

Assignment 4: Improving Efficiency with Sorting

Assigned: Wednesday, September 24, 2014

Due: Friday, October 3, 2014, 11:59 p.m.

Type: Individual

Problem Overview

This assignment will investigate an example *feature extraction* problem. Feature extraction is a subproblem of pattern recognition and is also used in areas such as statistical analysis, computer vision, and image processing. For example, an image processing problem may use a feature extraction algorithm to identify particular shapes or regions in a digitized image.

In this assignment, we're going to focus on a very simple feature extraction problem: Given a set of points in two-dimensional space, identify every subset of four or more points that are *collinear*. For example, given the set of points depicted in Figure 1, your program would detect the three groups of collinear points as depicted by the line segments in Figure 2.

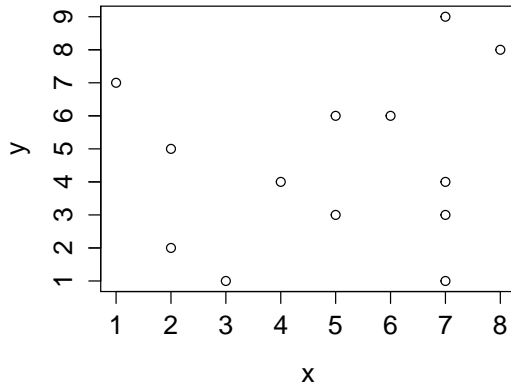


Figure 1: A set of 13 points.

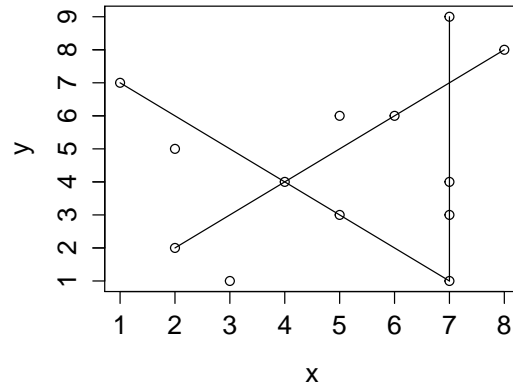


Figure 2: Three collinear groups identified.

As always, we want our solution to be useful at scale. For example, Figure 3 plots ~100,000 points and Figure 4 shows the 34 collinear groups identified by blue line segments. Each collinear group in Figure 4 is composed of far more than four points; four is just the minimum number of points to qualify for the collinear pattern that we're looking for. In the general problem statement we will refer to line segments instead of collinear groups, where each line segment must contain at least four points.

Problem Statement: Given a set of N distinct points in the plane, identify every line segment that connects a subset of four or more of the points. Each point will be specified by an (x, y) pair where x and y are int values in the range 0 to 32,767. For example, the thirteen points in Figures 1 and 2 are: (1, 7), (2, 2), (2, 5), (3, 1), (4, 4), (5, 3), (5, 6), (6, 6), (7, 1), (7, 3), (7, 4), (7, 9), (8, 8).

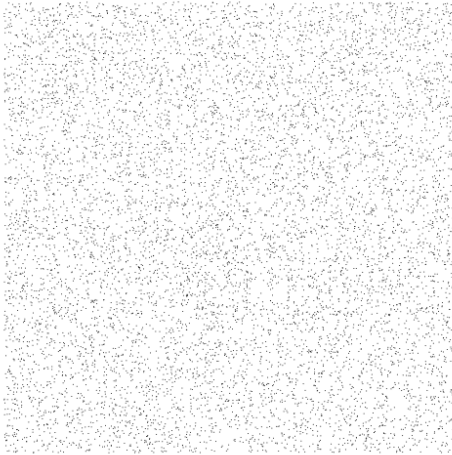


Figure 3: A set of ~100,000 points.



Figure 4: 34 collinear groups identified.

You must solve this problem in terms of the classes and methods described in the following sections.

