

Master Theorem

Given a recurrence of the form

$$T(n) = aT\left(\frac{n}{b}\right) + f(n),$$

for constants $a \geq 1$ and $b > 1$ with f asymptotically positive, the following statements are true:

- Case 1. If $f(n) = O(n^{\log_b a - \epsilon})$ for some $\epsilon > 0$, then $T(n) = \Theta(n^{\log_b a})$.
- Case 2. If $f(n) = \Theta(n^{\log_b a})$, then $T(n) = \Theta(n^{\log_b a} \log n)$.
- Case 3. If $f(n) = \Omega(n^{\log_b a + \epsilon})$ for some $\epsilon > 0$ (and $af\left(\frac{n}{b}\right) \leq cf(n)$ for some $c < 1$ for all n sufficiently large), then $T(n) = \Theta(f(n))$.

$$\begin{aligned} 1. \quad T(n) &= 2T(n/2) + O(n \log n) = O(n \log^2 n). \\ 2. \quad T(n) &= T(n/2) + O(1) = O(\log n) \text{ binary search.} \\ 3. \quad T(n) &= T(n/2) + O(n) = O(n) \text{ quicksort.} \\ 4. \quad T(n) &= 2T(n/2) + O(n) = O(n \log n) \text{ sorts.} \\ T(n) &= T(\sqrt{n}) + O(\log n) = \\ n^{\frac{1}{2^k}} &= 1 \quad \log n \\ k &= \log n \\ \frac{1}{2^k} &= \log n \\ k &= \log \log n \\ k &= O(\log \log n) \\ \frac{1}{4} &= O(\log n) \\ \vdots & \\ n^{\frac{1}{2^k}} &\rightarrow 2^k \log n = \underbrace{1}_{2^{k+1}} \underbrace{\log n}_{O(\log n)} \\ T(n) &= T(n/2) + O(n^k) = O(n^k) \end{aligned}$$

Algorithms:

Searching:

1. Binary Search $O(\log n)$

Functionality:

- If element is in the array, return index of element
- If element is not in array, return -1

Precondition: (fact that is true when the function begins) \Rightarrow loop for it to work correctly.

- Array is of size n
- Array is sorted

int search(A, key, n)

```

begin = 0
end = n-1
while begin < end do:
    mid = (begin + end)/2;
    if key <= A[mid]:
        then end = mid
    else:
        begin = mid + 1

return A[begin] == key
? begin
: -1

```

Postcondition: (fact that is true when the function ends) \Rightarrow show that the computation was done correctly.

$A[\begin{begin}] == key$

Invariant:

- relationship between variables that is always true

Loop Invariant

- relationship between variables that is true at the beginning (or end) of each iteration of a loop.

- $A[\begin{begin}] \leq key \leq A[\end]$. (Correctness)

(the key is in the range of the array).

- $(\end - \begin) \leq n/2^k$ in iteration k (Performance)

Not just for searching arrays

L eg. A complicated function.

The function that is always increasing.

L min value that satisfies the complicated function.

Logarithmic Time	$O(\log n) = O(\log \log n) <$
Sublinear Time	$O(n^{0.0000000}) <$
Linear Time	$O(20n) <$
Linearithmic Time	$O(n \log n) = O(\log n!) <$
Quadratic Time	$O(4n^2) <$
Polynomial Time	$O(n^{2.5}) <$
Exponential Time	$O(2^n) = O(2^{n+1}) <$
Factorial Time	$O((n-1)!) <$
!!!	$O(n^n)$

2. Peak (Local Max) Finding $O(\log n)$

Input: Some array $A[0 \dots n-1]$

Output: A local maximum in A $A[i-1] \leq A[i] \leq A[i+1] \leq A[i+2]$

Assume: $A[i-1] = A[n] = -\text{MAX_INT}$.

FindPeak(A, n)

```

if A[n/2+1] > A[n/2]
    then FindPeak(A[n/2+1 ... n], n/2) right half
else if A[n/2+1] > A[n/2]
    then FindPeak(A[1 ... n/2-1], n/2) left half.
else
    A[n/2] is a peak
    return n/2.

```

Invariant (key-property)

- If we recurse in the right half, then there exist a peak in the right half \Rightarrow The array must have a peak.
- Prove by Induction (Invariant)

Correctness

- There exists a peak in the range $[\begin{begin}, \end]$
- Every peak in $[\begin{begin}, \end]$ is a peak in $[1, n]$.

Sorting: Total Space = $O(n)$

Input: array $A[1 \dots n]$ of words/numbers (unsorted)
Output: array $B[1 \dots n]$ that is a permutation of A
 such that:
 $B[1] \leq B[2] \leq \dots \leq B[n]$ (sorted)

1. **BozoSort** ($A[1 \dots n]$) $O(n \cdot n!)$ in-place

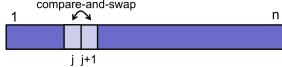
Repeat: \rightarrow not stable!

1. Choose a random permutation of the array A
2. If A is sorted, return A .

2. **BubbleSort** (A, n) $O(n) \rightarrow O(n^2)$ in-place

Repeat n times until no swaps:

- for $j = 1$ to $n-1$
 if $A[j] > A[j+1]$ \rightarrow only swap elems if they are different
 then swap ($A[j], A[j+1]$) \therefore stable



Best case:

(already sorted) $O(n)$

Avg case:

(random inputs) $O(n^2)$

Worst Case:

(worst input, n iterations) $O(n^2)$.

