



DEPARTAMENTO DE CIENCIAS DE LA COMPUTACIÓN

FACULTAD DE CIENCIAS FÍSICAS Y MATEMÁTICAS

UNIVERSIDAD DE CHILE

CC4102 - DISEÑO Y ANÁLISIS DE ALGORITMOS

DISEÑO Y ANÁLISIS DE ALGORITMOS

Intersección de Segmentos Ortogonales

Tarea 1

Profesor:

Pablo Barceló

Auxiliar:

Ariel Cáceres

Ayudantes:

Claudio Torres

Jaime Salas

Integrantes:

Juan Andrés Moreno

Tomás Perry

Fecha:

25 de mayo de 2017

Índice general

Índice general	1
1. Resumen	2
2. Introducción	3
2.1. Distribution Sweep	3
3. Metodología	4
3.1. Datos	4
3.1.1. Input	4
3.1.2. Parámetros	4
3.2. Algoritmos	4
3.2.1. Sort	4
3.2.2. Distribution Sweep	5
4. Resultados	6
5. Análisis	8
6. Conclusiones	9
7. Anexos	10
7.1. Códigos	10
7.1.1. Generación de Datos	10
7.1.2. Ordenamiento de Datos	12
7.1.3. Graficar resultados	18

Capítulo 1

Resumen

En este informe se trabaja el problema de “Intersección de segmentos ortogonales”, el cual es utilizado , mayormente , en computación geométrica. Debido a los grandes volúmenes de datos que se manejan en esta area resulta de interés la implementación eficiente de un algoritmo que opere en memoria secundaria. Se entiende por eficiente un algoritmo que optimice la cantidad de accesos a disco.

Primero se explicará el problema en detalle y la metodología ocupada para resolverlo, además de los factores a considerar para obtener una solución óptima. Luego se presentan los resultados. Se analizan estos tomando en cuenta lo que se desea optimizar, intentando identificar los aciertos y errores.

Por ultimo se presentan los resultados obtenidos para ambos algoritmos implementados y para todas las entradas pedidas por enunciado. Con estos, se analiza y se concluye sobre si se logro el orden de accesos a disco y tiempo de ejecución esperados.

Capítulo 2

Introducción

El problema de “Intersección de segmentos ortogonales” consiste en, a partir de un conjunto de segmentos verticales y horizontales, identificar todas las intersecciones entre estos, evidentemente siendo entre segmentos ortogonales. En particular, se trabaja con conjuntos suficientemente grandes para no ser procesables usando únicamente memoria principal, significando que se estudian por subconjuntos extraídos de memoria secundaria.

Para resolver el problema en esta tarea se propone implementar el algoritmo *Distribution Sweep*. Este requiere, entre otras cosas, generar listas de segmentos ordenadas por alguna coordenada. Dado que se trabajan con entradas de tamaño mayor a memoria principal, se tiene entonces que es necesaria la implementación de algún algoritmo de ordenamiento en memoria secundaria. Para esta tarea se decidió implementar Merge Sort.

2.1. Distribution Sweep

Es este el método que finalmente identifica las intersecciones, ocupando *slabs* para separar los segmentos en regiones y, recorriendo un *slab*, ver qué segmentos verticales intersectan a los segmentos horizontales. Este se explicará en detalle en la metodología.

Capítulo 3

Metodología

3.1. Datos

3.1.1. Input

Para estudiar el algoritmo, es necesario entregar una lista de segmentos, horizontales y verticales, definidos por la coordenada de inicio y de término. Para tener entradas suficientemente grandes para justificar la implementación de algoritmos que operen en memoria secundaria se generaron archivos de prueba con 2^i segmentos, donde $i \in \{9, 10, \dots, 21\}$.

Era necesario controlar una cantidad razonable de intersecciones entre segmentos. Para esto se decidió acotar los valores posibles que puede tomar una coordenada por $\pm 2^{\frac{i}{2}}$, con $i \in \{9, 10, \dots, 21\}$. Para el caso de distribuciones uniformes esto fue directo. Por otro lado, para el caso de las distribuciones normales se utilizó $\mu = 0$ y además se ocupó un valor de σ tal que el 90 % de los datos se encontraran en el rango $\pm 2^{\frac{i}{2}}$. Considerando que el 90 % de los datos se encuentran a $\pm 1,65\sigma$ del centro de la distribución, se tiene que con $\sigma = \frac{2^{\frac{i}{2}}}{1,65}$ se logra lo querido. Esto se implementó utilizando la librería *random* en Python. Se generaron listas con una distribución normal para la coordenada x , y repitiendo para distribución uniforme.

Además, se varió la proporción de segmentos verticales y horizontales, tal que

$$\frac{\#_{segmentos_horizontales}}{\#_{segmentos_verticales}} = \frac{\alpha}{(1 - \alpha)} \quad (3.1)$$

con $\alpha \in \{0,25, 0,5, 0,75\}$

3.1.2. Parámetros

Luego de generar los inputs de prueba (con script de Python adjunto a esta tarea) se comprobó que el tamaño de estos no era el necesario para justificar los accesos a memoria secundaria.

Definimos ciertos valores para regular los accesos a disco y verificar la eficiencia de nuestro algoritmo.

1. B: Tamaño del bloque de información que se trae por acceso a disco, en *bytes*.
2. M: Tamaño de memoria secundaria a ocupar, en *bytes*.

Además se cuenta con el parámetro:

1. N: Tamaño del input.

Con el fin de optimizar la forma en que se accede, se utilizan múltiplos de 4 para el *buffer*, asegurando así que no hayan bloques incompletos de información.