The Point class

You must create an immutable data type `Point` that represents a point in the plane. A shell of the `Point` class is provided for you, and its API is described below.

```
public class Point implements Comparable<Point> {
    public final Comparator<Point> SLOPE_ORDER;

    public Point(int x, int y)

    public void draw()
    public void drawTo(Point that)
    public double slopeTo(Point that)
    public int compareTo(Point that)
    public String toString()
}
```

The constructor, `draw` method, `drawTo` method, and `toString` method have been completed for you and must not be changed. You must add the bodies of the remaining methods yourself. You may add any number of private methods that you like, but you may not add any public method or constructor, nor may you change the signature of any public method or constructor.

The methods in the `Point` API that you must complete are described in more detail below.

The `compareTo` method.

This method must compare points by y -coordinates, breaking ties by x -coordinates. Thus, the invoking point (x_0, y_0) is less than the parameter point (x_1, y_1) if and only if either $y_0 < y_1$ or if $y_0 = y_1$ and $x_0 < x_1$.

For example, by this *natural order* of points, $(0, 1)$ is less than $(0, 2)$, $(7, 1)$ is less than $(5, 3)$, and $(3, 0)$ is less than $(4, 0)$.

The slopeTo method.

This method must return the slope between the invoking point (x_0, y_0) and the parameter point (x_1, y_1) , which is given by the formula:

$$\frac{(y_1 - y_0)}{(x_1 - x_0)}$$

For example, for the point (3, 3), the slope to (1, 1) is 1.0, the slope to (4, 5) is 2.0, and the slope to (5, 2) is -0.5.

Treat the slope of a horizontal line segment as positive zero¹; treat the slope of a vertical line segment as positive infinity¹; treat the slope of a degenerate line segment (between a point and itself) as negative infinity¹.

The SLOPE_ORDER Comparator

This field of the Point class must compare two points by the slopes they make with the invoking point (x_0, y_0) . Thus, the point (x_1, y_1) is less than the point (x_2, y_2) if and only if

$$\frac{(y_1 - y_0)}{(x_1 - x_0)} < \frac{(y_2 - y_0)}{(x_2 - x_0)}$$

Treat horizontal, vertical, and degenerate line segments the same as in the slopeTo() method.

For example, if the invoking point is (3, 3), then (5, 2) is less than (1, 1), and (1, 1) is less than (4, 5).

You will need to write a nested inner class that implements that slope-order behavior as a Comparator<Point>, and then set the SLOPE_ORDER field to be an instance of this class.

The Line Class

This class models a line segment as a set of points.

```
public class Line implements Comparable<Line>, Iterable<Point> {
    public Line()
    public Line(Collection<Point> l)
    public void add(Point p)
    public Point first()
    public Point last()
    public int length()
    public Iterator<Point> iterator()
    public int compareTo(Line that)
    public String toString()
}
```

The Constructors

The parameterless constructor creates a new line with no points. The second constructor creates a new line that contains all the distinct collinear points in the Collection parameter.

¹See the Java documentation of the Double class for a discussion of *positive zero*, *positive infinity*, and *negative infinity*.

The add method

The add method adds the given Point to the line, provided it is collinear with the existing points in the line and it isn't already present.

The first method

The first method returns the point on the line that is least with respect to natural ordering of Point. If the line has no points, first returns null.

The last method

The last method returns the point on the line that is greatest with respect to natural ordering of Point. If the line has no points, last returns null.

The length method

The length method returns the number of points on the line.

The iterator method

The iterator method returns an Iterator over the points on the line. The iteration order is natural order of Point.

The compareTo method

The compareTo method compares this line with the parameter line. The natural order of Line is based on the natural order of its first and last point. That is, given $line_1$ and $line_2$, $line_1 < line_2$ if $line_1.first < line_2.first$ or $line_1.first = line_2.first$ and $line_1.last < line_2.last$. Two lines $line_1$ and $line_2$ are equal if and only if $line_1.first = line_2.first$ and $line_1.last = line_2.last$.

