CSCI203 Algorithms and Data Structures

Multi-dimensional Search Trees

Lecturer: Dr. Xueqiao Liu

Room 3.117

Email: xueqiao@uow.edu.au

13/10/2024

Query Types

- Exact match query: Asks for the object(s) whose key matches query key exactly.
- Range query: Asks for the objects whose key lies in a specified query range (interval).
- Nearest-neighbor query: Asks for the objects whose key is "close" to query key.

Exact Match Query

Suppose that we store employee records in a database:

- Example:
 - key=ID: retrieve the record with ID=12345

Range Query

- Example:
 - key=Age: retrieve all records satisfying20 < Age < 50
 - key= #Children: retrieve all records satisfying 1 < #Children < 4</p>

ID Name Age Salary #Children

Nearest-Neighbor(s) (NN) Query

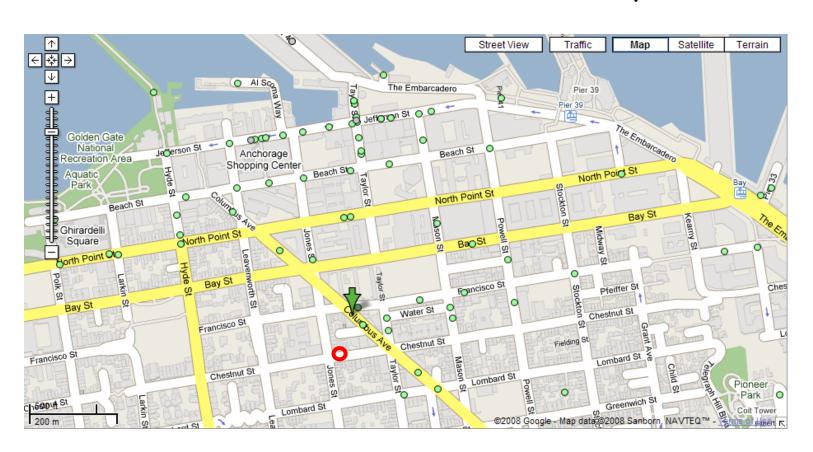
Example:

- key=Salary: retrieve the employee whose salary is closest to \$50,000 (i.e., 1-NN).
- key=Age: retrieve the 5 employees whose age is closest to 40 (i.e., k-NN, k=5).

ID Name Age Salary #Children

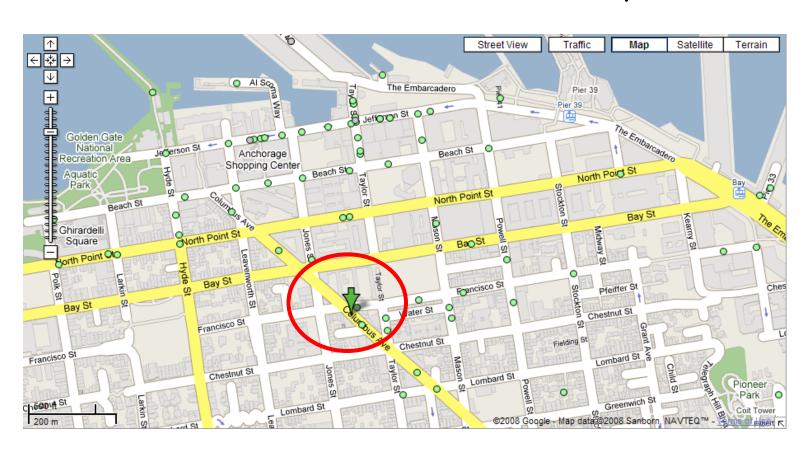
Nearest Neighbor(s) Query

What is the closest restaurant to my hotel?



Nearest Neighbor(s) Query...

Find the 4 closest restaurants to my hotel



Multi-dimensional Query

- In practice, queries might involve multidimensional keys.
 - key=(Name, Age): retrieve all records with Name="George" and "50 <= Age <= 70"</p>

ID Name Age Salary #Children

Nearest Neighbor Query in High Dimensions

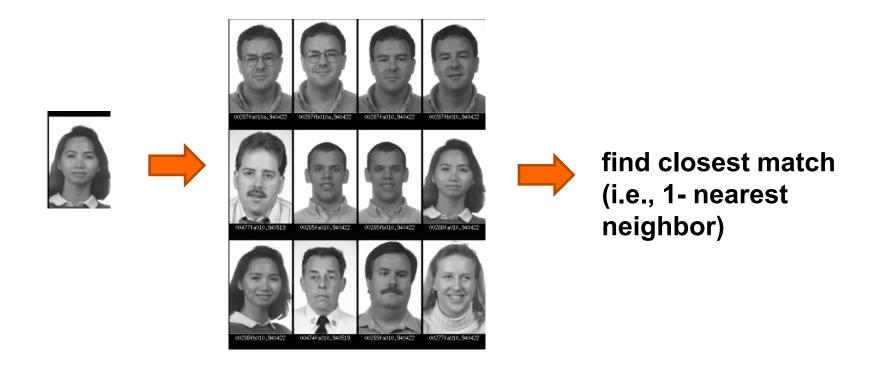
- Very important and practical problem!
 - Image retrieval



find K closest matches (i.e., K Nearest Neighbors)

Nearest Neighbor Query in High Dimensions

Face recognition



We will discuss ...