3.2. Algoritmos

3.2.1. Sort

Fue necesario diseñar un método de ordenamiento que funcionara en memoria secundaria, pudiendo elegir entre *MergeSort* y *DistributionSort*, optando por *MergeSort*, por la familiaridad con el algoritmo. Al igual que las particiones que se realizan en el *MergeSort* tradicional, en este algoritmo enfocado

en la optimización de accesos a disco, se utilizan $\Theta(m/n)$ “*runs*”, siendo cada uno una dirección a un archivo en memoria secundaria. Dado que cada uno de estos es de tamaño $\mathcal{O}(m)$ se pueden ordenar en memoria principal por separado.

Luego, se realiza el *merge* entre $\Theta(m)$ *runs* de igual tamaño y poniéndolos nuevamente en memoria secundaria. Evidentemente los *runs* dejarán de caber en memoria principal, por lo que no se traen completos, sino se ocupa un *buffer* de entrada para cada uno y un *buffer* de salida para el *run* obtenido, todos los anteriores de tamaño $\mathcal{O}(B)$.

3.2.2. Distribution Sweep

Como dicho antes, es este algoritmo el que identifica las intersecciones. Se identifican los siguientes pasos.

1. Generar una lista ordenada según X de los segmentos, la llamaremos *Lista X*.
2. Generar una lista ordenada según Y de los segmentos, la llamaremos *Lista Y*.
3. Dividir en k *slabs* según la *Lista X*, con una lista para cada uno.
4. Recorriendo la *Lista Y*, iterando en el valor de la coordenada Y :
 - Agrego los segmentos verticales a los *slabs* a los que pertenezcan, ocupando la lista de cada *slab*, la llamaremos l . Una vez el valor en que se va de la iteración es mayor al segmento, se remueve de l .
 - Si es un segmento horizontal, se recorren las listas de los *slabs* a los que el segmento pertenece, sin considerar el *slab* en el que inicia y en el que termina, y se registran todos los segmentos verticales que estan en las listas, guardando las intersecciones.
 - Se repite el algoritmo desde 3 en cada *slab* hasta que este quepa en memoria principal.

Se eligió un $B=4000$ y un $M=40000$.

Capítulo 4

Resultados

Debido a un mal estudio del tiempo y el trabajo necesario, solamente se trabajó el *MergeSort* con acceso a discos, quedando pendiente el algoritmo de *Distribution Sweep*. Se realizó el ordenamiento para los distintos tamaños de listas de segmentos aleatoriamente generados, variando las proporciones, y se midió el número de accesos a discos, comparando entre distribuciones y proporciones entre segmentos horizontales y verticales y una comparación entre todos. Cabe destacar que cada vez que se rellenaba un buffer se hacía un acceso a disco. Por otro lado, lo mismo ocurre al momento de escribir un buffer hacia su archivo correspondiente.

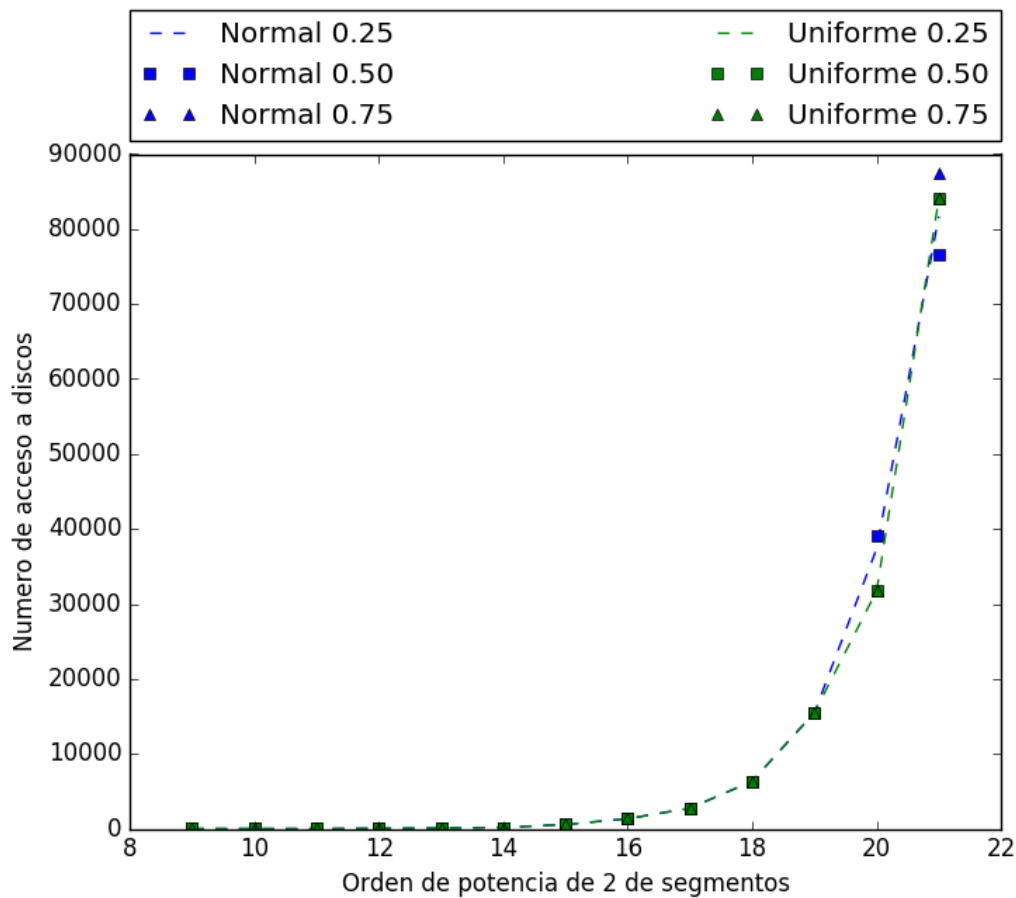
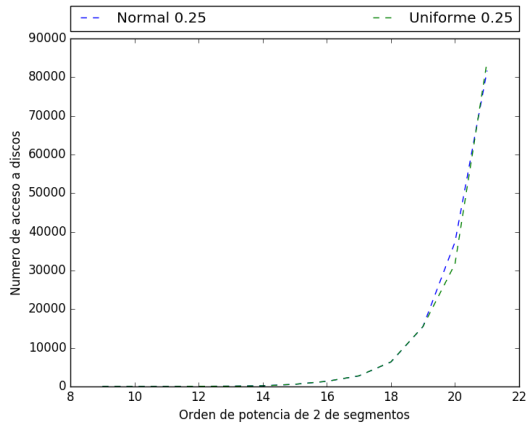
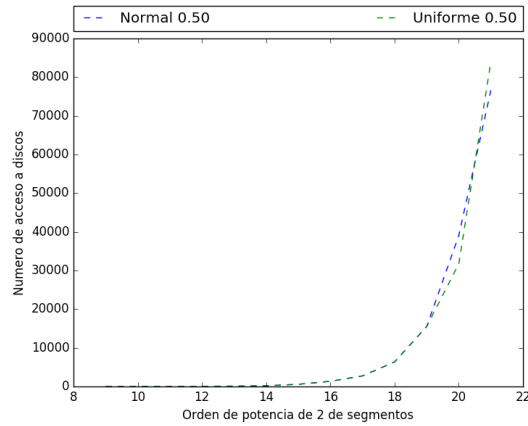