Loop Invariant:

At the end of iteration j ,
 the worst j items are correctly sorted
 to the first j position of the array.

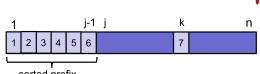
Correctness:

After n iterations \rightarrow sorted.

3. **SelectionSort** (A, n) $O(n^2)$ only. in-place

for $j = 1$ to $n-1$:

find min element $A[j]$ in $A[j \dots n]$ (Time: $(n-j)$)
 Swap ($A[j], A[k]$). \rightarrow not stable!



Loop invariant:

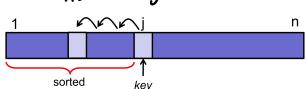
At the end of iteration j :
 the smallest j items are correctly sorted
 in the first j position of the array.

4. **Insertion Sort** (A, n) $O(n) \rightarrow O(n^2)$ in-place -

for $j = 2$ to n

Key: $A[j]$

[Insert key into sorted array $A[1 \dots j-1]$]
 \downarrow
 $i = j-1$ \quad Stable, do not swap identical elements.
 while $i > 0$ and $A[i] > \text{key}$ $\quad P$
 $A[i+1] = A[i]$
 $i = i-1$
 $A[i+1] = \text{key}$



Loop invariant:

At the end of iteration j :
 the first j items in the array
 are in sorted order.

Best Case:
 (already sorted) $O(n)$

Avg Case:
 (random permutation) $\sum_{j=1}^{n-1} O(j/2) = O(n^2)$

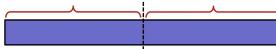
Worst Case:
 (inverse sorted) $O(n^2)$
 \downarrow on avg, a key in pos j needs to move back $j/2$ slots (the expectation).

Insertion vs Bubble \Rightarrow almost sorted array.

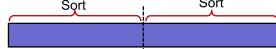
Note: Stability: preserves order of equal elements.

5. **MergeSort** (Divide-and-Conquer) $O(n \log n)$

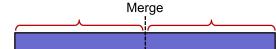
1. Divide: split array into 2 halves



2. Recurse: sort the 2 halves



3. Combine: merge the 2 sorted halves.



Space: $O(n \log n)$

iterative ver.: from bottom up.

$O(n)$ space ver.: use only 1 temp. array with merge.

$$S(n) = 2S(n/2) + O(1) = O(n)$$

MergeSort(A , begin, end, tempArray)

if $(\text{begin}=\text{end})$ then return;

else:

$$\text{mid} = \text{begin} + (\text{end}-\text{begin})/2$$

MergeSort(A , begin, mid, tempArray);

MergeSort(A , mid+1, end, tempArray);

Copy(tempArray, A , begin, end);

Not in-place!

MergeSort(A, n)

if $n=1$ then return; Base case: $O(1)$

else:

$X = \text{MergeSort}(A[1 \dots n/2], n/2)$; Recursive / $T(n/2)$

$Y = \text{MergeSort}(A[n/2+1 \dots n], n/2)$; Conquer $T(n/2)$

return Merge(X, Y, n/2); combine solutions $O(n)$

In each iter, move 1 elem to final list.
 After n iters, all n elems are in final list.
 Each iter takes $O(1)$ time to compare 2 elems and copy.

(mostly sorted) \Rightarrow slower than insertion, $O(n \log n)$

(small number of items) \Rightarrow slow, use insertion sort for $n < 1024 \Rightarrow$ due to caching / branch prediction.
 \Rightarrow base case = use a slower sort

Space usage:

- need extra space for merge

- merge copies data to new array:

↳ temporary array of size n .

Stability: Stable.

- MergeSort is stable if merge is stable.

- merge is stable if carefully implemented.

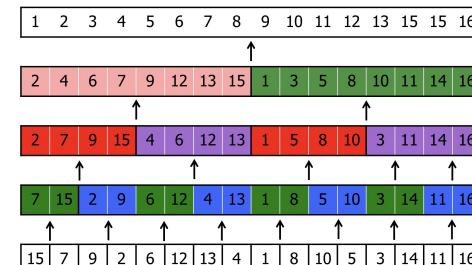
Merge:



1 2 7 9

Need temp array of size n .

Iterative version (doesn't require extra space)



Name	Best Case	Average Case	Worst Case	Extra Memory	Stable?
Bubble Sort	$O(n)$	$O(n^2)$	$O(n^2)$	$O(1)$	Yes
Selection Sort	$O(n^2)$	$O(n^2)$	$O(n^2)$	$O(1)$	No
Insertion Sort	$O(n)$	$O(n^2)$	$O(n^2)$	$O(1)$	Yes
Merge Sort	$O(n \log n)$	$O(n \log n)$	$O(n \log n)$	$O(n \log n)$	Yes
Quick Sort	$O(n \log n)$	$O(n \log n)$	$O(n^2)$	$O(n \log n)$	No

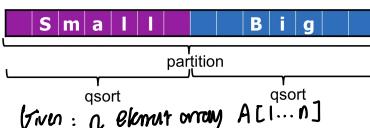
$O(n \log n)$
 in-place
 $\approx O(1)$.

