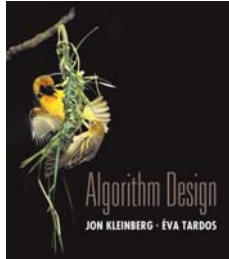---

# CS 580:  Algorithm Design and Analysis

Jeremiah Blocki
Purdue University
Spring 2018

**Announcement:** Homework 2 due on Tuesday, February 6th at 11:59PM

---

Greedy
Algorithms

**Algorithm Design**

JON KLEINBERG · ÉVA TARDOS

3

---

## Clustering

Clustering.  Given a set U of n objects labeled $p_1, ..., p_n$, classify into coherent groups.

photos, documents, micro-organisms

Distance function.  Numeric value specifying "closeness" of two objects.

number of corresponding pixels whose
intensities differ by some threshold

Fundamental problem.  Divide into clusters so that points in different clusters are far apart.
- Routing in mobile ad hoc networks.
- Identify patterns in gene expression.
- Document categorization for web search.
- Similarity searching in medical image databases
- Skycat:  cluster $10^9$ sky objects into stars, quasars, galaxies.

5

---

## Recap: Minimum Weight Spanning Trees

Cut Property: Minimum weight edge crossing a cut must be in the MST (assume edge weights are distinct)

Cycle Property: Maximum weight edge in a cycle must not be in the MST (assuming edge weights are distinct)

Prim's Algorithm
- **Repeatedly applies cut property to expand tree**
- **O(m log n) time with Binary Heap**
- **O(m+n log n) time with Fibonacci Heap**

Prim's Algorithm
- Consider edges in increasing order of weight
- For each edge we can either
  - Discard via Cycle Property, or
  - Add via Cut Property
- O(m log n) running time.

2

---

## 4.7  Clustering
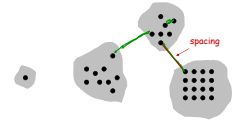


---

## Clustering of Maximum Spacing

k-clustering.  Divide objects into k non-empty groups.

Distance function.  Assume it satisfies several natural properties.
- $d(p_i, p_j) = 0$ iff $p_i = p_j$   (identity of indiscernibles)
- $d(p_i, p_j) \geq 0$           (nonnegativity)
- $d(p_i, p_j) = d(p_j, p_i)$       (symmetry)

Spacing.  Min distance between any pair of points in different clusters.

Clustering of maximum spacing.  Given an integer k, find a k-clustering of maximum spacing.



spacing

k = 4

6

---

2/1/2018

## Slide 7 — Greedy Clustering Algorithm
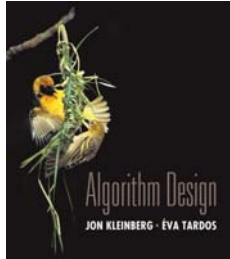
Single-link k-clustering algorithm.
- Form a graph on the vertex set U, corresponding to n clusters.
- Find the closest pair of objects such that each object is in a different cluster, and add an edge between them.
- Repeat n-k times until there are exactly k clusters.

Key observation. This procedure is precisely Kruskal's algorithm (except we stop when there are k connected components).

Remark. Equivalent to finding an MST and deleting the k-1 most expensive edges.
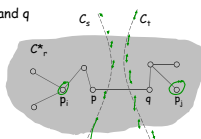
*Graph G has $O(n^2)$ edges*
*$O(n^2 \log n)$*

7

## Slide 9 — Divide and Conquer

Algorithm Design
JON KLEINBERG · ÉVA TARDOS

Slides by Kevin Wayne.
Copyright © 2005 Pearson-Addison Wesley.
All rights reserved.

9

## Slide — 5.1 Mergesort

### 5.1 Mergesort

## Slide 8 — Greedy Clustering Algorithm: Analysis

Theorem. Let $C^*$ denote the clustering $C^*_1, ..., C^*_k$ formed by deleting the k-1 most expensive edges of a MST. $C^*$ is a k-clustering of max spacing.

