



Lecture 10: Divide & Conquer; String Matching I

COMP526: Efficient Algorithms

Updated: November 5, 2024

Will Rosenbaum
University of Liverpool

Announcements

1. **NO QUIZ THIS WEEK!**
2. Programming Assignment Posted
 - Due Wednesday, 13 November
3. Attendance Code:

Meeting Goals

- Discuss more Divide & Conquer algorithms
 - Order Statistics
 - Majority
 - Closest Pair of Points
- Introduce the **String Matching** problem
 - Problem definition
 - Elementary algorithm

Divide & Conquer

Previously: Divide & Conquer Strategy

Generic Strategy

Given an algorithmic task:

1. Break the input into smaller instances of the task
2. Solve the smaller instances
 - this is typically recursive!
3. Combine smaller solutions to a solution to the whole task

Divide & Conquer Examples (so far):

- MERGESORT: divide an array *by index* to sort
 - $O(n \log n)$ time
- QUICKSORT: divide an array *by value* to sort
 - $O(n \log n)$ time
- BINARYSEARCH: divide a *sorted* array to search it
 - $O(\log n)$ time

Three More Problems

Problem 1. k -Selection:

- Given an array a of n numbers, find the k th largest number

Three More Problems

Problem 1. k -Selection:

- Given an array a of n numbers, find the k th largest number

Problem 2. Majority:

- Given an array a of n items, is there an item that is repeated more than $> n/2$ times?

Three More Problems

Problem 1. k -Selection:

- Given an array a of n numbers, find the k th largest number

Problem 2. Majority:

- Given an array a of n items, is there an item that is repeated more than $> n/2$ times?

Problem 3. Closest Points in the Plane

- Given n points p_1, p_2, \dots, p_n in the plane, which *pair* of points p_i, p_j are closest to one another?

k-Selection

k-Selection

Problem. Given an array a of n numbers, find the k th smallest number.

k-Selection

Problem. Given an array a of n numbers, find the k th smallest number.

Simple solution.

- sort a in $O(n \log n)$ time
- return $a[k]$

Can we do better?

k-Selection

Problem. Given an array a of n numbers, find the k th smallest number.

Simple solution.

- sort a in $O(n \log n)$ time
- return $a[k]$

Can we do better?

Modify QuickSort!

- Choose pivot p
- Perform split
- *only recurse on half that contains k th smallest value*
 - this will be the half that contains index k
- Random pivot selection
 $\implies O(n)$ expected time!

k-Selection

Problem. Given an array a of n numbers, find the k th smallest number.

Simple solution.

- sort a in $O(n \log n)$ time
- return $a[k]$

Can we do better?

Modify QuickSort!

- Choose pivot p
- Perform split
- *only recurse on half that contains k th smallest value*
 - this will be the half that contains index k
- Random pivot selection
 $\implies O(n)$ expected time!

```
1: procedure  
   QUICKSELECT( $a, \min, \max, k$ )  
2:   if  $\max - \min \leq 1$  then  
3:     return  $a[\min]$   
4:   end if  
5:    $p \leftarrow$  SELECTPIVOT( $a, \min, \max$ )  
6:    $j \leftarrow$  SPLIT( $a, \min, \max, p$ )  
7:   if  $j = k$  then  
8:     return  $a[k]$   
9:   else if  $j < k$  then  
10:    QUICKSELECT( $a, j + 1, \max, k$ )  
11:  else  
12:    QUICKSELECT( $a, \min, j - 1, k$ )  
13:  end if  
14: end procedure
```

k-Selection

Problem. Given an array a of n numbers, find the k th smallest number.

PollEverywhere Question

What is the *worst case* running time of QUICKSELECT on an array of n elements?



pollev.com/comp526

```
1: procedure  
   QUICKSELECT( $a, \min, \max, k$ )  
2:   if  $\max - \min \leq 1$  then  
3:     return  $a[\min]$   
4:   end if  
5:    $p \leftarrow$  SELECTPIVOT( $a, \min, \max$ )  
6:    $j \leftarrow$  SPLIT( $a, \min, \max, p$ )  
7:   if  $j = k$  then  
8:     return  $a[k]$   
9:   else if  $j < k$  then  
10:    QUICKSELECT( $a, j + 1, \max, k$ )  
11:  else  
12:    QUICKSELECT( $a, \min, j - 1, k$ )  
13:  end if  
14: end procedure
```

Deterministic k-Selection?

Question. Can we perform k -selection with a **worst case** $O(n)$ running time?

Deterministic k-Selection?

Question. Can we perform k -selection with a **worst case** $O(n)$ running time?

Idea. What if we can select better pivots?

- Suppose we can *guarantee* that our pivot is “good enough:”
 - rank of p is between cn and $(1 - c)n$ for $c > 0$
- How many recursive calls until we're done?

Deterministic k-Selection?

Question. Can we perform k -selection with a **worst case** $O(n)$ running time?

Idea. What if we can select better pivots?

- Suppose we can *guarantee* that our pivot is “good enough:”
 - rank of p is between cn and $(1 - c)n$ for $c > 0$
- How many recursive calls until we’re done?
 - each recursive call has size at most $(1 - 2c)n$
 - ℓ recursive calls \implies size at most $(1 - 2c)^\ell n$
 - \implies done after $\ell = O(\log n)$ levels of recursion
- What is overall running time?

Deterministic k-Selection?

Question. Can we perform k -selection with a **worst case** $O(n)$ running time?

Idea. What if we can select better pivots?