6. Quicksort ($A[1..n], n$) $O(n \log n)$

```

if n == 1 then return;
else:
    choose pivot index pIdx
    p = partition(A[1..n], pIdx) departs if partition's in-place
    x = Quicksort(A[1..p-1], p-1)
    y = Quicksort(A[p+1..n], n-p)

```



- Divide: Partition the array into 2 sub-arrays around a pivot x
such that elements in lower subarray $\leq x \leq$ elements in upper sub-array.



- Conquer: Recursively sort the 2 sub-arrays

- Combine: Trivial, do nothing.

Partition: $O(n)$

- Choose a pivot
- Find all elements smaller than the pivot
- Find all elements larger than the pivot

(Not in-place)

partition ($A[2..n], n, pivot$)

//assume no duplicates

```

B = new n element array
low = 1;
high = n;
for(i=2; i <= n; i++)
    if (A[i] < pivot)
        then B[low] = A[i];
        low++;
    else if (A[i] > pivot)
        then B[high] = A[i];
        high--;

```

$B[low] = pivot$

return <B, low>

Array B is partitioned around the pivot

Proof:

Invariants:

1. For every $i < low$: $B[i] < pivot$

2. For every $j > high$: $B[j] > pivot$.

In the end, every element from A is copied to B.

Thus: $B[i] = pivot$

By Invariants, B is partitioned around the pivot.

How to partition efficiently?

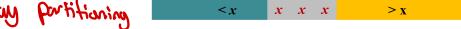
- Duplicates
- If you ignore duplicates, partitioning can be very slow!
- In-place partitioning
- One pass, use invariants to ensure correctness.
- Stability
- In-place partitioning isn't stable.
- Choosing a pivot (randomization)
- Deterministic pivots are generally bad.
- Randomization is an easy way to find a good pivot!

How to analyze?

- Divide-and-conquer recurrence
- It is sufficient to do a 90:10 split to get $O(n \log n)$ performance! *no need half-half*
- What is the probability that a random pivot yields a 90:10 split? Quite good! *high*
- Simplification: Paranoid QuickSort
- How many repetitions to find a good pivot, in expectation? $O(1)$
- Solve using recurrence: Linearity of expectation FTW.

Duplicator?

3-way partitioning

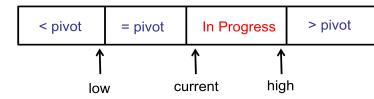


- 2 pass partitioning

- Regular partition
- pass duplicates.

- 1 pass partitioning

- more complicated
- maintain 4 regions of the array.



Choice of pivot?

Worst case doesn't matter \Rightarrow equality bias \Rightarrow Sorting the array takes n execution of partition.

- First element $A[1]$? eg. reverse sorted $\Rightarrow O(n^2)$
- Last element $A[n]$? eg. sorted $\Rightarrow O(n^2)$.
- Middle element $A[n/2]$?
- Median of $(A[1], A[n/2], A[n])$?

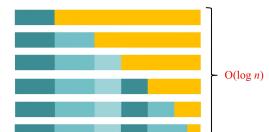
If array is split evenly

$$T(n) = T(n/2) + T(n/2) + cn \\ = 2T(n/2) + cn \\ = O(n \log n)$$

If array is split 1/10: $cn/10$

$$T(n) = T(n/10) + T(9n/10) + cn \\ \approx O(n \log n)$$

Assume larger part shrinks by at least 9/10 every iteration:



Choose pivot at random:

\Rightarrow running time is α , random variable.
 \Rightarrow for any input, there's a good probability of success.

Set of outcomes for X = $(e_1, e_2, e_3, \dots, e_n)$:

- $Pr(e_1 = p_1)$
- $Pr(e_2 = p_2)$
- ...
- $Pr(e_n = p_n)$

Linearity of expectation

$$E[A+B] = E[A] + E[B].$$

Expected outcome:
 $E[X] = e_1p_1 + e_2p_2 + \dots + e_np_n$

e.g. Flipping an (unfair) coin:
 $- Pr(\text{heads}) = p$
 $- Pr(\text{tails}) = (1-p)$

How many flips to get at least one head?

$$E[X] = (p)(1) + (1-p)(1 + E[X])$$

How many more flips to get a head?

Idea: If I flip "tails," the expected number of additional flips to get a "heads" is still $E[X]!!$

Forwards analysis:

Then we know:
 $E[T(e)] = E[T(e)] + E[T(n-e)] + E[\# \text{pivot chosen}(e)]$

$$\Leftrightarrow E[T(e)] = E[T(n-e)] + 2e$$

$$\approx O(n \log n)$$

Optimization:

Base case?

- Recurse all the way to single-element arrays.
- Switch to InsertionSort for small arrays.
- Halt recursion early, leaving small arrays unsorted. Then perform InsertionSort on entire array.

Relies on fact that
InsertionSort is very fast on
almost sorted arrays!

Double pivot.

Order Statistics $O(n)$

- Find the k^{th} smallest element in an unsorted array.

- L 1. Sort the array
- 2. return element number k
 $\Rightarrow O(n \log n)$

- L Only do min number of sorting
 - 1. Partition the array.
 - 2. Continue searching in the correct half.

