the mapper, calculating mean values of

all the points that belong to a certain center. The reducer returns the mean value (the centroid) as key.

```
hadoop jar /home/yifan/hw2/hw2.jar hw2.oneiteration 6
/home/yifan/hw2/data.txt /home/yifan/hw2/1
```
Screenshot of output:

part-r-00000 (~/hw2/11_0) - ge

Open ▾  ⊞

hadoop.log    ×       p

```
1858.7507,6405.8994
5371.5664,7100.6167
3951.7063,2406.9177
5935.647,540.1456
B235.479,7831.731
B514.726,3291.3528
```

2.Multi-iteration: We set the number of centers to 6. We made some changes in the DriverKmeans class. Each time a new map-reduce job starts, its input is the output from the last iteration. We set the number of iterations to 6.

```
hadoop jar /home/yifan/hw2/hw2.jar hw2.multiiteration 6
/home/yifan/hw2/data.txt /home/yifan/hw2/2
```
Screenshot of output:

part-r-00000 (~/hw2/2_5) - gedit

Open ▾  ⊞

```
1625.7435,7472.129
2168.8093,2512.1958
4936.062,7209.6343
6264.145,1555.0343
8327.882,3638.3225
8280.897,7966.4873
```

3. Same as above. We set the threshold to 100 and created a new variable which computes the difference between the latest centers and centers created from the last iteration. When the difference is under 100 or iteration comes to 6, the job will stop.

```
hadoop jar /home/yifan/hw2/hw2.jar hw2.problem3 6
/home/yifan/hw2/data.txt /home/yifan/hw2/3
```

```
part-r-00000 (~/hw2/3_5) - g
Open ▾    ⊞
    hadoop.log   ×         problem
1625.7435,7472.129
2168.8093,2512.1958
4936.062,7209.6343
6264.145,1555.0343
8327.882,3638.3225
8280.897,7966.4873
```

We set iteration to 21 and threshold to 3000.
```
hadoop jar /home/yifan/hw2/hw2.jar hw2.singleconvergence 6
/home/yifan/hw2/data.txt /home/yifan/hw2/22
```
Screenshot of output:

```
part-r-00000 (~/hw2/22_20) - gedit
Open ▾    ⊞
1722.6434,7731.0654,
1714.6415,2659.8892,
4980.459,6954.131,
5135.0366,2032.9302,
8245.471,7785.988,
8396.037,2759.5557,
KMeans clustering is convergent.
```

4.1 SingleReducer: When doing this optimization, we first set the number of reducers to 1 and then remove the code comparing the difference between latest center and centers created from the last iteration in DriverKmeans to the Reducer class.

```
hadoop jar /home/yifan/hw2/hw2.jar hw2.single1 6
/home/yifan/hw2/data.txt /home/yifan/hw2/4_1
```

```
1625.7435,7472.129
2168.8093,2512.1958
4936.062,7209.6343
6264.145,1555.0343
8327.882,3638.3225
8280.897,7966.4873
```

4.2 Combiner: When doing optimization, we want to use a combiner in front of the reducer to locally combine the same key. First we created a file that contains the initial k centers. We applied one mapper and reducer. In the mapper, the point in data.csv is read in line by line. We use the k center files to find the nearest center for each point, returning Key (center)-Pair (point that belongs to this center,1). The combiner takes the output of mapper, calculating sum values of all the points and the number of points that belong to a certain center, passing Key (center)-Pair(sum,and number of points) to reducer.

```
hadoop jar /home/yifan/hw2/hw2.jar hw2.problem3optimization 6
/home/yifan/hw2/data.txt /home/yifan/hw2/4_2
```



```
1625.802,7471.838
2169.2905,2512.884
4934.795,7210.5356
6261.5684,1556.0457
8328.411,3634.6682
8280.205,7963.113
KMeans clustering is not convergent
```