- Suppose we can *guarantee* that our pivot is “good enough:”
 - rank of p is between cn and $(1 - c)n$ for $c > 0$
- How many recursive calls until we're done?
 - each recursive call has size at most $(1 - 2c)n$
 - ℓ recursive calls \implies size at most $(1 - 2c)^\ell n$
 - \implies done after $\ell = O(\log n)$ levels of recursion
- What is overall running time?
 - $Cn + (1 - 2c)Cn + (1 - 2c)^2 Cn + \dots = O(n)$

Deterministic k-Selection?

Question. Can we perform k -selection with a **worst case** $O(n)$ running time?

Idea. What if we can select better pivots?

- Suppose we can *guarantee* that our pivot is “good enough:”
 - rank of p is between cn and $(1 - c)n$ for $c > 0$
- How many recursive calls until we're done?
 - each recursive call has size at most $(1 - 2c)n$
 - ℓ recursive calls \implies size at most $(1 - 2c)^\ell n$
 - \implies done after $\ell = O(\log n)$ levels of recursion
- What is overall running time?
 - $Cn + (1 - 2c)Cn + (1 - 2c)^2 Cn + \dots = O(n)$

But how can we find a good pivot *deterministically*?

- Need to find pivots close to the median...
- Median is (special case) of k selection!



Median of Medians Strategy

Strategy. To find a good pivot:

- Find a *smaller* set of values whose *median* is a good pivot
- Recursively find the median of the smaller set of values

Median of Medians Strategy

Strategy. To find a good pivot:

- Find a *smaller* set of values whose *median* is a good pivot
- Recursively find the median of the smaller set of values
- Consider blocks of size 5
 - sort each block
 - find the block median
- Claim: median of medians is a good pivot:

Median of Medians Strategy

Strategy. To find a good pivot:

- Find a *smaller* set of values whose *median* is a good pivot
- Recursively find the median of the smaller set of values
- Consider blocks of size 5
 - sort each block
 - find the block median
- Claim: median of medians is a good pivot:
 - at least $\frac{3}{10}$ -fraction is excluded

Median of Medians Strategy

Strategy. To find a good pivot:

- Find a *smaller* set of values whose *median* is a good pivot
- Recursively find the median of the smaller set of values
- Consider blocks of size 5
 - sort each block
 - find the block median
- Claim: median of medians is a good pivot:
 - at least $\frac{3}{10}$ -fraction is excluded

```
1: procedure SELECTPIVOT( $a, \ell, r$ )
2:    $m \leftarrow n/5$ 
3:   for  $i = 0, 1, \dots, m-1$  do
4:     SORT( $a[5i \dots 5i+4]$ )
5:     SWAP( $a, i, 5i+2$ )
6:   end for
7:   return QUICKSELECT( $a, 0, m, (m-1)/2$ )
8: end procedure
```

Median of Medians Strategy

Illustration:

```
1: procedure SELECTPIVOT( $a, \ell, r$ )
2:    $m \leftarrow n/5$ 
3:   for  $i = 0, 1, \dots, m-1$  do
4:     SORT( $a[5i \dots 5i+4]$ )
5:     SWAP( $a, i, 5i+2$ )
6:   end for
7:   return QUICKSELECT( $a, 0, m, (m-1)/2$ )
8: end procedure
9: procedure QUICKSELECT( $a, \ell, r, k$ )
10:  if  $r - \ell \leq 1$  return  $a[\ell]$ 
11:   $b \leftarrow$  SELECTPIVOT( $a, \ell, r$ )
12:   $j \leftarrow$  SPLIT( $a, \ell, r, a[b]$ )
13:  if  $j = k$  then
14:    return  $a[j]$ 
15:  else if  $j < k$  then
16:    QUICKSELECT( $a, j+1, r, k-j-1$ )
17:  else
18:    QUICKSELECT( $a, \ell, j, k$ )
19:  end if
20: end procedure
```


Median of Medians Strategy

Analysis.

Running time $T(n)$ satisfies

$$\begin{aligned}T(n) &\leq Cn + T\left(\frac{1}{5}n\right) + T\left(\frac{7}{10}n\right) \\ &\leq Cn + T\left(\frac{1}{5}n + \frac{7}{10}n\right) \\ &\leq Cn + T\left(\frac{9}{10}n\right)\end{aligned}$$

Therefore, $T(n) = O(n)$.

```
1: procedure SELECTPIVOT( $a, \ell, r$ )
2:    $m \leftarrow n/5$ 
3:   for  $i = 0, 1, \dots, m-1$  do
4:     SORT( $a[5i \dots 5i+4]$ )
5:     SWAP( $a, i, 5i+2$ )
6:   end for
7:   return QUICKSELECT( $a, 0, m, (m-1)/2$ )
8: end procedure
9: procedure QUICKSELECT( $a, \ell, r, k$ )
10:  if  $r - \ell \leq 1$  return  $a[\ell]$ 
11:   $b \leftarrow$  SELECTPIVOT( $a, \ell, r$ )
12:   $j \leftarrow$  SPLIT( $a, \ell, r, a[b]$ )
13:  if  $j = k$  then
14:    return  $a[j]$ 
15:  else if  $j < k$  then
16:    QUICKSELECT( $a, j+1, r, k-j-1$ )
17:  else
18:    QUICKSELECT( $a, 0, j, k$ )
19:  end if
20: end procedure
```

Median of Medians Strategy

Conclusion. The Median of Medians strategy allows us to

- solve k -selection in $O(n)$ time, worst case
- sort in $O(n \log n)$ time, worst case too
 - use k selection as a sub-routine for SELECTPIVOT in QUICKSORT