Pf. Let C denote some other clustering $C_1, ..., C_k$.
- The spacing of $C^*$ is the length $d^*$ of the $(k\text{-}1)^{st}$ most expensive edge.
- Let $p_i$, $p_j$ be in the same cluster in $C^*$, say $C^*_r$, but different clusters in C, say $C_s$ and $C_t$.
- Some edge (p, q) on $p_i$-$p_j$ path in $C^*_r$ spans two different clusters in C.
- All edges on $p_i$-$p_j$ path have length ≤ $d^*$ since Kruskal chose them.
- Spacing of C is ≤ $d^*$ since p and q are in different clusters. ∎

8

## Slide 10 — Divide-and-Conquer

Divide-and-conquer.
- Break up problem into several parts.
- Solve each part recursively.
- Combine solutions to sub-problems into overall solution.

Most common usage.
- Break up problem of size n into two equal parts of size ½n.
- Solve two parts recursively.
- Combine two solutions into overall solution in linear time.

Consequence.
- Brute force: $n^2$.
- Divide-and-conquer: n log n.

> Divide et impera.
> Veni, vidi, vici.
>     - Julius Caesar

10

## Slide 12 — Sorting

Sorting. Given n elements, rearrange in ascending order.

Applications.
- Sort a list of names.
- Organize an MP3 library.    obvious applications
- Display Google PageRank results.
- List RSS news items in reverse chronological order.

- Find the median.
- Find the closest pair.
- Binary search in a database.    problems become easy once items are in sorted order
- Identify statistical outliers.
- Find duplicates in a mailing list.

- Data compression.
- Computer graphics.
- Computational biology.
- Supply chain management.    non-obvious applications
- Book recommendations on Amazon.
- Load balancing on a parallel computer.
- . . .

12

Copyright 2000, Kevin Wayne

2

## Mergesort

**Mergesort.**
- Divide array into two halves.
- Recursively sort each half.
- Merge two halves to make sorted whole.

Jon von Neumann (1945)

$n^2$    $2(\frac{n}{2})^2 + n$

| A | L | G | O | R | I | T | H | M | S |
1  2  3  4  5       6

| A | L | G | O | R |   | I | T | H | M | S |    divide    $O(1)$

| A | G | L | O | R |   | H | I | M | S | T |    sort    $2T(n/2)$

| A | G | H | I | L | M | O | R | S | T |    merge    $O(n)$

$$T(n) = 2\,T\left(\frac{n}{2}\right) + O(n)$$

13

---

## Merging

**Merging.** Combine two pre-sorted lists into a sorted whole.

**How to merge efficiently?**
- Linear number of comparisons.
- Use temporary array.

05demo-merge.ppt

|   |   |   |   | O | R |   |   |   |   | M | S | T |

| A | G | H | I |   |   |   |   |   |

**Challenge for the bored.** In-place merge. [Kronrud, 1969]

using only a constant amount of extra storage

14

---

## A Useful Recurrence Relation

**Def.** $T(n)$ = number of comparisons to mergesort an input of size n.

**Mergesort recurrence.**

$$T(n) \le \begin{cases} 0 & \text{if } n=1 \\ \underbrace{T(\lceil n/2 \rceil)}_{\text{solve left half}} + \underbrace{T(\lfloor n/2 \rfloor)}_{\text{solve right half}} + \underbrace{n}_{\text{merging}} & \text{otherwise} \end{cases}$$

**Solution.** $T(n) = O(n \log_2 n)$.

**Assorted proofs.** We describe several ways to prove this recurrence. Initially we assume n is a power of 2 and replace ≤ with =.

15

---

## Proof by Recursion Tree



$$T(n) = \begin{cases} 0 & \text{if } n=1 \\ \underbrace{2T(n/2)}_{\text{sorting both halves}} + \underbrace{n}_{\text{merging}} & \text{otherwise} \end{cases}$$

merge steps

T(n)         $n$

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

T(n/4) T(n/4) T(n/4) T(n/4)   $4(n/4) = n$     $\log_2 n$

...