5.1(a) Convergence or Not: Using combiner, in the DriverKmeans class, if the difference between the latest center and centers created from the last iteration is below the threshold, we write "Convergence" into the file and otherwise, we write "not convergence."

```
hadoop jar /home/yifan/hw2/hw2.jar hw2.problem3optimization 6
/home/yifan/hw2/data.txt /home/yifan/hw2/5_1
```

```
1625.802,7471.838
2169.2905,2512.884
4934.795,7210.5356
6261.5684,1556.0457
8328.411,3634.6682
8280.205,7963.113
KMeans clustering is not convergent
```

5.2(a): For a singlereducer task, based on 4.1, in the DriverKmeans class, if the difference between the latest center and centers created from last iteration is below the threshold, we write "Convergence" into the file and otherwise, we write "not convergence."

```
hadoop jar /home/yifan/hw2/hw2.jar hw2.singleconvergence 6
/home/yifan/hw2/data.txt /home/yifan/hw2/5_3
```



```
1625.7435,7472.129,
2168.8093,2512.1958,
4936.062,7209.6343,
6264.145,1555.0343,
8327.882,3638.3225,
8280.897,7966.4873,
KMeans clustering is not convergent
```

We hypothesize that the reason our algorithm did not converge is because our synthetic dataset was built to simulate over 100 clusters, but we chose a k value that was too small. Or perhaps we need more rounds of iteration.

5.2(b):
For a singlereducer task, based on 4.1, in the Reducer class, we use a StringBuilder object to store all the clustered points for a certain center, returning Key (centers)-Pair (a set of points belongs to a center)

```
hadoop jar /home/yifan/hw2/hw2.jar hw2.optallpoint 6
/home/yifan/hw2/data.txt /home/yifan/hw2/5_4
```

1625.7325,7472.0737,

766.0,5867.0;2872.0,5332.0;1291.0,9585.0;754.0,7481.0;2788.0,8440.0;418.0,963
4.0;2998.0,8833.0;2016.0,9063.0;2429.0,5943.0;2162.0,8264.0;1144.0,8173.0;324
6.0,9739.0;2324.0,7562.0;974.0,8899.0;1028.0,8109.0;1562.0,7222.0;3262.0,8527
.0;2899.0,7554.0;2227.0,5725.0;1009.0,7039.0;207.0,9265.0;1678.0,9284.0;137.0
,9972.0;1867.0,6852.0;179.0,7362.0;631.0,9864.0;2433.0,6291.0;2473.0,5748.0;2
701.0,9700.0;2997.0,5378.0;1317.0,4868.0;1041.0,7160.0;1061.0,8187.0;2159.0,9
997.0;21.0,9794.0;3169.0,6074.0;414.0,9765.0;1000.0,8976.0;1345.0,8233.0;1912
.0,8264.0;487.0,5147.0;333.0,7370.0;3217.0,8360.0;123.0,5002.0;1449.0,6195.0;
2797.0,8395.0;495.0,9597.0;2512.0,9941.0;2862.0,9504.0;2623.0,5356.0;956.0,57
40.0;1939.0,6122.0;81.0,7840.0;1259.0,5253.0;330.0,9721.0;2170.0,9702.0;1792.
0,7043.0;618.0,7925.0;2989.0,8822.0;1070.0,9879.0;1513.0,5923.0;3146.0,7540.0
;1894.0,8748.0;2771.0,5051.0;2131.0,8430.0;2632.0,9179.0;2881.0,5638.0;2262.0
,5621.0;507.0,8481.0;638.0,8711.0;1273.0,9559.0;1649.0,5738.0;629.0,9165.0;68
6.0,8596.0;2817.0,5863.0;2440.0,9539.0;464.0,9175.0;746.0,8949.0;2955.0,9792.
0;2548.0,8029.0;3202.0,9787.0;1157.0,6720.0;2151.0,8060.0;461.0,8863.0;390.0,
9957.0;981.0,5014.0;490.0,9471.0;2641.0,8012.0;1683.0,8684.0;2393.0,6853.0;19
34.0,7017.0;1916.0,8982.0;298.0,6169.0;465.0,5511.0;392.0,5342.0;1487.0,9526.
0;633.0,6958.0;1908.0,7076.0;1584.0,8423.0;47.0,8926.0;1294.0,7369.0;101.0,97
32.0;250.0,7661.0;533.0,8896.0;2219.0,7245.0;2901.0,6297.0;2678.0,6679.0;213.
0,7985.0;603.0,8820.0;1337.0,8274.0;1959.0,5397.0;727.0,6964.0;1971.0,6190.0;
1937.0,5383.0;272.0,6859.0;1486.0,5200.0;1276.0,8890.0;2794.0,7152.0;3278.0,7
126.0;2581.0,8438.0;388.0,8484.0;219.0,8145.0;99.0,4797.0;2448.0,6914.0;116.0
,7189.0;1933.0,5790.0;1462.0,7524.0;1848.0,5763.0;833.0,7883.0;2041.0,5129.0;
1400.0,7399.0;2069.0,7050.0;1796.0,7726.0;2350.0,8913.0;2856.0,7154.0;2146.0,
6120.0;1405.0,7856.0;590.0,7053.0;135.0,8618.0;2346.0,9381.0;306.0,7887.0;29.
0,8215.0;663.0,5742.0;1283.0,5631.0;2308.0,6808.0;381.0,9226.0;2095.0,7053.0;
2442.0,6765.0;3004.0,9758.0;2323.0,9969.0;1035.0,6145.0;1821.0,6817.0;2683.0,
9060.0;2236.0,8712.0;973.0,5686.0;2860.0,7887.0;1351.0,9550.0;2787.0,7779.0;1
084.0,8263.0;773.0,7956.0;540.0,4919.0;3067.0,9462.0;1117.0,6389.0;259.0,6286
.0;2797.0,6187.0;435.0,5942.0;1949.0,9009.0;1774.0,8991.0;9.0,8199.0;2738.0,7
001.0;705.0,7007.0;1429.0,6691.0;1140.0,6209.0;3009.0,5717.0;2110.0,9642.0;59

