# Project Report for CSE 511: Data Processing at Scale, Project #2

Harsh Sanjaykumar Nagoriya
Arizona State University
Email: hnagoriy@asu.edu

Krushali Shah
Arizona State University
Email: kshah51@asu.edu

Savan Dixesh Doshi
Arizona State University
Email: sdoshi7@asu.edu

## I. INTRODUCTION

Geospatial data, often known as spatial data, is a resource with a geographic component.In other terms, the information in this type of data set often comprises 2 or more dimensional attributes such as graph coordinates. Two spatial query mechanisms shall be implemented in this case namely

1) Hot zone analysis
2) hot cell/hotspot analysis.

The level of geographical data collected has increased dramatically in recent years. The most significant characteristics of a cab company are the client and driver locations. As a result, they must strive to query and analyze such massive volumes of spatial data in order to achieve high performance and low latency. The 'hot cell analysis' applies spatial statistics to Spatio-temporal Big Data in order to identify statistically significant hot spots using Apache Spark.

## II. DEVELOPMENT APPROACH

In our project meetings, we began with a requirement analysis of Phase 2 and several ways we may utilize to solve the given problem statement. We also addressed the next steps, determining what needed to be done to complete our skeleton code. We compared our output with the offered result example after completing the Hot Zone and Hot Cell Analysis. We reviewed peer code after it was implemented, improved it, and tested it.

## III. REQUIREMENTS

As stated above, to compile the scaffolding code, we needed to have following resources in our machines.

### A. Machine requirements

We used the machines with following specifications in order to compile and test the code.

*1) MacBook Air:* 8-Core CPU, 7-Core GPU, 8GB Unified Memory, 256GB SSD Storage

*2) Dell Inspiron 15 Gaming 5577:* Intel Core i7-7700, NVIDIA GeForce GTX 1050 Ti, 8GB Unified Memory, 128GB SSD Storage

### B. Operating System requirements

Since the dependencies were developed in specified version of software and platform Ubuntu, we needed to install Ubuntu in our systems.

### C. Software requirements

We used the following software configurations in order to compile and test the code.

1) *Java:* Version 1.8.0_321
2) *Scala:* Version 2.11
3) *OpenJDK:* Version 8
4) *Spark:* Version 2.4.7
5) *Hadoop:* Version 2.7
6) *Gedit Text Editor:* Version 3.3
7) *Visual Studio Code:* Version 1.66

### D. Importing the solution

The provided answer JAR files for both phases can be imported into any system with enough RAM to perform spatial queries and complete the assessment, preferably one with enough RAM to run the software versions specified above. To complete the project, follow the steps listed below.

1) Go to command prompt and the project directory. Run sbt assembly.
2) The above step would create a jar file in target.scala-2.11 folder.
3) Copy jar file and paste it in the project root folder.
4) Run the command spark-submit with input files added as arguments
5) The solution for both functions is stored in respective csv files in respective folders.

So as to keep all these things working, we set proper environment variables.

## IV. WORK DONE

The Hot Zone analysis and the Range Join Query from Phase 1 had the same requirements. There would be a set of Rectangles R and a set of points S. The first step in implementing this query would be to discover all the (Point, Rectangle) pairs where the point is inside the rectangle, and then use a GROUP BY clause to aggregate the output for each rectangle to determine its hotness. We went over the offered template to understand the provided inputs, pre-processing, and desired output in order to fully understand the task for Hot Cell Analysis. The template included cells with two spatial and one temporal feature, such as latitude, longitude, and month day. To calculate

$$\sum_{j=1}^{n} W_{i,j}$$

(which essentially is the number of neighbors for cell ) we have used the geometric location of the cell to calculate its number of neighbors. There will be 27 count for number of neighbors if the dimension on boundary is zero. Similarly, the neighbor count for cells 1, 2, and 3 would be 18,12, and 8, respectively. Is then the

$$\sum_{j=1}^{n} W_{i,j}$$

count of trips in neighboring cells as the weight for each neighbor is equal. Now that we have all of the relevant tables, we can construct Getis-ord statistics and look at the numerator and denominator. The top 50 cells are then returned as a scala dataframe after sorting the result in descending order based on the G—score value.

## V. SOLUTION VERIFICATION

As a consequence, a Data-Frame with 50 cells and three locations for the cells with the greatest Getis-Ord statistics is constructed. We debugged the code, revealing intermediate data frames from the previous part's calculations as well as statistics like Mean and Standard Deviation. To make the JAR file, we used the **sbt clean assembly** command. To launch the JAR file that we used to check the findings, type **spark-submit**. As a result, a CSV file is created in the desired output directory. The input test data for both Hot Cell Analysis and Hot Zone Analysis was provided in CSV format. We compared our output to the actual CSV-formatted result file. **Fig 1** shows the obtained output after running the JAR file.



Fig. 1.   Output

## VI. CONCLUSION

To create our data, we used Apache Spark and Scala. The built solution was utilized to extract critical and decisive data from the provided dataset, which could later be used to make strategic decisions. We've also included a few custom spatial queries in the given solution. The solution provides the client with critical geographic data that can assist them make business strategy decisions and give better service to their customers.

## VII. CONTRIBUTIONS

We, all three students, have worked equally and all together with no conflicts of interest.

## VIII. ACKNOWLEDGEMENT

We, authors, would like to thank Dr. Zhichao Cao (Instructor), Akkamahadevi Hanni (Teaching Assistant), Rithvik Chokkam (Grad Service Assistant) Arizona State University for their constant encouragement towards the realization of this work.

## REFERENCES

[1] Spark SQL - Quick Guide, Available at: https://www.tutorialspoint.com/spark_sql/spark_sql_quick_guide.htm
[2] Downloads — Apache Spark, https://spark.apache.org/downloads.html
[3] Jia Yu, "CSE512-Project-Hotspot-Analysis-Template", https://github.com/jiayuasu/CSE512-Project-Hotspot-Analysis-Template
[4] Jia Yu, "CSE512-Project-Phase2-Requirement", https://github.com/jiayuasu/CSE512-Project-Phase2-Template
[5] Download & install apache-spark on MacOS (Big Sur, Monterey, Catalina, Mojave) via Homebrew / brew, Available at: https://www.youtube.com/watch?v=6c1uP_UbuBg