$\text{QuickSelect}(A[1..n], n, k)$

```
if (n==1) then return A[1];
else Choose random pivot index pIdx;
    P = partition(A[1..n], n, pIdx)
    if (k==p) then return A[p];
    else if (k < p) then
        return Select(A[1..p-1], k) → only reuse once! → not sorting
    else if (k > p) then
        return Select(A[p+1..n], k-p) → partitioning right and left
                                         are not exactly the same.
```

Recurrence: $T(n) = T(\frac{n}{2}) + O(n) \approx O(n)$

Trees

Binary Tree : either (a) empty

(b) a node pointing to 2 binary trees

```
public class TreeNode {
    private TreeNode leftTree;
    private TreeNode rightTree;
    private longType key;
    private valueType values;
    ...
}
```

property:

all in left sub-tree $<$ key $<$ all in right sub-tree.

height

- Number of edges on longest path from root to leaf.
- $h(v)=0$
- $h(v)=\max(h(v.\text{left}), h(v.\text{right}))+1$.
- $h(\text{null})=-1$.

```
- public int height() {
    int leftHeight = -1;
    int rightHeight = -1;
    if (leftTree != null)
        leftHeight = leftTree.height();
    if (rightTree != null)
        rightHeight = rightTree.height();
    return max(leftHeight, rightHeight) + 1;
    ...
}
```

max of subtrees.

Search (max/min) $O(h)$

```
public TreeNode searchMax() {
    if (rightTree == null)
        return rightTree.searchMax();
    ...
}
```

else return this; \Rightarrow key is here!

Worst Case: $O(n) \Rightarrow$ e.g. linear tree

Search (any key) Avg Case: $O(h) \Rightarrow O(\log n)$.

```
public TreeNode search(int queryKey) {
    if (queryKey < this.key) {
        if (leftTree != null)
            return leftTree.search(queryKey);
        else return null;
    }
    else if (queryKey > this.key) {
        if (rightTree != null)
            return rightTree.search(queryKey);
        else return null;
    }
    else return this;
}
```

insert $O(h)$

```
public void insert(int insKey, int insValue) {
    if (insKey < this.key) {
        if (leftTree != null)
            leftTree.insert(insKey, insValue);
        else leftTree = new TreeNode(insKey, insValue);
    }
    else if (insKey > key) {
        if (rightTree != null)
            rightTree.insert(insKey, insValue);
        else rightTree = new TreeNode(insKey, insValue);
    }
    else return;  $\Rightarrow$  key already in tree.
}
```

performance of tree algorithms depends on shape.

- \hookrightarrow order of insertion.
- \hookrightarrow each order does not yield a unique shape.
- \hookrightarrow insert keys in a random order \Rightarrow balanced.

\hookrightarrow can reconstruct a BST using only-order-traversal.

in-Order-traversal (V) $O(n) \Rightarrow$ visit each node at most once

```
- left-subtree
  self
  right-subtree
- public void in-order-traversal() {
    pre-order → // Traverse left subtree
    if (leftTree != null)
        leftTree.in-order-traversal();
    in-order → visit(this);
    // Traverse right subtree
    if (rightTree != null)
        rightTree.in-order-traversal();
    post-order →
}
```

successor $O(h)$

- If you search for a key not in tree
 - \hookrightarrow either find predecessor or successor.

- Successor(key)

1. Search for key in the tree
2. If (result > key), then return result.
3. If (result == key), then search for successor of result.

```
public TreeNode successor() {
    if (rightTree != null)
        return rightTree.searchMin();

    TreeNode parent = parentTree;
    TreeNode child = this;
    while ((parent != null) && (child == parent.rightTree))
        child = parent;
        parent = child.parentTree;
    }
    return parent;
}
```

delete $O(h)$

- 3 cases
 - No children
 - remove v
 - 1 child
 - remove v
 - connect child(v) to parent(v)
 - 2 children.
 - x = successor(v) \parallel swap x w/ v.
 - delete(x)
 - remove v
 - connect x to left(v), right(v), parent(v).
- Successor of deleted node has at most 1 child.

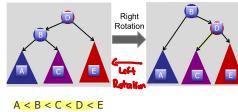
Balanced Trees

- Operations take $O(h)$ time
- $\log(n)-1 \leq h \leq n$

* A BST is balanced if $h = O(\log n)$

- On a balanced BST: all operations run in $O(\log n)$ time
- How to get a balanced tree?
 - define a good property of a tree
 - Show that if the good property holds, then the tree is balanced.
 - After every insert/delete, make sure the good property still holds. If not, fix it.

Tree Rotations \Rightarrow can create every possible tree shape!



$A < B < C < D < E$

Rotations maintain ordering of keys.

\Rightarrow Maintains BST property.

Right Rotation: The root of the subtree moves right. (Craq, left child)

Left Rotation: The root of the subtree moves left. (Craq, right child)

right-rotate(v)

\Rightarrow Assume v has left := null

v.left = w.left

w.parent = v

v.parent = w

w.left = v.right

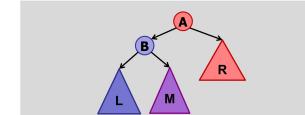
w.right = v



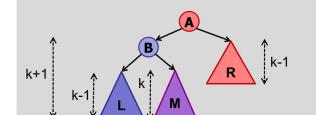
A is Left-heavy if left sub-tree has larger height than right sub-tree.