2168.8752,2512.389,

73.0,3921.0;2434.0,1076.0;3702.0,3100.0;2185.0,757.0;4529.0,3376.0;3007.0,484
8.0;4406.0,3944.0;3319.0,1135.0;2064.0,1939.0;1036.0,3481.0;954.0,4794.0;479.
0,122.0;3587.0,67.0;439.0,4032.0;623.0,3248.0;847.0,3449.0;2225.0,126.0;3872.
0,3854.0;1018.0,4196.0;1141.0,2319.0;57.0,144.0;3626.0,125.0;4270.0,2811.0;37
82.0,2651.0;4731.0,4093.0;353.0,835.0;2574.0,1526.0;3928.0,2897.0;3889.0,4689
.0;1771.0,1415.0;3856.0,2321.0;2593.0,680.0;1549.0,316.0;1390.0,2851.0;557.0,
1737.0;2903.0,1701.0;3796.0,3054.0;1405.0,1584.0;419.0,1966.0;3813.0,4684.0;3
57.0,3435.0;1144.0,4536.0;832.0,4721.0;3346.0,3972.0;838.0,3809.0;3754.0,4477
.0;216.0,346.0;2333.0,2958.0;92.0,3167.0;1214.0,2684.0;3552.0,4594.0;3638.0,6
57.0;4035.0,1750.0;2167.0,1129.0;1219.0,3345.0;1841.0,2497.0;1822.0,664.0;233
6.0,354.0;3752.0,276.0;1878.0,1243.0;1804.0,2849.0;2778.0,1857.0;1124.0,1992.
0;3226.0,316.0;2372.0,1165.0;3337.0,2389.0;1008.0,4067.0;2517.0,2645.0;432.0,
2909.0;1045.0,2861.0;1845.0,2271.0;1535.0,2535.0;2666.0,3724.0;3641.0,220.0;1
977.0,3776.0;2099.0,2233.0;519.0,3301.0;1162.0,998.0;2908.0,3955.0;2171.0,398
7.0;47.0,345.0;3483.0,323.0;2254.0,2103.0;3245.0,3934.0;2974.0,2149.0;3998.0,
2893.0;3143.0,1786.0;4280.0,4207.0;1546.0,2970.0;2981.0,1205.0;4167.0,1702.0;
429.0,2532.0;3674.0,1258.0;401.0,4742.0;2734.0,3676.0;498.0,1325.0;2006.0,316

