

Announcements:

HW1 available tonight, due Jan 20, 11:59p.

Today: 1. intro to algorithm analysis – code 2. review of asymptotics

Algorithm Analysis – example 1:

```

1 int index of first q (vector<int> arr, int q){
2     for (int i = 0; i < arr.size(); i++)
3         if (arr[i] == q)
4             return i;
5     return -1;
6 }
    
```

dynamic (size can change)

• vector - standard template library array implementation

• c++ by default passes by value

What does it do? (a local copy is made when called by parameter)

0	1	2	3	4	5	6	7
32	11	73	21	29	86	81	92

Algorithmic Analysis – running time

```

1 int _____ (vector<int> & arr, int q){
2     for (int i = 0; i < arr.size(); i++)
3         if (arr[i] == q)
4             return i;
5     return -1;
6 }
    
```

Comparison may take a long time.
eg. comparing pictures

How long does it take? It depends

What does it depend on?

1. size of array
2. clock speed of hardware
3. location of q in the array

What should we count? steps, lines of code, operations.

Algorithmic Analysis Discussion

```

1 int _____ (vector<int> & arr, int q){
2     for (int i = 0; i < arr.size(); i++)
3         if (arr[i] == q)
4             return i;
5     return -1;
6 }
    
```

0	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15
32	11	73	21	29	86	81	92	57	61	64	15	79	44	7	45

What's the best case running time? q is at arr[0]

What's the worst case running time? q is not in the array
longest possible running time

What's the average case running time?
requires a probability distribution over the inputs.

Algorithmic Analysis Discussion

```

1 int _____ (vector<int> & arr, int q){
2     for (int i = 0; i < arr.size(); i++)
3         if (arr[i] == q)
4             return i;
5     return -1;
6 }
    
```

How many lines are executed in the worst case?

$T(n) = 2n + 1$ (let $n = \text{arr.size}()$)

Discussion:

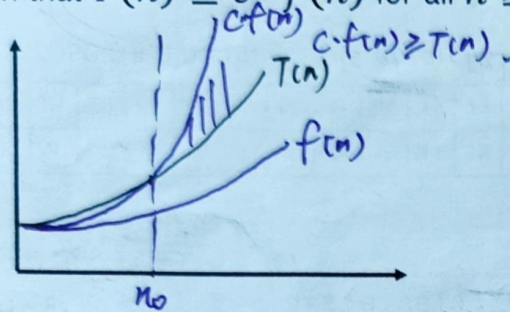
concrete discussion will always lead to in consequential disagreement

$T(n) = cn + d$ where c, d are constants

$T(n) \in O(n)$

Aside on Asymptotics

Defn: $T(n) = O(f(n))$ if there are positive constants c and n_0 such that $T(n) \leq c \cdot f(n)$ for all $n \geq n_0$.



We typically use $f(n)$ to describe the asymptotic upper bound on the worst case running time of an algorithm

Asymptotic definitions

- ▶ $T(n) \in O(f(n))$ if there are positive constants c and n_0 such that $T(n) \leq cf(n)$ for all $n \geq n_0$.
 - ▶ $T(n) \in \Omega(f(n))$ if there are positive constants c and n_0 such that $T(n) \geq cf(n)$ for all $n \geq n_0$ *asymptotic lower bound*
 - ▶ $T(n) \in \Theta(f(n))$ if $T(n) \in O(f(n))$ and $T(n) \in \Omega(f(n))$.
 - ▶ $T(n) \in o(f(n))$ if for any positive constant c , there exists n_0 such that $T(n) < cf(n)$ for all $n \geq n_0$.
 - ▶ $T(n) \in \omega(f(n))$ if for any positive constant c , there exists n_0 such that $T(n) > cf(n)$ for all $n \geq n_0$.
- not seen in the course*

The reminder page:

Typical growth rates in order (increasing order of complexity)