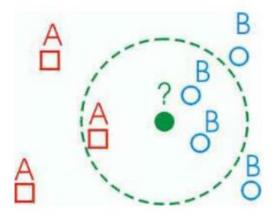
Range trees

▶ KD-trees

Interpreting Queries Geometrically

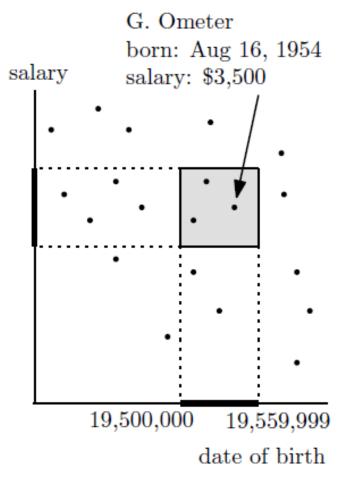
Multi-dimensional keys can be thought as "points" in high dimensional spaces.

Queries about records -> Queries about points



Example 1- Range Search in 2D

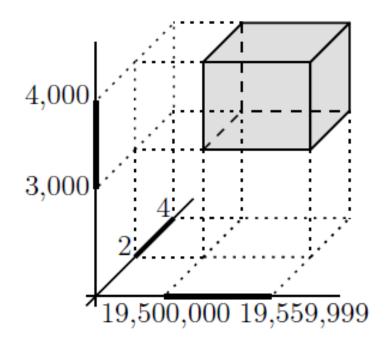
A database query may ask for all employees with age between a_1 and a_2 , and salary between s_1 and s_2



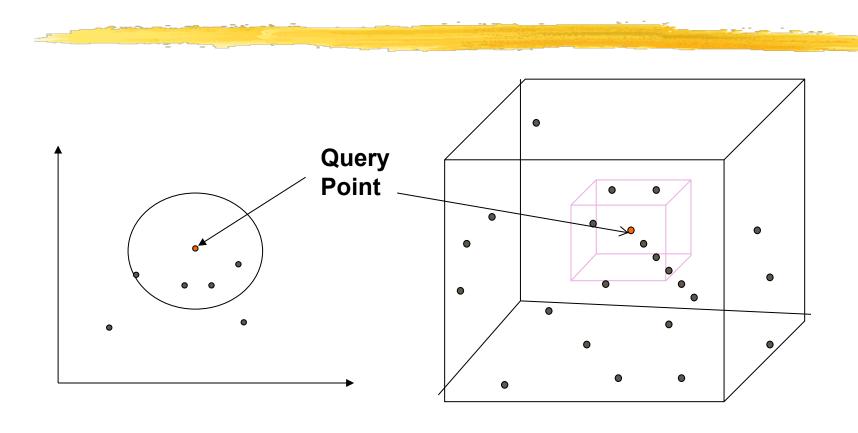
age = $10,000 \times year + 100 \times month + day$

Example 2 - Range Search in 3D

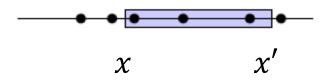
Example of a 3-dimensional (orthogonal) range query: children in [2,4], salary in [3000,4000], date of birth in [19,500,000,19,559,999]



Example 3 - Nearest Neighbors Search



- Data
 - $P = \{p_1, p_2, \dots, p_n\}$ in 1D space (a set of real numbers)
- Query
 - Which points are in the interval [x, x']



13/10/2024

Range: [x, x']

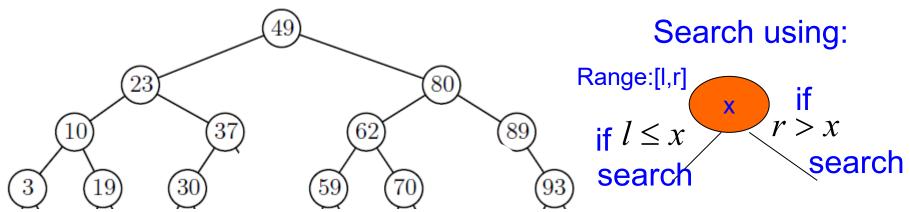
Data structure 1: Sorted Array

- A= 3 9 27 28 29 98 141 187 200 201 202 999
- Query: Search for x & x' in A by binary searchO(logn)
 Output all points between them. O(k)
 Total O(logn+k)

Example: retrieve all points in [25, 90]

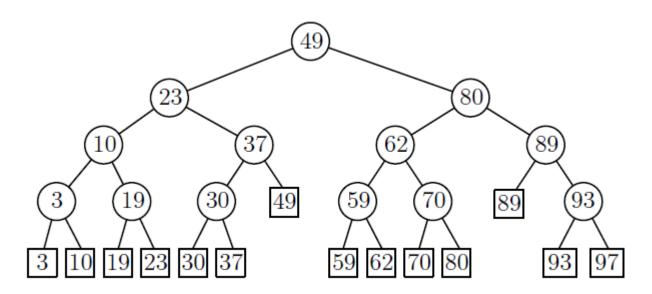
Does not generalize well to high dimensions.

- Data Structure 2: BST
 - Search using binary search property.
 - Some subtrees are eliminated during search.

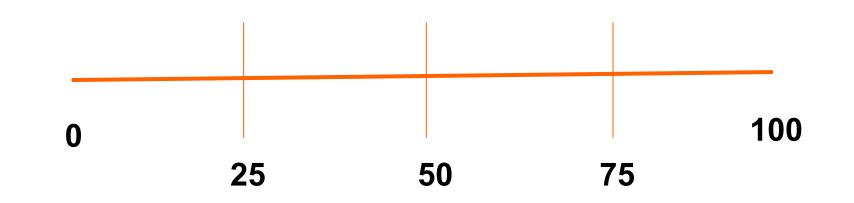


Example: retrieve all points in [25, 90]

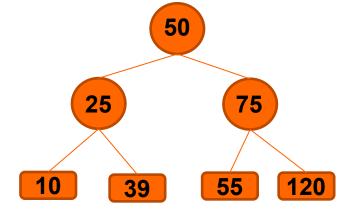
- Data Structure 3: BST with data stored in leaves
 - Internal nodes store splitting values (i.e., not necessarily same as data).
 - Data points are stored in the leaf nodes.



