

Data Structures

K-d Tree

CS 225

September 29, 2023

Brad Solomon & G Carl Evans



UNIVERSITY OF
ILLINOIS
URBANA - CHAMPAIGN

Department of Computer Science

MP_Lists Plagiarism Report

Significant increase in plagiarism

Still processing all the FAIR cases

Remember course policies!

MP_Mosaic Extra Credit Extension

Today's lecture will 'review' several key concepts

Concepts may be new to some, extra credit is extended

Extra credit deadline: Wednesday

Learning Objectives

Discuss (one) extension beyond BST

Introduce lambda functions in C++

Finish AVL proof and introduce B-Trees

Summary of Balanced BST

AVL Trees

- Max height: $1.44 * \lg(n)$

- Rotations:

 - Zero rotations on find

 - One rotation on insert

 - $O(h) == O(\lg(n))$ rotations on remove

Range-based Searches

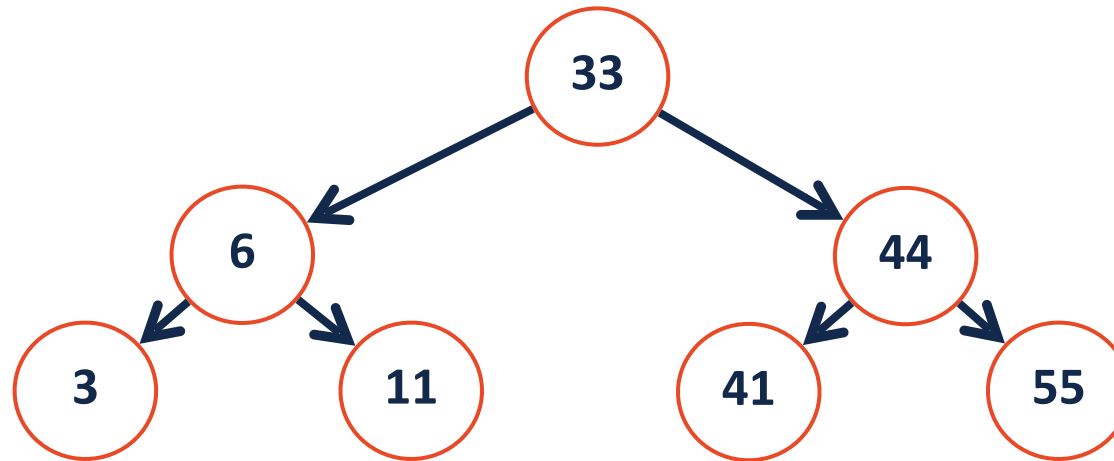
Balanced BSTs are useful structures for range-based and nearest-neighbor searches.

Q: Consider points in 1D: $\mathbf{p} = \{p_1, p_2, \dots, p_n\}$.
...what points fall in $[11, 42]$?



Range-based Searches

Q: Consider points in 1D: $\mathbf{p} = \{p_1, p_2, \dots, p_n\}$.
...what points fall in $[11, 42]$?



Red-Black Trees in C++

C++ provides us a balanced BST as part of the standard library:

```
std::map<K, V> map;
```

```
V & std::map<K, V>::operator[] ( const K & )
```

```
std::map<K, V>::erase ( const K & )
```


Red-Black Trees in C++



C++ provides us a balanced BST as part of the standard library:

```
iterator std::map<K, V>::lower_bound( const K & );
```

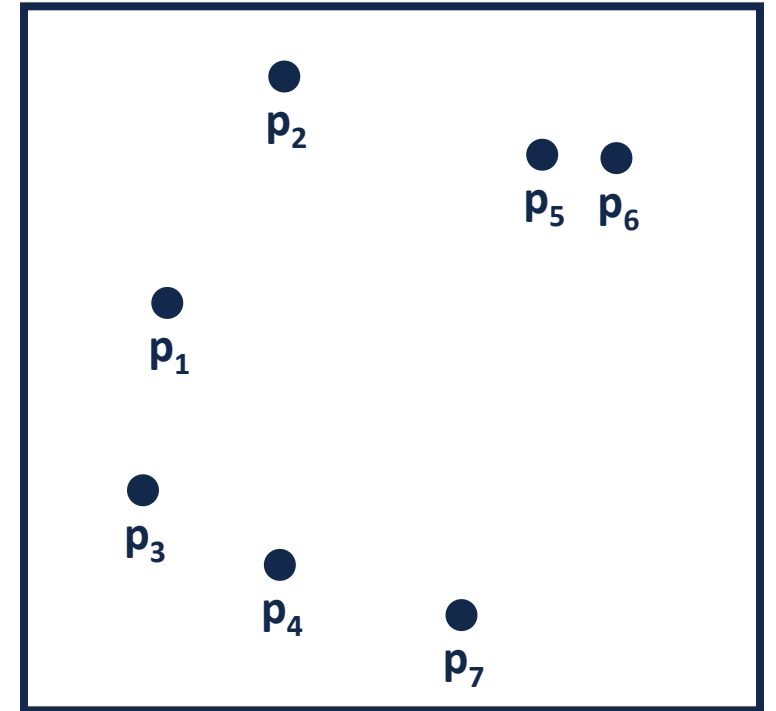
```
iterator std::map<K, V>::upper_bound( const K & );
```

Range-based Searches

Consider points in 2D: $\mathbf{p} = \{p_1, p_2, \dots, p_n\}$.

Q: What points are in the rectangle:
[$(x_1, y_1), (x_2, y_2)$]?

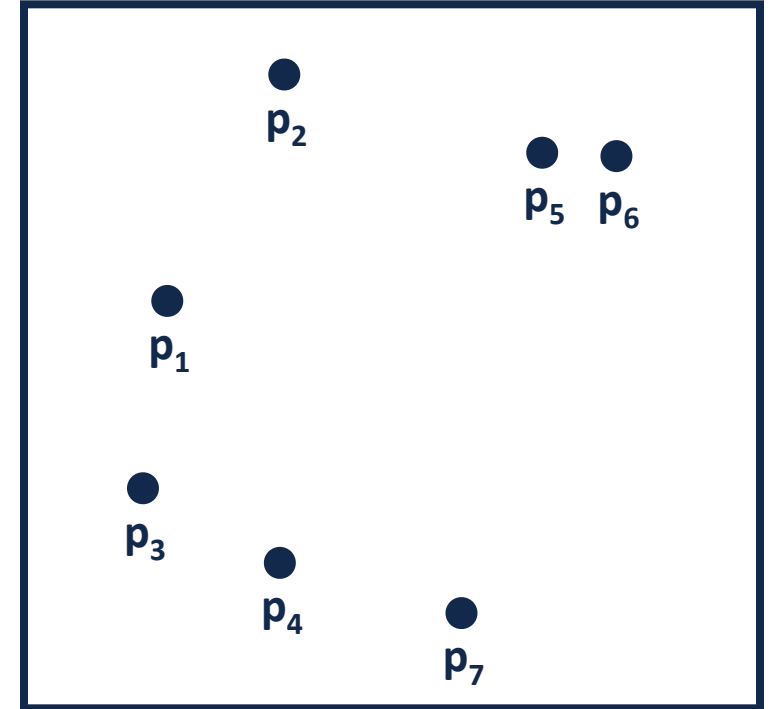
Q: What is the nearest point to (x_1, y_1) ?



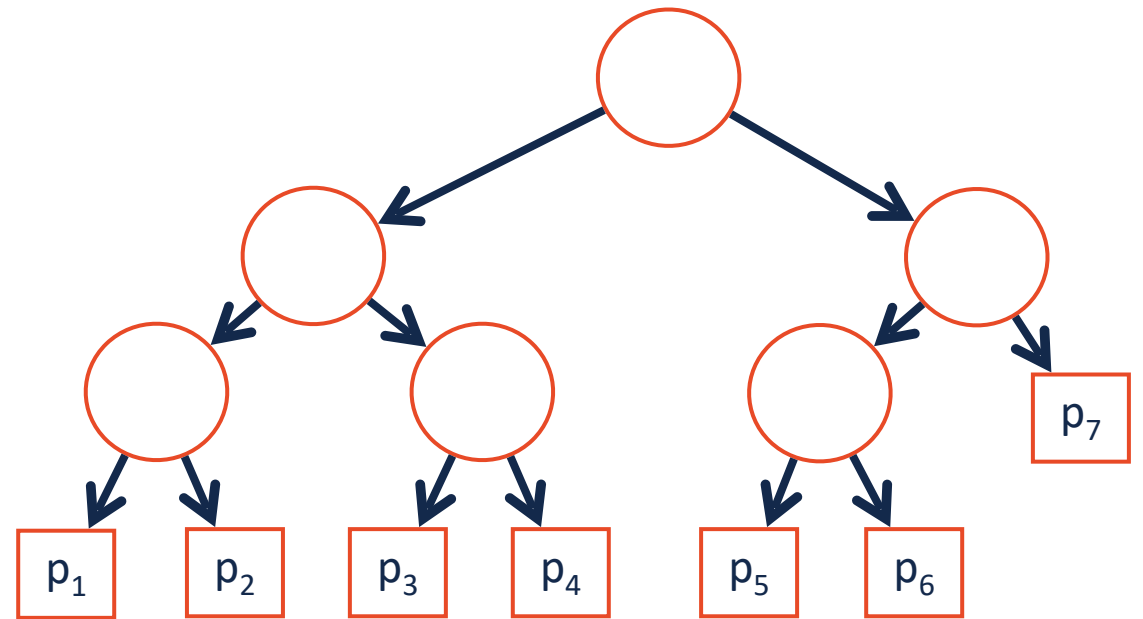
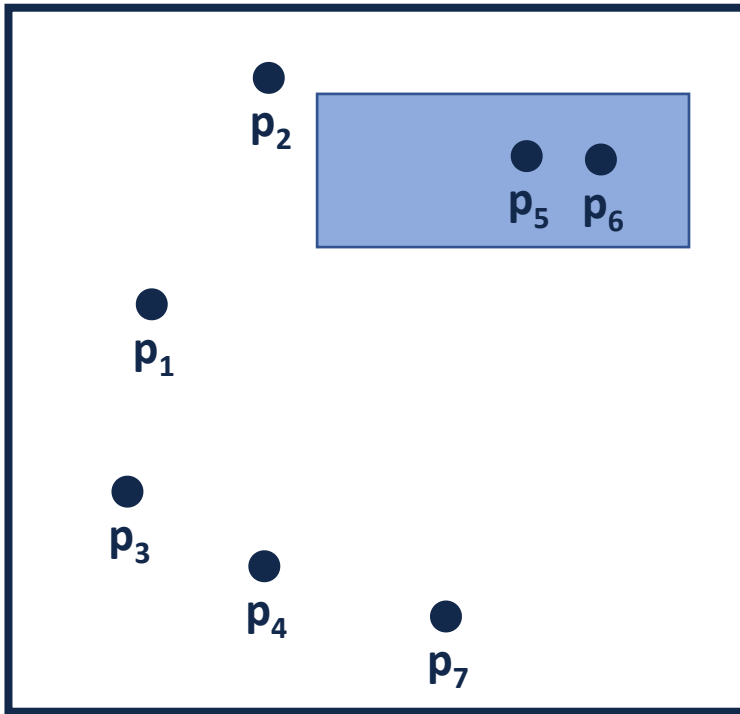
Range-based Searches