T(n / 2ᵏ)    $2^k (n / 2^k) = n$

...

T(2) T(2) T(2) T(2) T(2) T(2) T(2) T(2)   $n/2\,(2) = n$

$n \log_2 n$

16

---

## Proof by Telescoping

**Claim.** If $T(n)$ satisfies this recurrence, then $T(n) = n \log_2 n$.

assumes n is a power of 2

$$T(n) = \begin{cases} 0 & \text{if } n=1 \\ \underbrace{2T(n/2)}_{\text{sorting both halves}} + \underbrace{n}_{\text{merging}} & \text{otherwise} \end{cases}$$

**Pf.** For n > 1:

$$\frac{T(n)}{n} = \frac{2T(n/2)}{n} + 1$$

$$= \frac{T(n/2)}{n/2} + 1$$

$T(n) = n \log_2 n$

$$= \frac{T(n/4)}{n/4} + 1 + 1$$

...

$$= \underbrace{\frac{T(n/n)}{n/n}}_{} + \underbrace{1 + \cdots + 1}_{\log_2 n}$$

$$= \log_2 n$$

17

---

## Proof by Induction

**Claim.** If $T(n)$ satisfies this recurrence, then $T(n) = n \log_2 n$.

Guess   $T(n) \le C n^d$

assumes n is a power of 2

$$T(n) = \begin{cases} 0 & \text{if } n=1 \\ \underbrace{2T(n/2)}_{\text{sorting both halves}} + \underbrace{n}_{\text{merging}} & \text{otherwise} \end{cases}$$

**Pf.** (by induction on n)   $T(n) = T\left(\frac{n}{2}\right) + T\left(\frac{2n}{3}\right) + n$
- Base case: n = 1.
- Inductive hypothesis: $T(n) = n \log_2 n$.
- Goal: show that $T(2n) = 2n \log_2 (2n)$.

$$\begin{aligned} T(2n) &= 2T(n) + 2n \\ &= 2n \log_2 n + 2n \\ &= 2n(\log_2 (2n) - 1) + 2n \\ &= 2n \log_2 (2n) \end{aligned}$$

18

**Analysis of Mergesort Recurrence**

Claim. If $T(n)$ satisfies the following recurrence, then $T(n) \leq n \lceil \lg n \rceil$.

$$T(n) \leq \begin{cases} 0 & \text{if } n=1 \\ T(\lceil n/2 \rceil) + T(\lfloor n/2 \rfloor) + n & \text{otherwise} \end{cases}$$

solve left half — solve right half — merging — $\log_2 n$

Pf. (by induction on n)
- Base case: $n = 1$.
- Define $n_1 = \lfloor n/2 \rfloor$, $n_2 = \lceil n/2 \rceil$.
- Induction step: assume true for $1, 2, \ldots, n-1$.

$$\begin{aligned} T(n) &\leq T(n_1) + T(n_2) + n \\ &\leq n_1 \lceil \lg n_1 \rceil + n_2 \lceil \lg n_2 \rceil + n \\ &\leq n_1 \lceil \lg n_2 \rceil + n_2 \lceil \lg n_2 \rceil + n \\ &= n \lceil \lg n_2 \rceil + n \\ &\leq n(\lceil \lg n \rceil - 1) + n \\ &= n \lceil \lg n \rceil \end{aligned}$$

$$\begin{aligned} n_2 &= \lceil n/2 \rceil \\ &\leq \lceil 2^{\lceil \lg n \rceil}/2 \rceil \\ &= 2^{\lceil \lg n \rceil}/2 \\ \Rightarrow \lg n_2 &\leq \lceil \lg n \rceil - 1 \end{aligned}$$

19

---

**More General Analysis**

$n^c$

$a(n/b)^c = n^c(a/b^c)$ $\log_b n$

$n^c(a/b^c)^2$

$T(n/b^k)$ ... $n^c(a/b^c)^k$ ...