- Constant: $O(1)$
- Logarithmic: $O(\log n)$
- Poly-log: $O((\log n)^k)$
- Linear: $O(n)$
- Log-linear: $O(n \log n)$
- Superlinear: $O(n^{1+c})$ (c is a constant > 0)
- Quadratic: $O(n^2)$
- Cubic: $O(n^3)$
- Polynomial: $O(n^k)$ (k is a constant) "tractable"
- Exponential: $O(c^n)$ (c is a constant > 0) "intractable"

if $f(n) \leq O(n \log n)$
then $f(n) \leq O(n^{1+c})$

A: $T(n) = \omega(n^2)$

B: $T(n) = O(n \log n)$

since $A \cap B = \emptyset$

any conclusion will be true $\omega(n^2)$

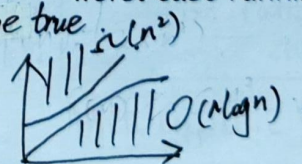
Algorithmic Analysis Discussion

```

1 int _____ (vector<int> & arr, int q){
2     for (int i = 0; i < arr.size(); i++)
3         if (arr[i] == q)
4             return i;
5     return -1;
6 }
    
```

$T(n) = \underline{\hspace{2cm}}$

Choose which definition we should use for describing worst case running time? ($o, O, \omega, \Omega, \theta, \Theta$)



Announcements:

HW1 available, due Jan 20, 11:59p.

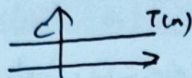
Warm-up:

```

1 int return first (vector<int> & arr, int q){
2     for (int i = 0; i < arr.size(); i++)
3         if (i == 0)
4             return arr[i];
5     return -1;
6 }
    
```

Good name?

Running time:



$O(n)$ And this $T(n)$ is not $\omega(n)$
 $T(n) = c$ so $T(n) = \theta(1)$
 by definition of θ : $T(n) \leq c \cdot 1$
 $\wedge T(n) \geq c \cdot 1$
 $\Rightarrow T(n) = \theta(1)$ by θ : $T(n) = \theta(1)$

Algorithmic Analysis Discussion

```

1 int _firstLoc_ (vector<int> & arr, int q){
2     for (int i = 0; i < arr.size(); i++)
3         if (arr[i] == q)
4             return i;
5     return -1;
6 }
    
```

$T(n) = cn + d$

Choose which definition we should use for describing worst case running time? (o , O , ω , Ω , θ , Θ)

since $T(n) \leq cn + d$. $T(n) \geq cn + d$.
 $T(n) = \theta(n)$

Runtime analysis guidelines:

always worst case

- Single operations, constant time
- Consecutive operations additive
- Conditionals: constant for Boolean eval + branch time
- Loops: sum of loop body times *And, or, not. worst branch (longest time)*
- Function call: time for function *Cannot assume that the body of for loop always takes the same amounts of time*

Some puzzles:

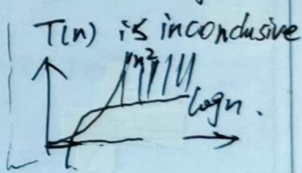
Typical growth rates in order

- Constant: $O(1)$
- Logarithmic: $O(\log n)$
- Poly-log: $O((\log n)^k)$
- Linear: $O(n)$
- Log-linear: $O(n \log n)$
- Superlinear: $O(n^{1+c})$ (c is a constant > 0)
- Quadratic: $O(n^2)$
- Cubic: $O(n^3)$
- Polynomial: $O(n^k)$ (k is a constant)
- Exponential: $O(c^n)$ (c is a constant > 0) "intractable"

Tractable

if we take $k > 0$ $k \rightarrow \infty$.
 n^c $c \rightarrow 0$.
 we still have $n^c > (\log n)^k$

if $T(n) = O(n^2)$ and $T(n) = \omega(\log n)$



Announcements:

HW1 available, due Jan 20, 11:59p.