Consider points in 2D: $\mathbf{p} = \{p_1, p_2, \dots, p_n\}$.

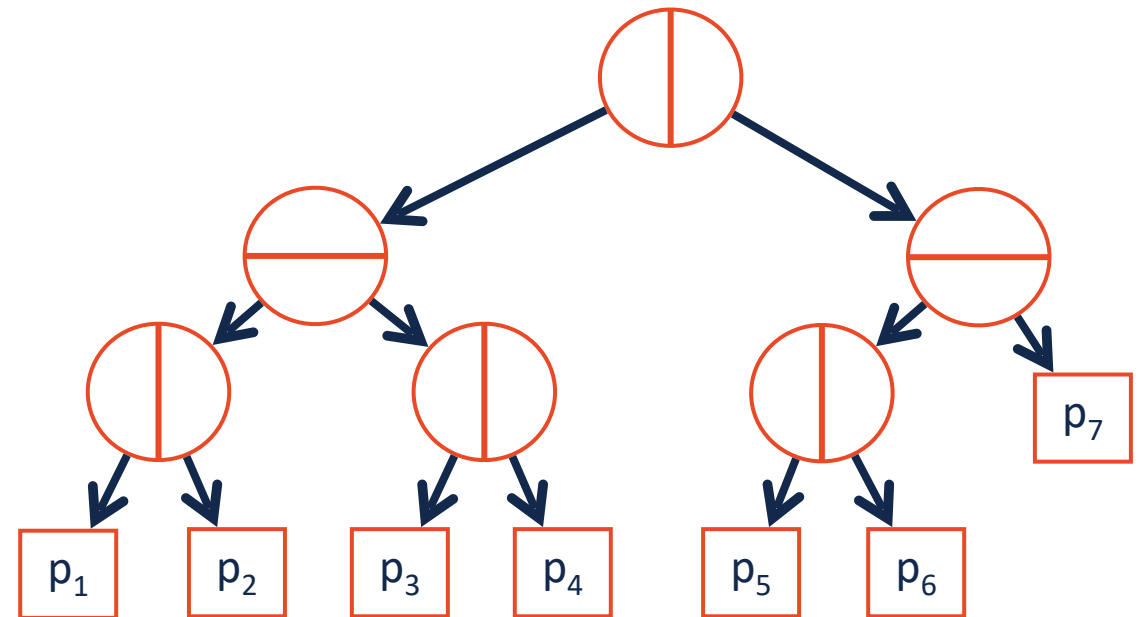
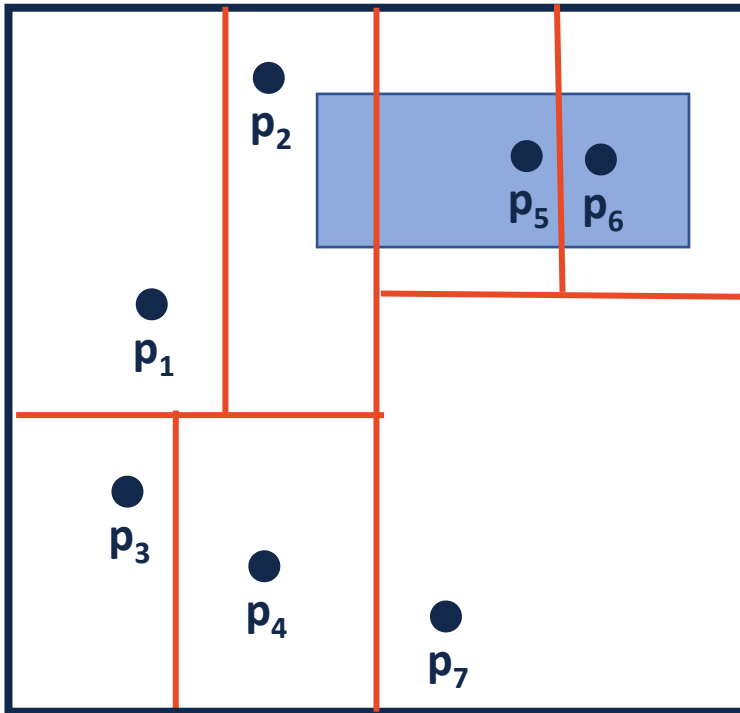
Tree construction:



Range-based Searches



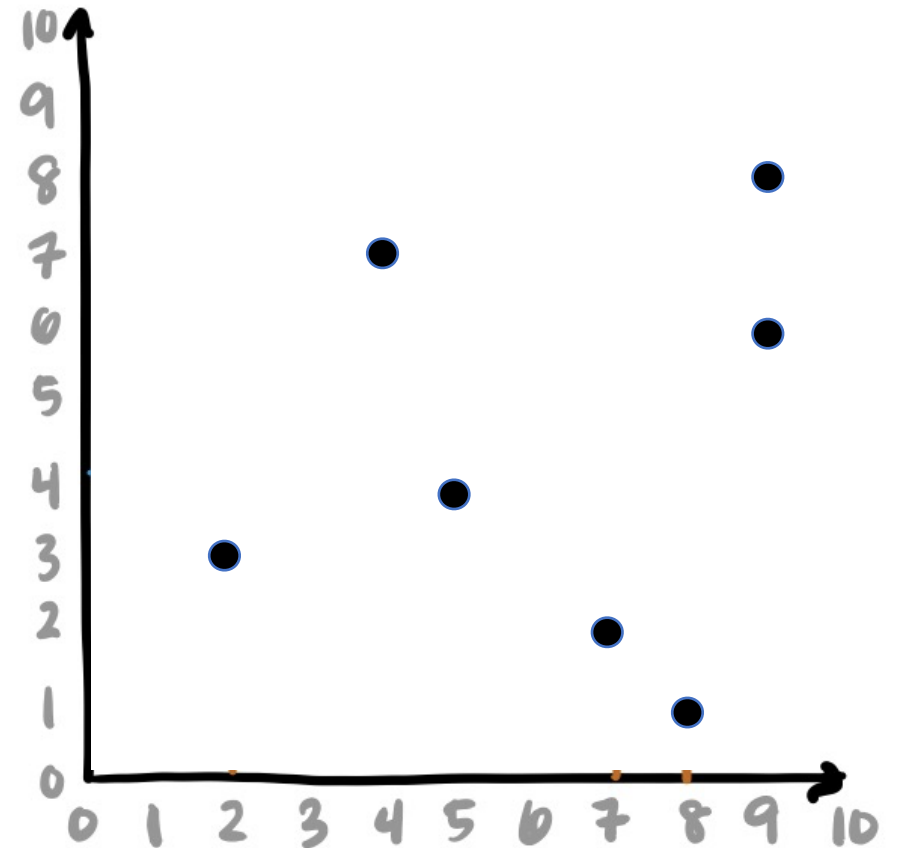
Range-based Searches



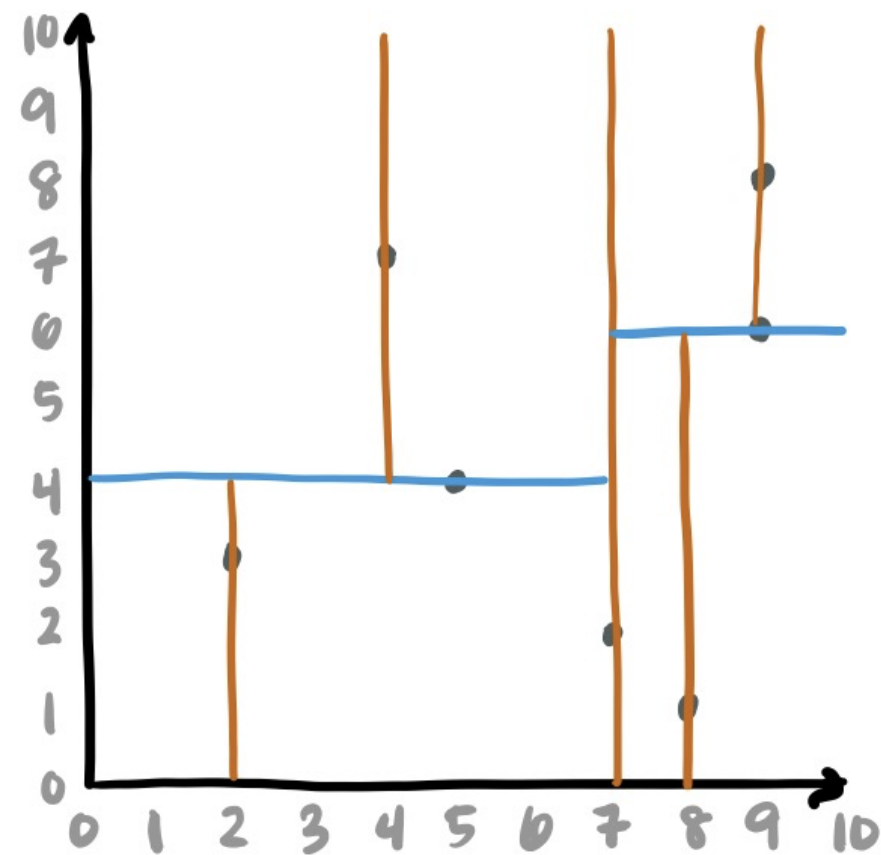
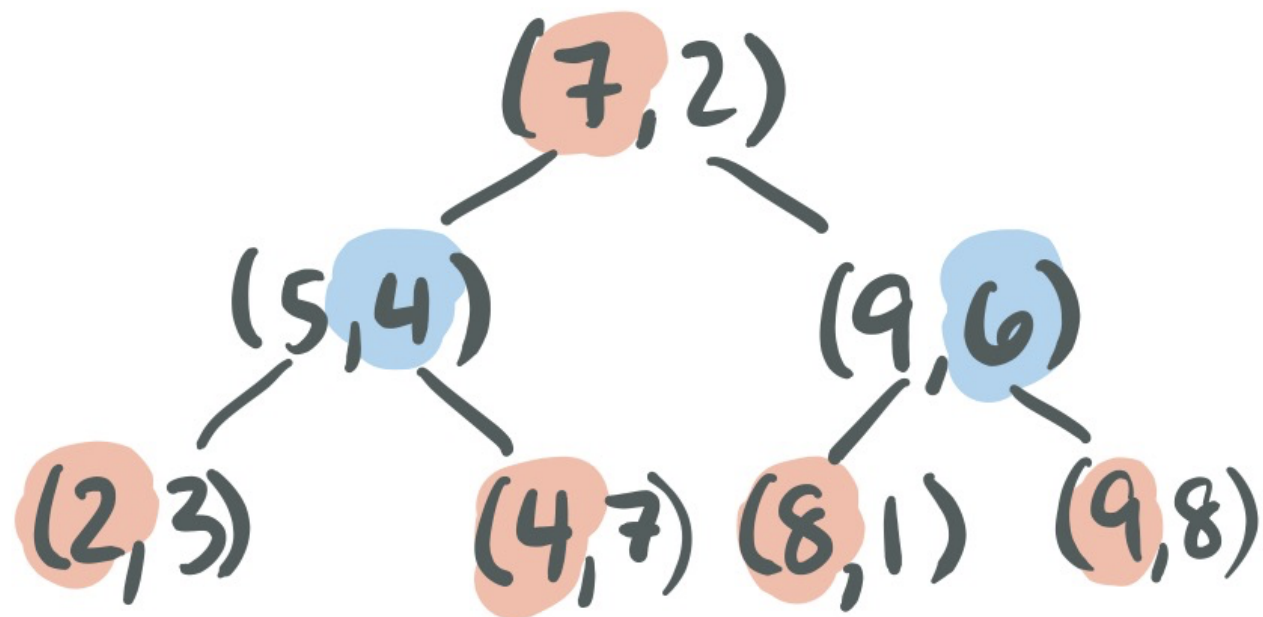
Nearest Neighbor: k-d tree

A **k-d tree** is similar but splits on points:

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



Nearest Neighbor: k-d tree



Nearest Neighbor: k-d tree

This construction seems easy conceptually but...

1. Review, understand, and use **quickselect**
2. Review, understand, and use **lambda functions**

Functions as arguments

Consider the function from Excel
`COUNTIF(range, criteria)`

	A	B	C
1	1		
2	102		
3	105		
4	4		
5	5		
6	27		
7	41		
8	-7		
9	999		
10	1		
11			

Functions as arguments

Countif.hpp

```
10 template <typename Iter, typename Pred>
11 int Countif(Iter begin, Iter end, Pred pred) {
12     int count = 0;
13     auto cur = begin;
14
15     while(cur != end) {
16         if(pred(*cur))
17             ++count;
18         ++cur;
19     }
20
21     return count;
22 }
```

Lambda Functions in C++

main.cpp

```
1 bool isNegative(int num) { return (num < 0); }
2
3 class IsNegative {
4 public:
5     bool operator() (int num) { return (num < 0); }
6 };
7
8 int main() {
9     std::vector<int> numbers = {1, 102, 105, 4, 5, 27, 41, -7, 999};
10
11     auto isnegl = [](int num) { return (num < 0); };
12     auto isnegfp = isNegative;
13     auto isnegfunctor = IsNegative();
14
15     cout << "There are " << Countif(numbers.begin(), numbers.end(), _____)
16         << " negative numbers" << std::endl;
17
```

Lambda Functions in C++

[] () { }

Lambda Functions in C++

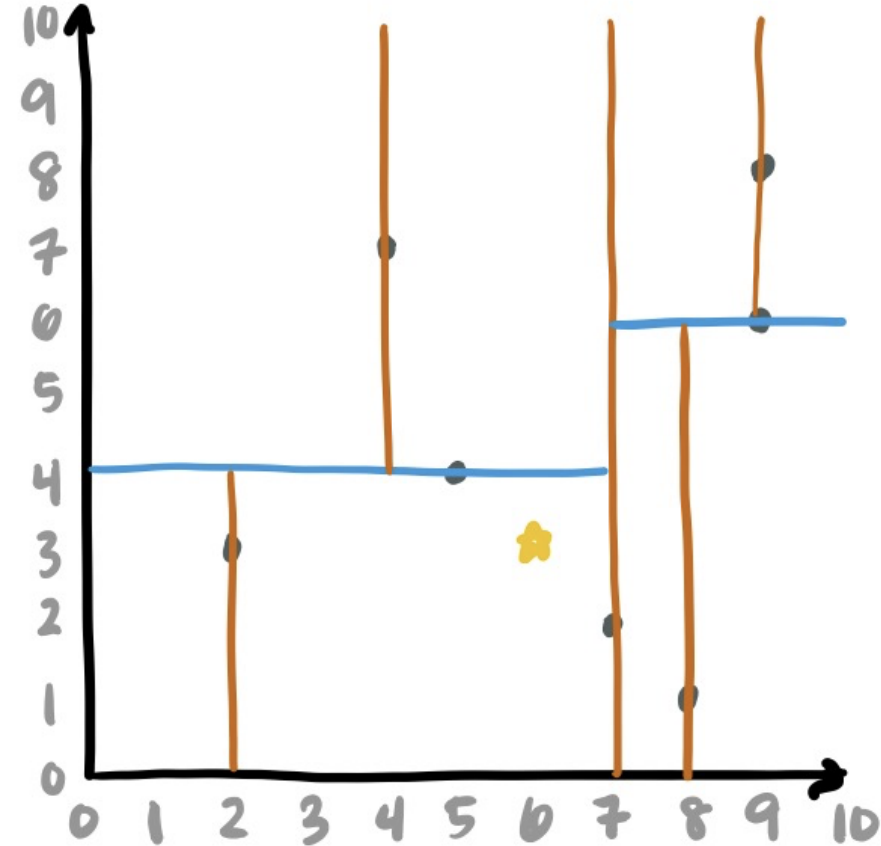
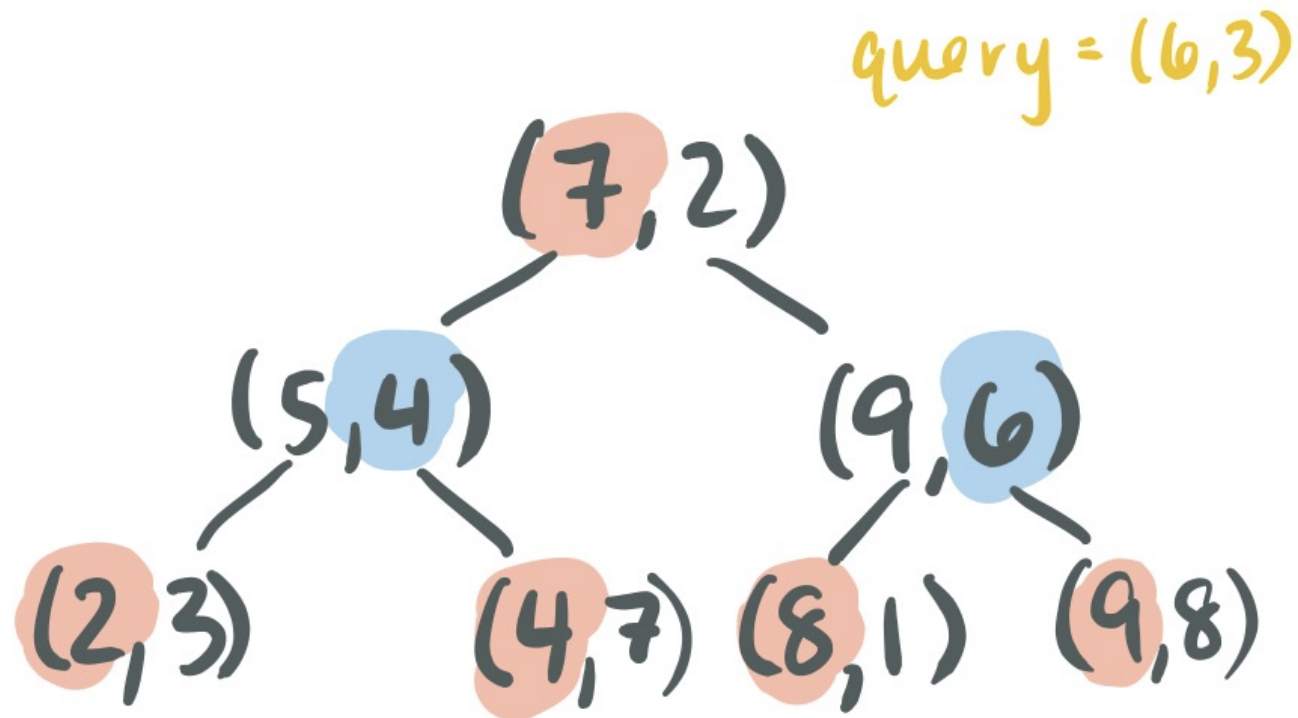


main.cpp

```
29     int big;
30
31
32     std::cout << "How big is big? ";
33     std::cin >> big;
34
35     auto isbig = [big](int num) { return (num >= big); };
36
37     std::cout << "There are " << Countif(numbers.begin(), numbers.end(), isbig)
38         << " big numbers" << std::endl;
}
```