Figura 4.1: Comparación de todos los *mergesorts*.

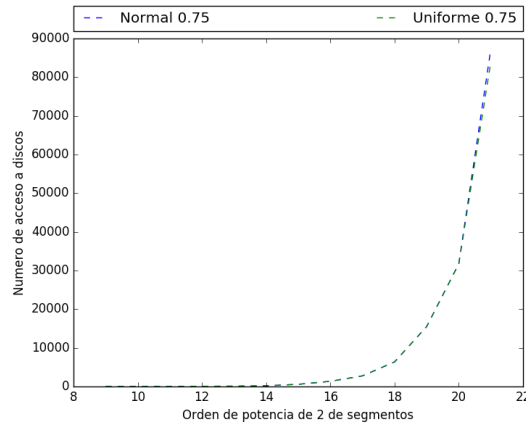
A partir de la Figura 4.1, se observan comportamientos similares entre el ordenamiento de listas con distribución normal como uniforme, teniendo un crecimiento similar, pero se observa que el acceso a discos es igual para todas las distribuciones uniformes de igual tamaño, sin importar la proporción entre segmentos verticales y horizontales, mientras que el caso de la distribución normal se observa que una mayor cantidad de segmentos verticales es mejor.



(a) Comparación para $\alpha = 0,25$



(b) Comparación para $\alpha = 0,50$



(c) Comparación para $\alpha = 0,75$

Figura 4.2: Comparaciones entre distribuciones, según su proporción entre segmentos horizontales y verticales.

A partir de la Figura 4.2, se puede observar que el ordenamiento de una distribución es superior a la otra dependiendo de la proporción entre segmentos verticales y horizontales, creciendo el número de accesos para las distribuciones normales y manteniéndose similar para las distribuciones uniformes. Por simplicidad, no se incluyen los resultados obtenidos al ordenar según la coordenada y , pues estos eran iguales.

Capítulo 5

Análisis

A partir de los datos obtenidos¹, se calcula que los accesos a disco son $\Theta(n \log_m(n))$, con $n = \frac{N}{B}$, $m = \frac{M}{B}$, con una constante cercana a 60, teniendo una constante alta, pero cumpliendo con el orden teórico.

La constante podría deberse a una mala elección de tamaño de bloque y/o memoria secundaria a ocupar, cambiando la base del logaritmo.

La diferencia obtenida entre ordenar una distribución normal y una uniforme podría deber a optimizaciones en el *Branch Predictor* del procesador, siendo mas eficiente en distribuciones uniformes.

¹Adjuntos en el archivo 'results.csv'

Capítulo 6

Conclusiones

A partir de la experiencia, es necesario destacar el *cuello de botella* que se genera al requerir accesos a disco, ralentizando considerablemente la velocidad de cómputo y disminuyendo el uso útil de la CPU. Se señala el impacto que tiene la distribución sobre el ordenamiento, pudiendo deberse a optimizaciones que se hacen a nivel de la arquitectura del computador.

Se rescata la influencia de la proporción entre segmentos sobre el número de accesos a disco, pudiendo observar un beneficio en tener una proporción equitativa.