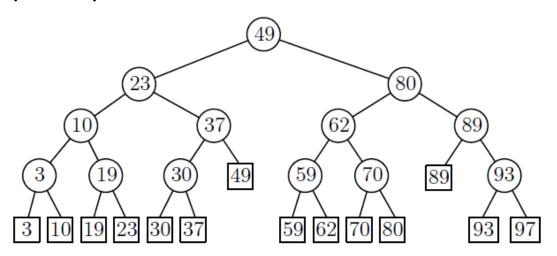
BST with data stored in leaves



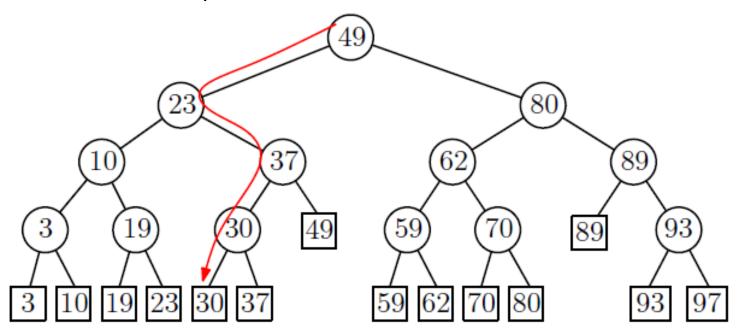
Data: 10, 39, 55, 120



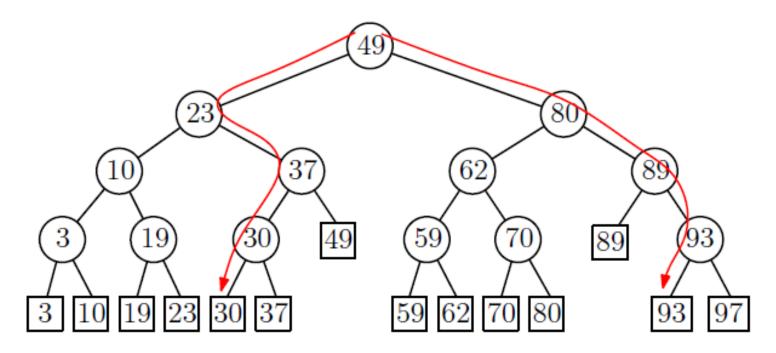
- Retrieving data in [x, x']
 - Perform binary search twice, once using x and the other using x^{\prime}
 - Suppose binary search ends at leaves l and l'
 - The points in [x, x'] are the ones stored between l and l' plus, possibly, the points stored in l and l'



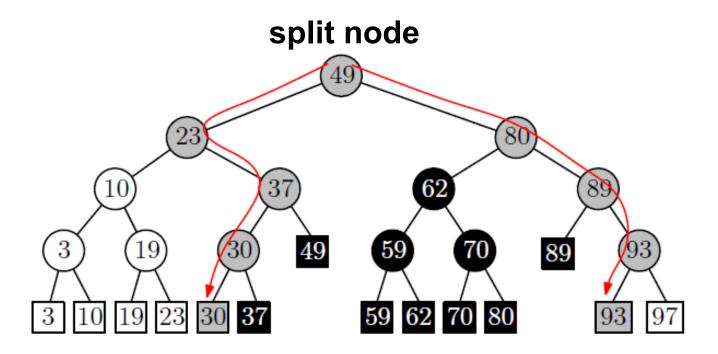
- **Example:** retrieve all points in [25,90]
 - The search path for 25 is:



> The search for 90 is:



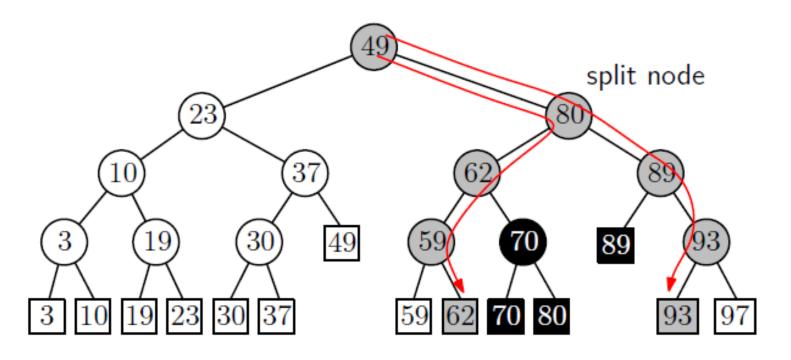
Examine the leaves in the sub-trees between the two traversing paths from the root.

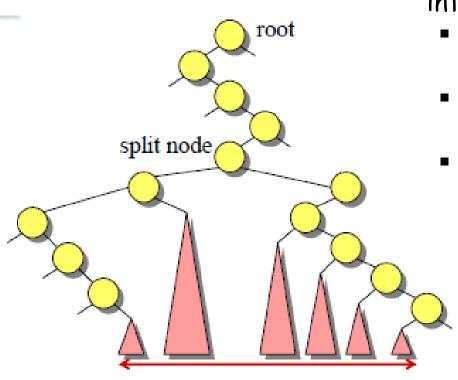


retrieve all points in [25, 90]

1D Range Search - Another Example

A 1-dimensional range query with [61, 90]

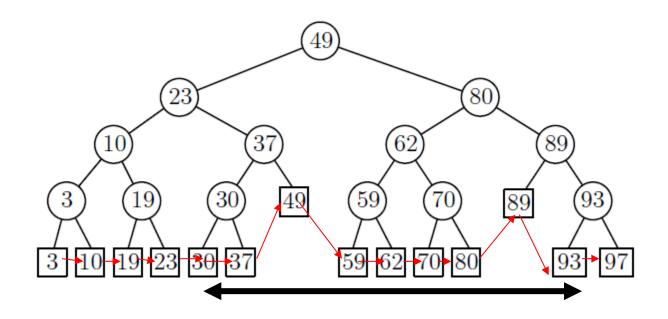




- How do we find the leaves of interest?
 - Find split node (i.e., node where the paths to x and x' split).
 - Left turn: report leaves in right subtrees
 - Right turn: report leaves in left substrees