Note. Randomized variants tend to be more efficient in practice.

```
1: procedure SELECTPIVOT( $a, \ell, r$ )
2:    $m \leftarrow n/5$ 
3:   for  $i = 0, 1, \dots, m-1$  do
4:     SORT( $a[5i \dots 5i+4]$ )
5:     SWAP( $a, i, 5i+2$ )
6:   end for
7:   return QUICKSELECT( $a, 0, m, (m-1)/2$ )
8: end procedure
9: procedure QUICKSELECT( $a, \ell, r, k$ )
10:  if  $r - \ell \leq 1$  return  $a[\ell]$ 
11:   $b \leftarrow$  SELECTPIVOT( $a, \ell, r$ )
12:   $j \leftarrow$  SPLIT( $a, \ell, r, a[b]$ )
13:  if  $j = k$  then
14:    return  $a[j]$ 
15:  else if  $j < k$  then
16:    QUICKSELECT( $a, j+1, r, k-j-1$ )
17:  else
18:    QUICKSELECT( $a, \ell, j, k$ )
19:  end if
20: end procedure
```

Majority

Majority

Problem 2. Majority:

- Given an array a of n items, is there an item that is repeated more than $n/2$ times?

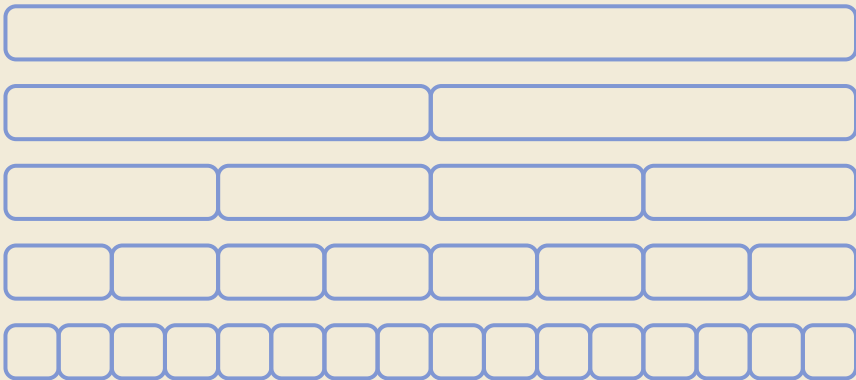
Naive Solution

- Iterate over elements and compare each element to all others to see if occurs at least $n/2$ times
- Takes $\Theta(n^2)$ time

Observation. If a value m is a majority, then m must either be a majority in $a[0 \dots n/2]$ or $a[n/2 + 1 \dots n - 1]$ as well.

- Split a in half
- Recursively find candidate majority m_ℓ and m_r for halves
- Check to see if either is a majority

Divide & Conquer Majority Illustration



1	2	1	1	2	3	3	1	2	1	1	2	3	3	3	2	3	2	1	1	1	3	1	2	1	1	3	1	1	2	1	1
0	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19	20	21	22	23	24	25	26	27	28	29	30	31

Divide & Conquer Majority in Code

```
1: procedure ISMAJORITY( $a, \ell, r, v$ )
2:    $count \leftarrow 0$ 
3:   for  $i = \ell, \ell + 1, \dots, r$  do
4:     if  $a[i] = v$  then
5:        $count \leftarrow count + 1$ 
6:     end if
7:   end for
8:   return  $count > (r - \ell + 1) / 2$ 
9: end procedure
10: procedure MAJORITY( $a, \ell, r$ )
11:   if  $\ell - r < 1$  return  $a[\ell]$ 
12:    $j \leftarrow (r - \ell) / 2$ 
13:    $v_\ell \leftarrow$  MAJORITY( $a, \ell, j$ )
14:    $v_r \leftarrow$  MAJORITY( $a, j + 1, r$ )
15:   if ISMAJORITY( $a, \ell, r, v_\ell$ ) then
16:     return  $v_\ell$ 
17:   else if ISMAJORITY( $a, \ell, r, v_r$ ) then
18:     return  $v_r$ 
19:   end if
20:   return  $\perp$ 
21: end procedure
```

Divide & Conquer Majority in Code

PollEverywhere Question

What is the *worst case* running time of MAJORITY on an array of n elements?



pollev.com/comp526

```
1: procedure ISMAJORITY( $a, \ell, r, v$ )
2:    $count \leftarrow 0$ 
3:   for  $i = \ell, \ell + 1, \dots, r$  do
4:     if  $a[i] = v$  then
5:        $count \leftarrow count + 1$ 
6:     end if
7:   end for
8:   return  $count > (r - \ell + 1) / 2$ 
9: end procedure
10: procedure MAJORITY( $a, \ell, r$ )
11:   if  $\ell - r < 1$  return  $a[\ell]$ 
12:    $j \leftarrow (r - \ell) / 2$ 
13:    $v_\ell \leftarrow$  MAJORITY( $a, \ell, j$ )
14:    $v_r \leftarrow$  MAJORITY( $a, j + 1, r$ )
15:   if ISMAJORITY( $a, \ell, r, v_\ell$ ) then
16:     return  $v_\ell$ 
17:   else if ISMAJORITY( $a, \ell, r, v_r$ ) then
18:     return  $v_r$ 
19:   end if
20:   return  $\perp$ 
21: end procedure
```

Divide & Conquer Majority in Code

Analysis.

- Almost identical to MERGESORT
- Each call to ISMAJORITY(a, ℓ, r, v) takes time $\Theta(\ell - r)$
- Running time $T(n)$ satisfies $T(n) \leq 2T(n/2) + \Theta(n)$
- Solve recursion \implies done!

```
1: procedure ISMAJORITY( $a, \ell, r, v$ )
2:    $count \leftarrow 0$ 
3:   for  $i = \ell, \ell + 1, \dots, r$  do
4:     if  $a[i] = v$  then
5:        $count \leftarrow count + 1$ 
6:     end if
7:   end for
8:   return  $count > (r - \ell + 1) / 2$ 
9: end procedure
10: procedure MAJORITY( $a, \ell, r$ )
11:   if  $\ell - r < 1$  return  $a[\ell]$ 
12:    $j \leftarrow (r - \ell) / 2$ 
13:    $v_\ell \leftarrow$  MAJORITY( $a, \ell, j$ )
14:    $v_r \leftarrow$  MAJORITY( $a, j + 1, r$ )
15:   if ISMAJORITY( $a, \ell, r, v_\ell$ ) then
16:     return  $v_\ell$ 
17:   else if ISMAJORITY( $a, \ell, r, v_r$ ) then
18:     return  $v_r$ 
19:   end if
20:   return  $\perp$ 
21: end procedure
```


Divide & Conquer Majority in Code

Analysis.