Capítulo 7

Anexos

7.1. Códigos

7.1.1. Generación de Datos

```
1 import numpy as np
2 import random
3 import pandas as pd
4 import math
5
6 alpha = [0.25, 0.5, 0.75]
7 N = list(range(9,22))
8 distributions = ['uniform', 'normal']
9 file_names = []
10
11 def generate_files():
12     for n in N:
13         for a in alpha:
14             for d in distributions:
15                 f_name = 'input-2-{}-{}-{}.csv'.format(str(n),d,a)
16                 segments = generate_segments(n, d, a)
17                 df = pd.DataFrame(segments)
18                 df.to_csv(f_name, index=False, header=False)
19                 print("Finished generating input of size 2^{}, alpha = {} and distribution = {}\n".format(n, a, d))
20
21 def generate_segments(N, distribution, alpha):
22     vs_qty = int((2**N)*(1-alpha))
23     hs_qty = int((2**N)*alpha)
24
25     vs_ctr = 0
26     hs_ctr = 0
27
28     ##FFFDUsing sets make checking if a segment is in the generated set faster.
29     segments = set()
30     bound = math.floor(2**(N/2))
31     ##
32     ## Vertical Segments generation:
33     while True:
34         ##
35         ##
36         ## First we generate x coordinates with its corresponding distribution:
37         if(distribution == 'uniform'):
38             x_1 = x_2 = math.floor(np.random.uniform(-bound, bound))
39         elif(distribution == 'normal'):
40             ## 1.0 : The idea here is to have 95% of the points between -2^(N/2) and 2^(N/2) .
41             ## Still not sure if it will behave well, we could need to change this.
42             sigma = bound / 1.65
43             x_1 = x_2 = math.floor(np.random.normal(0, sigma))
44
45         ##
46         ##
47
48         ##FFFDNow we generate the y coordinates:
49         ## 1.1 : Not sure if we need to check that y_1 != y_2, could this really not happen?
50         y_1 = math.floor(np.random.uniform(-bound, bound))
51         y_2 = math.floor(np.random.uniform(-bound, bound))
52         while(y_1 == y_2):
53             y_1 = math.floor(np.random.uniform(-bound, bound))
54             y_2 = math.floor(np.random.uniform(-bound, bound))
55
56         ##
57
58         if((x_1, y_1, x_2, y_2) not in segments):
59             segments.add((x_1, y_1, x_2, y_2))
60             vs_ctr += 1
61         if(vs_ctr == vs_qty):
62             break
63
64     ## Horizontal Segments generation:
65     while True:
66
67         ## 1.2 : IDEM 1.1
68         x_1 = math.floor(np.random.uniform(-bound, bound))
69         x_2 = math.floor(np.random.uniform(-bound, bound))
70         while(x_1 == x_2):
71             x_1 = math.floor(np.random.uniform(-bound, bound))
72             x_2 = math.floor(np.random.uniform(-bound, bound))
73
74         y_1 = y_2 = math.floor(np.random.uniform(-bound, bound))
```

```
75         if ((x_1, y_1, x_2, y_2) not in segments):
76             segments.add((x_1, y_1, x_2, y_2))
77             hs_ctr += 1
78         if (hs_ctr == hs_qty):
79             break
80
81     segments = list(segments)
82     ###FFFD13 : If we don't shuffle, then the segments list would end up with all the
83     ### vertical segments first and then all the horizontal ones.
84     np.random.shuffle(segments)
85     return segments
86
87 generate_files()
```