Warm-up:

```

1 int find_min (vector<string> & candy, int a){
2     int retloc = a;
3     string voi = candy[a];
4     for (int i = a+1; i < candy.size(); i++)
5         if (candy[i] < voi){
6             retloc = i;
7             voi = candy[i];
8         }
9     return retloc;
10 }
    
```

keep track of the same thing

changed simultaneously

n-1 iterations

n-1 iterations in worst case

retloc holds voi's location

voi holds min in candy[a..i]

Good name?

Does it work?

Let $C_1 = \min(d, c)$
 $C_2 = \max(d, c)$

Running time:

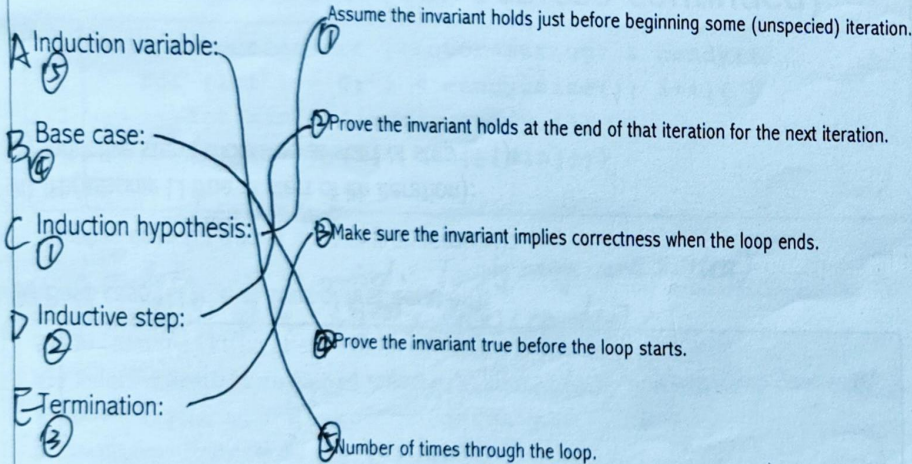
$$T(n) = d + c(n-1) = O(n)$$

$$T(n) \leq (C_2 + d)n \text{ for } n \geq C_2$$

$$n \geq C_2 \geq d - c$$

$$T(n) \geq (C_1 - 1)n \text{ for } n \geq C_1$$

Proving a (given) loop invariant--



Proving correctness is only one benefit of loop invariants, they are also a natural way to think about your program!

Correctness continued:

```

1 int find_min (vector<string> & candy, int a){
2     int retloc = a;
3     string voi = candy[a];
4     for (int i = a+1; i < candy.size(); i++)
5         if (candy[i] < voi){
6             retloc = i;
7             voi = candy[i];
8         }
9     return retloc;
10 }
    
```

Correctness? (argued formally via induction) *stopped at i.* "loop invariant"

Thm: for any i in $\{a+1, \dots, n\}$, $\text{findMin}(\text{candy}, a)$ gives index of min value in $\text{candy}[a..i]$

iterative variable is $i \in \{a+1, \dots, n-1\}$.

Base case: $i = a+1$ $\text{retloc} = a$ $\text{voi} = \text{candy}[a]$, represents the min of $\text{candy}[a..a]$ is true

IH: retloc and voi holds min value in $\text{candy}[a, \dots, i-1]$.

start iteration i , if $\text{candy}[i] \geq \text{voi}$, we do nothing, $\text{voi}, \text{retloc}$ represent min in $\text{candy}[a..i]$

if $\text{candy}[i] < \text{voi}$, then voi and retloc are updated. it is still true that retloc and voi represent the minimum in $\text{candy}[a, \dots, i]$

Another Example

```

1 void _____ (vector<string> & candy){
2     for (int i = 0; i < candy.size(); i++){
3         int min = findMin(candy, i);
4         string temp = candy[i];
5         candy[i] = candy[min];
6         candy[min] = temp;
7     }
    
```

Functionality?