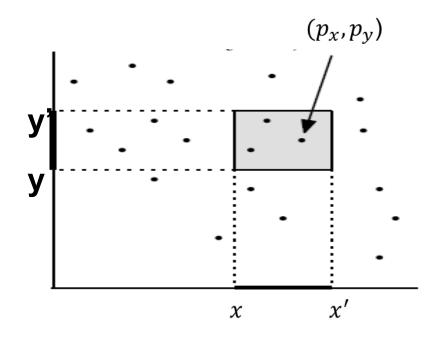
O(logn + k) time where k is the number of items reported.

Speed-up search by keeping the leaves in sorted order using a linked-list.

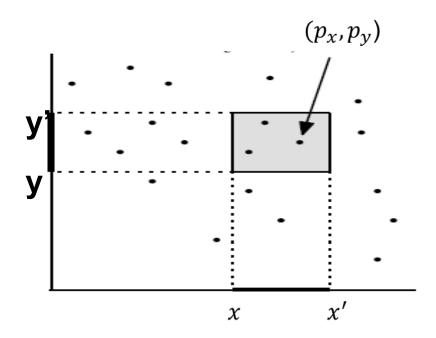


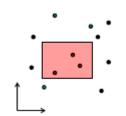
- \square A 2D range query asks for the points inside a query rectangle $[x, x'] \times [y, y']$
 - \Box A point (p_x, p_y) lies in this rectangle if and only if:

$$p_x \in [x, x']$$
 and $p_y \in [y, y']$

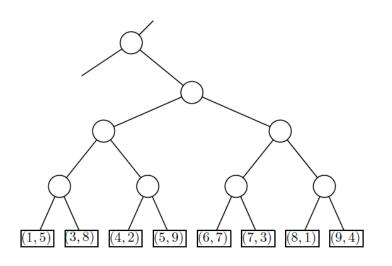


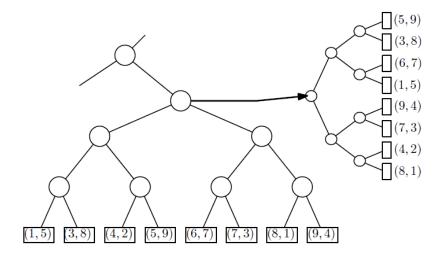
- A 2D range query can be decomposed in two 1D range queries:
 - One on the x-coordinate of the points.
 - The other on the y-coordinates of the points.

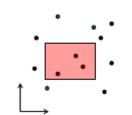


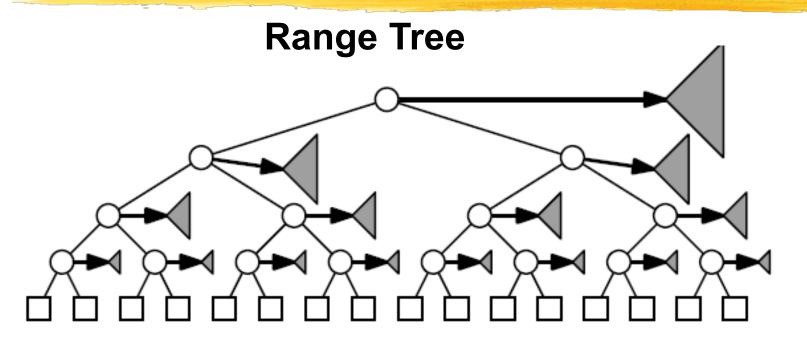


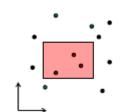
- Store a *primary 1D range tree* for all the points based on *x-coordinate*.
- For each node, store a **secondary 1D range tree** based on *y-coordinate*.



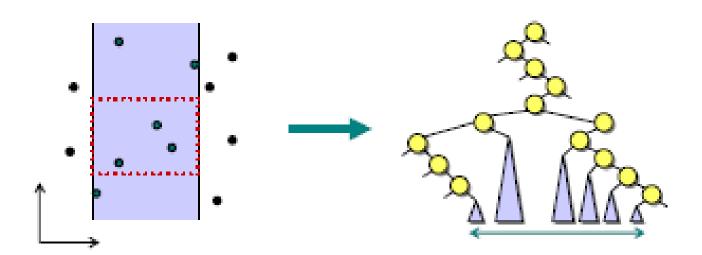




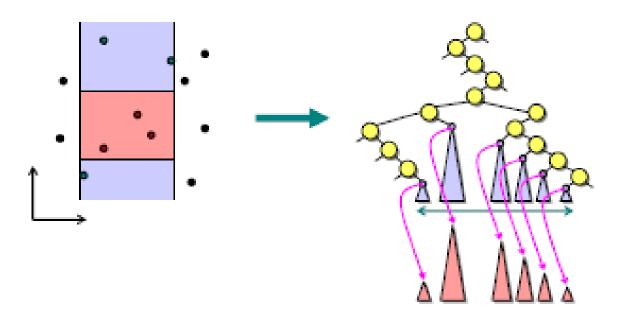




- Search using the x-coordinate only.
- ▶ How to restrict to points with proper y-coordinate?



Recursively search within each subtree using the y-coordinate.



