Closest Pair Algorithm

https://cs.pomona.edu/classes/cs140/

Outline

Topics and Learning Objectives

- Learn more about Divide and Conquer paradigm
- Learn about the closest-pair problem and its O(n lg n) algorithm
 - Gain experience analyzing the run time of algorithms
 - Gain experience proving the correctness of algorithms

<u>Exercise</u>

• Closest Pair

Extra Resources

• Algorithms Illuminated: Part 1: Chapter 3

Closest Pair Problem

- Input: P, a set of n points that lie in a (two-dimensional) plane
- <u>Output</u>: a pair of points (p, q) that are the "closest"
 Distance is measured using Euclidean distance:

$$d(p, q) = sqrt((p_x - q_x)^2 + (p_y - q_y)^2)$$

• Assumptions: None

Closest Pair Problem



- What is the brute force method for this search?
- What is the asymptotic running time of the brute force method?

Input

p1 p2 p3 p4 p5 p6 p7

One-dimensional closest pair

How would you find the closest two points?

- Sort by position : O(n lg n) p6 p4 p1 p3 p5 p7 p2
- Return the closest two using a linear scan : O(n)
- Total time : $O(n \lg n) + O(n) = O(n \lg n)$

Any problems using this approach for the two-dimensional case?

- <u>Sorting does not generalize to higher dimensions!</u>
- How do you sort the points?





y

1. Which two are closest on the y-axis?



1. Which two are closest on the y-axis?

y

12



y

1. Which two are closest on the y-axis?



1. Which two are closest on the y-axis?

2. Which two are closest on the x-axis?



1. Which two are closest on the y-axis?

2. Which two are closest on the x-axis?



Y

1. Which two are closest on the y-axis?

2. Which two are closest on the x-axis?

3. Which two are closest?



Y

1. Which two are closest on the y-axis?

2. Which two are closest on the x-axis?

3. Which two are closest?

Closet Pair—Two-Dimensions

- 1. Create a copy of the points (we now have two separate copies of P)
 - 1. Sort by x-coordinate
 - 2. Sort other by y-coordinate



Now we know we can't do better than O(n lg n)

Sorted by x coordinate

Px : [p0(1,10), p1(2,8), p5(3,5), p3(5,7), p2(7,3), p4(8,4), p7(9,1), p6(10,9)]

Sorted by y coordinate

Py : [p7(9,1), p2(7,3), p4(8,4), p5(3,5), p3(5,7), p1(2,8), p6(10,9), p0(1,10)]



Closet Pair—Two-Dimensions

- 1. Create a copy of the points (we now have two separate copies of P)
 - 1. Sort by x-coordinate
 - 2. Sort other by y-coordinate



- Can we still end up with a O(n lg n) algorithm for finding the closest pair?
- Does the closeness of two points on one axis matter?

1. FUNCTION FindClosestPair(points)

- 2. points_x = copy_and_sort_by_x(points)
- 3. points_y = copy_and_sort_by_y(points)
- 4. **RETURN** ClosestPair(points_x, points_y)

Closet Pair—Two-Dimensions

- 1. Create a copy of the points (we now have two separate copies of P)
 - 1. Sort by x-coordinate
 - 2. Sort other by y-coordinate



- Can we still end up with a O(n lg n) algorithm for finding the closest pair?
- Does the closeness of two points on one axis matter?

2. Apply the Divide-and-Conquer method

Divide-and-Conquer

- 1. **DIVIDE** into smaller subproblems
- 2. CONQUER the subproblems via recursive calls
- 3. COMBINE solutions from the subproblems

• How would you divide the problems?



Y

- 1. Which two are closest on the y-axis?
- 2. Which two are closest on the x-axis?
- 3. Which two are closest?
- 4. How would you divide the search space?



V

```
1. Which two are closest on the y-axis?
```

- 2. Which two are closest on the x-axis?
- 3. Which two are closest?
- 4. How would you divide the search space?