$n^c(a/b^c)^{\log_b n}$

$$T(n) \leq \begin{cases} 1 & \text{if } n=1 \\ a \times T\left(\frac{n}{b}\right) + n^c & \text{otherwise} \end{cases}$$

Case 1: $\left(\frac{a}{b^c}\right) = 1$

$$T(n) \leq \sum_{i=0}^{\log_b n} n^c \left(\frac{a}{b^c}\right)^i = n^c \log_b n$$

$$\sum n^c \left(\frac{a}{b^c}\right)^i = n^c \sum \gamma^i$$

21

---

**More General Analysis**

$n^c$

$a(n/b)^c = n^c(a/b^c)$ $\log_b n$

$n^c(a/b^c)^2$

$T(n/b^k)$ ... $n^c(a/b^c)^k$ ...

$n^c(a/b^c)^{\log_b n}$

$$T(n) \leq \begin{cases} 1 & \text{if } n=1 \\ a \times T\left(\frac{n}{b}\right) + n^c & \text{otherwise} \end{cases}$$

Case 3: $\left(\frac{a}{b^c}\right) > 1$

$$T(n) \leq \sum_{i=0}^{\log_b n} n^c \left(\frac{a}{b^c}\right)^i = \Theta(n^{\log_b a})$$

$\log_b a \geq \log_b b^c = c$

23

---

**More General Analysis**

$n^c$

$a(n/b)^c = n^c(a/b^c)$

$n^c(a/b^c)^2$ $\log_b n$

$T(n/b^k)$ ... $n^c(a/b^c)^k$ ...

$n^c(a/b^c)^{\log_b n}$

\# subproblems
size of subproblem

$$T(n) \leq \begin{cases} 1 & \text{if } n=1 \\ a \times T\left(\frac{n}{b}\right) + n^c & \text{otherwise} \end{cases}$$

merge

$$T(n) \leq \sum_{i=0}^{\log_b n} n^c \left(\frac{a}{b^c}\right)^i$$

20

---

**More General Analysis**

$n^c$

$a(n/b)^c = n^c(a/b^c)$

$n^c(a/b^c)^2$ $\log_b n$

$T(n/b^k)$ ... $n^c(a/b^c)^k$ ...

$n^c(a/b^c)^{\log_b n}$

$$T(n) \leq \begin{cases} 1 & \text{if } n=1 \\ a \times T\left(\frac{n}{b}\right) + n^c & \text{otherwise} \end{cases}$$

Case 2: $\left(\frac{a}{b^c}\right) < 1$

$$T(n) \leq \sum_{i=0}^{\log_b n} n^c \left(\frac{a}{b^c}\right)^i = \Theta(n^c)$$

$$n^c \frac{(1-\gamma)}{(1-\gamma)} \sum_{i=0}^{\infty} \gamma^i = \frac{n^c}{(1-\gamma)}$$

22

---

**5.3 Counting Inversions**

## Counting Inversions

Music site tries to match your song preferences with others.
- You rank n songs.
- Music site consults database to find people with similar tastes.

Similarity metric: number of inversions between two rankings.
- My rank: 1, 2, ..., n.
- Your rank: $a_1, a_2, ..., a_n$.
- Songs i and j inverted if i < j, but $a_i > a_j$.

*Songs*

| | A | B | C | D | E |
|---|---|---|---|---|---|
| Me | 1 | 2 | 3 | 4 | 5 |
| You | 1 | 3 | 4 | 2 | 5 |

Inversions
3-2, 4-2

Brute force: check all $\Theta(n^2)$ pairs i and j.

25

---

## Applications

Applications.
- Voting theory.
- Collaborative filtering.
- Measuring the "sortedness" of an array.
- Sensitivity analysis of Google's ranking function.
- Rank aggregation for meta-searching on the Web.
- Nonparametric statistics (e.g., Kendall's Tau distance).