Same as above, except we have multiple reducers.

```
hadoop jar /home/yifan/hw2/hw2.jar hw2.problem3opt4 6
/home/yifan/hw2/datas.txt /home/yifan/hw2/5_2
```
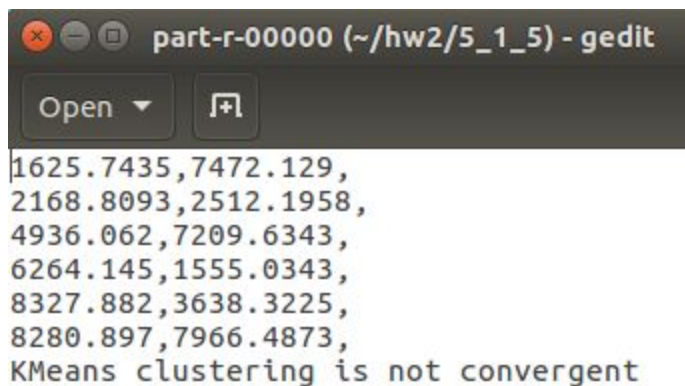
```
1625.7325,7472.0737,
766.0,5867.0;2872.0,5332.0;1291.0,9585.0;754.0,7481.0;2788.0,8440.0;418.0,963
4.0;2998.0,8833.0;2016.0,9063.0;2429.0,5943.0;2162.0,8264.0;1144.0,8173.0;324
6.0,9739.0;2324.0,7562.0;974.0,8899.0;1028.0,8109.0;1562.0,7222.0;3262.0,8527
.0;2899.0,7554.0;2227.0,5725.0;1009.0,7039.0;207.0,9265.0;1678.0,9284.0;137.0
,9972.0;1867.0,6852.0;179.0,7362.0;631.0,9864.0;2433.0,6291.0;2473.0,5748.0;2
701.0,9700.0;2997.0,5378.0;1317.0,4868.0;1041.0,7160.0;1061.0,8187.0;2159.0,9
997.0;21.0,9794.0;3169.0,6074.0;414.0,9765.0;1000.0,8976.0;1345.0,8233.0;1912
.0,8264.0;487.0,5147.0;333.0,7370.0;3217.0,8360.0;123.0,5002.0;1449.0,6195.0;
2797.0,8395.0;495.0,9597.0;2512.0,9941.0;2862.0,9504.0;2623.0,5356.0;956.0,57
40.0;1939.0,6122.0;81.0,7840.0;1259.0,5253.0;330.0,9721.0;2170.0,9702.0;1792.
0,7043.0;618.0,7925.0;2989.0,8822.0;1070.0,9879.0;1513.0,5923.0;3146.0,7540.0
;1894.0,8748.0;2771.0,5051.0;2131.0,8430.0;2632.0,9179.0;2881.0,5638.0;2262.0
,5621.0;507.0,8481.0;638.0,8711.0;1273.0,9559.0;1649.0,5738.0;629.0,9165.0;68
6.0,8596.0;2817.0,5863.0;2440.0,9539.0;464.0,9175.0;746.0,8949.0;2955.0,9792.
0;2548.0,8029.0;3202.0,9787.0;1157.0,6720.0;2151.0,8060.0;461.0,8863.0;390.0,
9957.0;981.0,5014.0;490.0,9471.0;2641.0,8012.0;1683.0,8684.0;2393.0,6853.0;19
34.0,7017.0;1916.0,8982.0;298.0,6169.0;465.0,5511.0;392.0,5342.0;1487.0,9526.
0;633.0,6958.0;1908.0,7076.0;1584.0,8423.0;47.0,8926.0;1294.0,7369.0;101.0,97
32.0;250.0,7661.0;533.0,8896.0;2219.0,7245.0;2901.0,6297.0;2678.0,6679.0;213.
0,7985.0;603.0,8820.0;1337.0,8274.0;1959.0,5397.0;727.0,6964.0;1971.0,6190.0;
1937.0,5383.0;272.0,6859.0;1486.0,5200.0;1276.0,8890.0;2794.0,7152.0;3278.0,7
126.0;2581.0,8438.0;388.0,8484.0;219.0,8145.0;99.0,4797.0;2448.0,6914.0;116.0
,7189.0;1933.0,5790.0;1462.0,7524.0;1848.0,5763.0;833.0,7883.0;2041.0,5129.0;
1400.0,7399.0;2069.0,7050.0;1796.0,7726.0;2350.0,8913.0;2856.0,7154.0;2146.0,
6120.0;1405.0,7856.0;590.0,7053.0;135.0,8618.0;2346.0,9381.0;306.0,7887.0;29.
0,8215.0;663.0,5742.0;1283.0,5631.0;2308.0,6808.0;381.0,9226.0;2095.0,7053.0;
2442.0,6765.0;3004.0,9758.0;2323.0,9969.0;1035.0,6145.0;1821.0,6817.0;2683.0,
```