7.1.2. Ordenamiento de Datos

```
1  /*
2  * For this script to run correctly one must have a folder in the same directory filled with the input files.
3  * The name of this folder is passed in as the 3 argument.
4  * Usage: Main B M InputsFolder CoordinateToSortBy
5  * Last test : 21/05/2017
6  * Arguments : 40000 400000 Inputs x
7  * This means, B = 4000 bytes | M = 40000 | inputs folder = Input
8  *
9  * */
10
11 import java.io.IOException;
12 import java.nio.ByteBuffer;
13 import java.nio.channels.FileChannel;
14 import java.nio.file.Files;
15 import java.nio.file.Path;
16 import java.nio.file.Paths;
17 import java.util.*;
18 import java.util.stream.Collectors;
19 import java.util.stream.Stream;
20
21 import static java.nio.file.StandardOpenOption.APPEND;
22 import static java.nio.file.StandardOpenOption.CREATE;
23
24 public class Main {
25
26
27     static int B;
28     static int M;
29     static int m;
30     static Random nameGenerator = new Random();
31     static int numberOfDiskIO = 0;
32
33     public static void main(String[] args) throws IOException{
34
35         // We receive as an argument the directory where we previously stored all the generated inputs.
36         Path inputsFolderPath = Paths.get(args[2]);
37         // Path to output file where we will write the following:
38         // Input-Name.csv , Number of disk IO's.
39         Path resultPath = Paths.get("resultsy.csv");
40
41         B = Integer.parseInt(args[0]);
42         M = Integer.parseInt(args[1]);
43         m = (int) Math.ceil(M/B);
44         char coordToSort = 'y';
45
46
47         // Here we store the paths to all input files.
48         List<Path> inputPaths;
49         try (Stream<Path> paths = Files.walk(inputsFolderPath)) {
50             inputPaths = paths.filter(Files::isRegularFile).collect(Collectors.toList());
51         }
52
53         // We iterate through the input files, calling mergeSort on each one of them and counting
54         // the number of disk IO's.
55         for (Path aPath: inputPaths) {
56             numberOfDiskIO = 0;
57             List<Path> sortedFilePath = externalMergeSort(aPath, coordToSort);
58             // If after the MergeSort procedure, specifically after the merge part we have more than one
59             // sorted run, it means something went wrong with the merge.
60             if(sortedFilePath.size() != 1){
61                 System.out.println("ERROR: Merge ended with more than one resulting sorted run.");
62                 break;
63             }
64             // Construct the line to be written to the file.
65             String line = String.format(aPath.getFileName().toString()+" "+ "%d"+System.lineSeparator(), ←
66                                     numberOfDiskIO);
67             // Write to it.
68             Files.write(resultPath, line.getBytes(), CREATE, APPEND);
69         }
70     }
71
72     /*
73     * External MergeSort procedure. Given a path to the input file to be sorted and the coordinate to it by←
74     * it
75     * proceeds to run the External MergeSort procedure. This complies of the following:
76     * 1. Phase one: Generate n/m sorted runs and write them to secondary memory.
77     * 2. Phase two: At each step merge m sorted runs into one.
78     * Repeat this last step until there is only one sorted run.
79     *
80     * Path inputPath : Path to the input file to be sorted.
81     * char coordToSort : coordinate to sort the file by. (can be 'x' or 'y')
```

```

82 public static List<Path> externalMergeSort(Path inputPath, char coordToSort) throws IOException{
83     FileChannel anInput = FileChannel.open(inputPath);
84     List<Path> phaseOnePaths = phaseOne(anInput, coordToSort);
85     phaseOnePaths = phaseTwo(phaseOnePaths, coordToSort);
86     return phaseOnePaths;
87 }
88
89 /*
90 * Phase one of the external mergeSort algorithm. It takes an input file and produces n/m sorted runs ←
91 * which are
92 * stored in secondary memory for the second phase of the algorithm (Merge). These runs are sorted by the ←
93 * coordinate
94 * indicated by the argument coordToSort.
95 *
96 * FileChannel input : input file.
97 * char coordToSort : coordinate to sort the file by. (can be 'x' or 'y')
98 *
99 * Returns : List of paths to the files where each sorted run was written to.
100 */
101 public static List<Path> phaseOne(FileChannel input, char coordToSort) throws IOException{
102     List<Path> sortedRuns = new ArrayList<>();
103     List<String> linesToWrite;
104     List<int[]> segmentsRun;
105
106     while(Math.abs(input.position() - input.size()) > 1){
107
108         segmentsRun = getSegmentsRun(input);
109         sortSegmentsRun(segmentsRun, coordToSort);
110
111         linesToWrite = segmentsRunToString(segmentsRun);
112
113         Path path = generateSegmentsRunName();
114         sortedRuns.add(path);
115
116         writeSortedRunToSM(linesToWrite, path);
117     }
118
119     return sortedRuns;
120 }
121
122 /*
123 * Phase two of the external mergeSort algorithm. It takes a list of paths to sorted runs and performs ←
124 * phase two steps
125 * until there is only one path in the lists of paths. The amount of paths decreases because at each step ←
126 * of the phase two
127 * we merge O(m) sorted runs into one large sorted run.
128 *
129 * List<Path> phaseOneOutput : Phase one output, ie, a list of paths to the sorted run generated by ←
130 * phaseOne method.
131 * char coordToSort : coordinate to sort the file by. (can be 'x' or 'y')
132 */
133 public static List<Path> phaseTwo(List<Path> phaseOneOutput, char coordToSort) throws IOException{
134     List<Path> output = phaseOneOutput;
135     while(output.size() != 1){
136         output = phaseTwoStep(output, coordToSort);
137     }
138
139     return output;
140 }
141
142 /*
143 * A phase two step consists basically in merging O(m) sorted runs into a single one in secondary memory.
144 * This is done by taking the first m sorted runs in the sortedRuns list and continuously selecting the ←
145 * minimum
146 * between each one of them and storing it in an output buffer. Once this buffer reaches a size of B we ←
147 * flush it
148 * to secondary memory.
149 *
150 * List<Path> sortedRuns : List containing paths to each sorted run generated previously by another phase ←
151 * two step
152 * or the phase one. The idea here is to take m of these and after the merge is completed we delete them.
153 * char coordToSort : coordinate to sort the file by. (can be 'x' or 'y')
154 */
155 public static List<Path> phaseTwoStep(List<Path> sortedRuns, char coordToSort) throws IOException{
156     List<ArrayDeque> inputBuffers = new ArrayList<>();
157     List<FileChannel> inputFiles = new ArrayList<>();
158     /* Path to file where we will store our generated merged run */
159     Path outputFilePath = generateSegmentsRunName();
160
161     /* This for loop basically initializes the input buffers and files. */
162     for (Path pathToInput: sortedRuns) {
163         /* We need to break from the loop when we have m - 1 input buffers (O(m)) */
164         if(inputBuffers.size() == m - 1){
165             break;
166         }
167     }
168     /* FileChannel input to allow us to read from the file from disk whenever its corresponding ←
169     inputbuffer is empty. */

```

```

161     FileChannel input = FileChannel.open(pathToInput);
162     /* The next two lines are the actual stuffing of the input buffer. */
163     List<int[]> segments = segmentsRunToInt(readBlockFromFile(input, B));
164     ArrayDeque segmentsQueue = new ArrayDeque(segments);
165     /* Now we add the input buffer to the input buffers list and the FileChannel to the files list. ←
        */
166     inputBuffers.add(segmentsQueue);
167     inputFiles.add(input);
168 }
169
170 /* We need to keep track of the size of the output buffer so we know when to perform a write to disk ←
        */
171 int outputBufferSize = 0;
172 List<String> outputBuffer = new ArrayList<>();
173
174 /* Iterate until we have consumed all the input buffers. */
175 while(inputBuffers.size() > 0){
176
177     /* Before anything, we need to check if the output buffer is full (= size ~ B).
178     * If it is we need to perform a write to disk operation.
179     */
180
181     if(outputBufferSize < B){
182         /* We get the index associated with the input buffer which contains the minimum segment. */
183         int minLocation = getInputBuffersMinLocation(inputBuffers, new SegmentComparator(coordToSort ←
            ));
184
185         /* Typical parsing. */
186         int [] segment = (int[]) inputBuffers.get(minLocation).removeFirst();
187         String minSegment = String.format("%d,%d,%d,%d", segment[0], segment[1], segment[2], segment ←
            [3]);
188
189         outputBuffer.add(minSegment);
190         outputBufferSize += minSegment.getBytes().length;
191
192         /* This if checks whether the input buffer from which we took the minimum is now empty or ←
            not. */
193         if(inputBuffers.get(minLocation).size() == 0){
194             /* If there is anything left in the corresponding input file then we load a chunk of ←
            data to the input buffer. */
195             if(Math.abs(inputFiles.get(minLocation).position() - inputFiles.get(minLocation).size()) ←
            > 1){
196                 inputBuffers.get(minLocation).addAll(segmentsRunToInt(readBlockFromFile(inputFiles. ←
            get(minLocation), B)));
197             }
198             /* Else, it means that the input buffer is empty and the corresponding file has nothing ←
            new to read from. Thus
199             * we have reached the end of that sorted run and we can delete it. */
200             else{
201                 inputBuffers.remove(minLocation);
202                 Files.delete(sortedRuns.get(minLocation));
203                 inputFiles.remove(minLocation);
204                 sortedRuns.remove(minLocation);
205             }
206         }
207     }
208 } else{
209     writeSortedRunToSM(outputBuffer, outputFilePath);
210     outputBuffer.clear();
211     outputBufferSize = 0;
212 }
213 }
214
215 /* Before returning we write the remaining segments in the output buffer to the output file. */
216 writeSortedRunToSM(outputBuffer, outputFilePath);
217
218 /* We add the newly created sorted run -its path- (product of merging O(m) sorted runs) to the lists ←
    of sorted runs. */
219 sortedRuns.add(outputFilePath);
220
221 return sortedRuns;
222 }
223
224 /*
225 * Takes a list of queues where each one of them represents an input buffer of a sorted run. This method
226 * peeks the first elements of each buffer and returns the index of the list the contains the minimum.
227 *
228 * List<ArrayDeque> inputBuffers : List of queues representing input buffers.
229 * Comparator segmentComparator : segmentComparator class instance used to compare each segment and ←
    decide the minimum.
230 */
231 public static int getInputBuffersMinLocation(List<ArrayDeque> inputBuffers, Comparator segmentComparator ←
    ){
232     int [] min = {Integer.MAX_VALUE, Integer.MAX_VALUE, Integer.MAX_VALUE, Integer.MAX_VALUE};
233     int index = 0;

