Name: Sidharth Dinesh
ASU ID: 1225238352
Course: CSE 511 Data Processing at Scale

# Individual Report for Project - 2

## Purpose of Project:

Purpose of Project-2 is to understand and perform Spatial Data Analysis (Hot Spot). We accomplish this by completing two different hot spot analysis tasks:

1. Hot Zone Analysis: Using Scala and Spark, this analysis will count the points falling within the area of rectangular datasets. Each rectangle will have a corresponding count associated, and the "hotter" it is, the more dense it is, i.e. more points. The objective of the task is to find how "hot" each rectangle is.
2. Hot Cell Analysis: This analysis will use certain spatial statistical parameters on spatio-temporal big data using Spark and Scala to locate statistically significant geographical hot spots.

## Hot Zone Analysis:

We modify the "ST_Contains" function, which returns "true" value in case the points are in the areal bounds of the rectangles, otherwise returns "false". By performing a JOIN operation on the rectangles dataset and points dataset and using WHERE condition giving us boolean values, we filter the results. GROUP BY and ORDER BY helps us to count and sort the Scala Dataframe.

```scala
HotzoneUtils.scala
1    package cse512
2
3    object HotzoneUtils {
4
5      def ST_Contains(queryRectangle: String, pointString: String ): Boolean = {
6        if (queryRectangle == null || queryRectangle.isEmpty() || pointString == null || pointString.isEmpty())
7            return false
8
9          val rectangleArr = queryRectangle.split(",")
10         var x1 = rectangleArr(0).toDouble
11         var y1 = rectangleArr(1).toDouble
12         var x2 = rectangleArr(2).toDouble
13         var y2 = rectangleArr(3).toDouble
14
15         val pointArr = pointString.split(",")
16         var x = pointArr(0).toDouble
17         var y = pointArr(1).toDouble
18
19         if (x >= x1 && x <= x2 && y >= y1 && y <= y2)
20             return true
21         else if (x >= x2 && x <= x1 && y >= y2 && y <= y1)
22             return true
23         else
24             return false
25     }
26
27     // YOU NEED TO CHANGE THIS PART
28
29   }
30
```

```scala
HotzoneAnalysis.scala
32          // YOU NEED TO CHANGE THIS PART
33          val countDataframe = joinDf.groupBy("rectangle").count()
34              val outputDataframe = countDataframe.sort("rectangle").coalesce(1)
35
36          outputDataframe
```

# Hot Cell Analysis:

We need to calculate the z-score as part of Hot Cell Analysis. For each cell in a space-time cube, it has 26 neighbors, i.e. total weights in a cell are usually 27. To perform transformation of datasets, we use a highly popular column-based function of Spark SQL, which is: "User-Defined Function" (UDF). To calculate weights, we have it set according to the cell boundaries in a way that weights are equal to:

   a. 18 if the point lies on the x OR y OR z boundary.
   b. 12 if the point lies on the (x AND y) OR (y AND z) OR (x AND z) boundary.
   c. 8 if the point lies on the (x AND y AND z) boundary.

We calculate the Mean and Standard deviation using the below formulas:

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

where $x_j$ is the attribute value for cell $j$, $w_{i,j}$ is the spatial weight between cell $i$ and $j$, $n$ is equal to the total number of cells,

Finally using withColumn method, we calculate z-score using:

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

```scala
HotcellUtils.scala
55      def isAdjacentCubeFn(x1: Int, y1: Int, z1: Int,
56                           x2: Int, y2: Int, z2: Int): Boolean = {
57        if (math.abs(x1 - x2) > 1
58          || math.abs(y1 - y2) > 1
59          || math.abs(z1 - z2) > 1) false else true
60      }
61
62      def gScoreFn(mean: Double, std: Double, numCells: Double)
63                  (sum: Long, adjacentCount: Long): Double = {
64        val ad = adjacentCount.toDouble
65        (sum - mean * ad) / (
66          std * math.sqrt(
67            (numCells * ad - (ad * ad)) / (numCells - 1.0))
68          )
69      }
70
71      val is_adjacent_cube: UserDefinedFunction =
72        udf[Boolean, Int, Int, Int, Int, Int, Int](isAdjacentCubeFn)
73
74      def g_score(mean: Double, std: Double, numCells: Double): UserDefinedFunction =
75        udf[Double, Long, Long](gScoreFn(mean, std, numCells))
76      }
```

```scala
53      val COL_X = "x"
54      val COL_Y = "y"
55      val COL_Z = "z"
56      val COL_COUNT = "count"
57      val COL_SUM = "sum"
58      val COL_NOAC = "num_of_adj_cells"
59      val COL_GSCORE = "g_score"
60      val FIRST_DF = "first_df"
61      val SECOND_DF = "second_df"
62
63      /**
64       * Filter the coordinates to the area
65       * under observation and count the number
66       * of pickups for the given cell.
67       */
68      val selectedCellsCountDf = pickupInfo
69        .filter(
70          col(COL_X) >= minX and col(COL_X) <= maxX
71            and col(COL_Y) >= minY and col(COL_Y) <= maxY
72            and col(COL_Z) >= minZ and col(COL_Z) <= maxZ
73        )
74        .groupBy(COL_X, COL_Y, COL_Z)
75        .count() // This is the aggregation step for the groupBy and not the "count" action on the [[DataFrame]].
76        .cache()
77
78      val (total, sqrTotal) = selectedCellsCountDf.select(COL_COUNT)
79        .rdd
80        .aggregate((0.0, 0.0))(
81          (acc: (Double, Double), r) => {
82            val curr = r.getLong(0).toDouble
83            (
84              acc._1 + curr,
85              acc._2 + (curr * curr)
86            )
87          },
88          (x: (Double, Double), y: (Double, Double)) => (
89            x._1 + y._1,
90            x._2 + y._2
91          )
92        )
93
94      val mean = total / numCells
95      val std = math.sqrt((sqrTotal / numCells) - (mean * mean))
96
97      log.info(s"Mean: $mean")
98      log.info(s"Standard Deviation: $std")
```

```scala
100       /**
101        * Cross join to get all the cells and the filter
102        * the cells which are adjacent to the each other.
103        *
104        * NOTE: This cross join operation is one of the expensive operation
105        * which joins millions of cells together just to filter upto 26 surrounding
106        * cells. Future work can be done to index these cells and find out a way to
107        * join only the necessary cells.
108        */
109       val adjacentDf = selectedCellsCountDf.as(FIRST_DF)
110         .crossJoin(selectedCellsCountDf.as(SECOND_DF))
111         .filter(
112           HotcellUtils.is_adjacent_cube(
113             col(s"$FIRST_DF.$COL_X"), col(s"$FIRST_DF.$COL_Y"),
114             col(s"$FIRST_DF.$COL_Z"), col(s"$SECOND_DF.$COL_X"),
115             col(s"$SECOND_DF.$COL_Y"), col(s"$SECOND_DF.$COL_Z")
116           )
117         )
118         .select(
119           col(s"$FIRST_DF.$COL_X"), col(s"$FIRST_DF.$COL_Y"),
120           col(s"$FIRST_DF.$COL_Z"), col(s"$SECOND_DF.$COL_COUNT")
121         )
122         .groupBy(COL_X, COL_Y, COL_Z)
123         .agg(sum(COL_COUNT) as COL_SUM, count(COL_COUNT) as COL_NOAC)
124
125       /**
126        * Sorting after repartition may not be required.
127        * But we have added the line as we the order is maintained
128        * when there is only one part file to write, and sorting
129        * 50 elements is not computationally heavy in the given case.
130        */
131       val scoredDf = adjacentDf.withColumn(COL_GSCORE,
132         HotcellUtils.g_score(mean, std, numCells)(col(COL_SUM), col(COL_NOAC)))
133         .sort(desc(COL_GSCORE))
134         .limit(50)
135         .repartition(1)
136         .sort(desc(COL_GSCORE))
137         .select(COL_X, COL_Y, COL_Z, COL_GSCORE)
138
139       scoredDf.select(COL_X, COL_Y, COL_Z)
```