Performance:

| Job | Running Time |
| --- | --- |
| One iteration | 60s |
| Multiple iteration | 337s |
| Not convergent (iteration: 6) | 345s |
| convergent (iteration: 20) | 1258s |
| Single Reducer | 344s |
| With Combiner | 215s |
| Whether convergent (Single Reducer) | 330s |

| Write out final clustered points (Single Reducer) | 475s |
|---|---|
| Write out final clustered points (Multiple Reducer) | 500s |

Findings: we need more running time if we set more iterations. Single Reducer is a little faster than multiple reducer because we have less output and input. We can calculate difference in the reducer node for a single reducer. Using the combiner will reduce running time because combiner provides local combining, decreasing tasks for shuffling. Writing out all final points needs more time than just indicating whether it is convergent because there is more write out and read in cost.

# Team Member Background and Contributions

Although our group is quite diverse in terms of demographics and program of study, we were surprised to find that we are quite similar in terms of background skills. Although this class was our first time using Hadoop, all three of us were comfortable coding in Java. Thus, our general strategy for this project was that each of us selected 1-2 of the problems that were most interesting to us and solved them. We shared our solutions via GitHub, read each other's solutions, and provided feedback until we were satisfied with three final solutions.

Problem 1: Allison and Abhinav both completed this problem independently. Abhinav chose to solve the problem with a small rectangles dataset and large points dataset, while Allison also explored other solutions. Xiaoting reviewed their codes.

Problem 2: Allison and Abhinav both completed this problem independently. Their solutions were very similar, and they discussed the pros and cons of the slight differences in their code before submitting a final solution.

Problem 3: Xiaoting completed this problem in its entirety. Allison decided to generate synthetic data that had "clusters," i.e. not uniformly distributed, to test the code on a more meaningful dataset. Abhinav re-ran her code to verify performance. The team reviewed Xiaoting's code and provided comments, which she applied herself.

Report: Allison wrote the documentation for Problem 1, Abhinav wrote Problem 2, and Xiaoting wrote Problem 3. Everyone read and edited the document to prepare it for final submission.

The skills we all acquired while working on this project include implementation of a broadcast join, handling different input format types, and more familiarity with the Hadoop framework.  Our team worked together effectively to complete all problems and organize our results.