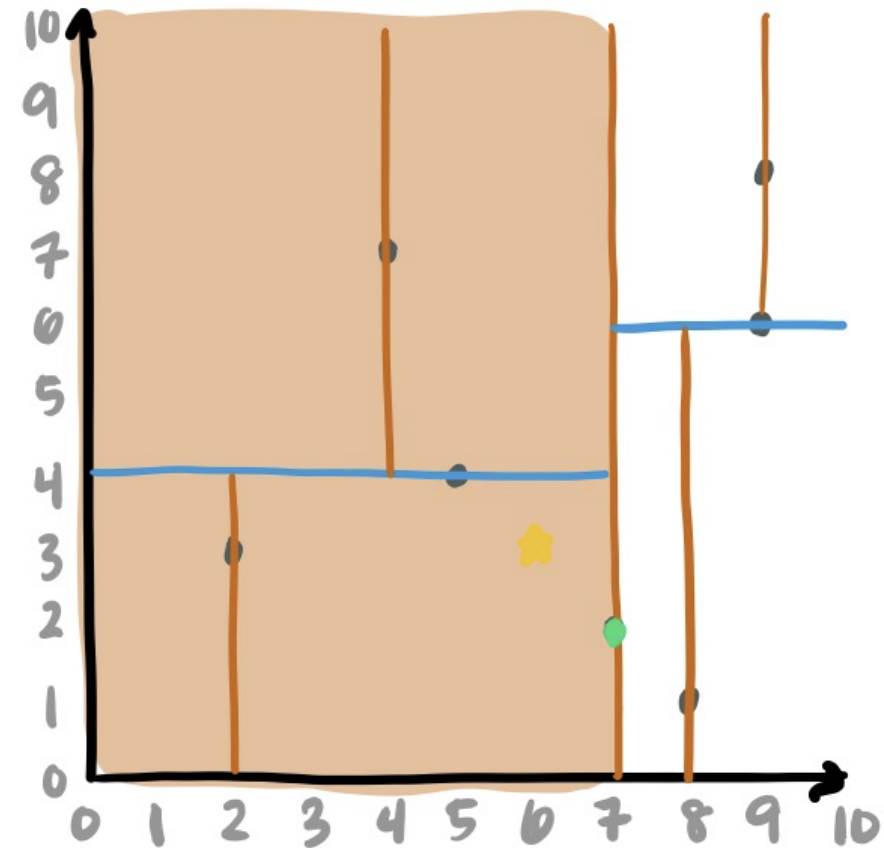
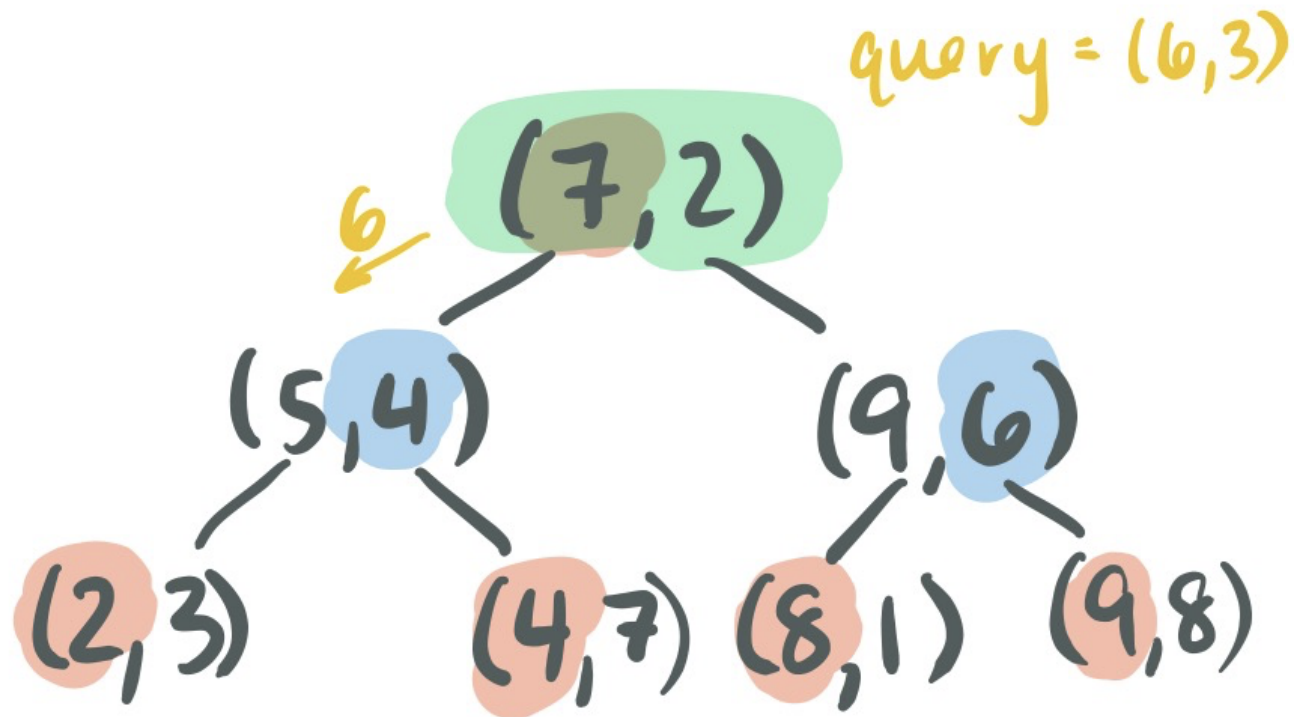
Nearest Neighbor: k-d tree

When querying a k-d tree, it acts like a BST* at first...



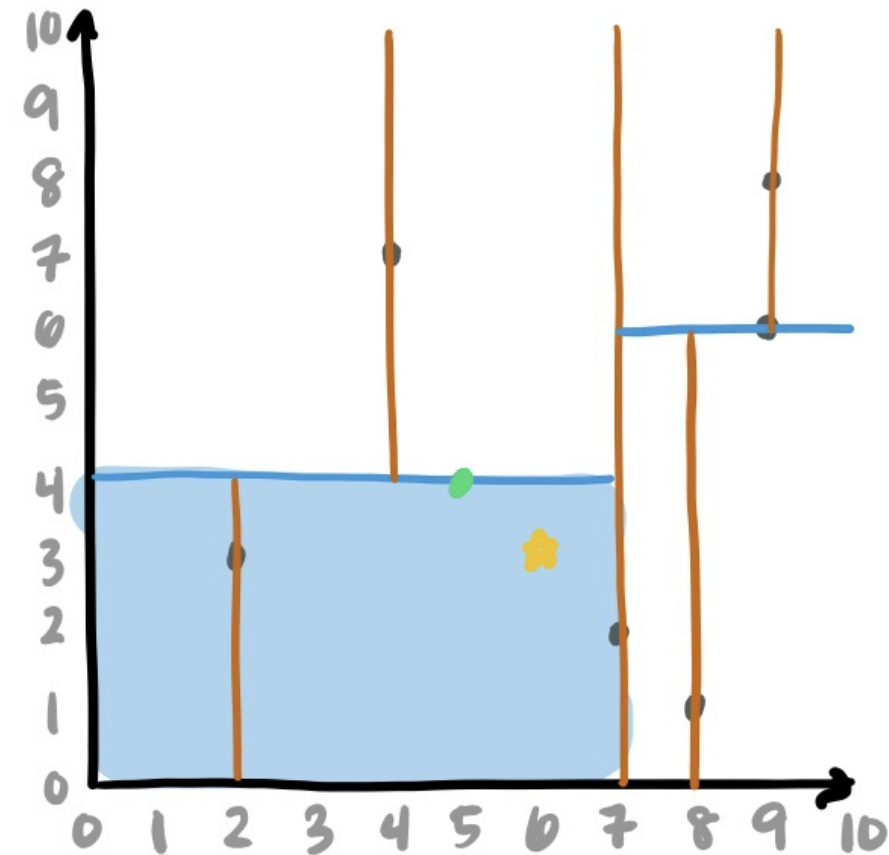
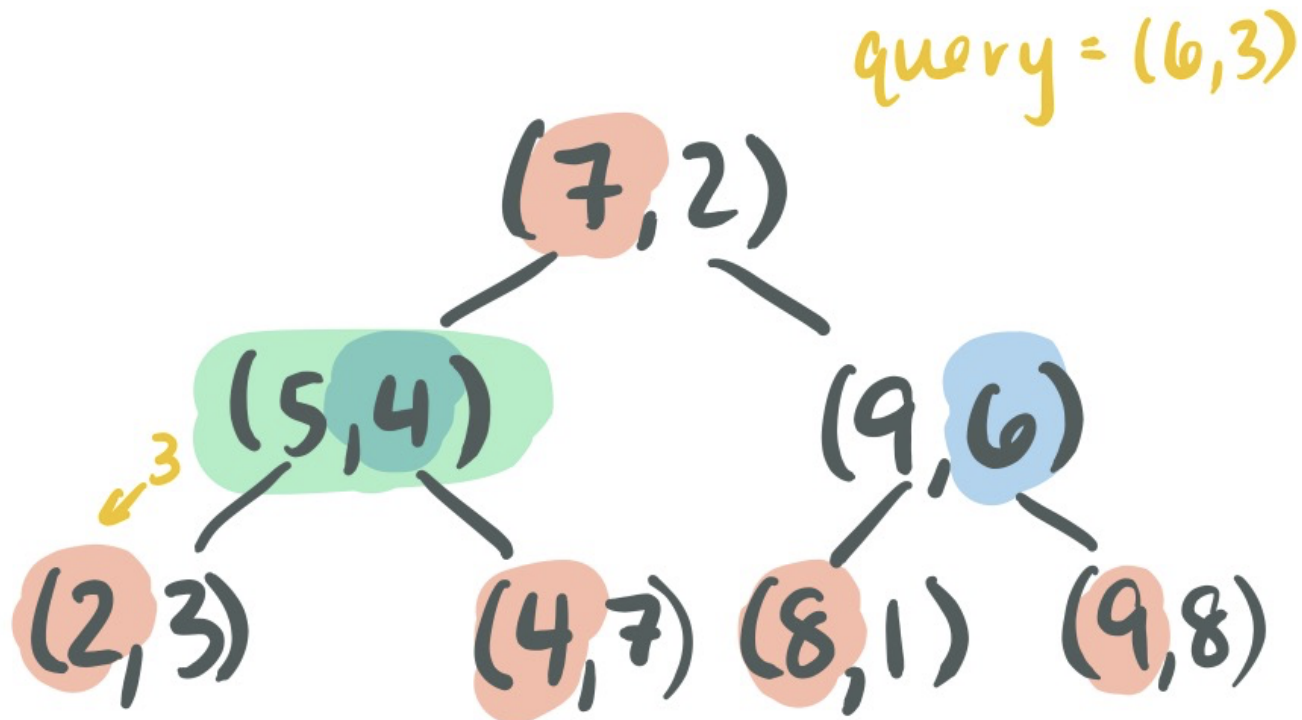
Nearest Neighbor: k-d tree