- Almost identical to MERGESORT
- Each call to ISMAJORITY(a, ℓ, r, v) takes time $\Theta(\ell - r)$
- Running time $T(n)$ satisfies $T(n) \leq 2T(n/2) + \Theta(n)$
- Solve recursion \implies done!

Challenge. Devise an algorithm that finds the majority in $\Theta(n)$ time (worst case). (Hint: don't use Divide & Conquer)

```
1: procedure ISMAJORITY( $a, \ell, r, v$ )
2:    $count \leftarrow 0$ 
3:   for  $i = \ell, \ell + 1, \dots, r$  do
4:     if  $a[i] = v$  then
5:        $count \leftarrow count + 1$ 
6:     end if
7:   end for
8:   return  $count > (r - \ell + 1) / 2$ 
9: end procedure
10: procedure MAJORITY( $a, \ell, r$ )
11:   if  $\ell - r < 1$  return  $a[\ell]$ 
12:    $j \leftarrow (r - \ell) / 2$ 
13:    $v_\ell \leftarrow$  MAJORITY( $a, \ell, j$ )
14:    $v_r \leftarrow$  MAJORITY( $a, j + 1, r$ )
15:   if ISMAJORITY( $a, \ell, r, v_\ell$ ) then
16:     return  $v_\ell$ 
17:   else if ISMAJORITY( $a, \ell, r, v_r$ ) then
18:     return  $v_r$ 
19:   end if
20:   return  $\perp$ 
21: end procedure
```

Closest Points in the Plane

Closest Points in the Plane

Problem 3. Given n points p_1, p_2, \dots, p_n in the plane, which *pair* of points p_i, p_j are closest to one another?



Closest Points in the Plane

Problem 3. Given n points p_1, p_2, \dots, p_n in the plane, which *pair* of points p_i, p_j are closest to one another?

Naive Strategy suggested by

GenAI:

- Compute distances between all pairs of points

```
1: procedure NAIVEMINDIST( $p$ )
2:    $d \leftarrow \infty$ 
3:   for  $i = 1, 2, \dots, n - 1$  do
4:     for  $j = 0, 1, \dots, i - 1$  do
5:       if DIST( $p[i], p[j]$ ) <  $d$  then
6:          $d \leftarrow$  DIST( $p[i], p[j]$ )
7:       end if
8:     end for
9:   end for
10:  return  $d$ 
11: end procedure
```

Closest Points in the Plane

Problem 3. Given n points p_1, p_2, \dots, p_n in the plane, which *pair* of points p_i, p_j are closest to one another?

PollEverywhere Question

What is the worst case running time of NAIVEMINDIST on a set of n points in the plane?



pollev.com/comp526

```
1: procedure NAIVEMINDIST( $p$ )
2:    $d \leftarrow \infty$ 
3:   for  $i = 1, 2, \dots, n - 1$  do
4:     for  $j = 0, 1, \dots, i - 1$  do
5:       if DIST( $p[i], p[j]$ ) <  $d$  then
6:          $d \leftarrow$  DIST( $p[i], p[j]$ )
7:       end if
8:     end for
9:   end for
10:  return  $d$ 
11: end procedure
```

Closest Points in the Plane

Problem 3. Given n points p_1, p_2, \dots, p_n in the plane, which *pair* of points p_i, p_j are closest to one another?

Naive Strategy suggested by

GenAI:

- Compute distances between all pairs of points

Question. How could we use **Divide & Conquer** to improve on this running time?

```
1: procedure NAIVEMINDIST( $p$ )
2:    $d \leftarrow \infty$ 
3:   for  $i = 1, 2, \dots, n - 1$  do
4:     for  $j = 0, 1, \dots, i - 1$  do
5:       if DIST( $p[i], p[j]$ ) <  $d$  then
6:          $d \leftarrow$  DIST( $p[i], p[j]$ )
7:       end if
8:     end for
9:   end for
10:  return  $d$ 
11: end procedure
```