A is Right-heavy if right sub-tree has larger height than left sub-tree.

Note: Rotation reduces root height by 1, and everything higher in tree remains unchanged.



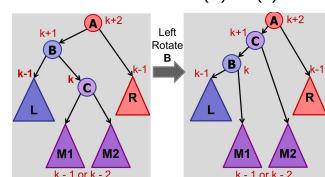
Assume A is the lowest node in the tree violating balance property.



Assume A is the lowest node in the tree violating balance property.

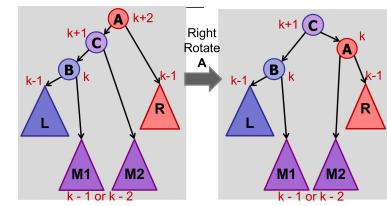
Case 3: B is right-heavy : $h(L) = h(M) - 1$

$$h(R) = h(L)$$



Left-rotate B

After left-rotate B: A and C still out of balance.



Summary:

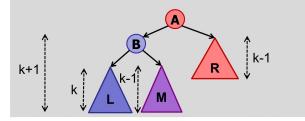
If v is out of balance and left heavy:

1. v.left is balanced: right-rotate(v)

2. v.left is left-heavy: right-rotate(v)

3. v.left is right-heavy: left-rotate(v.left)

right-rotate(v)



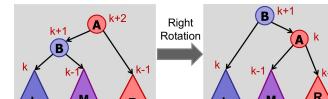
Assume A is the lowest node in the tree violating balance property.

Case 1: B is equi-height : $h(L) = h(M)$

$$h(R) = h(L) - 1$$

Case 2: B is left-heavy : $h(L) = h(M) + 1$

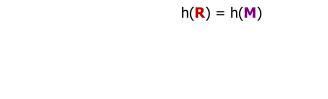
$$h(R) = h(M)$$



right-rotate:

Case 2: B is left-heavy: $h(L) = h(M) + 1$

$$h(R) = h(M)$$



right-rotate:

Case 2: B is left-heavy: $h(L) = h(M) + 1$

$$h(R) = h(M)$$

Trees are a good way to store, summarize and search dynamic data.

Operations that create a data structure
- build (process)

Operations that modify the structure:

- insert

- delete

- Query Operations

- search, select etc.

Notes:

- Every height-balanced tree is balanced.
- Not every balanced tree is height-balanced.

AVL Trees

Step 0: Augment

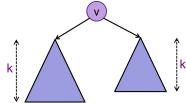
- In every node v, store heights: can be difference in height from parent etc.
- $v.height = h(v)$

On insert & delete update height:

```
insert(v)
  if (v.left == null)
    v.left.height = 0
  else
    v.left.height = v.left.height + 1
  height = max(left.height, right.height) + 1
```

Step 1: Define Balance Condition / Invariant*

- A node is height-balanced if: $|v.left.height - v.right.height| \leq 1$



A binary search tree is height-balanced if every node in the tree is height-balanced.

A height-balanced tree with a node has at most height $h < 2\log(n)$.

$$\begin{aligned} \leftrightarrow \frac{n}{2} &< \log(n) \\ \leftrightarrow 2^{\lfloor \log(n) \rfloor} &< n \\ \leftrightarrow 2^{\lfloor \log(n) \rfloor} &< n \end{aligned}$$

A height-balanced tree with height h has at least $n > 2^{\lfloor \log(n) \rfloor}$ nodes.

Proof:

Let n_h be the min num of nodes in a height-balanced tree of height h .

Show $n_h > 2^{\lfloor \log(n) \rfloor}$

$$n_h > n_{h-1}$$

$$n_h \geq 1 + n_{h-1} + n_h$$

note:
no. of nodes
ratio of 2 to 1
at node $\frac{1}{2}n$

$$\therefore n_h \geq 2 \cdot n_{h-2} \geq 2^2 n_{h-4}$$

$$\geq 4 \cdot n_{h-4} \geq 2^3 n_{h-6}$$

$$\geq 8 \cdot n_{h-6} \geq 2^5 n_{h-8}$$

$$\geq \dots$$

$$\geq 2^h n_0$$

= base case. Note:

$$\therefore n_h \geq 2^h n_0$$

$$\geq 2^{h-1} n_0$$

$$\geq 2^{h-2} n_0$$

$$\geq \dots$$

$$\geq 2^0 n_0$$

= base case. Note:

$$\therefore n_h \geq 2^{\lfloor \log(n) \rfloor}$$

$$\geq 2^{\lfloor \log(n) \rfloor}$$

<math display

(a, b) -trees

$$\lfloor \frac{2}{b+1} \leq a \leq \frac{b+1}{2} \rfloor.$$

$\lfloor a$: min no. of children in an internal node of a tree.

$\lfloor b$: max no. of children in an internal node of a tree.

Binary trees	(a, b) -trees
Each node has at most 2 children	Each node can have more than 2 children
Each node stores exactly one key	Each node can store multiple keys

Rules:

1. (a, b) -child policy

$$\begin{aligned} \text{root: } & 1 \leq \text{keys} \leq b-1 \\ & 2 \leq \text{children} \leq b \end{aligned}$$

no. of children
always one more than keys.

$$\begin{aligned} \text{internal: } & a-1 \leq \text{keys} \leq b-1 \\ & a \leq \text{children} \leq b \end{aligned}$$

$$\begin{aligned} \text{leaf: } & a-1 \leq \text{keys} \leq b-1 \\ & \text{no children!} \end{aligned}$$

2. Key ranges

A non-leaf node must have one more child than its number of keys.