Running Time?

Correctness?

1) iterative variable:

2) loop invariant:

Announcements:

HW1 due today 11:59p. PA1 due 02/03, 11:59p. Quizzes

Warm-up: Weird Mystery Function

```

1 void slide (vector<string> & candy, int loc){
2 // assumes candy[0:loc-1] is sorted, loc valid
3 string temp = candy[loc]; temp=C.
4 int j = loc; j=5
5 while (j > 0 && candy[j-1] > temp) {
6 candy[j] = candy[j-1];
7 j--; }
8 candy[j] = temp;
9 }

```

loc

0	1	2	3	4	5	6	7
A	D	E	H	L	C	J	R

largest loc for this candy is 5
A D E H L | J R

Good name?

Does it work?

Running time: $O(n)$ if standalone
 $\Theta(\text{loc})$ since loc determined by insertion sort.
 Assume candy[0:loc-1] is sorted.

Sort values candy[0...loc].

Insertion Sort

```

1 vector<string> insertionSort (vector<string> & candy){
2 for (int i = 1; i < candy.size(); i++) {
3     slide (candy, i);
4 return candy; }

```

3) Base case: $i=1$

candy[0, ..., i-1] = candy[0...0] is sorted.

4) IH: at start of iteration i , candy[0, ..., i-1] is sorted.

5) Inductive step:

at start of iteration i , candy[0, ..., i-1] is sorted.

by correctness of slide, after line 3, candy[0, ..., i] is sorted.

LI restored for iteration $i+1$

6) Termination: $i=n$.

LI gives candy[0, ..., n-1] is sorted.

You write Insertion Sort...

```

1 void insertionSort (vector<string> & candy){
2     for (int i=1; i < candy.size(); i++){
3         slide (candy, i);
4     }
5 }

```

0	1	2	3
M	A	T	H

Functionality? insertion sort.

Running Time?

$$\sum_{i=1}^n i = \frac{(n-1)n}{2} = \Theta(n^2)$$

Correctness?

- Iterative variable: $i \in \{1, \dots, n\}$ ← not happen
- Loop invariant: candy[0...i-1] is sorted.

Linear Sorts, recap

We have learned and analyzed selection sort and insertion sort.

Which is better?

- Asymptotically? the same $\Theta(n^2)$.
- Empirically? (with data, practically) selection
- What if list is already sorted? ~~sort~~ insertion
- What if list is almost sorted? selection
- What if list is in reverse order? selection

<https://www.toptal.com/developers/sorting-algorithms>

Something NEW!!

Make at least 3 observations about this code:

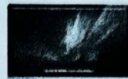
```
1 template <class LIT>
2 struct Node{
3     LIT data;
4     Node * next;
5     Node(LIT nData, Node * next=NULL): data(nData) {}
6 };
```

struct is the same as class

Syntax for default

- recursive struct
- line 1 parameterizes type info
- struct is class with public default access
- line 5 is a constructor

Switching gears... for line 4:



Configure your iMac 27-inch

Use the options below to build the system of your dreams

Memory

More memory (RAM) increases performance and enables your computer to perform faster and better. Choose additional 1066MHz DDR3 memory for your iMac.

The more memory your computer has, the more programs you can run simultaneously, and the better performance you get from your computer.

- Select the standard memory configuration to support day-to-day tasks such as email, word processing, and web browsing as well as more complex tasks such as editing photos, creating illustrations, and building presentations.
- Upgrade your memory to enjoy greater performance for more advanced computing tasks such as video editing and DVD authoring.

Your iMac uses one of the following configurations without wasting clock

- 4GB 1066MHz DDR3
- 8GB 1066MHz DDR3
- 16GB 1066MHz DDR3
- 32GB 1066MHz DDR3



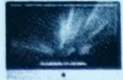
Announcements:

PA1 due 02/03, 11:59p. Quizzes ongoing.