Closest Points: Divide & Conquer

Step 1. split the array according to x -coordinate



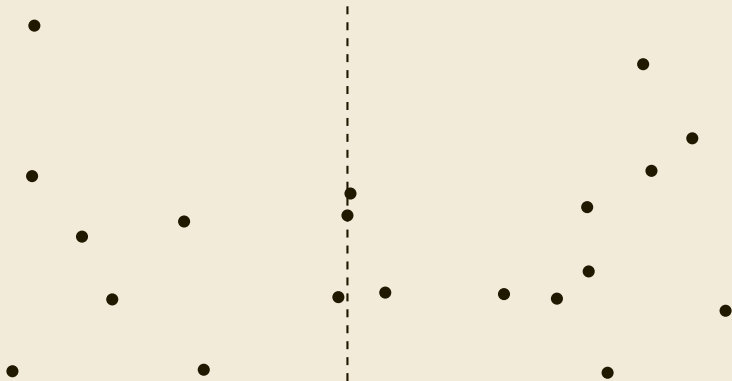
Closest Points: Divide & Conquer

Step 1a. sort the array by x coordinate



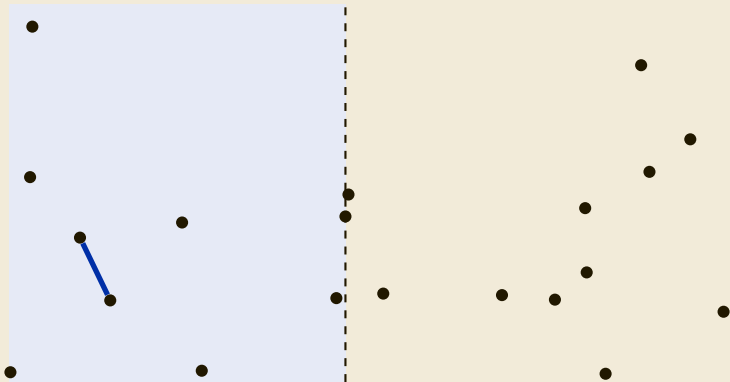
Closest Points: Divide & Conquer

Step 1b. find median according to x coordinate, p_m



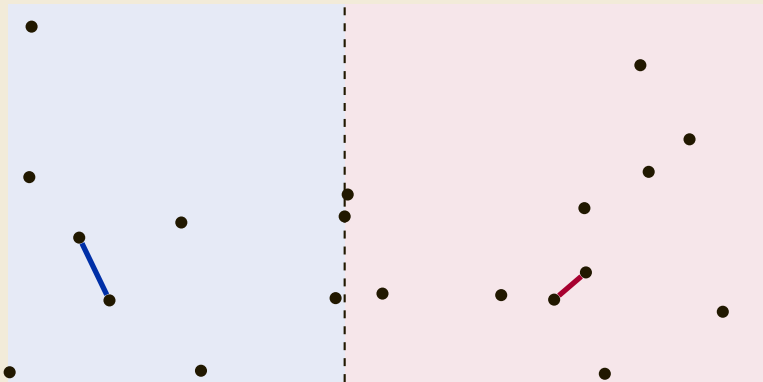
Closest Points: Divide & Conquer

Step 2a. (recursively) solve the problem for left half



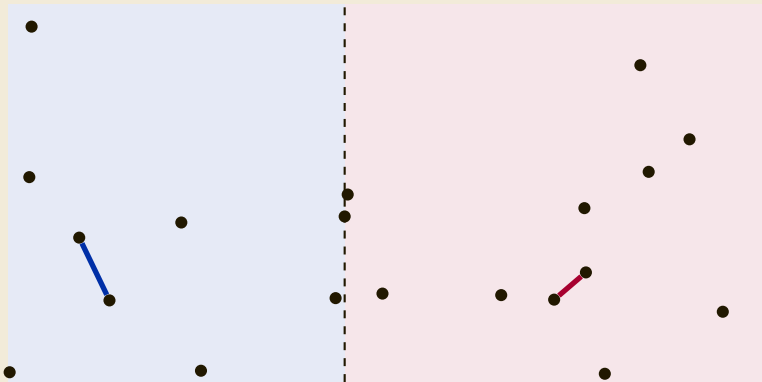
Closest Points: Divide & Conquer

Step 2b. (recursively) solve the problem for right half



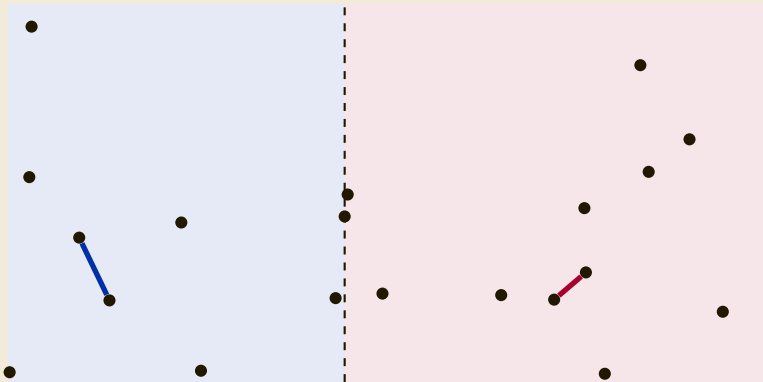
Closest Points: Divide & Conquer

Step 3. merge solutions together



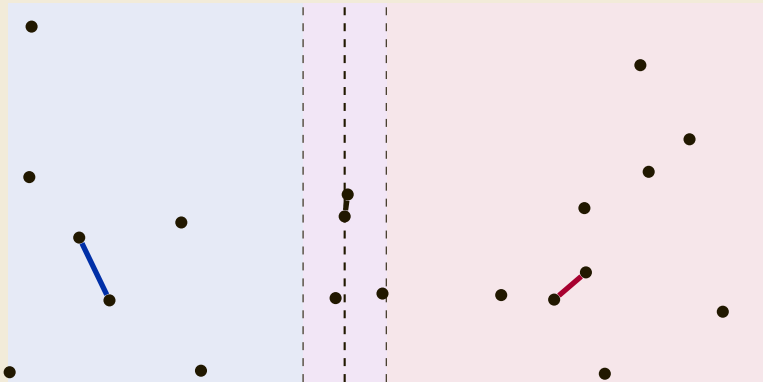
Closest Points: Divide & Conquer

Step 3. merge solutions together ... but **how?**



Closest Points: Divide & Conquer

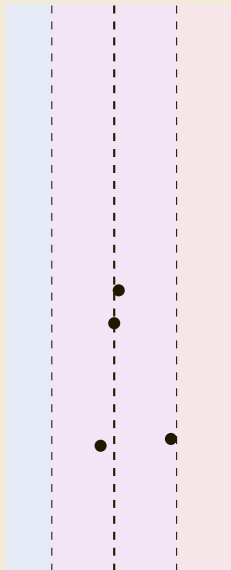
Critical Analysis. What happens in the middle strip?