\lfloor to ensure all value ranges due to its keys are covered in the addressed.

3. Leaf depth

All leaf nodes must all be at the same depth (from root).

Note: (a, b) -trees in general are good for searching big data.

\hookrightarrow more cache memory utilized.

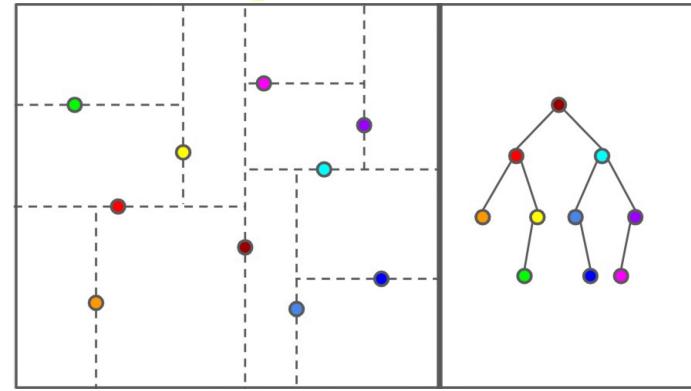
kd-trees

A kd-tree is another simple way to store geometric data in a tree. Let's think about 2-dimensional data points, i.e., points (x, y) in the plane. The basic idea behind a kd-tree is that each node represents a rectangle of the plane. A node has two children which divide the rectangle into two pieces, either vertically or horizontally.

For example, some node v in the tree may split the space vertically around the line $x = 10$: all the points with x -coordinates ≤ 10 go to the left child, and all the points with x -coordinates > 10 go to the right child.

Typically, a kd-tree will alternate splitting the space horizontally and vertically. For example, nodes at even levels split the space vertically and nodes at odd levels split the space horizontally. This helps to ensure that the data is well divided, no matter which dimension is more important.

All the points are stored at the leaves. When you have a region with only one node, instead of dividing further, simply create a leaf.



B-trees

Simply $(B, 2B)$ -trees.

\lfloor Operations: • Split

Conditions:

- 1. node contains $\geq 2a$ keys because after splitting,
LHS has $\geq a-1$ keys.
RHS has $\geq a$ keys

- 2. node's parent contains $\leq b-2$ keys.

- 1. pick median key V_m as split key.
- 2. split node x into LHS and RHS using V_m .
- 3. create a new node y .
- 4. Move LHS split from x to y .

- dealing with root root splitting
- 5. Create a new empty node e
 - 6. insert V_m into e
 - 7. promote e to new root node
 - 8. Assign y and z to be left and right child of e
 - 9. Assign previous subtree t_m associated with V_m to be final child of y .

• Insertion $O(b \log n)$

- 1. Start at node $x = \text{root}$

- 2. Repeat until $x = \text{leaf}$:

- 2.1 If x contains $b-1$ keys,
- 2.1.1 split x into y (LHS) and z (RHS) using split key V_m .
- 2.1.2 Extract V_m returned by the split operation
- If $x \leq V_m$, set $x = y$
- Else, set $x = z$.

- 2.2 Search for a in the keylist of x
- If a is larger than all the keys, set a to be the final child
- Else, look for first key $v_i \geq a$ and set x to be t_i (subtree w/ key range $\leq v_i$)

- 3. Add x to w 's keylist.

Scenario	Operation	Algorithm
• y and z are siblings • Have $< b-1$ keys together	merge(y, z)	<ol style="list-style-type: none"> In parent, delete key v (the key separating the siblings) Add v to the keylist of y In y, merge in z's keylist and treelist Remove z
• y and z are siblings • Have $\geq b-1$ keys together	share(y, z)	<ol style="list-style-type: none"> merge(y, z) Split newly combined node using the regular split operation

• Deletion

- Start at node $w = \text{root}$
- Repeat until w contains x or w is a leaf:
 - If w contains $a-1$ keys and z is a sibling of w
 - merge(w, z)
 - Set w to be the newly merged node
 - Else if w and z together contain $\geq b-1$ keys,
 - share(w, z)
 - Set w to be the newly split node whose key range contains x
- Search for x in the keylist of w
 - If x is larger than all the keys, set w to be the final child
 - Else, look for the first key $v_i \geq x$ and set w to be t_i (subtree w/ key range $\leq v_i$)
- Check for x in w 's keylist,
 - If it exists, remove it and return success
 - Else, return failure

Augmented Trees

Basic Methodology:

1. Choose underlying data structure
(tree, hash table, linked list, stack etc.)
2. Determine additional information
3. Modify data structure to maintain additional info
when the structure changes
(subject to insert/delete etc.)
4. Develop new operations.

Order Statistics.

Inputs: A set of integers.

Output: select(k)

↳ the k th item in the set.



preprocess: $O(n \log n)$

select: $O(1)$

↳ tradeoff: how many items to select?