Memory



MLLOXO.



Configure your iMac 27-inch

Use the options below to build the system of your dreams

Memory

More memory (RAM) increases performance and enables your computer to perform faster and better. Choose additional 16GB (8GB) memory for your Mac.

The more memory you get, the better the situation with opening all your applications.

Upgrade your memory and drive performance.

Your Mac uses less energy when you're not using it.

Without wasting class

16GB (8GB) DDR3

16GB (8GB) DDR3

16GB (8GB) DDR3

16GB (8GB) DDR3



Variables and Pointers:

Stack

every variable has

loc	name	val	type
a1b	p	mem addr	int*
a2c	y		int
a3b	tyler	int	Node
a4c	x	15	int

Backtrack to variables:

int x; ← before assigning, val is garbage
x = 15;

Node<int> tyler; // const invoked
int y;

Special variable type:
(pointer == memory address)

int *p;

P = x; is a compiler error (type mismatch)

P = &x; (& is memory address operator, gives the memory address of x)

Pointers and dynamic memory:

Stack

Heap

```
int x;
int *p;
p = &x;
cout << *p << endl;
```

* unary operator:
follows pointer and
returns the target.

```
p = new int;
*p = 42;
```

loc	name	val	type
a3b	p	b12	int*
a4c	x	15	int

loc	name	val	type
b12		42	int

Fun and games with pointers: (warm-up)

int * p, q;

What type is q? int

int *p;

int x;

p = &x;

*p = 6; x = 6

cout << x;

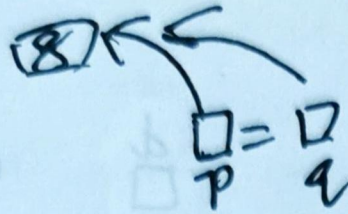
What is output? 6

cout << p;

What is output? address of x

Write a statement whose output is the value of x, using variable p: cout << *p;

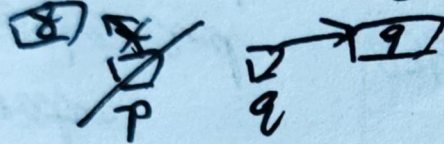

```
int *p, *q;
p = new int;
q = p;
*q = 8;
```



```
cout << *p;
```

What is output? 8 ✓

```
q = new int;
*q = 9;
```



```
p = NULL;
```

Do you like this? X memory leak

delete q; give the memory location at q back to the system.

```
q = NULL;
```

Do you like this? ✓

pull back the point.

if *q=30, runtime error, seg fault

Memory leak:

occurs when all pointers to some heap memory are removed.

Deleting a null pointer: nothing happens

Dereferencing a null pointer: seg fault, runtime error.

Announcements: PA1 due 02/03, 11:59p. Quizzes ongoing.

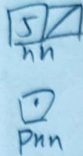
```

1 template <class LIT>
  struct Node{
2     LIT data;
3     Node * next;
4     Node(LIT newData):data(newData),next(NULL){};
  };
  
```

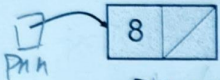
What is the result of each statement?

```

int main(){
  Node<int> nn(5);
  Node<int> * pnn; }
  
```



Write code that would result in each of these memory configurations (in order):



```

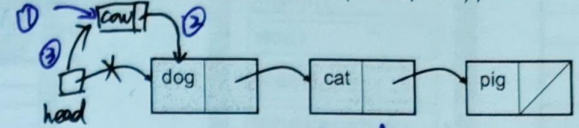
pnn = new Node<int>(8);
Node<int>* n1 = new Node<int>(6);
Node<int>* tmp = pnn;
  
```



```

pnn = n1;
pnn->next = tmp;
Or Node* t = new Node(6);
t->next = pnn;
pnn = t;
  