When querying a k-d tree, it acts like a BST* at first...



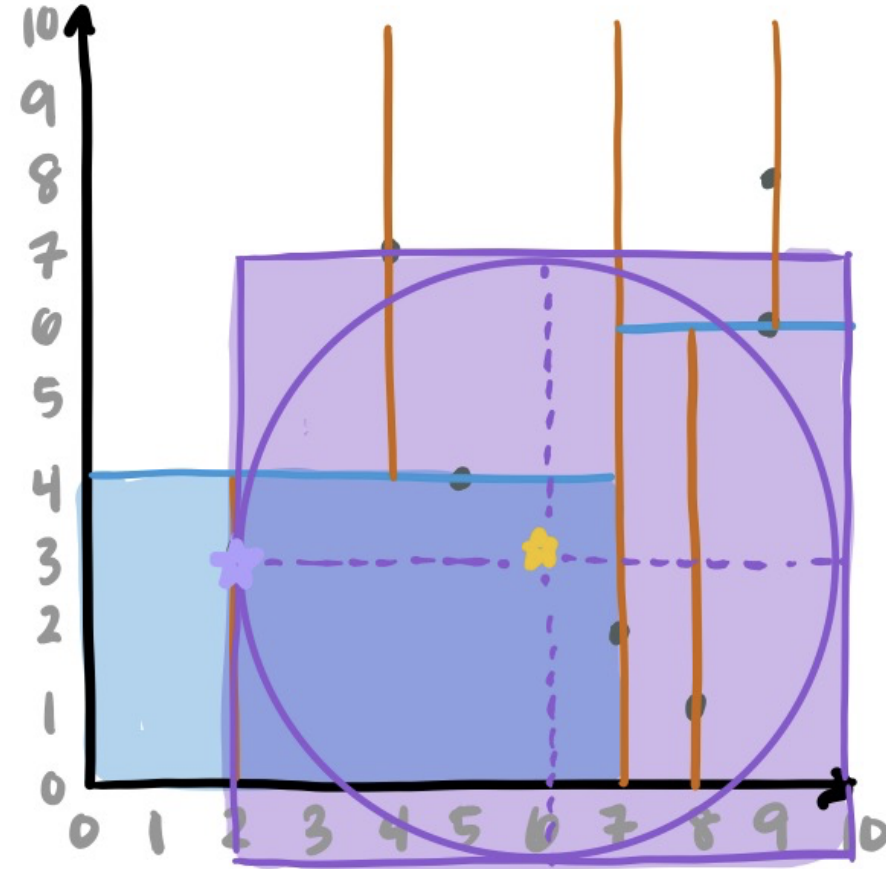
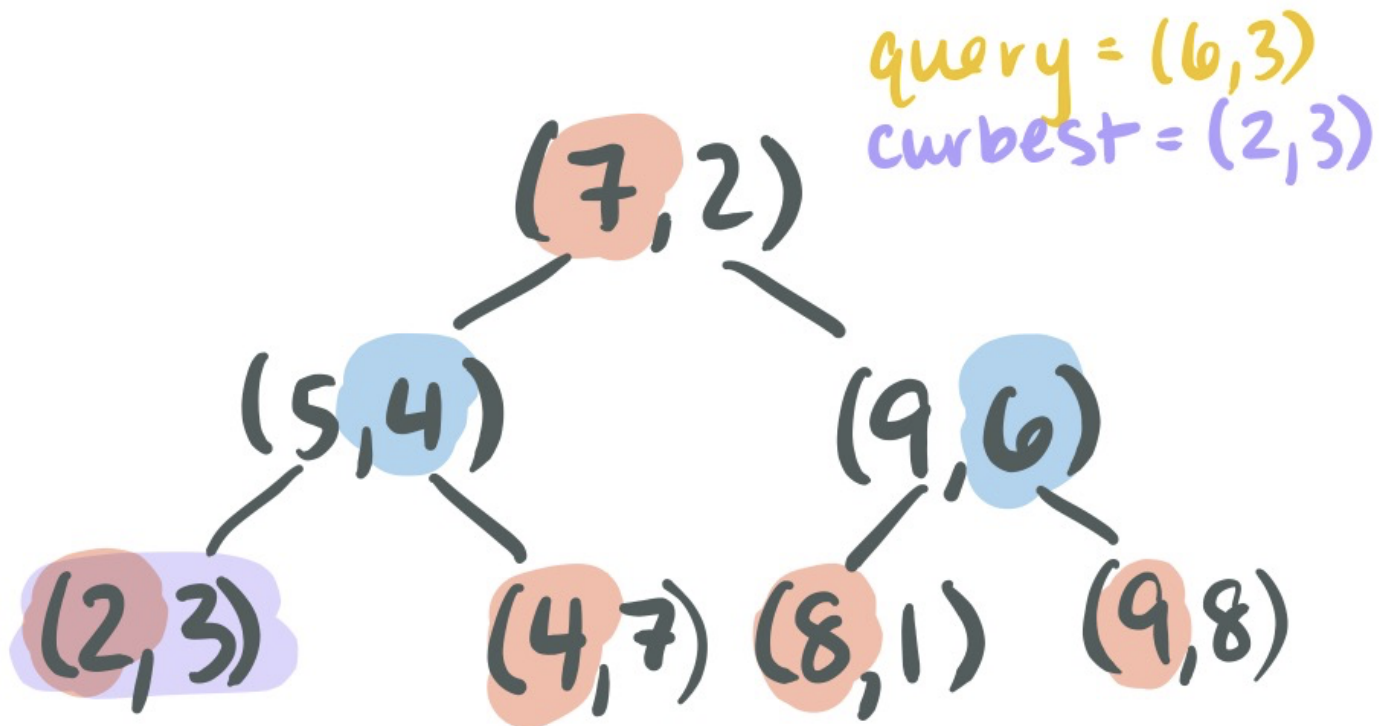
Nearest Neighbor: k-d tree

When querying a k-d tree, it acts like a BST* at first...