Analysis of the Middle Strip

Suppose:

- d_ℓ is minimal distance on the left
- d_r is minimal distance on the right
- $\delta = \min \{d_\ell, d_r\}$
- x_m is the median x -coordinate among points

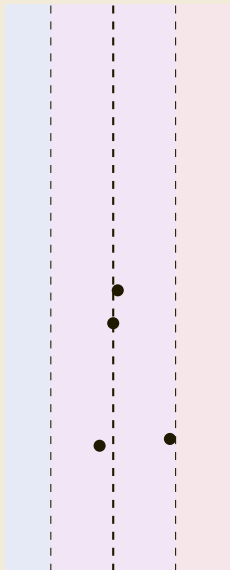


Analysis of the Middle Strip

Suppose:

- d_ℓ is minimal distance on the left
- d_r is minimal distance on the right
- $\delta = \min \{d_\ell, d_r\}$
- x_m is the median x -coordinate among points

Claim 1. If p is in left half and q is on right half have with $\text{DIST}(p_i, p_j) < \delta$, then $x_m - \delta < x_i \leq x_m$ and $x_m \leq x_j \leq x_m + \delta$.



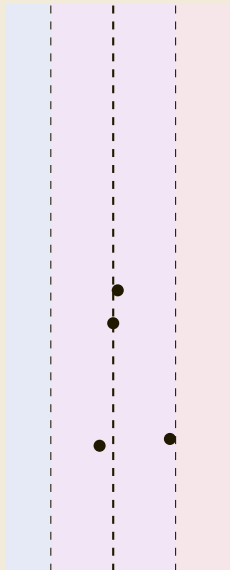
Analysis of the Middle Strip

Suppose:

- d_ℓ is minimal distance on the left
- d_r is minimal distance on the right
- $\delta = \min \{d_\ell, d_r\}$
- x_m is the median x -coordinate among points

Claim 1. If p is in left half and q is on right have with $\text{DIST}(p_i, p_j) < \delta$, then $x_m - \delta < x_i \leq x_m$ and $x_m \leq x_j \leq x_m + \delta$.

Claim 2. With p as above, there are at most 8 points q on the right side with $\text{DIST}(p, q) \leq \delta$.



Analysis of the Middle Strip

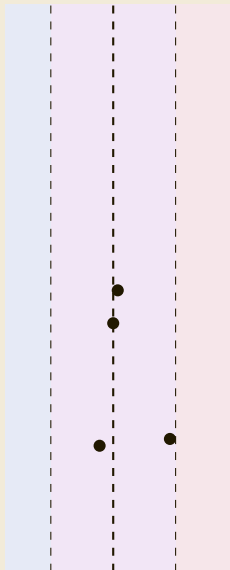
Suppose:

- d_ℓ is minimal distance on the left
- d_r is minimal distance on the right
- $\delta = \min \{d_\ell, d_r\}$
- x_m is the median x -coordinate among points

Claim 1. If p is in left half and q is on right have with $\text{DIST}(p_i, p_j) < \delta$, then $x_m - \delta < x_i \leq x_m$ and $x_m \leq x_j \leq x_m + \delta$.

Claim 2. With p as above, there are at most 8 points q on the right side with $\text{DIST}(p, q) \leq \delta$.

Consequence. We only need to make $O(n)$ further distance computations to compute overall minimum distance.



Putting it Together

Algorithm Sketch. Find the closest pair of points among p_1, p_2, \dots, p_n in the plane:

1. Sort points by x -coordinate, x_m is the median value.
2. Recursively sort left and right halves.
3. Set δ to be the minimum distance on either half.
4. Consider points within distance δ of median line, and compute distances across the halves.
 - this can be done in $O(n)$ time
5. Report the smallest distance found.

Putting it Together

Algorithm Sketch. Find the closest pair of points among p_1, p_2, \dots, p_n in the plane:

1. Sort points by x -coordinate, x_m is the median value.
2. Recursively sort left and right halves.
3. Set δ to be the minimum distance on either half.
4. Consider points within distance δ of median line, and compute distances across the halves.
 - this can be done in $O(n)$ time
5. Report the smallest distance found.

Running time analysis.

- Preprocessing takes $O(n \log n)$ to sort the points.
- The main algorithm running time satisfies the recursion
$$T(n) \leq 2T(n/2) + O(n)$$

⇒ overall running time is $O(n \log n)$.

Concluding Thoughts

Divide & Conquer is a powerful algorithm design strategy.

Efficiency improvement over naive solutions:

- Sorting $\Theta(n^2) \rightarrow \Theta(n \log n)$
- k -Selection $\Theta(n^2) \rightarrow \Theta(n)$
- Majority $\Theta(n^2) \rightarrow \Theta(n \log n)$
- Closest points in the plane $\Theta(n^2) \rightarrow \Theta(n \log n)$

Concluding Thoughts

Divide & Conquer is a powerful algorithm design strategy.

Efficiency improvement over naive solutions:

- Sorting $\Theta(n^2) \rightarrow \Theta(n \log n)$
- k -Selection $\Theta(n^2) \rightarrow \Theta(n)$
- Majority $\Theta(n^2) \rightarrow \Theta(n \log n)$
- Closest points in the plane $\Theta(n^2) \rightarrow \Theta(n \log n)$

Other applications:

- Matrix multiplication (Strassen's algorithm):
 $\Theta(n^3) \rightarrow \Theta(n^{\log_2 7 + o(1)}) \approx \Theta(n^{2.807})$
- Integer multiplication: $\Theta(B^2) \rightarrow \Theta(B^{\log_2 3}) \rightarrow \Theta(B \log B)$
- Fast Fourier Transform: $\Theta(n^2) \rightarrow \Theta(n \log n)$

Concluding Thoughts

Divide & Conquer is a powerful algorithm design strategy.

