

Juliana Porto

February 28, 2024

DS503 Big Data Management

Project 2

| Question | Status        | Comment   |
|----------|---------------|---|
| Q1       | Fully Working | <p><b>Step 1: Creating Datasets</b></p> <p>I created 2 datasets that are each over 100MB as outlined in the Project 2 Q1 Part1 instructions with the points in Points.txt and the rectangles in Rectangles.txt. I also created 2 smaller datasets with 10 points and 10 rectangles for testing purposes.</p> <p><i>Assumptions:</i></p> <ul style="list-style-type: none"><li>• All values are integers</li><li>• No rectangles can extend past the boundaries</li></ul> <p><b>Source Code:</b><br/>Project2/Q1_SpatialJoin/Step1_CreatingDatasets/Main.java</p> <p><b>Points Dataset:</b><br/>Project2/Q1_SpatialJoin/ Step1_CreatingDatasets/ Points.txt</p> <p><b>Rectangles Dataset:</b><br/>Project2/Q1_SpatialJoin/ Step1_CreatingDatasets/ Rectangles.txt</p> <p><b>Test Datasets:</b><br/>.../ Test_Points.txt<br/>.../ Test_Rectangles.txt</p> <hr/> <p><b>Step 2: Spatial Join</b></p> <p><i>Plan:</i></p> <p>Point Map Logic:</p> <ol style="list-style-type: none"><li>1. For a given record, do parsing and extract fields</li><li>2. Single output (one for each point):<ol style="list-style-type: none"><li>a. Key = x value, y value</li><li>b. Value = "Point", x value, y value</li></ol></li></ol> <p>Rectangle Map Logic:</p> <ol style="list-style-type: none"><li>1. For a given record, do parsing and extract fields</li><li>2. Get max x &amp; y of rectangle:<ol style="list-style-type: none"><li>a. max_x = bottomLeft_x + width</li></ol></li></ol> |

|  |  |  |
|--|--|--|
|  |  | <p>b. <code>max_y = bottomLeft_y + height</code></p> <p>3. Get all possible points in rectangle from knowing min &amp; max x &amp; y</p> <p>4. Multiple outputs (one for each point in rectangle)</p>  |
|  |  | <p>a. Key = x value, y value</p> <p>b. Value = "Rectangle", <code>bottomLeft_x</code>, <code>bottomLeft_y</code>, <code>height</code>, <code>width</code></p> <p>Reduce Logic:</p> <ol style="list-style-type: none"> <li>1. Separate array based on dataset source, either "Point" or "Rectangle"</li> <li>2. Join records by x and y values</li> <li>3. Print points within window <ol style="list-style-type: none"> <li>a. If window is defined as an argument, use provided parameters</li> <li>b. Otherwise, use default values for entire grid: 0, 0, 10,000, 10,000</li> </ol> </li> <li>4. Output: <ol style="list-style-type: none"> <li>c. key = NULL;</li> <li>d. value = <code>bottomLeft_x</code>, <code>bottomLeft_y</code>, <code>height_width</code>, (x, y)</li> </ol> </li> </ol> <p><b>Source Code:</b><br/> Project2/Q1_SpatialJoin/Step2_SpatialJoin/SpatialJoin.java</p> <p><b>Output with Small Test Dataset &amp; no W():</b><br/> Project2/Q1_SpatialJoin/Step2_SpatialJoin/TestDataOutputQ1/part-r-00000</p> <p><b>Output with Small Test Dataset &amp; W(3,4,7,8):</b><br/> Project2/Q1_SpatialJoin/Step2_SpatialJoin/TestDataOutputQ1_window/part-r-00000</p> <p><b>Additional Notes:</b></p> <ul style="list-style-type: none"> <li>• Single map-reduce job implemented for the spatial join operation (no points lost).</li> <li>• Optional input parameter W(x1, y1, x2, y2) that indicates a spatial window (rectangle) of interest within which we want to report the joined objects. If W is given, then any rectangle that is entirely outside W and any point that is outside W should be skipped. If W is omitted, then the entire two sets should be joined.</li> </ul> |

|    |               |   |
|----|---------------|---|
|    |               | <ul style="list-style-type: none"> <li>○ <b>Note:</b> I needed to put W(x1, y1, x2, y2) in quotes when running in the terminal so that it didn't think the parenthesis were special characters.</li> </ul>  |
| Q2 | Fully Working | <p><b>Step 1: Dataset</b></p> <p>Used the point dataset created in Q1. The initial points in the HDFS file are totally in random order, and there is no specific organization. I also created 2 smaller datasets with a small number of points in specific areas for testing purposes.</p>  |
|    |               | <p><i>Assumptions:</i></p> <ul style="list-style-type: none"> <li>• All values are integers</li> </ul> <p><b>Points Dataset:</b><br/>Project2/Q2_OutlierDetection/ Step1_CreatingDataset/ Points.txt</p> <p><b>Test Point Datasets:</b><br/>.../Test_OutlierPoints.txt<br/>.../Test_OutlierPoints0.txt</p> <p><b>Step 2: Outlier</b></p> <p><i>Plan:</i></p> <p>Point Map Logic:</p> <ol style="list-style-type: none"> <li>1. For a given record, do parsing and extract fields</li> <li>2. Assign point a current grid cell within the space             <ol style="list-style-type: none"> <li>a. Make sure to round up if not an int</li> </ol> </li> <li>3. Assign point's neighboring grid cells             <ol style="list-style-type: none"> <li>a. Cell above</li> <li>b. Cell below</li> <li>c. Cell left</li> <li>d. Cell right</li> <li>e. Cell above left</li> <li>f. Cell above right</li> <li>g. Cell below right</li> <li>h. Cell below left</li> </ol> </li> <li>4. If any grid cells are equal to 0 or above the max x &amp; y value / r then do not print; otherwise...</li> <li>5. Multiple outputs (one for each defined grid cell per point):             <ol style="list-style-type: none"> <li>c. Key = grid cell</li> <li>d. Value = x value, y value</li> </ol> </li> </ol> <p>Reduce Logic:</p> |