Query cost: O(log²n+k)

Range Search in d dimensions

1D query time: O(logn + k)

2D query time: $O(log^2n + k)$

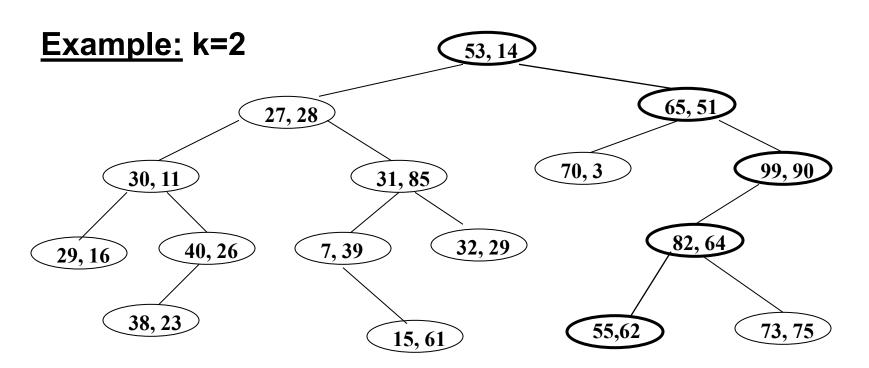
d dimensions:

Query time: $O(k + \log^d n)$ to report k points.

Space: $O(n \log^{d-1} n)$

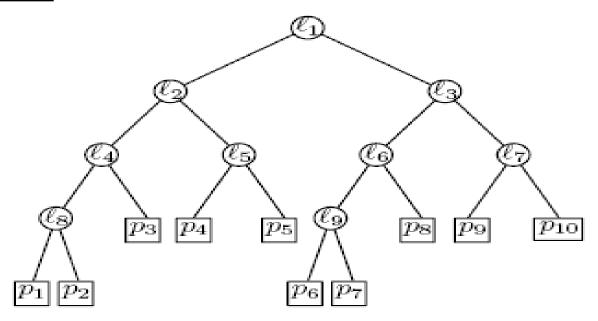
KD Tree

 A binary search tree where every node is a k-dimensional point.



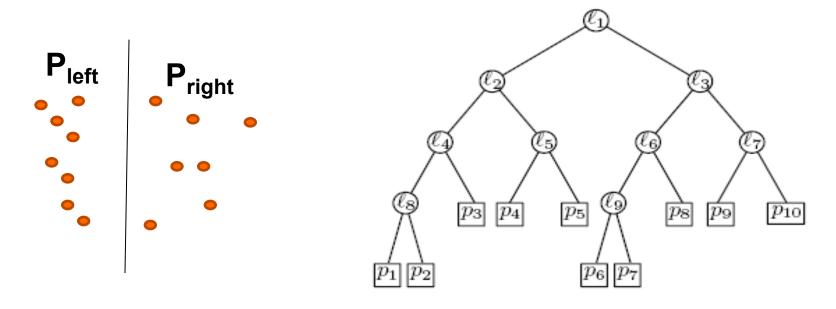
KD Tree...

Example: data stored at the leaves



KD Tree...

- Every node (except leaves) represents a hyperplane that divides the space into two parts.
- Points to the left (right) of this hyperplane represent the left (right) sub-tree of that node



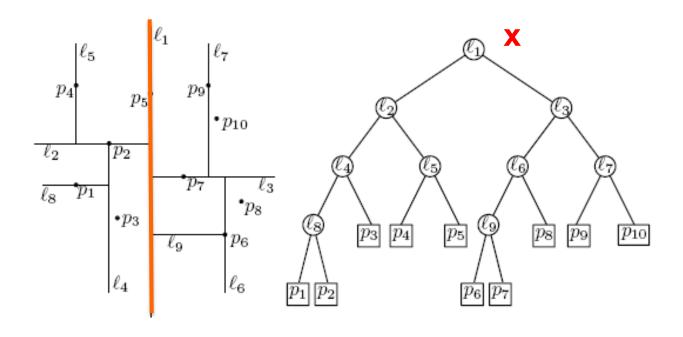
KD Tree...

As we move down the tree, we divide the space along alternating (but not always) axis-aligned hyperplanes:

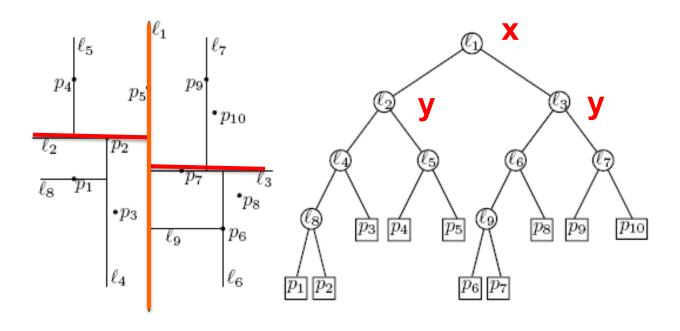
Split by x-coordinate: split by a vertical line that has (ideally) half the points left or on, and half right.

Split by y-coordinate: split by a horizontal line that has (ideally) half the points below or on and half above.

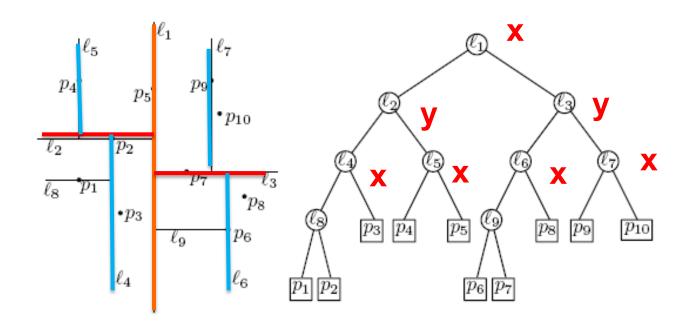
Split by x-coordinate: split by a vertical line that has approximately half the points left or on, and half right.