```

```

237     for (Deque segmentsDeque: inputBuffers) {
238         int [] segment = (int[]) segmentsDeque.getFirst();
239         if(segmentComparator.compare(min, segment) == 1){
240             min = segment;
241             index = inputBuffers.indexOf(segmentsDeque);
242         }
243     }
244
245     return index;
246 }
247
248 /*
249 * Takes a List of segments, where each segment is an array of integers of the form [x-1, y-1, x-2, y-2]
250 * and returns a List of Strings where each component is a string of the form "x-1,y-1,x-2,y-2".
251 *
252 * List<int[]> segmentsRun : List of segments to "convert".
253 * Returns : List of segments in String representation explained above.
254 */
255 public static List<String> segmentsRunToString(List<int[]> segmentsRun){
256     List<String> linesToWrite = new ArrayList<>();
257     for (int[] segment: segmentsRun) {
258         String line = String.format("%d,%d,%d,%d", segment[0], segment[1], segment[2], segment[3]);
259         linesToWrite.add(line);
260     }
261     return linesToWrite;
262 }
263
264 /*
265 * Takes a List of segments, where each segment is a string of the form "x-1,y-1,x-2,y-2"
266 * and returns a List of int[] where each component of the list is an array of ints of the form [x-1, y-1↵
267     , x-2, y-2].
268 *
269 * List<String> segmentsRun : List of segments to "convert".
270 * Returns : List of segments in int[] representation explained above.
271 */
272 public static List<int[]> segmentsRunToInt(List<String> segmentsRun){
273     List<int[]> segments = new ArrayList<>();
274     for (String line: segmentsRun) {
275         int [] segment = Stream.of(line.split(",")).mapToInt(Integer::parseInt).toArray();
276         segments.add(segment);
277     }
278     return segments;
279 }
280
281 /*
282 * Generates a name for storing a sorted run in secondary memory.
283 * To generate random names it used nameGenerator, which is an instance of Random.
284 * Returns : Path containing name where a sorted run will be stored in secondary memory.
285 */
286 public static Path generateSegmentsRunName(){
287     return Paths.get(String.format("%d.csv", nameGenerator.nextInt()));
288 }
289
290 /*
291 * Given an input File of size N, we wish to take chunks of size B from it and start the sorting process.
292 * This method returns a buffer filled with the segments inside the chunk of size B of the input file.
293 * FileChannel input : Input file. Using a FileChannel allows us read B bytes from the file.
294 * int B : Size of chunks we read from the file (bytes).
295 *
296 * Returns : List<String> where each component is a String representing a segment. Each String looks like↵
297     "x-1,y-1,x-2,y-2"
298 */
299 public static List<String> readBlockFromFile(FileChannel input, int B) throws IOException {
300
301     /*
302     * A big problem in reading text files by a certain amount of bytes is dealing with line cuts.
303     * It could happen that we wish to read X bytes from the file and that means reading a non integer
304     * amount of lines. In this case we need to rewind the file pointer to the last \n read. So we can ↵
305         start
306     * from there the next time.
307     */
308
309     numberOfDiskIO++;
310
311     boolean EOBBuffer = false;
312     // Here we create our byte buffer with a fixed size B.
313     ByteBuffer byteBuffer = ByteBuffer.allocate(B);
314     // We read from the file. Starting from the last position we were in.
315     input.read(byteBuffer);
316     // To read from the buffer we need to move the pointer position to the beggining.
317     byteBuffer.flip();
318
319     // Lines is the list in which we will store the segments read from the chunk of size B.
320     List<String> lines = new ArrayList<>();
321     String lineString = "";
322     char c;
323
324     // We read from the buffer while it still has bytes to read.