This is the median x-value This is not the average x-value

```
1.
                                                              FUNCTION FindClosestPair(points)
                                                          2.
                                                                points_x = copy_and_sort_by_x(points)
1.
      FUNCTION ClosestPair(px, py)
                                                          3.
                                                                points_y = copy_and_sort_by_y(points)
2.
        n = px.length
                                                                RETURN ClosestPair(points x, points y)
                                                          4.
3.
        # What is the base case?
4.
        IF n == 2
5.
          RETURN px[0], px[1], dist(px[0], px[1])
6.
7.
8.
9.
        # What are the recursive cases?
10.
        pl, ql, dl = ClosestPair(left_px, left_py)
11.
12.
                                                How do we create these arrays?
13.
14.
        pr, qr, dr = ClosestPair(right_px, right_py)
```



```
1. FUNCTION ClosestPair(px, py)
```

```
2. n = px.length
```

```
3. IF n == 2
```

5.

9.

```
4. RETURN px[0], px[1], dist(px[0], px[1])
```

```
6. left_px = px[0 .. < n//2]
```

- 7. left_py = [p FOR p IN py IF p.x < px[n//2].x]
- 8. pl, ql, dl = ClosestPair(left_px, left_py)

```
10. right_px = px[n//2 ..< n]
```

- 11. right_py = [p FOR p IN py IF p.x \ge px[n//2].x]
- 12. pr, qr, dr = ClosestPair(right_px, right_py)

What is the running time of these operations?

Median x value



Any problems with our current approach?

```
1. FUNCTION ClosestPair(px, py)
```

```
2. n = px.length
```

```
3. IF n == 2
```

```
4. RETURN px[0], px[1], dist(px[0], px[1])
```

```
6. left_px = px[0 .. < n//2]
```

- 7. left_py = [p FOR p IN py IF p.x < px[n//2].x]
- 8. pl, ql, dl = ClosestPair(left_px, left_py)

```
10. right_px = px[n//2 ..< n]
```

```
11. right_py = [p FOR p IN py IF p.x \ge px[n//2].x]
```

```
12. pr, qr, dr = ClosestPair(right_px, right_py)
```

```
13.
```

5.

9.

```
14. d = min(dl, dr)
```

- 15. ps, qs, ds = ClosestSplitPair(px, py, d)
- 16. 17. **RETI**
 - **RETURN** Closest(pl, ql, dl, pr, qr, dr, ps, qs, ds)

What time complexity does this process need such that the overall algorithm runs in O(n lg n)? **Hint: think about Merge Sort.**

Exercise Question 1

 What must be the running time of ClosestSplitPair if the ClosestPair algorithm is to have a running time of O(n lg n)?

```
FUNCTION ClosestPair (px, py)
   n = px.length
   IF n == 2
      RETURN px[0], px[1], dist(px[0], px[1])
   left px = px[0 ... < n//2]
   left py = [p FOR p IN py IF p.x < px[n//2].x]
   pl, ql, dl = ClosestPair(left px, left py)
   right px = px[n//2 ... < n]
   right py = [p FOR p IN py IF p.x \geq px[n//2].x]
   pr, qr, dr = ClosestPair(right px, right py)
   d = min(dl, dr)
   ps, qs, ds = ClosestSplitPair(px, py, d)
   RETURN Closest(pl, ql, dl, pr, qr, dr, ps, qs, ds)
```

Merge Sort and It's Recurrence

FUNCTION RecursiveFunction(some_input) IF base_case: # Usually O(1) RETURN base_case_work(some_input)

Two recursive calls, each with half the data
one = RecursiveFunction(some_input.first_half)
two = RecursiveFunction(some_input.second_half)

Combine results from recursive calls (usually O(n))
one_and_two = Combine(one, two)

RETURN one_and_two

```
1. FUNCTION ClosestPair(px, py)
```

```
2. n = px.length
```

```
3. IF n == 2
```

```
4. RETURN px[0], px[1], dist(px[0], px[1])
```

```
6. left_px = px[0 .. < n//2]
```

- 7. left_py = [p FOR p IN py IF p.x < px[n//2].x]
- 8. pl, ql, dl = ClosestPair(left_px, left_py)

```
10. right_px = px[n//2 ..< n]
```

```
11. right_py = [p FOR p IN py IF p.x \ge px[n//2].x]
```

```
12. pr, qr, dr = ClosestPair(right_px, right_py)
```

```
13.
```

5.