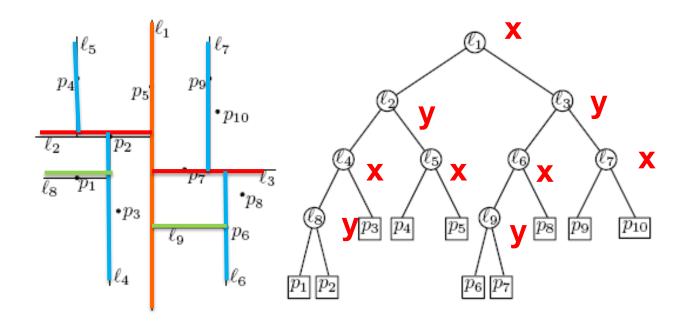
Split by y-coordinate: split by a horizontal line that has half the points below or on and half above.



Split by x-coordinate: split by a vertical line that has half the points left or on, and half right.

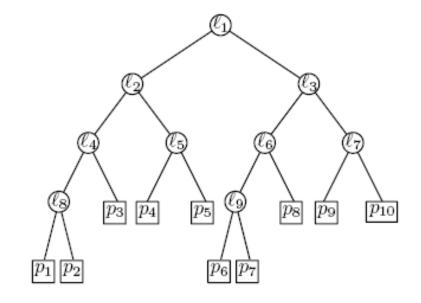


<u>Split by y-coordinate</u>: split by a horizontal line that has half the points below or on and half above.



Node Structure

- > A KD-tree node has 5 fields
 - Splitting axis
 - Splitting value
 - Data
 - Left pointer
 - Right pointer

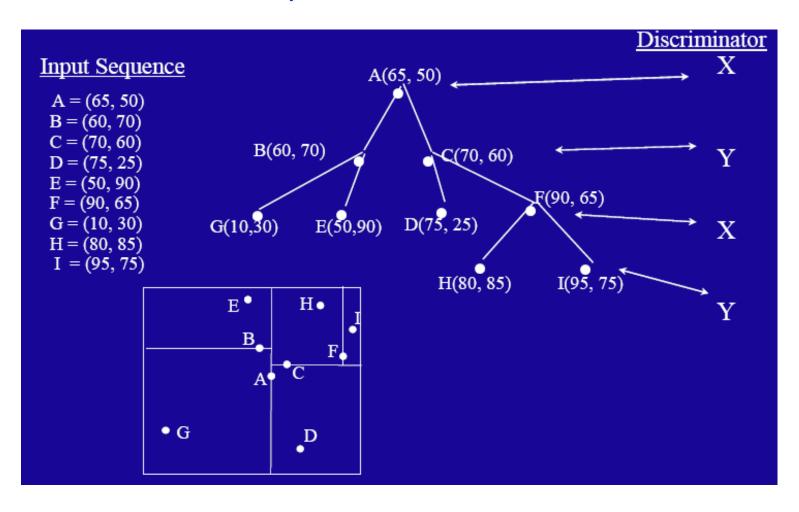


Splitting Strategies

- Divide based on order of point insertion
 - Assumes that points are given one at a time.
- Divide by finding median
 - Assumes all the points are available ahead of time.
- Divide perpendicular to the axis with widest spread
 - Split axes might not alternate

... and more!

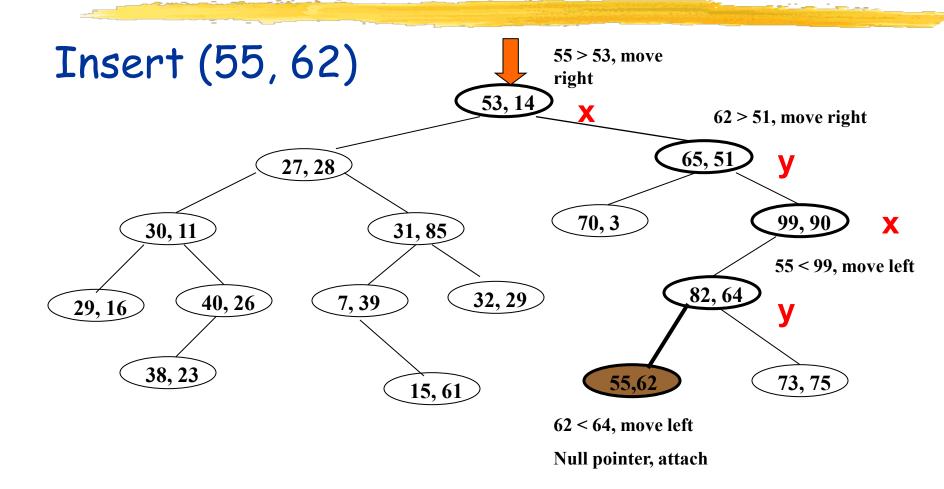
Example - using order of point insertion (data stored at nodes)



KD Tree...

- Let's discuss
 - Insert
 - Delete
 - Search

Insert new data



Delete data

- Suppose we need to remove p = (a, b)
 - Find node t which contains p
 - If t is a leaf node, replace it by null
 - Otherwise, find a replacement node r = (c, d) see below!
 - Replace (a, b) by (c, d)
 - Remove r
- Finding the replacement r = (c, d)
 - If t has a right child, use the successor*
 - Otherwise, use node with minimum value* in the left subtree

)

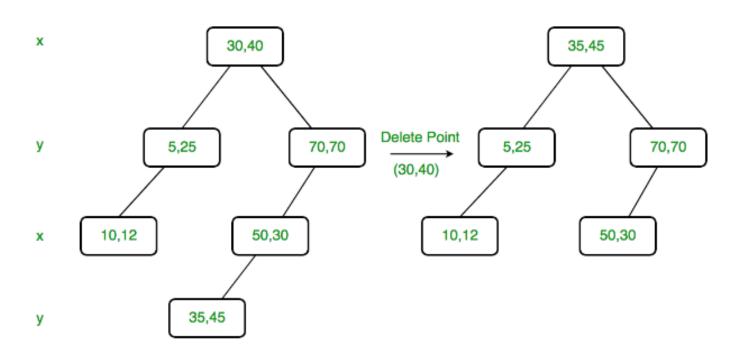
Delete data

- How to find the replacement (c,d) for (a,b)?
 - *(depending on what axis the node discriminates)
- 1. If (a,b) has right child as not NULL
 - 1) Find minimum of current node's dimension in right subtree.
 - 2) Replace the node with above found minimum and recursively delete minimum in right subtree.