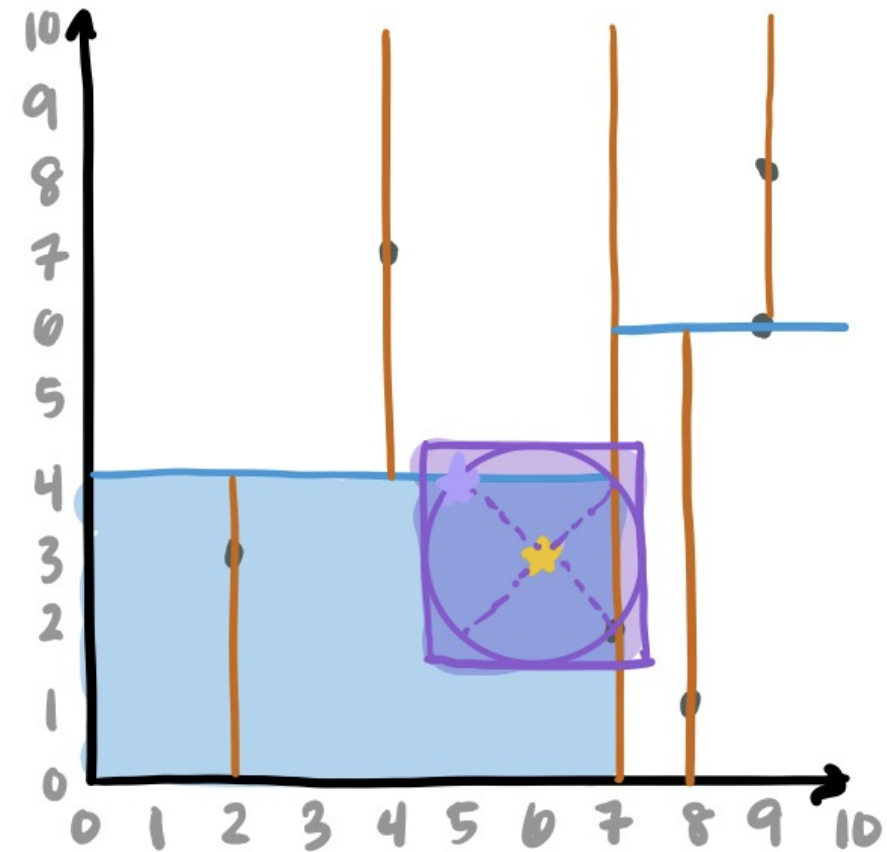
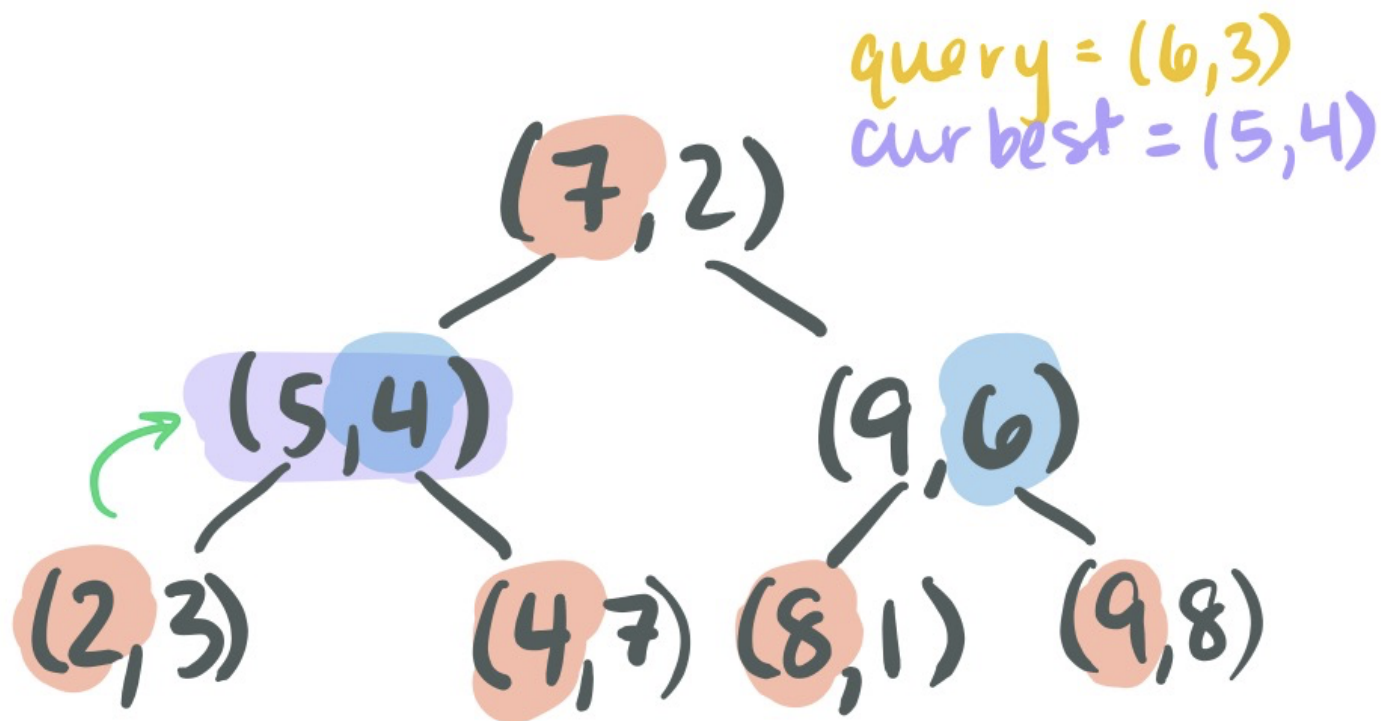
Nearest Neighbor: k-d tree

When querying a k-d tree, it acts like a BST* at first...

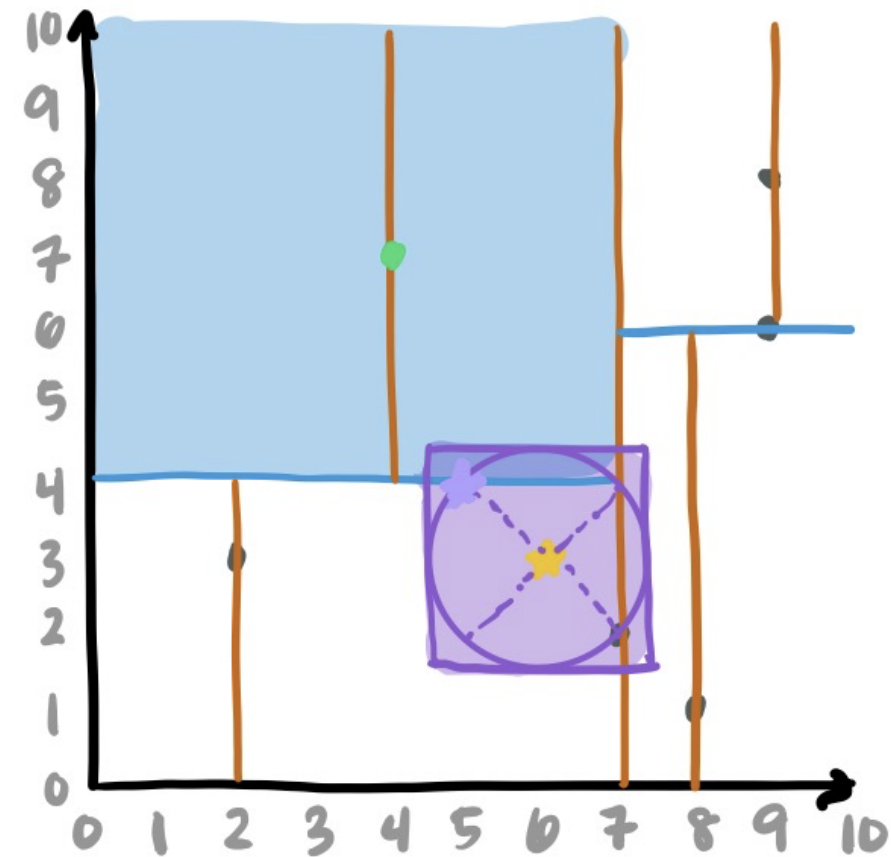
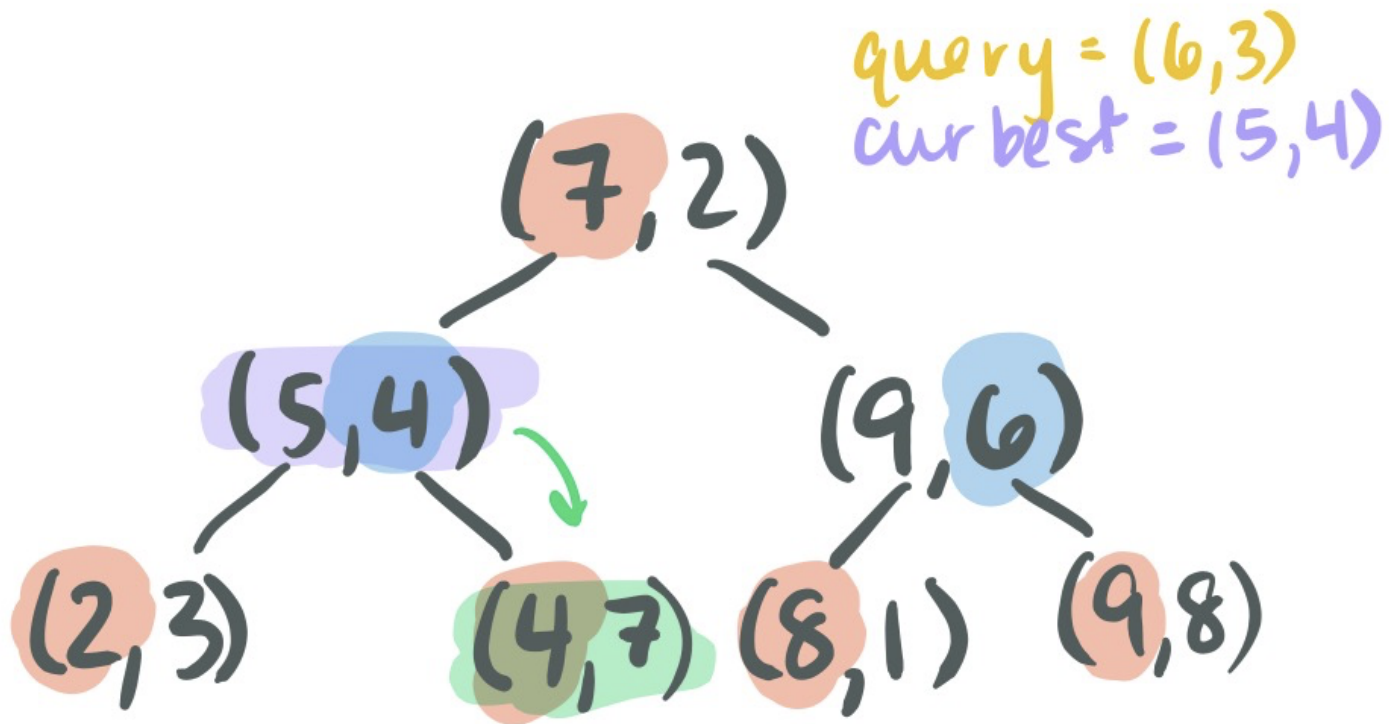