Efficiency improvement over naive solutions:

- Sorting $\Theta(n^2) \rightarrow \Theta(n \log n)$
- k -Selection $\Theta(n^2) \rightarrow \Theta(n)$
- Majority $\Theta(n^2) \rightarrow \Theta(n \log n)$
- Closest points in the plane $\Theta(n^2) \rightarrow \Theta(n \log n)$

Other applications:

- Matrix multiplication (Strassen's algorithm):
 $\Theta(n^3) \rightarrow \Theta(n^{\log_2 7 + o(1)}) \approx \Theta(n^{2.807})$
- Integer multiplication: $\Theta(B^2) \rightarrow \Theta(B^{\log_2 3}) \rightarrow \Theta(B \log B)$
- Fast Fourier Transform: $\Theta(n^2) \rightarrow \Theta(n \log n)$

Other considerations:

- Practical because of **parallelism**

String Matching

String Matching: Motivation

Fundamental Problems. Given a (large) **text** T and (small) **pattern** P :

- Determine if T contains the pattern P .
- Find the *first occurrence* of P in T (if any)
- Find the number of occurrences of P in T

String Matching: Motivation

Fundamental Problems. Given a (large) **text** T and (small) **pattern** P :

- Determine if T contains the pattern P .
- Find the *first occurrence* of P in T (if any)
- Find the number of occurrences of P in T

Example applications.

- Search on your computer: Ctrl + F

String Matching: Motivation

Fundamental Problems. Given a (large) **text** T and (small) **pattern** P :

- Determine if T contains the pattern P .
- Find the *first occurrence* of P in T (if any)
- Find the number of occurrences of P in T

Example applications.

- Search on your computer: Ctrl + F
- Bioinformatics:
 - does a DNA sequence (T) contain a particular gene (P)?

String Matching: Motivation

Fundamental Problems. Given a (large) **text** T and (small) **pattern** P :

- Determine if T contains the pattern P .
- Find the *first occurrence* of P in T (if any)
- Find the number of occurrences of P in T

Example applications.

- Search on your computer: Ctrl + F
- Bioinformatics:
 - does a DNA sequence (T) contain a particular gene (P)?
- Computer virus detection
 - does your hard drive store a known program?

String Matching: Motivation

Fundamental Problems. Given a (large) **text** T and (small) **pattern** P :

- Determine if T contains the pattern P .
- Find the *first occurrence* of P in T (if any)
- Find the number of occurrences of P in T

Example applications.

- Search on your computer: Ctrl + F
- Bioinformatics:
 - does a DNA sequence (T) contain a particular gene (P)?
- Computer virus detection
 - does your hard drive store a known program?
- (Counter) Espionage
 - does a data transmission contain the phrase “ATTACK AT DAWN?”

String Matching: Motivation

Fundamental Problems. Given a (large) **text** T and (small) **pattern** P :

- Determine if T contains the pattern P .
- Find the *first occurrence* of P in T (if any)
- Find the number of occurrences of P in T

Example applications.

- Search on your computer: Ctrl + F
- Bioinformatics:
 - does a DNA sequence (T) contain a particular gene (P)?
- Computer virus detection
 - does your hard drive store a known program?
- (Counter) Espionage
 - does a data transmission contain the phrase “ATTACK AT DAWN?”

Interesting parameters. $|T|$ is large ($\sim 1\text{B}$), $|P|$ is relatively small ($\sim 1\text{K}$)

Making Things Precise

Notation

- Σ is a finite **alphabet** or set of **characters**, $\sigma = |\Sigma|$
 - $\Sigma = \{0, 1\}$ is binary alphabet
 - $\Sigma = \{A, B, \dots\}$ is Roman alphabet
 - $\Sigma = \dots$ e.g., ASCII, Unicode,
- $\Sigma^n = \Sigma \times \Sigma \times \dots \times \Sigma = \{(c_1, c_2, \dots, c_n) \mid \text{each } c_i \in \Sigma\} =$ strings of exactly n characters
- $\Sigma^* = \bigcup_{n=0}^{\infty} \Sigma^n =$ all *finite* strings
- $\Sigma^+ = \bigcup_{n=0}^{\infty} \Sigma^n =$ all *nonempty* (finite) strings
- $\varepsilon \in \Sigma^0$ is the **empty string**
- for $S \in \Sigma^n$, $S[i]$ is i th character of S
- for $S, T \in \Sigma^*$, ST is the **concatenation** of S and T
- for $S \in \Sigma^n$, $S[i..j] = S[i]S[i+1] \dots S[j]$ is a **substring**
 - $S[0..j]$ is a **prefix**, $S[j..n-1]$ is a **suffix**
 - $S[i..j] = S[i..j-1] \implies S = S[0..n]$

The String Matching Problem

Input:

- A **text** $T \in \Sigma^*$ of length n
- A **pattern** $P \in \Sigma^*$ of length m (typically $m \ll n$)

Output:

- The index of the **first occurrence** of P in T , or -1 if T does not contain P as a substring:
 - $\min \{i \mid T[i, i + m) = P\}$

Example.

- $T = 10110011011101$
- $P_1 = 1101$

The String Matching Problem

Input:

- A **text** $T \in \Sigma^*$ of length n
- A **pattern** $P \in \Sigma^*$ of length m (typically $m \ll n$)

Output:

- The index of the **first occurrence** of P in T , or -1 if T does not contain P as a substring:
 - $\min \{i \mid T[i, i + m) = P\}$

Example.