```

Example 1: insertAtFront<farmAnimal>(head, cow);



```

void insertAtFront(Node* curr, LIT e) {
  Node<LIT> * t = new Node<LIT>(e);
  t->next = curr;
  curr = t;
}
  
```

pass by reference

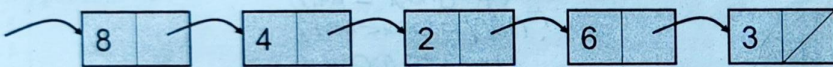
Running time?
 $\Theta(1)$

```

1 struct Node{
2   LIT data;
3   Node * next;
4   Node(LIT newData):data(newData),next(NULL){};
5 };
  
```

Example 2:

8 4 2 6 3
outcome: 3 6 2 4 8



```

void printReverse(Node* curr) {
  if (curr != NULL) {
    printReverse(curr->next);
    cout << curr->data;
  }
}
  
```

$T(n)$: running time on list of length n
 Running time?
 $T(0) = C$
 $T(n) = C + T(n-1)$
 $= \Theta(n)$

```

1 struct Node{
2   LIT data;
3   Node * next;
4   Node(LIT newData):data(newData),next(NULL){};
5 };
  
```

Example 3: 8 4 2 6 3 1 \rightarrow 3 2 8



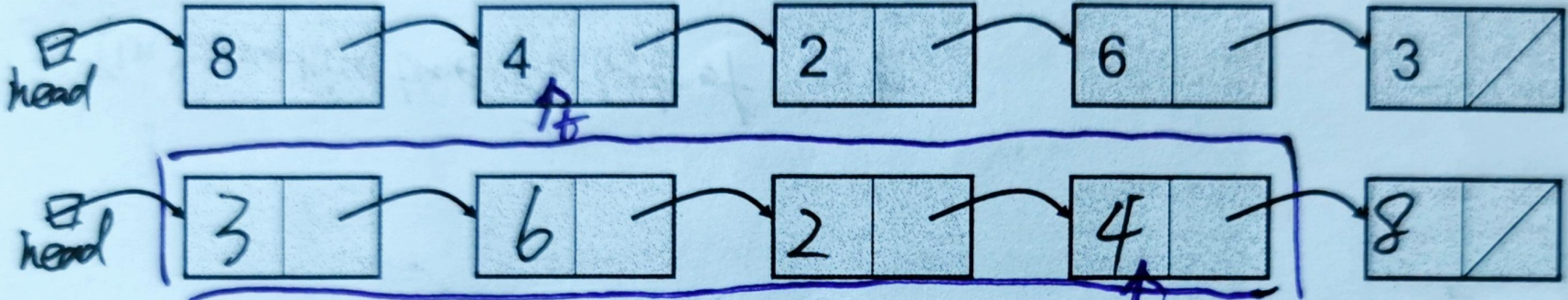
```

void printReverseODDS(Node* curr) {
  // no items in list
  if (curr == NULL) {
    return;
  }
  // one item in list
  else if (curr->next == NULL) { cout << curr->data; }
  // > 1 item in list
  else { printReverseODDS(curr->next->next);
         cout << curr->data; }
}
  
```

Running time? $T(0) = C$ $T(1) = C$
 $T(n) = C + T(n-2) = \Theta(n)$

Example 4:

$head = reverse(head);$



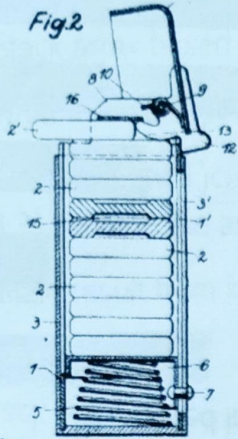
```
node * reverse(Node * curr) {  
    if (curr != NULL && curr->next != NULL) {  
        Node * temp = curr->next; // one element  
        Node * revRest = reverse(curr->next);  
        temp->next = curr;  
        curr->next = NULL;  
        curr = revRest;  
    }  
    return curr;  
}
```