Implement a data structure that supports:

- insert (C int keys)

- delete (C int keys)

and also:

- select (C int k)



Solution 1:

Basic structure: sorted array A.

insert(int item): add item to sorted array A.

select(int k): return $A[k]$

Solution 2:

Basic structure: unsorted array A.

insert(int item): add item to end of array A.

select(int k): run QuickSelect(k)

Solution 1 > Solution 2

iff there are more selects than inserts.

	Insert	Select
Solution 1: Sorted Array	$O(n)$	$O(1)$
Solution 2: Unsorted Array	$O(1)$	$O(n)$

Solution 3: Use a (balanced) tree

- Use rank (C index of the element in the sorted array)
 - ↳ do not store ranks in every node \Rightarrow just update all nodes per insert.
 - ↳ Store size (weight) of sub-tree in every node
 - $w(\text{leaf}) = 1$
 - $w(v) = w(v.\text{left}) + w(v.\text{right}) + 1$.
 - "rank in sub-tree" = $\text{left.weight} + 1$.

↳ select (k)

```
rank = m.left.weight + 1;
if (k == rank) then
  return v;
else if (k < rank) then
  return m.left.select(k);
else if (k > rank) then
  return m.right.select(k - rank);
```

· Examples:

- Find the l th tallest student in the class. \rightarrow select(k) : find the node with rank k .
- determine the percentile of Johnny's height.
↳ following the l th percentile or g th percentile? \rightarrow rank(k)

· rank (node)

```
rank = node.left.weight + 1;
while (node != null) do
  if (node is left child) then
    do nothing;
  else if (node is right child) then
    rank = rank + parent.left.weight + 1; // Take into account
                                              // that parent is smaller
                                              // than node.
  node = node.parent;
return rank;
```

Basic Methodology:

1. Choose underlying data structure

AVL tree

2. Determine additional information $\text{Weight of each node}$

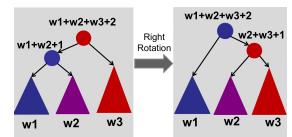
3. Modify data structure to maintain additional info
when the structure changes
update weight as needed.

4. Develop new operations.

select and rank.

Maintain weight during insertions:

- maintain weight during rotations

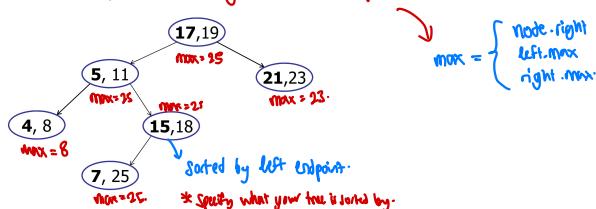


Interval Queries

- Solutions where space is (near) linear in the number of things stored.

- operations are logarithmic in number of things stored.

- given an interval, (a, b) , alignment: minimum endpoint in subtree



- note: maintain max during rotation \rightarrow note that are inserted in bottom.

interval-search(x): find interval containing x \rightarrow $O(\log n)$.

```
C = root;
while (C != null and x is not in C.interval) do
  if (C.left = null) then
    C = C.right;
  else if (x > C.left.max) then
    C = C.right;
  else C = C.left;
return C.interval;
  
```

- if search goes right,
then there is no overlap in the
left subtree.
- if a search goes left, then move to go left.
- if a search goes left and fails, then
key < every interval in right subtree.

List all intervals that overlap at that point: $O(K \log n)$

All-Overlap algorithm:

\downarrow
no. of overlapping intervals.

Repeat until no more intervals:

- search for a interval
- add to list
- delete interval.

\rightarrow add & delete cheap.

Repeat for all intervals on list:

- add interval back to tree.

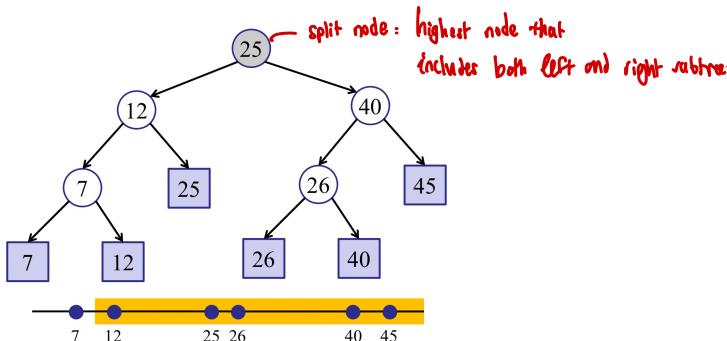
Final algorithm: $O(K + \log n)$.

aka orthogonal range searching.

Orthogonal Range Searching

One dimension

- use a binary tree
- Store all pts in the leaves of the tree.
- ↳ interval not store only copies.
- Each internal node v stores the MAX of any leaf in the left sub-tree.



Algorithm:

- find the split node: $v = \text{findsplit}(low, high)$
- Do left traversal: $\text{LeftTraversal}(v.\text{left}, low, high)$
- Do right traversal: $\text{RightTraversal}(v.\text{right}, low, high)$.