26

---

## Counting Inversions: Divide-and-Conquer

Divide-and-conquer.

| 1 | 5 | 4 | 8 | 10 | 2 | 6 | 9 | 12 | 11 | 3 | 7 |
|---|---|---|---|---|---|---|---|---|---|---|---|

27

---

## Counting Inversions: Divide-and-Conquer

Divide-and-conquer.
- Divide: separate list into two pieces.

| 1 | 5 | 4 | 8 | 10 | 2 | 6 | 9 | 12 | 11 | 3 | 7 |
|---|---|---|---|---|---|---|---|---|---|---|---|

Divide: O(1).

| 1 | 5 | 4 | 8 | 10 | 2 | | 6 | 9 | 12 | 11 | 3 | 7 |

28

---

## Counting Inversions: Divide-and-Conquer

Divide-and-conquer.
- Divide: separate list into two pieces.
- Conquer: recursively count inversions in each half.

| 1 | 5 | 4 | 8 | 10 | 2 | 6 | 9 | 12 | 11 | 3 | 7 |
|---|---|---|---|---|---|---|---|---|---|---|---|

Divide: O(1).

| 1 | 5 | 4 | 8 | 10 | 2 | | 6 | 9 | 12 | 11 | 3 | 7 |

Conquer: 2T(n / 2)

5 blue-blue inversions         8 green-green inversions

5-4, 5-2, 4-2, 8-2, 10-2        6-3, 9-3, 9-7, 12-3, 12-7, 12-11, 11-3, 11-7

29

---

## Counting Inversions: Divide-and-Conquer

Divide-and-conquer.
- Divide: separate list into two pieces.
- Conquer: recursively count inversions in each half.
- Combine: count inversions where $a_i$ and $a_j$ are in different halves, and return sum of three quantities.

| 1 | 5 | 4 | 8 | 10 | 2 | 6 | 9 | 12 | 11 | 3 | 7 |
|---|---|---|---|---|---|---|---|---|---|---|---|

Divide: O(1).

| 1 | 5 | 4 | 8 | 10 | 2 | | 6 | 9 | 12 | 11 | 3 | 7 |

Conquer: 2T(n / 2)

5 blue-blue inversions         8 green-green inversions

9 blue-green inversions
5-3, 4-3, 8-6, 8-3, 8-7, 10-6, 10-9, 10-3, 10-7

Combine: ???

Total = 5 + 8 + 9 = 22.

30

## Counting Inversions: Combine

**Combine:** count blue-green inversions

- Assume each half is **sorted**.
- Count inversions where $a_i$ and $a_j$ are in different halves.
- **Merge** two sorted halves into sorted whole.

*to maintain sorted invariant*

▷ play

| 3 | 7 | 10 | 14 | 18 | 19 |  | 2 | 11 | 16 | 17 | 23 | 25 |
|---|---|----|----|----|----|--|---|----|----|----|----|----|
|   |   |    |    |    |    |  | 6 | 3  | 2  | 2  | 0  | 0  |

13 blue-green inversions: 6 + 3 + 2 + 2 + 0 + 0          Count: O(n)

| 2 | 3 | 7 | 10 | 11 | 14 | 16 | 17 | 18 | 19 | 23 | 25 |
|---|---|---|----|----|----|----|----|----|----|----|----|

Merge: O(n)

$$T(n) \leq T(\lfloor n/2 \rfloor) + T(\lceil n/2 \rceil) + O(n) \Rightarrow T(n) = O(n \log n)$$

---

# 5.4 Closest Pair of Points

---

## Closest Pair of Points: First Attempt

**Divide.** Sub-divide region into 4 quadrants.



---

## Counting Inversions: Implementation

**Pre-condition.** [Merge-and-Count] A and B are sorted.
**Post-condition.** [Sort-and-Count] L is sorted.

```
Sort-and-Count(L) {
    if list L has one element
        return 0 and the list L

    Divide the list into two halves A and B
    (r_A, A) ← Sort-and-Count(A)
    (r_B, B) ← Sort-and-Count(B)
    (r , L) ← Merge-and-Count(A, B)

    return r = r_A + r_B + r and the sorted list L
}
```

---

## Closest Pair of Points

**Closest pair.** Given n points in the plane, find a pair with smallest Euclidean distance between them.

**Fundamental geometric primitive.**
- Graphics, computer vision, geographic information systems, molecular modeling, air traffic control.
- Special case of nearest neighbor, Euclidean MST, Voronoi.