Running time? $O(n)$

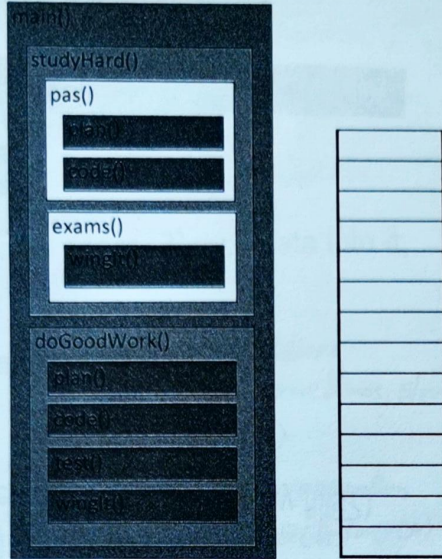
Announcements: PA1, 2/3. Exam1, 2/5.

Stacks

Fig.2



if left then push;
if right remove and match;
{ } { () } { () } { () } { () }



$$45 + 72 - * 3 - 6 /$$

Abstract Data Type: description of the functionality of a structure (interface)

Stack ADT:

```

1 template<class LIT>
2 class Stack{
3 public:
4     Stack(); // ~Stack, ctor, op=?
5     bool empty() const;
6     void push(const LIT & e); insert.
7     LIT pop(); remove, return data.
8 private:
9
10    ?
11
12 };
    
```

push(3)
push(8)
push(4)
pop() → 4
pop() → 8
push(6)
pop() → 6
push(2)
pop() → 2
pop() → 3.

Push 1, 2, 3 in order. What output is impossible?
3 | 2 is impossible.

Stack linked memory implementation:

```

1 template<class LIT>
2 class Stack{
3 public:
4     Stack(); // ~Stack, ctor, op=?
5     bool empty() const;
6     void push(const LIT & e);
7     LIT pop();
8 private:
9     struct Node {
10        LIT data;
11        Node * next;};
12    Node * top;
13    int size;
14 };
    
```

Running time.
 $T(\text{pop}) = T(\text{push}) = \Theta(1)$



```

1 template<class LIT>
2 LIT Stack<LIT>::pop(){
3     assert(!empty());
4     LIT retval = top->data;
5     top = top->next; <Node* tmp = top;
6     return retval; <delete tmp;
7 }                                     size--;
                                        tmp = NULL;
    
```

insert At Front-

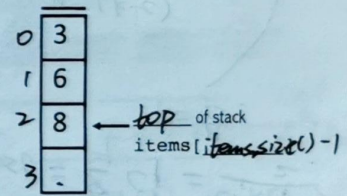
```

1 template<class LIT>
2 void Stack<LIT>::push(LIT d) {
3     Node* newN = new Node(d);
4     newN->next = top;
5     top = newN; size++;
6 }
    
```

Stack array based implementation:

```

1 template<class LIT>
2 class Stack{
3 public:
4     Stack(); // ~Stack, ctor, op=?
5     bool empty() const;
6     void push(const LIT & e);
7     LIT pop();
8 private:
9     vector<LIT> items;
10 };
    
```



```

1 template<class LIT>
2 LIT Stack<LIT>::pop(){
3     assert(!empty());
4     LIT retval = items[items.size()-1];
5     return retval;
6 }
    
```

$T(\text{pop}) = \Theta(1)$

```

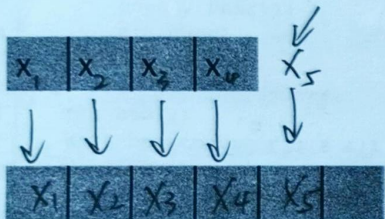
1 template<class LIT>
2 void Stack<LIT>::push(LIT d){
3
4     items.push_back(d);
5
6 }
    
```

Stack array based implementation: (what if array fills?)