9.

```
14. d = min(dl, dr)
```

- 15. ps, qs, ds = ClosestSplitPair(px, py, d)
- 16.

```
17. RETURN Closest(pl, ql, dl, pr, qr, dr, ps, qs, ds)
```

How do we find the closest pair that splits the two sides?



- In ClosestSplitPair we only need to check for pairs that are closer than those found in the recursive calls to ClosestPair
- <u>This is easier (faster) than trying to find the closest split pair without any extra information!</u>

d = min[d(pl, ql), d(pr, qr)]

```
FUNCTION ClosestSplitPair(px, py, d)
n = px.length
x_median = px[n//2].x
middle_py = [p FOR p IN py IF x_median - d < p.x < x_median + d]</pre>
```

```
closest_d = INFINITY, closest_p = closest_q = NONE
FOR i IN [0 ..< middle_py.length - 1]
FOR j IN [1 ..= min(7, middle_py.length - i)]
p = middle_py[i], q = middle_py[i + j]
IF dist(p, q) < closest_d
closest_d = dist(p, q)
closest_p = p, closest_q = q</pre>
```

```
RETURN closest_p, closest_q, closest_d
```

Exercise Question 2

What is the running time of the nested for-loop (looping over j)?

```
FUNCTION ClosestSplitPair(px, py, d)
   n = px.length
   x median = px[n//2].x
   middle py = [p FOR p IN py IF x median - d < p.x < x median + d]
   closest d = INFINITY, closest p = closest q = NONE
   FOR i IN [0 .. < middle py.length - 1]
      FOR j IN [1 ..= min(7, middle py.length - i)]
         p = middle py[i], q = middle py[i + j]
         IF dist(p, q) < closest d</pre>
            closest d = dist(p, q)
            closest p = p, closest q = q
   RETURN closest p, closest q, closest d
```

Loop Unrolling

...

```
FOR j IN [1 ..= min(7, middle_py.length - i)]
p = middle_py[i], q = middle_py[i + j]
IF dist(p, q) < closest_d
        closest_d = dist(p, q)
        closest_p = p, closest q = q</pre>
```

```
IF dist(middle_py[i], middle_py[i + 1]) < closest_d
closest_d = dist(middle_py[i], middle_py[i + 1])
closest_p = middle_py[i]
closest_q = middle_py[i + 1]
IF dist(middle_py[i], middle_py[i + 2]) < closest_d
closest_d = dist(middle_py[i], middle_py[i + 2])
closest_p = middle_py[i]
closest_p = middle_py[i + 2]
```

```
FUNCTION ClosestSplitPair(px, py, d)
n = px.length
x_median = px[n//2].x
middle_py = [p FOR p IN py
IF x_median - d < p.x < x_median + d]
closest_d = INFINITY, closest_p = closest_q = NONE
FOR i IN [0 ..< middle_py.length - 1]
FOR j IN [1 ..= min(7, middle_py.length - i)]
p = middle_py[i], q = middle_py[i + j]
IF dist(p, q) < closest_d
closest_d = dist(p, q)</pre>
```

```
closest_p = p, closest_q = q
```

RETURN closest_p, closest_q, closest_d



Theorem for correctness of ClosestPair

Theorem:

Provided a set of n points called P, the ClosestPair algorithm find the closest pair of points according to their pairwise Euclidean distances.
ClosestPair finds the closest pair

Let $p \in left$, $q \in right$ be a split pair with d(p, q) < dThen

- A. p and $q \in middle_py$, and
- B. p and q are at most 7 positions apart in middle_py

If the claim is true:

<u>Corollary 1</u>: If the closest pair of P is in a split pair, then our ClosestSplitPair procedure finds it.