- $T = 10110011011101$
- $P_1 = 1101$
 - Output: $i \leftarrow 6$

The String Matching Problem

Input:

- A **text** $T \in \Sigma^*$ of length n
- A **pattern** $P \in \Sigma^*$ of length m (typically $m \ll n$)

Output:

- The index of the **first occurrence** of P in T , or -1 if T does not contain P as a substring:
 - $\min \{i \mid T[i, i + m) = P\}$

Example.

- $T = 10110011011101$
- $P_1 = 1101$
 - Output: $i \leftarrow 6$
- $P_2 = 000$

The String Matching Problem

Input:

- A **text** $T \in \Sigma^*$ of length n
- A **pattern** $P \in \Sigma^*$ of length m (typically $m \ll n$)

Output:

- The index of the **first occurrence** of P in T , or -1 if T does not contain P as a substring:
 - $\min \{i \mid T[i, i + m) = P\}$

Example.

- $T = 10110011011101$
- $P_1 = 1101$
 - Output: $i \leftarrow 6$
- $P_2 = 000$
 - Output: $i \leftarrow -1$

Brute Force Matching

Brute Force Matching

Guess an index i where a match might occur

- Possible guesses $i = 0, 1, \dots, n - m - 1$

Check if match at i :

- is $T[i, i + m) = P$?
- verify each character individually

Cost = number of comparisons made

Brute Force Matching

Guess an index i where a match might occur

- Possible guesses $i = 0, 1, \dots, n - m - 1$

Check if match at i :

- is $T[i, i + m) = P$?
- verify each character individually

```
1: procedure VERIFYMATCH( $T, P, i$ )
2:    $j \leftarrow 0$ 
3:   while  $j < m$  do
4:     if  $T[i + j] \neq P[j]$  then
5:       return FALSE
6:     end if
7:      $j \leftarrow j + 1$ 
8:   end while
9:   return TRUE
10: end procedure
```

Cost = number of comparisons made

Brute Force Matching

Guess an index i where a match might occur

- Possible guesses $i = 0, 1, \dots, n - m - 1$

Check if match at i :

- is $T[i, i + m) = P$?
- verify each character individually

```
1: procedure VERIFYMATCH( $T, P, i$ )
2:    $j \leftarrow 0$ 
3:   while  $j < m$  do
4:     if  $T[i + j) \neq P[j)$  then
5:       return FALSE
6:     end if
7:      $j \leftarrow j + 1$ 
8:   end while
9:   return TRUE
10: end procedure
```

Cost = number of comparisons made

PollEverywhere Question

What are the worst case and best case running times of VERIFYMATCH?



pollev.com/comp526

Brute Force Matching

Guess an index i where a match might occur

- Possible guesses $i = 0, 1, \dots, n - m - 1$

Check if match at i :

- is $T[i, i + m) = P$?
- verify each character individually

Best and Worst Cases:

```
1: procedure VERIFYMATCH( $T, P, i$ )
2:    $j \leftarrow 0$ 
3:   while  $j < m$  do
4:     if  $T[i + j] \neq P[j]$  then
5:       return FALSE
6:     end if
7:      $j \leftarrow j + 1$ 
8:   end while
9:   return TRUE
10: end procedure
```

Cost = number of comparisons made

Brute Force Matching

Guess an index i where a match might occur

- Possible guesses $i = 0, 1, \dots, n - m - 1$

Check if match at i :

- is $T[i, i + m) = P$?
- verify each character individually

Cost = number of comparisons made

Brute force. Guess and check every value

$i = 0, 1, \dots, n - m - 1$

- Worst case running time is $\Theta(nm)$
 - What is example has cost $\Omega(nm)$?
- Best case cost is $\Theta(m)$

Brute Force Example

Example

- $T = abbbababbab$
- $P = abba$

0	1	2	3	4	5	6	7	8	9	10
a	b	b	b	a	b	a	b	b	a	b

procedure

BRUTEFORCEMATCH(T, P)

for $i = 0, 1, \dots, n - m - 1$ **do**

if VERIFYMATCH(T, P, i) **then**

return i

end if

end for

return -1

end procedure

Brute Force Efficiency

The **worst case** complexity of brute force search is $\Theta(nm)$...
...but when is this **actually** achieved?

Brute Force Efficiency

The **worst case** complexity of brute force search is $\Theta(nm)$...
...but when is this **actually** achieved?

Example. Consider the case where P contains *no repeated characters*.

Brute Force Efficiency

The **worst case** complexity of brute force search is $\Theta(nm)$...
...but when is this **actually** achieved?

Example. Consider the case where P contains *no repeated characters*.

- Claim: brute force search running time is now $O(n)$
 - In fact, at most $2n$ comparisons made!
 - *Why?*

Brute Force Efficiency

The **worst case** complexity of brute force search is $\Theta(nm)$...
...but when is this **actually** achieved?

Example. Consider the case where P contains *no repeated characters*.

- Claim: brute force search running time is now $O(n)$
 - In fact, at most $2n$ comparisons made!
 - *Why?*
- Which of these comparisons were unnecessary?
 - How can you search with fewer comparisons?

Brute Force Efficiency

The **worst case** complexity of brute force search is $\Theta(nm)$...
...but when is this **actually** achieved?

Example. Consider the case where P contains *no repeated characters*.

- Claim: brute force search running time is now $O(n)$
 - In fact, at most $2n$ comparisons made!
 - *Why?*
- Which of these comparisons were unnecessary?
 - How can you search with fewer comparisons?

More generally: How can we use results of *previous comparisons* to avoid making unnecessary comparisons in the future?

For Next Time

Consider How could we improve upon BRUTEFORCEMATCH

- How can we use information about *previous matches* in order to avoid doing some *future checks*?

Scratch Notes