```



```

322 while(byteBuffer.hasRemaining()){
323
324     // We get the first char from the buffer and check if it is a end of line.
325     c = (char)byteBuffer.get();
326     if(c == System.lineSeparator().charAt(0) && !byteBuffer.hasRemaining()){
327         EOBuffer = true;
328     }
329     // This case happens when the last chunk read left us at the end of a line.
330     if(c == System.lineSeparator().charAt(0) && byteBuffer.hasRemaining()){
331         c = (char)byteBuffer.get();
332     }
333     while(c != System.lineSeparator().charAt(0)){
334         lineString += c;
335         /*
336          * While we don't see an end of line we keep on reading the line. Mind the fact that if we ←
337          * are in a line
338          * cutted by the read method we would never see an end of line.
339          * For this, if the bytebuffer does not have any more bytes to get and we haven't reached and←
340          * end of line
341          * it means we need to stop reading.
342          */
343         if(!byteBuffer.hasRemaining()){
344             EOBuffer = true;
345             break;
346         }
347         c = (char)byteBuffer.get();
348     }
349     if(EOBuffer){
350         break;
351     }
352     /*
353     * This one-liner simply gets the lineString of the form "x1,y1,x2,y2", strips it by ',' and ←
354     * converts each
355     * component to an integer, then an array of integers.
356     */
357     //System.out.println(lineString);
358     //segments = Stream.of(lineString.split(",")).mapToInt(Integer::parseInt).toArray();
359     lines.add(lineString);
360     lineString = "";
361 }
362
363 /*
364 * Now if we ended up in the middle of a line we need to calculate how much we need to rewinding the ←
365 * input file
366 * position, so we can read the whole line the next time.
367 */
368 int fileRewind = 0;
369 bufferPosition = byteBuffer.position()-1;
370 char lastCharRead = (char)byteBuffer.get(bufferPosition);
371 if(lastCharRead != System.lineSeparator().charAt(0)){
372     bufferPosition--;
373     lastCharRead = (char)byteBuffer.get(bufferPosition);
374     while(lastCharRead != System.lineSeparator().charAt(0)){
375         fileRewind++;
376         bufferPosition--;
377         lastCharRead = (char)byteBuffer.get(bufferPosition);
378     }
379 }
380 // Finally, we update de input file's position to the end of the last line we were able to read ←
381 // completely.
382 input.position(input.position() - fileRewind - 1);
383
384 return lines;
385 }
386
387 /*
388 * Given segmentsRun with each component inside of it of the form [x-1, y-1, x-2, y-2]
389 * it uses the SegmentComparator to sort the list given a coordinate (coordToSort)
390 *
391 * int [] segmentsRun : a run of unordered segments.
392 * char coordToSort : coordinate to sort by segments. (Values can be 'x' or 'y')
393 *
394 * Its void since the sort is done in-place.
395 */
396 public static void sortSegmentsRun(List<int[]> segmentsRun, char coordToSort){
397     Collections.sort(segmentsRun, new SegmentComparator(coordToSort));
398 }
399
400 /*
401 * Reads m times blocks of size B from the input. Each block is a set of lines in the input.
402 * It puts them all together in the segmentRun List.
403 *
404 * FileChannel input: Input file.
405 *
406 * Returns: List<int[]> of segments. Each component of the list a segment represented in an array of ints
407 * of the form [x-1, y-1, x-2, y-2]

```

```

405  * */
406  public static List<int[]> getSegmentsRun(FileChannel input) throws IOException{
407      List<String> blockReadFromFile;
408      List<int[]> segmentsRun = new ArrayList<>();
409
410      for (int i = 0; i < m; i++) {
411          blockReadFromFile = readBlockFromFile(input, B);
412          for (String line: blockReadFromFile) {
413              int [] segment = Stream.of(line.split(",")).mapToInt(Integer::parseInt).toArray();
414              segmentsRun.add(segment);
415          }
416      }
417
418      return segmentsRun;
419  }
420
421  /*
422  * Gets a sorted run of segments and writes it to secondary memory in the specified path.
423  * List<String> sortedRun : List of sorted segments (single run).
424  * Path path : Path to file in which the sorted run is going to be written.
425  * */
426  public static void writeSortedRunToSM(List<String> sortedRun, Path path) throws IOException{
427      numberOfDiskIO++;
428      Files.write(path, sortedRun, CREATE, APPEND);
429  }
430
431  /*
432  * Comparison function for sorting the random sample in getPivots.
433  * Constructor: SegmentComparator
434  * */
435  static private class SegmentComparator implements Comparator<int[]> {
436
437      char cts;
438      /*
439      * Constructor for the class Comparator.
440      * coordToSort : coordinate by which we are sorting the file (x or y).
441      * */
442      private SegmentComparator(char coordToSort){
443          cts = coordToSort;
444      }
445
446      /*
447      * segmentOne : First segment involved in comparison.
448      * segmentTwo : Second segment involved in comparison.
449      * */
450      public int compare(int[] segmentOne, int[] segmentTwo){
451          // Notice that the homework pdf states that there a some difference in sorting by x or y.
452          if (cts == 'x'){
453              if (Integer.compare(segmentOne[0], segmentTwo[0]) == 0){
454                  return Integer.compare(segmentOne[2], segmentTwo[2]);
455              }else{
456                  return Integer.compare(segmentOne[0], segmentTwo[0]);
457              }
458          }
459
460          // Here we need to check if a segment is vertical or horizontal to decide ties between segments.
461          if (cts == 'y'){
462              // Case 1: The two segments both have the same y coordinates, ie:
463              // S1.y1 = S2.y1 and S1.y2 = S2.y2
464              if (Integer.compare(segmentOne[1], segmentTwo[1]) == 0 && Integer.compare(segmentOne[3], ←
465                  segmentTwo[3]) == 0){
466
467                  // Note that a segment is vertical if both of its y coordinates are the same.
468                  // So if Segment one is vertical, we claim it larger by comparison.
469                  if(Integer.compare(segmentOne[0], segmentOne[2]) == 0){
470                      return 1;
471                  }else{
472                      return -1;
473                  }
474              }
475              // Case 2: The two segments differ in at least one y coordinate.
476          }else{
477              if (Integer.compare(segmentOne[1], segmentTwo[1]) == 0){
478                  return Integer.compare(segmentOne[3], segmentTwo[3]);
479              }else{
480                  return Integer.compare(segmentOne[1], segmentTwo[1]);
481              }
482          }
483      }
484
485      return 0;
486  }
487
488  }
489
490  }
    
```