<u>Corollary 2</u>: ClosestPair is correct and runs in O(n lg n) since it has the same recursion tree as merge sort

Proof—Part A



Otherwise, p and q would not be the closest pair with d(p, q) < d

Proof—Part A



Otherwise, p and q would not be the closest pair with d(p, q) < d

ClosestPair finds the closest pair

Let $p \in left$, $q \in right$ be a split pair with d(p, q) < dThen

A. $p and q \in middle py, and$

B. p and q are at most 7 positions apart in middle py

If the claim is true:

<u>Corollary 1</u>: If the closest pair of P is in a split pair, then our ClosestSplitPair procedure finds it.

<u>Corollary 2</u>: ClosestPair is correct and runs in O(n lg n) since it has the same recursion tree as merge sort



































Proof—Part B

p and q are at most 7 positions apart in middle_py



<u>Lemma 1</u>: All points of middle_py with a y-coordinate between those of p and q lie within those 8 boxes.

Proof:

- 1. First, recall that the y-coordinate of p, q differs by less than d.
- Second, by definition of middle_py, all have an x-coordinate between x_median += d.

Proof—Part B

p and q are at most 7 positions apart in middle py



<u>Lemma 1</u>: All points of middle_py with a y-coordinate between those of p and q lie within those 8 boxes.

<u>Lemma 2</u>: At most one point of P can be in each box.

<u>Proof</u>: By contradiction. Suppose points a and b lie in the same box. Then

- 1. a and b are either both in L or both in R
- 2. $d(a, b) \le d/2 \operatorname{sqrt}(2) \le d$

This is a contradiction! How did we define d?



ClosestPair finds the closest pair

Let $p \in left$, $q \in right$ be a split pair with d(p, q) < dThen

- A. p and $q \in middle_py$, and
- B. p and q are at most 7 positions apart in middle_py

If the claim is true:

<u>Corollary 1</u>: If the closest pair of P is in a split pair, then our ClosestSplitPair procedure finds it.

<u>Corollary 2</u>: ClosestPair is correct and runs in O(n lg n) since it has the same recursion tree as merge sort

Closest Pair

- 1. Copy P and <u>sort</u> one copy by x and the other copy by y in O(n lg n)
- 2. Divide P into a left and right in O(n)
- 3. Conquer by recursively searching left and right
- 4. Look for the closest pair in middle_py in O(n)
 - Must filter by x
 - And scan through middle_py by looking at adjacent points

T(n) **FUNCTION** ClosestPair(px, py)

0(1) = px.length
0(1) n == 2
0(1) RETURN px[0], px[1], dist(px[0], px[1])

```
O(n) t_px = px[0 ..< n//2]
O(n) t_py = [p FOR p IN py IF p.x < px[n//2].x]
T(n/2) ql, dl = ClosestPair(left_px, left_py)
```

```
O(n) ht_px = px[n//2 ... < n]
O(n) ht_py = [p FOR p IN py IF p.x \ge px[n//2].x]
T(n/2) qr, dr = ClosestPair(right_px, right_py)
```

O(1) = min(dl, dr) O(n) , qs, ds = ClosestSplitPair(px, py, d)

O(1) TURN Closest(pl, ql, dl, pr, qr, dr, ps, qs, ds)

$$T(n) = 2 T(n/2) + O(n)$$

= O(n lg n)

$$T(n) = 2 T(n/2) + O(n)$$

= O(n lg n)

T(n/2)ft_sorted = MergeSort(array[0 ..< n//2]) T(n/2)ght_sorted = MergeSort(array[n//2 ..< n])

O(n) ray_sorted = Merge(left_sorted, right_sorted)

O(1) ETURN array_sorted

T(n) = 2 T(n/2) + O(n)= O(n lg n)

T(n) FUNCTION RecursiveFunction(some_input)
O(1) base_case:
 # Usually O(1)
O(1) RETURN base_case_work(some_input)

Two recursive calls, each with half the data
T(n/2) ne = RecursiveFunction(some_input.first_half)
T(n/2) vo = RecursiveFunction(some_input.second_half)

Combine results from recursive calls (usually O(n))
O(n) ne_and_two = Combine(one, two)

O(1) ETURN one_and_two

Supplementary slides showing an example execution.


































Closest is Split

Closest on Right



Closest on Left



Closest on Left





Closest is Split

Closest on Right



Closest is Split



Closest is Split