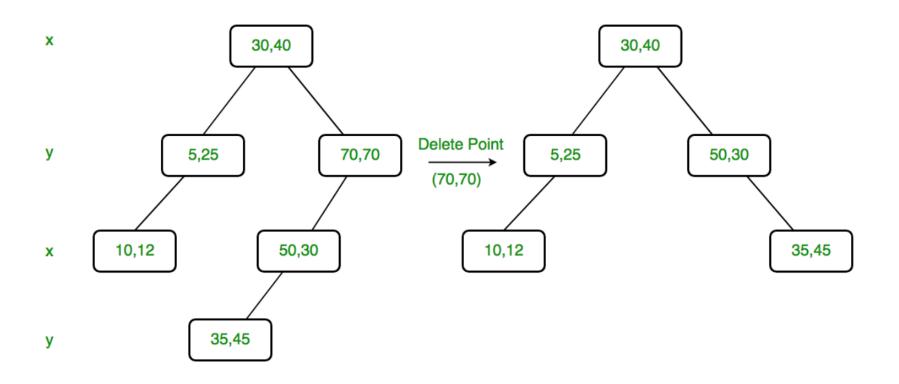
Delete Data

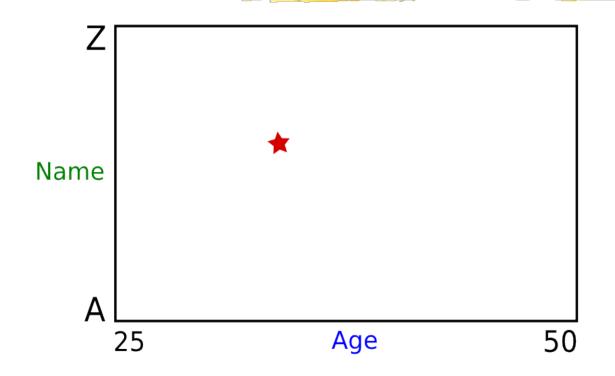
- ▶ How to find the replacement (c,d) for (a,b)?
 - *(depending on what axis the node discriminates)
- 2. Else If node to be deleted has left child as not NULL
 - 1) Find minimum of current node's dimension in left subtree
 - 2) Replace the node with above found minimum and recursively delete minimum in left subtree.
 - 3) Make new left subtree as right child of current node.