*fast closest pair inspired fast algorithms for these problems*

**Brute force.** Check all pairs of points p and q with $\Theta(n^2)$ comparisons.

**1-D version.** O(n log n) easy if points are on a line.

**Assumption.** No two points have same x coordinate.
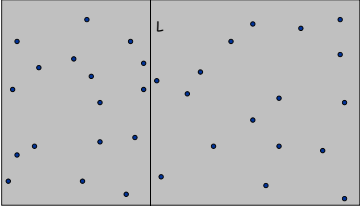
*to make presentation cleaner*

---

## Closest Pair of Points: First Attempt

**Divide.** Sub-divide region into 4 quadrants.
**Obstacle.** Impossible to ensure n/4 points in each piece.



---

**Slide 37**

Closest Pair of Points

Algorithm.
- **Divide**: draw vertical line L so that roughly ½n points on each side.

37

**Slide 38**
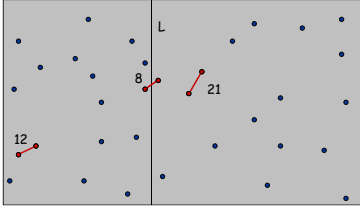
Closest Pair of Points

Algorithm.
- Divide: draw vertical line L so that roughly ½n points on each side.
- **Conquer**: find closest pair in each side recursively.

21

12

38

**Slide 39**

Closest Pair of Points

Algorithm.
- Divide: draw vertical line L so that roughly ½n points on each side.
- Conquer: find closest pair in each side recursively.
- **Combine**: find closest pair with one point in each side. ← seems like $\Theta(n^2)$
- Return best of 3 solutions.

8

21

12

39

**Slide 40**

Closest Pair of Points

Find closest pair with one point in each side, assuming that distance < δ.

21

$\delta = \min(12, 21)$

12

40

**Slide 41**

Closest Pair of Points

Find closest pair with one point in each side, assuming that distance < δ.
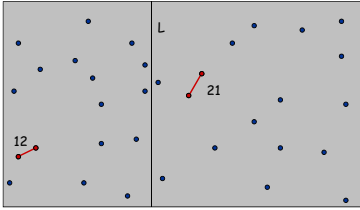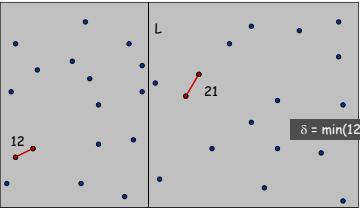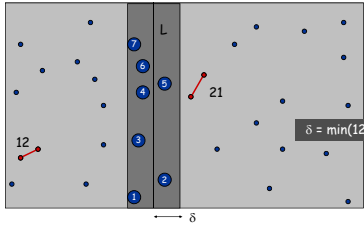- Observation: only need to consider points within δ of line L.

21

$\delta = \min(12, 21)$

12

δ

41

**Slide 42**

Closest Pair of Points

Find closest pair with one point in each side, assuming that distance < δ.
- Observation: only need to consider points within δ of line L.
- Sort points in 2δ-strip by their y coordinate.

7

6

5

4

21

3

$\delta = \min(12, 21)$

12

2

1

δ

42

7

## Closest Pair of Points

Find closest pair with one point in each side, assuming that distance < δ.
- Observation: only need to consider points within δ of line L.
- Sort points in 2δ-strip by their y coordinate.
- Only check distances of those within 11 positions in sorted list!