Nearest Neighbor: k-d tree

Backtracking: start recursing backwards -- store "best" possibility as you trace back

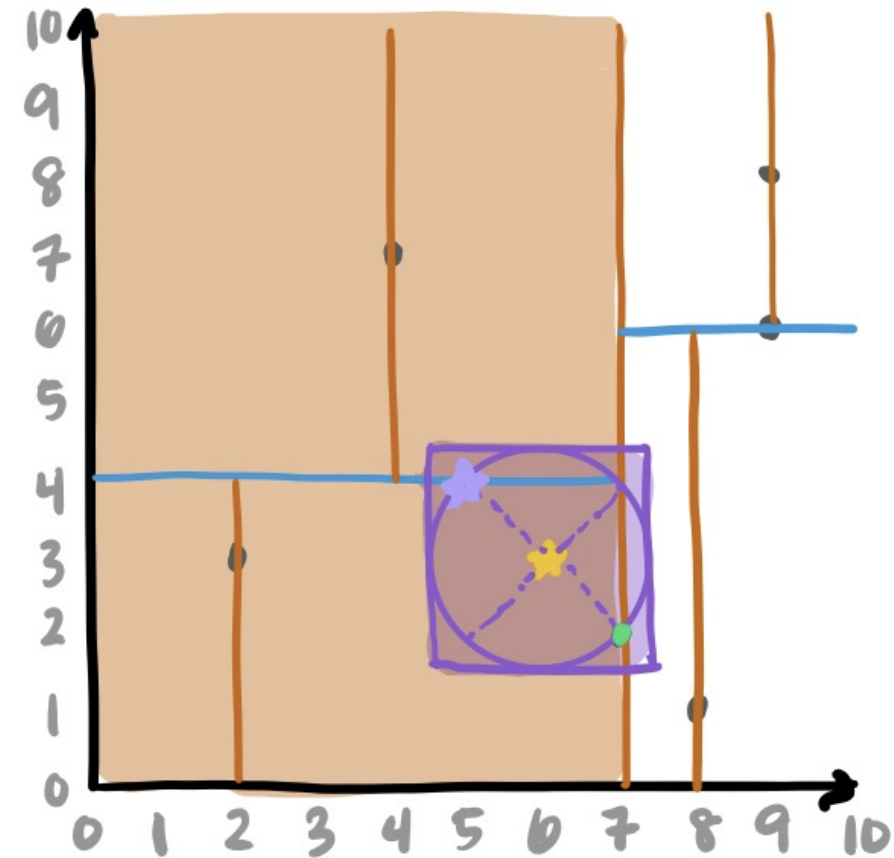
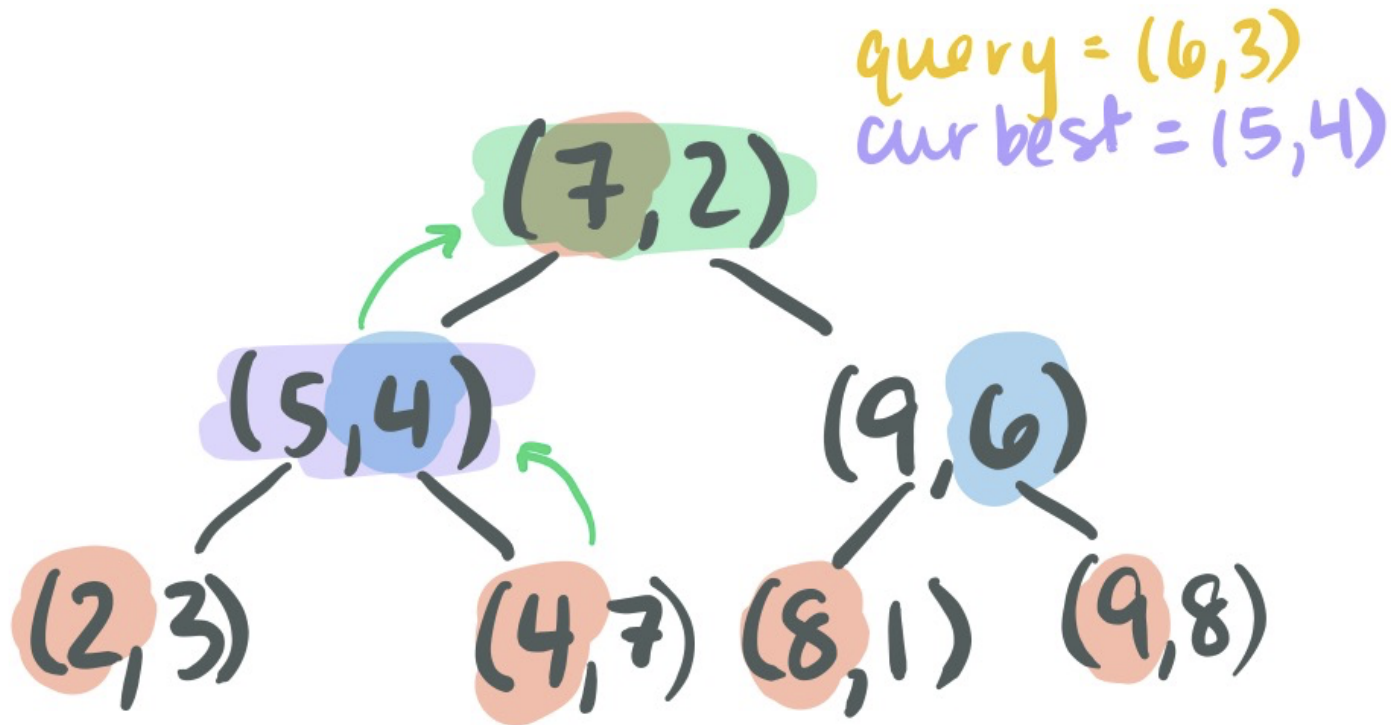


Nearest Neighbor: k-d tree

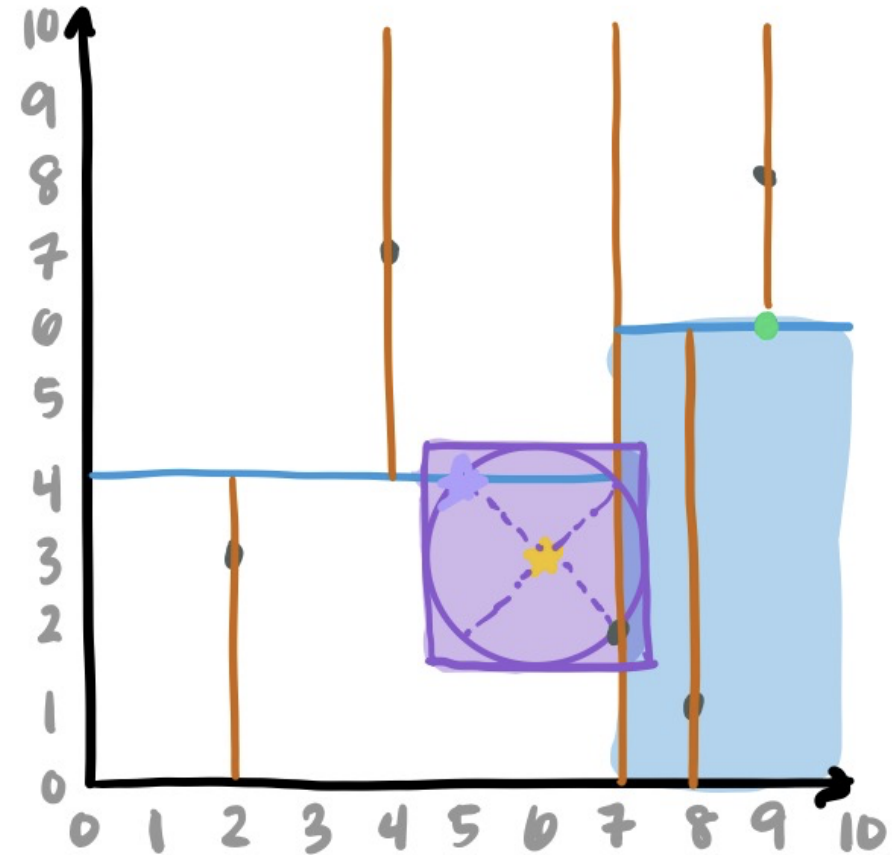
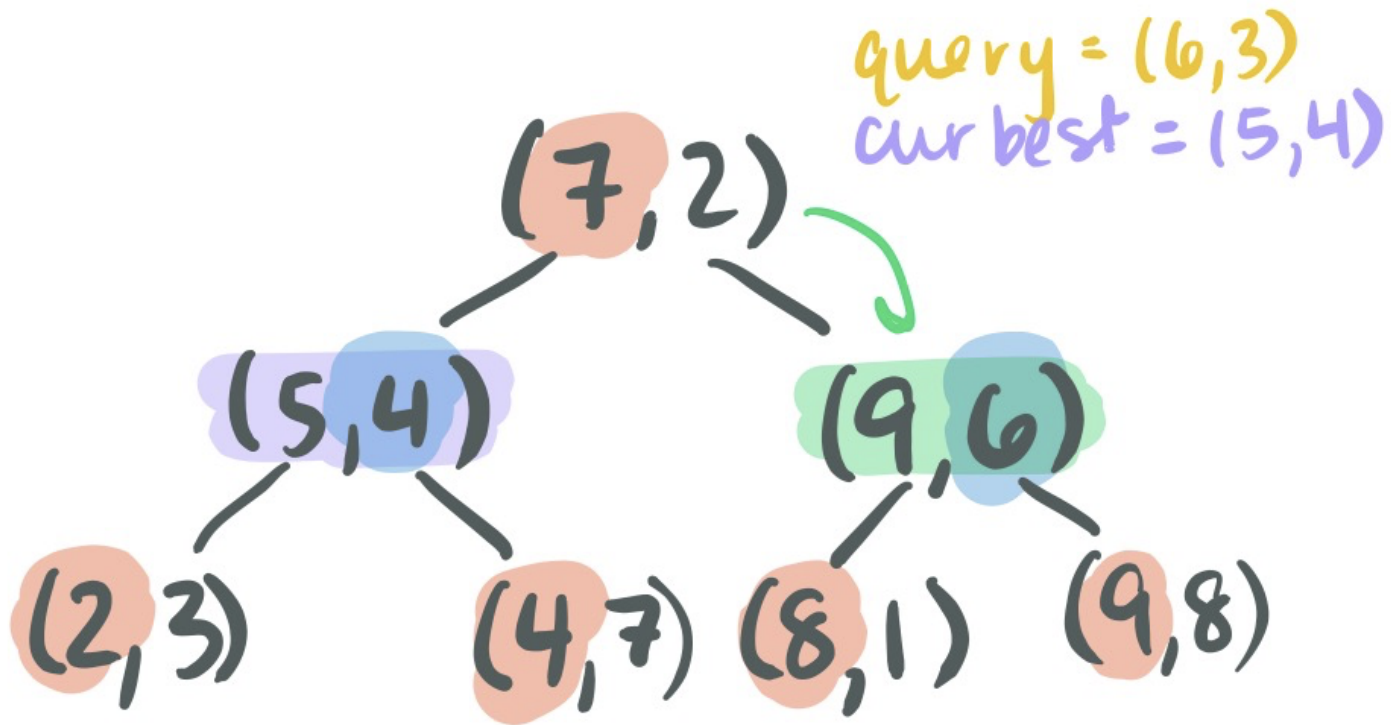


