# PROJECT 2 REPORT

## Purpose of Project:

### *Why did we implement hot zone analysis?*

ST_Contains function essentially checks if a point is inside rectangular coordinates. The data input are split and corners are found based on which we give a True/False value back.

In this function, we work with some data only. We do a "rangeJoin" query on the rectangle and point dataset. We then count the number of points inside each rectangle and then sort it based on value. If a rectangle has the max points then we can label it as the "hottest" among all rectangles. We also have a function called run_hot_zone_analysis which will sort the data based on point and rectangle Path it will give an output that gives the rectangle and no of points in that particular rectangle.

### *Why did we implement hot cell analysis?*

In this, we process spatial data of New York City taxi trips. We calculate Getis-Ord coordinate statistics to help us in the ST function. Getis-ord Gi is calculated as:

$$G_i^* = \frac{\sum\limits_{j=1}^{n} w_{i,j} x_j - \bar{X} \sum\limits_{j=1}^{n} w_{i,j}}{S \sqrt{\frac{[n \sum\limits_{j=1}^{n} w_{i,j}^2 x_j - (\sum\limits_{j=1}^{n} w_{i,j}^2)]}{n-1}}}$$

$$\bar{X} = \frac{\sum\limits_{j=1}^{n} x_j}{n}$$

$$S = \sqrt{\frac{\sum\limits_{j=1}^{n} x_j^2}{n} - (\bar{X})^2}$$

For reach cell we calculate the impact of current cell to its adjacent ones. The Z score is then calculated for each cell and ordered by largest to smallest in the final csv output. We use the concepts of "Spacetime cube" and "SpatialWeight" to calculate scores.

## Reflection:

What I particularly enjoyed in this project was the implementation of the function ST_Contains where we had a data point and a rectangle and we needed to find whether the point was inside it. We do range and range join queries inside this function to find all points inside the rectangle and range join to find points and rectangles where the point is inside it. Spatial function, ST Within calculates whether two points are within range for given points and distances. The most enjoyable part was working with apache spark, I thoroughly enjoyed the project right from the setup of JDK, Spark and Microsoft C++ distributions to working with scala files.

## Lessons learned:

In this project, we got to explore a lot of new technologies and implement them for hot cell and hot zone analysis. We use scala which supports OOPs and functional programming and runs on JVM. We use Java development kit to run scala. Apache spark is used for large-scale data processing as in the case of yellow taxi 2009 01.csv file. The main point to do this project was to execute spatial queries on large amount of data. Along the way we learned a lot about concepts like Getis-Ord statistic. Scala allows for immutable codes which allows for concurrent transactions and can be run on Java runtime env.

## Code:

The code is written in scala and then compiler into a jar file which is then run on inputs to generate outputs which are then compared to sample outputs.



HotCellUtils.scala



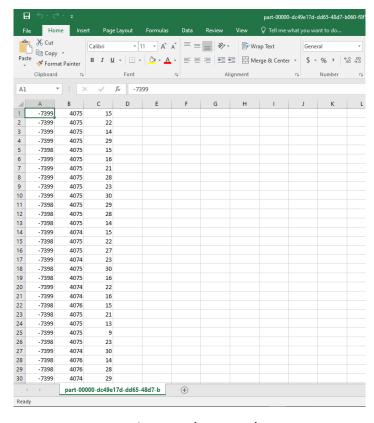HotCellAnalysis.scala



HotZoneUtils.scala



HotZoneAnalysis.scala

# DATA PROCESSING AT SCALE - CSE 511
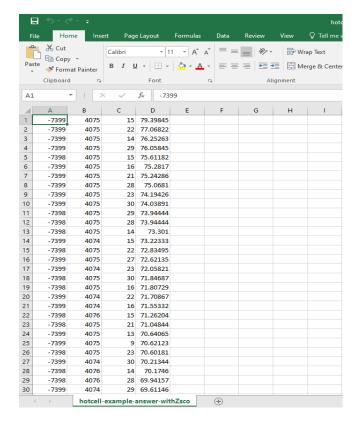
Submitted By: Aman Peshin (1225476655)
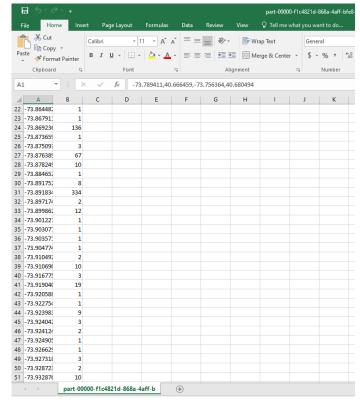
## OUTPUT: HotCell and HotZone Analysis

### Executed Output(HotCell)



### Sample Output(HotCell)



### Executed Output(HotZone)



### Sample Output(HotZone)