δ = min(12, 21)

---

## Closest Pair Algorithm

```
Closest-Pair(p₁, …, pₙ) {
    Compute separation line L such that half the points      O(n log n)
    are on one side and half on the other side.

    δ₁ = Closest-Pair(left half)                             2T(n / 2)
    δ₂ = Closest-Pair(right half)
    δ  = min(δ₁, δ₂)

    Delete all points further than δ from separation line L  O(n)

    Sort remaining points by y-coordinate.                   O(n log n)

    Scan points in y-order and compare distance between      O(n)
    each point and next 11 neighbors. If any of these
    distances is less than δ, update δ.

    return δ.
}
```

---

## MST Algorithms: Theory

Deterministic comparison based algorithms.
- $O(m \log n)$              [Jarník, Prim, Dijkstra, Kruskal, Boruvka]
- $O(m \log \log n)$.        [Cheriton-Tarjan 1976, Yao 1975]
- $O(m \, \beta(m, n))$.     [Fredman-Tarjan 1987]
- $O(m \log \beta(m, n))$.   [Gabow-Galil-Spencer-Tarjan 1986]
- $O(m \, \alpha(m, n))$.    [Chazelle 2000]

Holy grail.  $O(m)$.

Notable.
- $O(m)$ randomized.         [Karger-Klein-Tarjan 1995]
- $O(m)$ verification.       [Dixon-Rauch-Tarjan 1992]

Euclidean.
- 2-d: $O(n \log n)$.        compute MST of edges in Delaunay
- k-d: $O(k \, n^2)$.        dense Prim

---

## Closest Pair of Points



Def.  Let $s_i$ be the point in the 2δ-strip, with the $i^{th}$ smallest y-coordinate.

Claim.  If $|i - j| \geq 12$, then the distance between $s_i$ and $s_j$ is at least δ.
Pf.
- No two points lie in same $\frac{1}{2}$δ-by-$\frac{1}{2}$δ box.
- Two points at least 2 rows apart have distance $\geq 2(\frac{1}{2}\delta)$.  •

Fact.  Still true if we replace 12 with 7.

---

## Closest Pair of Points: Analysis

Running time.

$$T(n) \leq 2T(n/2) + O(n \log n) \;\Rightarrow\; T(n) = O(n \log^2 n)$$

Q.  Can we achieve O(n log n)?

A.  Yes. Don't sort points in strip from scratch each time.
- Each recursive returns two lists: all points sorted by y coordinate, and all points sorted by x coordinate.
- Sort by merging two pre-sorted lists.

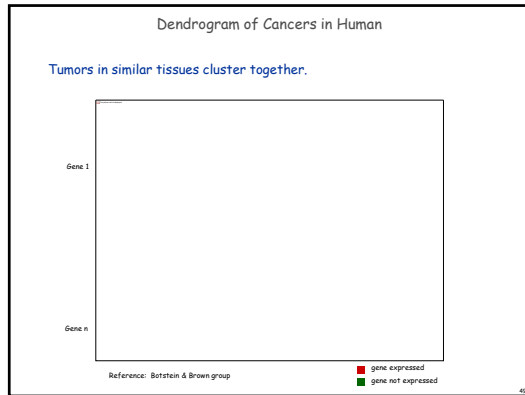$$T(n) \leq 2T(n/2) + O(n) \;\Rightarrow\; T(n) = O(n \log n)$$

---

## Dendrogram

Dendrogram.  Scientific visualization of hypothetical sequence of evolutionary events.
- Leaves = genes.
- Internal nodes = hypothetical ancestors.



Reference: http://www.biostat.wisc.edu/bmi576/fall-2003/lecture13.pdf

## Dendrogram of Cancers in Human

Tumors in similar tissues cluster together.

Gene 1

Gene n

Reference: Botstein & Brown group

■ gene expressed
■ gene not expressed

49

## Extra Slides