Nearest Neighbor: k-d tree

On ties, use `smallerDimVal` to determine which point remains `curBest`

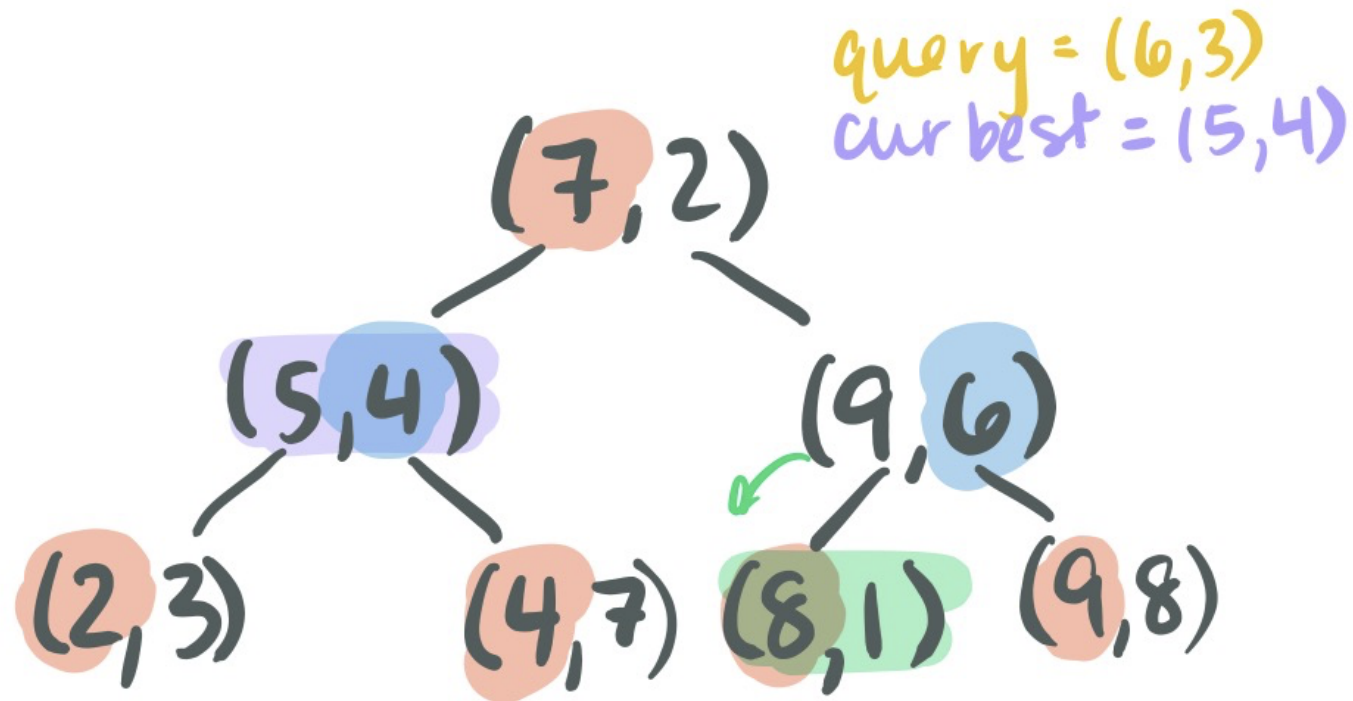


Nearest Neighbor: k-d tree

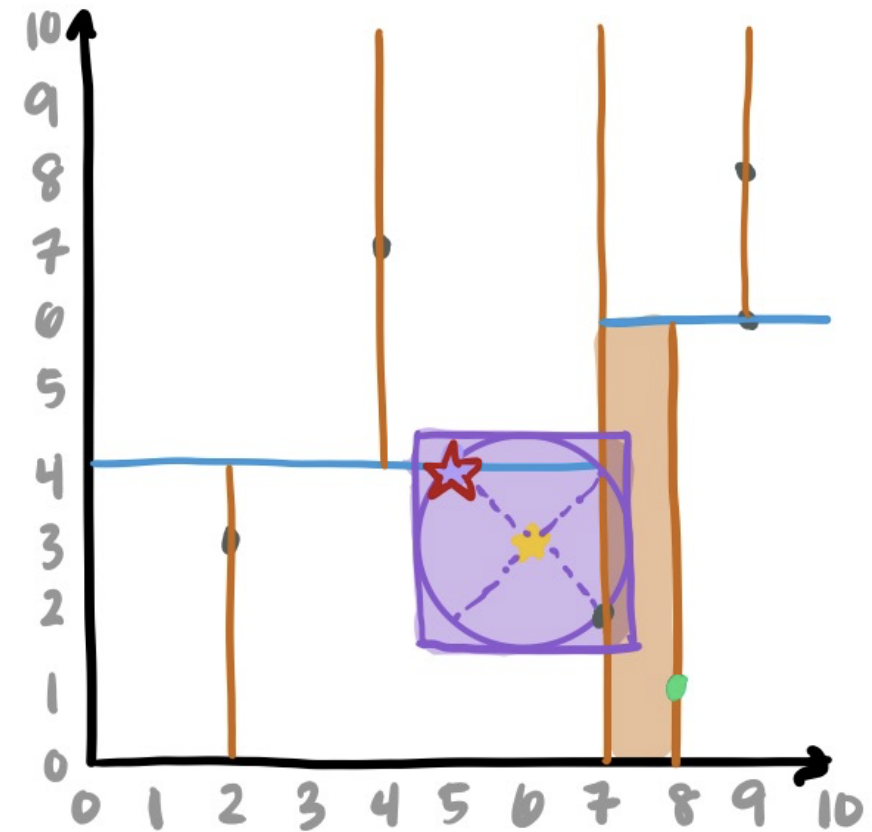




Nearest Neighbor: k-d tree



BEST: (5,4)



Nearest Neighbor: k-d tree

Final tips:

The mp_mosaic writeup is long. **READ IT**

The suggestions in the writeup should be followed carefully

Plan of Action

Since our goal is to find the lower bound on n given h , we can begin by defining a function given h which describes the smallest number of nodes in an AVL tree of height h :

$N(h)$ = minimum number of nodes in an AVL tree of height h

Simplify the Recurrence

$$N(h) = 1 + N(h - 1) + N(h - 2)$$

State a Theorem

Theorem: An AVL tree of height h has at least _____.

Proof by Induction:

I. Consider an AVL tree and let h denote its height.

II. Base Case: _____

An AVL tree of height _____ has at least _____ nodes.

Prove a Theorem

III. Base Case: _____

An AVL tree of height _____ has at least _____ nodes.

Prove a Theorem

IV. Induction Case: _____

Assume for all heights $i < h$, $N(i) \geq 2^{i/2}$. Prove that $N(h) \geq 2^{h/2}$

Prove a Theorem



V. Using a proof by induction, we have shown that:

...and inverting:

AVL Runtime Proof

An upper-bound on the height of an AVL tree is **$O(\lg(n))$** :

$N(h)$:= Minimum # of nodes in an AVL tree of height h

$$N(h) = 1 + N(h-1) + N(h-2)$$

$$> 1 + 2^{(h-1)/2} + 2^{(h-2)/2}$$

$$> 2 \times 2^{(h-2)/2} = 2^{(h-2)/2+1} = 2^{h/2}$$

Theorem #1:

Every AVL tree of height h has at least $2^{h/2}$ nodes.

Summary of Balanced BST

AVL Trees

- Max height: $1.44 * \lg(n)$
- Rotations:
 - Zero rotations on find
 - One rotation on insert
 - $O(h) == O(\lg(n))$ rotations on remove

Red-Black Trees

- Max height: $2 * \lg(n)$
- Constant number of rotations on insert (max 2), remove (max 3).

Summary of Balanced BST

Pros:

- Running Time:
 - Improvement Over:
- Great for specific applications:

Next Week: Considering hardware limitations

Can we always fit our data in main memory?

Where else can we keep our data?

Does this match our assumption that all memory lookups are $O(1)$?