7.1.3. Graficar resultados

```
1 import csv
2 import matplotlib.pyplot as plt
3
4
5 # Auxiliary function to get the second value of a tuple
6 def getKey(item):
7     return item[1]
8
9
10 # Plotter for the results
11 def plot(coord):
12     data = open('results'+coord+'.csv', 'rb')
13     reader = csv.reader(data, delimiter=',')
14     uniform25 = list()
15     normal25 = list()
16     uniform50 = list()
17     normal50 = list()
18     uniform75 = list()
19     normal75 = list()
20
21     # Separate the results in corresponding distribution and proportion
22     for row in reader:
23         file = row[0]
24         value = row[1]
25         att = file.split('-')
26         if att[3] == "normal":
27             if att[4] == "0.25.csv":
28                 normal25.append([int(value), int(att[2])])
29             elif att[4] == "0.5.csv":
30                 normal50.append([int(value), int(att[2])])
31             else:
32                 normal75.append([int(value), int(att[2])])
33         else:
34             if att[4] == "0.25.csv":
35                 uniform25.append([int(value), int(att[2])])
36             elif att[4] == "0.5.csv":
37                 uniform50.append([int(value), int(att[2])])
38             else:
39                 uniform75.append([int(value), int(att[2])])
40
41     # Remove repeated values
42     old_uniform25 = uniform25
43     uniform25 = []
44     old_uniform50 = uniform50
45     uniform50 = []
46     old_uniform75 = uniform75
47     uniform75 = []
48
49     old_normal25 = normal25
50     normal25 = []
51     old_normal50 = normal50
52     normal50 = []
53     old_normal75 = normal75
54     normal75 = []
55
56     for x in old_uniform25:
57         if x not in uniform25:
58             uniform25.append(x)
59
60     for x in old_uniform50:
61         if x not in uniform50:
62             uniform50.append(x)
63
64     for x in old_uniform75:
65         if x not in uniform75:
66             uniform75.append(x)
67
68     for x in old_normal25:
69         if x not in normal25:
70             normal25.append(x)
71
72     for x in old_normal50:
73         if x not in normal50:
74             normal50.append(x)
75
76     for x in old_normal75:
77         if x not in normal75:
78             normal75.append(x)
79
80     # Sort values and only save accesses
81     uniform25 = [x[0] for x in sorted(uniform25, key=getKey)]
82     uniform50 = [x[0] for x in sorted(uniform50, key=getKey)]
83     uniform75 = [x[0] for x in sorted(uniform75, key=getKey)]
```

```
84 normal25 = [x[0] for x in sorted(normal25, key=getKey)]
85 normal50 = [x[0] for x in sorted(normal50, key=getKey)]
86 normal75 = [x[0] for x in sorted(normal75, key=getKey)]
87 numbers = range(9, 22)
88
89 # Plot comparing all the tests
90 plt.plot(numbers, normal25, 'b—', label='Normal 0.25')
91 plt.plot(numbers, normal50, 'bs', label='Normal 0.50')
92 plt.plot(numbers, normal75, 'b^', label='Normal 0.75')
93 plt.plot(numbers, uniform25, 'g—', label='Uniforme 0.25')
94 plt.plot(numbers, uniform50, 'gs', label='Uniforme 0.50')
95 plt.plot(numbers, uniform75, 'g^', label='Uniforme 0.75')
96 plt.legend(bbox_to_anchor=(0., 1.02, 1., .102), loc=3, ncol=2, mode="expand", borderaxespad=0.)
97 plt.xlabel('Orden de potencia de 2 de segmentos')
98 plt.ylabel('Numero de acceso a discos')
99 plt.savefig('all'+coord+'.png', bbox_inches='tight')
100 plt.clf()
101
102 # Plot comparing test with alpha = 0.25
103 plt.plot(numbers, normal25, 'b—', label='Normal 0.25')
104 plt.plot(numbers, uniform25, 'g—', label='Uniforme 0.25')
105 plt.legend(bbox_to_anchor=(0., 1.02, 1., .102), loc=3, ncol=2, mode="expand", borderaxespad=0.)
106 plt.xlabel('Orden de potencia de 2 de segmentos')
107 plt.ylabel('Numero de acceso a discos')
108 plt.savefig('25'+coord+'.png')
109 plt.clf()
110
111 # Plot comparing test with alpha = 0.50
112 plt.plot(numbers, normal50, 'b—', label='Normal 0.50')
113 plt.plot(numbers, uniform50, 'g—', label='Uniforme 0.50')
114 plt.legend(bbox_to_anchor=(0., 1.02, 1., .102), loc=3, ncol=2, mode="expand", borderaxespad=0.)
115 plt.xlabel('Orden de potencia de 2 de segmentos')
116 plt.ylabel('Numero de acceso a discos')
117 plt.savefig('50'+coord+'.png')
118 plt.clf()
119
120 # Plot comparing test with alpha = 0.75
121 plt.plot(numbers, normal75, 'b—', label='Normal 0.75')
122 plt.plot(numbers, uniform75, 'g—', label='Uniforme 0.75')
123 plt.legend(bbox_to_anchor=(0., 1.02, 1., .102), loc=3, ncol=2, mode="expand", borderaxespad=0.)
124 plt.xlabel('Orden de potencia de 2 de segmentos')
125 plt.ylabel('Numero de acceso a discos')
126 plt.savefig('75'+coord+'.png')
127 plt.clf()
128
129 # Plot for x and y tests
130 plot('x')
131 plot('y')
```