The Extractor Class

This class provides methods that allow clients to read in an input file of Point data and find all lines segments of four or more collinear points in that data. The API of this class is described below, and as in the case of the Point class above, you may not modify this API in any way.

```
public class Extractor {
    public Extractor(String filename)
    public Extractor(Collection<Point> c)
    public void drawPoints()
    public void drawLines()
    public Set<Line> getLinesBrute()
    public Set<Line> getLinesFast()
}
```

The Constructors

The first constructor for the `Extractor` class takes a single parameter of type `String`. This parameter is a filename for a file of `Point` data formatted as follows: The first line of the file contains a single `int` value N that is the number of lines `Point` data that follow. Each of the following N lines contains two `int` values separated by one or more blanks. The first `int` is the x value of a `Point` and the second `int` is the y value of a `Point`. There may be lines of text past these first $N + 1$ lines of data, but they should be ignored.

A sample input file is shown below.

```
5
11000 11000
12000 10000
13000 10000
14000 10000
15000 10000
```

Instantiating an `Extractor` object with this data file would ensure that five distinct instances of the `Point` class are stored in a suitable data structure inside the new `Extractor` object.

The second constructor takes a `Collection` of points and creates an `Extractor` for this data.

The drawPoints method

This method uses the `StdDraw` class in the provided `stdlib.jar` file to graphically display all the points associated with this `Extractor` object. Specifically, this method should iterate over all the `Point` objects in the `Extractor` object's internal data structure and invoke each `Point`'s `draw` method.

This method is only to help you visualize the data that your program is processing. It will not be invoked by the grading program, and thus it will not affect your grade in any way.

The drawLines method

This method uses the `StdDraw` class in the provided `stdlib.jar` file to graphically display all the lines already identified by this `Extractor` object, if any. If no lines have yet been identified, this method will still raise a graphics window with no lines drawn. If lines have been identified, then this method can use the `drawTo` method of the `Point` object at one end of the line segment to draw a line from that `Point` to the `Point` object at the other end of the line segment.

This method is only to help you visualize the data that your program is processing. It will not be invoked by the grading program, and thus it will not affect your grade in any way.

The getLinesBrute method

This method implements the straight-forward, *brute force* approach to extracting the feature that we're interested in. Since *any* combination of four distinct points that are collinear qualify as our feature, we could generate *all* combinations of four distinct points and check each one to see if those four points are collinear. So, we could describe this brute force solution as a *combinatoric* approach to the problem: We're generating the combination of N things taken four at a time, and each time we generate a new combination, we're testing it based on our feature criteria (collinearity).

For example, let's name the points in the given sample input file p_1 through p_5 , as shown below.

```
5
11000 11000 ( $p_1$ )
12000 10000 ( $p_2$ )
```

13000 10000 (p_3)
 14000 10000 (p_4)
 15000 10000 (p_5)

The table below shows the combinations that our code would generate, along with the result of testing each combination for collinearity. Note that to check if four points p , q , r , and s are collinear, check whether the slope between p and q , between p and r , and between p and s are all equal.

Combination	Collinear?
p_1, p_2, p_3, p_4	no
p_1, p_2, p_3, p_5	no
p_1, p_2, p_4, p_5	no
p_1, p_3, p_4, p_5	no
p_2, p_3, p_4, p_5	yes

The advantage of this approach is that it's fairly simple to code. Since the number of points being selected out of the set of N total points is fixed at four, then four nested for loops should do the trick. But of course, that's also the problem with this approach: four nested loops each dependent on N will have $O(N^4)$ time complexity.

In fact, since this is a combinatoric solution we can calculate exactly how many combinations of four points will be computed for a given N .

$$\binom{N}{4} = \frac{N!}{4!(N-4)!}$$

For our example above, this would give $\frac{5!}{4!} = 5$. If the input had 10 points, the brute force solution would have to test $\frac{10!}{4! \times 6!} = \frac{10 \times 9 \times 8 \times 7}{4 \times 3 \times 2} = 210$ different combinations of four points. For $N = 20$, the brute force solution would generate and test 4845 combinations of four points. For $N = 1000$, over 41 billion combinations of four points would be generated and tested by our program. You can see how this escalates very quickly and the brute force solution becomes infeasible to apply for even moderately large values of N .