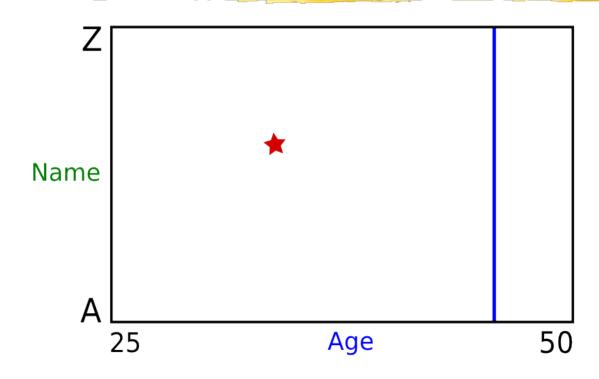
Delete Data



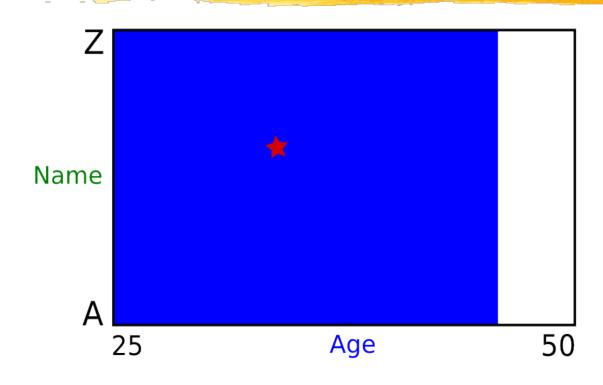
Delete Data



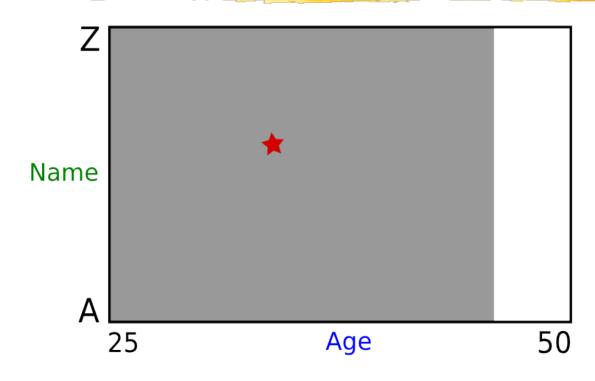




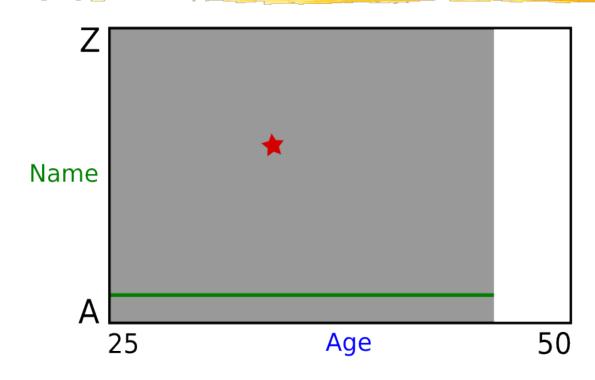
Current node's key: Age=45



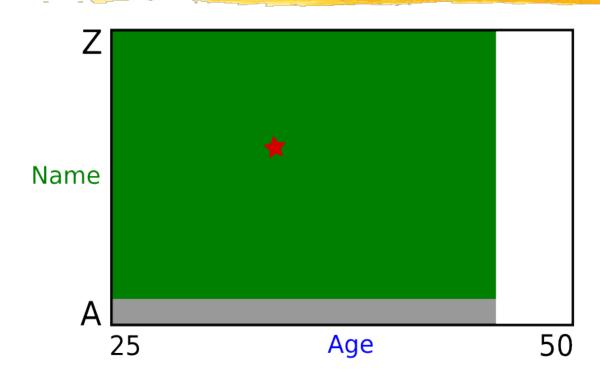
Current node's key: Age=45



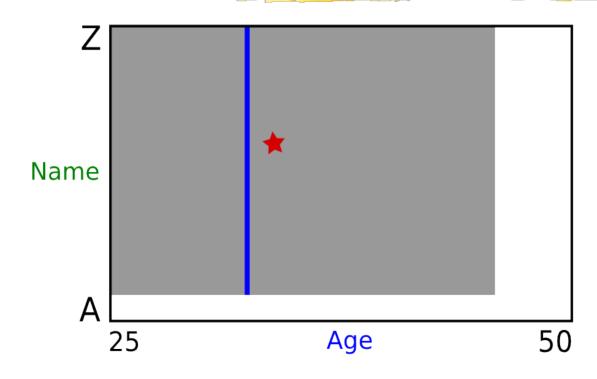
Current node's key: Age=45



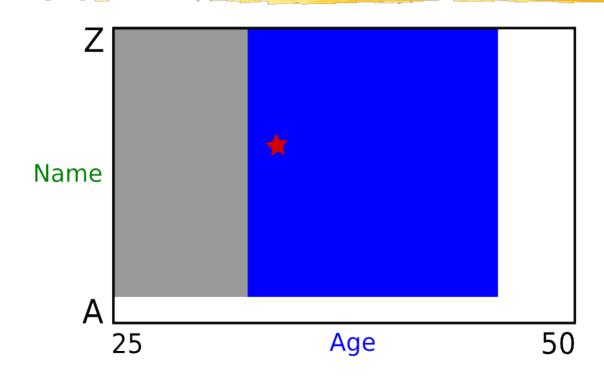
Current node's key: Name=B



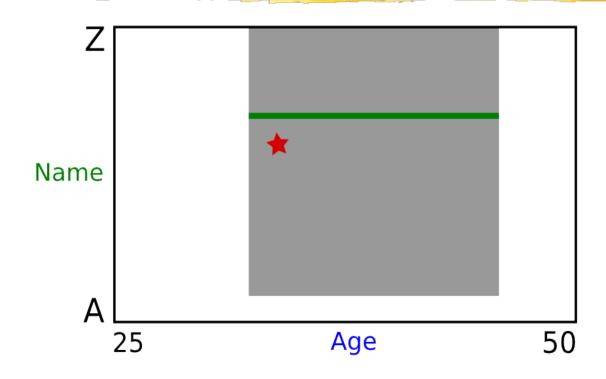
Current node's key: Name=B



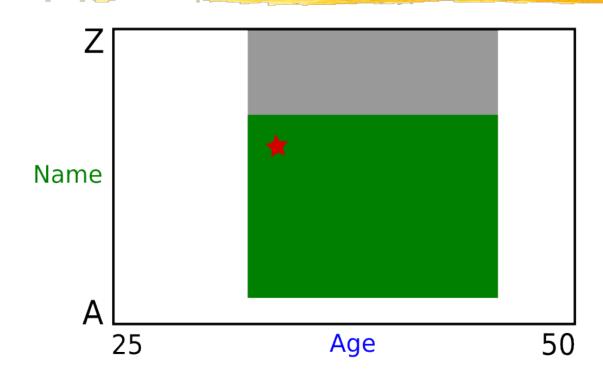
Current node's key: Age=30



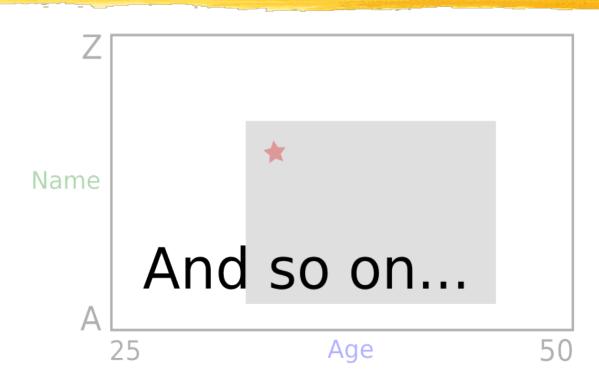
Current node's key: Age=30



Current node's key: Name=R



Current node's key: Name=R



Current node's key: Name=R

KD Tree - Complexity

- **Construction** $O(dn \log n)$
 - Sort points in each dimension: $O(dn \log n)$
 - Determine splitting line (median finding): O(dn)
- Space requirements:
 - O(n)
- Query requirements:
 - KD tree: $O(n^{1-\frac{1}{d}} + k)$

O(n+k) as d increases!

Related References

- Slides are modified based on the slides created by Dr George Bebis from Department of Computer Science & Engineering University of Nevada (UNR) Reno
- https://en.wikipedia.org/wiki/K-d_tree