$\text{findsplit}(low, high) \quad O(1 \log n)$

```

v = root;
done = false;
while !done {
    if (high <= v.key) then v = v.left;
    else if (low > v.key) then v = v.right;
    else (done = true);
}
return v;
    ↑
    found split
  
```

$\text{LeftTraversal}(v, low, high) \quad O(1 \log n)$

```

if (low <= v.key) {
    ← all-leaf-traversal(v.right);
    LeftTraversal(v.left, low, high);
}
else {
    ← LeftTraversal(v.right, low, high);
    ↓ return right.
}
  
```

$\text{RightTraversal}(v, low, high) \quad O(1 \log n)$

```

if (v.key <= high) {
    ← all-leaf-traversal(v.left);
    RightTraversal(v.right, low, high);
}
else {
    RightTraversal(v.left, low, high);
    ↓ return left.
}
  
```

Count Points?

- ↳ augment the tree.
- ↳ keep a count of the number of nodes in each subtree.
- ↳ Instead of walking the entire sub-tree, just remember the count.

$\text{LeftTraversal}(v, low, high)$

```

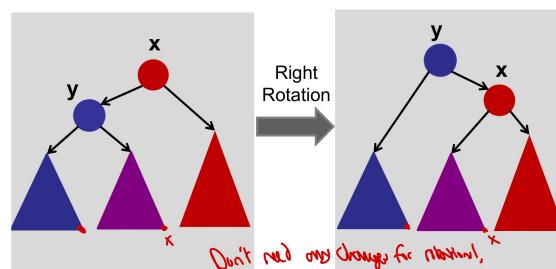
if (low <= v.key) {
    all-leaf-traversal(v.right);
    total += v.right.count; ← keep track of the count
    LeftTraversal(v.left, low, high);
}
else {
    LeftTraversal(v.right, low, high);
}
  
```

Invariant:

The search interval for a left-traversal at node v including the max item in the subtree rooted at v

All-leaf-traversal $\rightarrow O(k)$, $k = \text{num of items found}$

Rotations:



\therefore Query Time Complexity = $O(k + \log n)$
 \downarrow
 k is the no. of pts. found.

Preprocessing (Build tree) Time Complexity = $O(n \log n)$

Total space complexity = $O(n)$.

↳ each pt stored once
 $+ O(n)$ guide nodes.

Symbol Table

Chaining

- no order. - no predecessor / successor
 - search by name
 - no comparison
 - supports insert, search, delete, contains, filter.
 $O(1)$ $O(n)$
 - Sorting
 - Dictionary? \Rightarrow AVLTree $O(n \log n)$
 - Symbol Table? $O(n^2)$

L Direct Access Tables

0	null
1	null
2	item1
3	null
4	Seth
5	item3
6	null
7	null
8	item2
9	null

How to get Index?

↳ Hashing.

Hash Functions

- ↳ There are a lot of possible keys, but very small table.

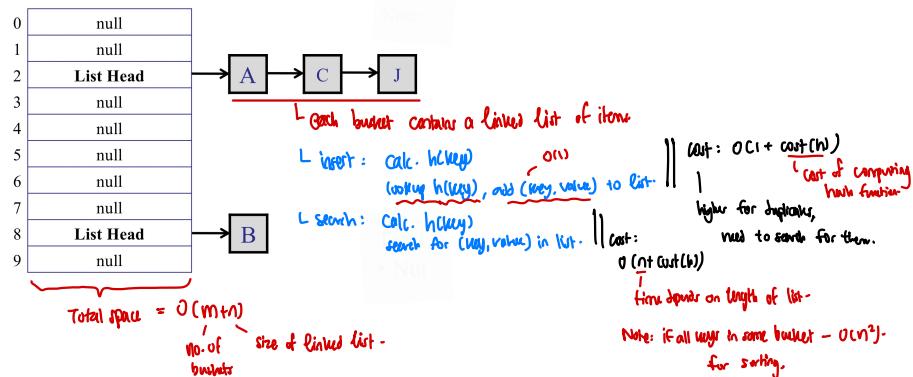
↳ Define $b: U \rightarrow \{1 \dots m\}$

- Store key k in bucket $h(k)$
 - Time: $O(\underbrace{\text{time to compute}}_{\approx O(1 + \text{smith})} + \text{time to access bucket})$

↳ Collision

- 2 distinct keys k_1 and k_2 collide if :

$$h(k_1) = h(k_2)$$
 - Unsolvable by pigeon hole principle.
 - Solved by chaining & open addressing.



Simple Uniform Hashing Assumption

- each key is equally likely to map to any bucket (random bucket)
 - keys are mapped independently.
 - as long as there is enough buckets, there won't be too many keys in one bucket
 - Define: load (hash table) = n/m
 - = average number of items buckets.
 - Expected search time = $1 + \text{expected no. of items per bucket.}$

Note: randomly choose bucket
would increase searching later slow

- rule: $P(\text{choose 1 bin}) = 1/m$

$$\begin{aligned} \text{- Expected search time} &= 1 + n/m \quad (\text{assume } m=2n) \\ &= O(1) \end{aligned}$$

- Worst Case search time = $\Theta(n)$

- Worst case insertion time = $O(1)$

- Max cost to insert n items in hash-table: $O(\log n) \approx \Theta(\log n / \log \log n)$ -
 \hookrightarrow max nr. of items in 1 budget.

Table size $\approx m = \Theta(n)$

Since Expected search time : $O(|tn|m)$

Don't know how many keys we need to store.

Group table?

↳ Choose new size mm

↳ choose your best function

1 b c

Compute & copy to new table.

$$m_2 > 2n.$$

$\approx o(n)$.
initialising a
table of size
 $m_1 \times m_2$