The getLinesFast method

A fundamental property of sorting is that it brings duplicates together. We can make use of this and solve the problem much faster if we use sorting as part of our solution. Collinear points have the same slope with respect to each other, and thus are duplicates with respect to SLOPE_ORDER.

To see if point p is part of a group of four or more collinear points we can do the following.

1. Sort the N points with respect to the slope that they make with p .
2. Scan the sorted points to find all groups of three or more consecutive points having the same slope to p . Each such group is collinear with p and is thus, together with p , part of a line segment of at least four points.
3. Repeat for the remaining $N - 1$ points.

Here's an example of sorting the points in Figure 1 with respect to the slope they make with $(7, 1)$.

(7, 1)	(6, 6)	(5, 6)	(1, 7)	(4, 4)	(5, 3)	(2, 5)	(2, 2)	(3, 1)	(8, 8)	(7, 3)	(7, 4)	(7, 9)

Note how this sorting brings together the points of the two line segments that contain (7, 1) (underlined in gray).

How much faster is this sort-and-scan approach? We can sort in $O(N \log N)$ time and the subsequent scan is $O(N)$. We have to perform these operations for all N points, so the total cost of this *sort and scan* approach is $N \times (N \log N + N)$ which is $O(N^2 \log N)$. This is a significant asymptotic improvement since $O(N^2 \log N) \prec O(N^4)$, and the clock-time difference is dramatic. Problem sizes that are infeasible for the brute force solution are solved quickly (or at least in reasonable amount of time) by this sort-and-scan solution.

There is an extra benefit of this approach: We are not limited to identifying four-point segments. We can now identify *maximal* line segments of four or more collinear points.

Notes and other requirements

Here are a couple of extra requirements plus a few things to keep in mind.

- **Start this one early.** There's more reading, thinking, and up-front understanding to take care of on this assignment. Read this handout carefully. Ask questions of your TA and of me. Ask questions on Piazza. Start early and be proactive.
- You've been provided with the file `stdlib.jar`, which contains easy to use input/output and drawing capabilities along with other things. This file is provided by Princeton University under the GNU General Public License. Its documentation can be found here: <http://introcs.cs.princeton.edu/java/stdlib>. *You are not required to use this file at all.* If you choose to do so, you will have to add the path to this jar file to your CLASSPATH.
- You've been provided with shells of the `Point`, `Line`, and `Extractor` classes. *Don't modify the things that have been done for you.*
- Start by completing the remaining methods of the `Point` class. There's no *point*² in attempting the `List` or `Extractor` class before this is complete and correct.
- Complete the `List` class after `Point` but before `Extractor`.
- When you begin work on the `Extractor` class, start with the `getLinesBrute` method. This is shorter and easier to get correct quickly. After you can produce correct output with the `getLinesBrute` method, turn your attention to `getLinesFast`.
- In `getLinesFast`, do not print subsegments of lines containing five or more collinear points. For example, if the line segment $p \rightarrow q \rightarrow r \rightarrow s \rightarrow t$ exists in the data, identify it but not any four-point subsegment such as $p \rightarrow q \rightarrow r \rightarrow s$.

Assignment Submission

You must turn in the following three files to Web-CAT for grading: `Point.java`, `Line.java`, and `Extractor.java`. While not required, it is strongly recommended that you also submit your own test cases. Note the following rules regarding your Web-CAT submissions:

²Sorry; couldn't resist. Here's more: www.punoftheday.com

- Separate submissions for Point and Line are available on Web-CAT. You can submit an unlimited number of times and these submissions are not counted toward your grade. This is strictly to help you make sure Point and Line work correctly before you make submissions that include the Extractor class.
- The feedback hints for failed test cases will be unavailable beginning one day prior to the assignment due date.
- You can submit to Web-CAT no more than eight times for this assignment.
- The *last* submission that you make to Web-CAT will be used to determine your grade on the assignment, even if its score is lower than that of an earlier submission.
- Submissions made within the 24 hour period after the published deadline will be assessed a late penalty of 15 points.
- No submissions will be accepted more than 24 hours after the published deadline.