|    |                   |  |
|----|-------------------|--|
|    |                   | <ol style="list-style-type: none"> <li>1. Input: all points connected to a certain grid cell <ol style="list-style-type: none"> <li>a. Gives us all points in that grid and its neighboring grids</li> </ol> </li> <li>2. For each point, count number of neighbors within radius r</li> <li>3. If number of neighbors are less than k, print as outlier</li> <li>4. Track non-outliers to identify unique outliers</li> <li>5. Output: <ol style="list-style-type: none"> <li>a. key = NULL</li> <li>b. value = x, y</li> </ol> </li> </ol> <p><b>Source Code:</b><br/> Project2/Q2_OutlierDetection/Step2_ReportingOutliers/<br/> OutlierDetection.java</p> <p><b>Test Outputs with Small Datasets:</b><br/> .../TestDataOutputQ2/part-r-00000<br/> .../TestData0OutputQ2/part-r-00000</p> |
|    |                   | <p><b>Additional Notes:</b></p> <ul style="list-style-type: none"> <li>• Single map-reduce job implemented for the outlier detection operation (no points lost).</li> <li>• Program takes in two mandatory parameters r and k.</li> <li>• Mapper has <i>max</i> variable to define x- and y-axis maximums and handle boundary restrictions.</li> <li>• Source referenced to understand partitioning and supporting (AKA neighboring) areas:<br/> <a href="https://web.cs.wpi.edu/~meltabakh/Publications/DOD-ICDE-2017.pdf">https://web.cs.wpi.edu/~meltabakh/Publications/DOD-ICDE-2017.pdf</a></li> </ul>  |
| Q3 | Partially Working | <p><b>Step 1: Dataset</b></p> <p>Used the point dataset created in Q1. The initial points in the HDFS file are totally in random order, and there is no specific organization.</p> <p>Within the program I also created a list of centroid points that contains K initial seed points. I tried to make this save to HDFS but was unsuccessful.</p> <p><i>Assumptions:</i></p> <ul style="list-style-type: none"> <li>• Initial centroid values are integers, but then become doubles with averaging</li> </ul> <p><b>Points Dataset:</b><br/> Project2/Q3_Clustering/Step1_CreatingDataset/Points.txt</p>  |

|  |  |  |
|--|--|--|
|  |  | <p><b>Test Dataset:</b><br/>.../ Test0_Points.txt</p> <p><b>Step 2: Clustering</b></p> <p><i>Plan:</i></p> <p>Map Logic:</p> <ol style="list-style-type: none"> <li>1. Generate all centroids, do parsing and extract fields</li> <li>2. For a given point record, do parsing and extract fields</li> <li>3. Assign each point to the closest center (k) <ol style="list-style-type: none"> <li>a. <math>D = \sqrt{(x_2 - x_1)^2 + (y_2 - y_1)^2}</math></li> </ol> </li> <li>4. Single output (one for each point): <ol style="list-style-type: none"> <li>a. Key = closest centroid x, closest centroid y</li> <li>b. Value = x value, y value</li> </ol> </li> </ol> <p>Reduce Logic:</p> <ol style="list-style-type: none"> <li>1. Recompute all centroids <ol style="list-style-type: none"> <li>a. x=average of all x-values</li> <li>b. y=average of all y-values</li> </ol> </li> <li>2. Output: <ol style="list-style-type: none"> <li>a. key = NULL</li> </ol> </li> </ol> |
|  |  | <ol style="list-style-type: none"> <li>b. value = new centroid x, new centroid y</li> </ol> <p>Iterative Check:</p> <ol style="list-style-type: none"> <li>1. Check if iterations have met max (6) <ol style="list-style-type: none"> <li>a. If yes, then stop</li> </ol> </li> <li>2. If not at max, check if new centroids are the same as the previous <ol style="list-style-type: none"> <li>a. If same, then stop</li> </ol> </li> <li>3. If not same, run another iteration to get new centroids and do everything again</li> </ol> <p><b>Source Code:</b><br/>Project2/Q3_Clusering/Step2_Clustering/Clustering.java</p> <p><b>Output with Small Test Dataset &amp; k=3:</b><br/>Project2/Q3_Clustering/Step2_Clustering/TestOutputQ3_k3/<br/>iteration_0/part-r-00000<br/>.../ iteration_1/part-r-00000<br/>.../ iteration_2/part-r-00000<br/>.../ iteration_3/part-r-00000<br/>.../ iteration_4/part-r-00000<br/>.../ iteration_5/part-r-00000</p>                          |

|  |  |  |
|--|--|--|
|  |  | <p><b>Additional Notes:</b></p> <ul style="list-style-type: none"><li>• Logic for stopping if new and old centroids the same not implemented</li></ul> |
|--|--|--|