How do vectors work? STL implementation of a dynamic resizing array

What does it mean for an array to be "dynamic"? detects when it needs more space and grabs it

Analysis holds for array based implementations of Lists, Stacks, Queues, Heaps...

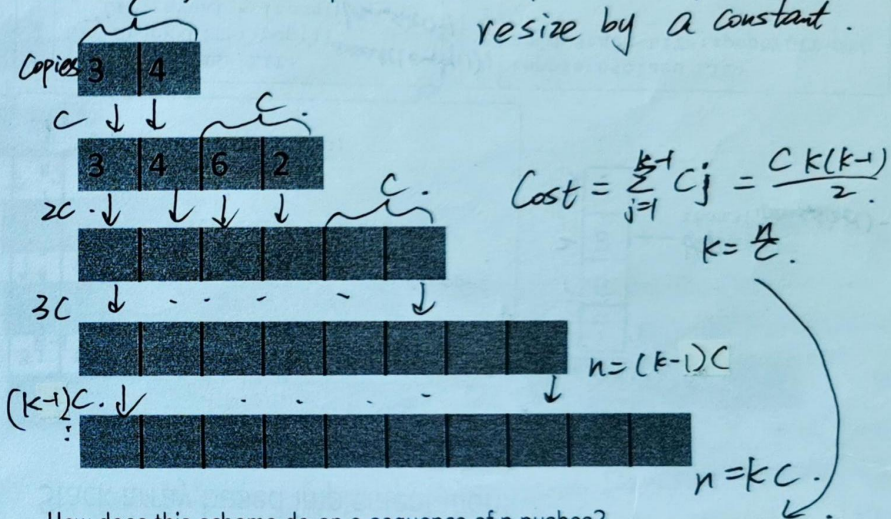


General Idea: upon an insert (push), if the array is full, create a larger space and copy the data into it.

Main question: What's the resizing scheme? We examine 2...

Stack array based implementation: (what if array fills?)

resize by a constant.



How does this scheme do on a sequence of n pushes?

$$\text{Cost} = C \cdot \frac{\frac{n}{C}(\frac{n}{C}-1)}{2} = \Theta(n^2)$$

over n additions \Rightarrow avg of $\Theta(n)$ per addition.

Announcements: PA1, 2/3. Exam1, 2/5.

Array Resizing fin...

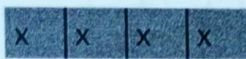
... Stack array based implementation:

... What if the array fills?

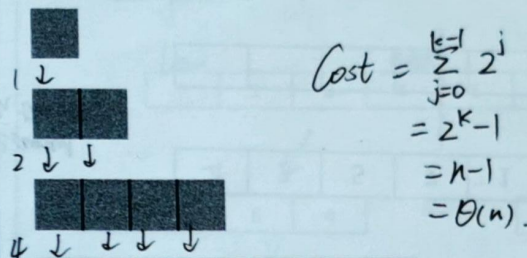
... General idea: 1) create larger space, 2) copy data into it, and 3) rename the array.

... Main result from last time? *analysis over n additions sometimes fast, sometimes slow (amortized analysis)*

$\Theta(n^2)$ cost over n operation
 \Rightarrow an average of $\Theta(\frac{n^2}{n}) = \Theta(n)$ per operation.



Another resizing strategy: (current size) $\times 2$ when array fills



when popping, pop to $\frac{1}{2}n$, then shrink to $\frac{1}{4}n$
 if we grow by an arbitrary constant, the base changes $n = C \cdot k$

How does this scheme do on a sequence of n pushes? $\Theta(n)$ over n pushes.
 $\Theta(\frac{n}{n}) = \Theta(1)$ per push.

$n = 2^k$ we resize $\log_2 n$ times.

Stack Summary:

Linked list based implementation of a stack:

Constant time push and pop.

Array based implementation of a stack:

$\Theta(1)$ time pop, if not too empty
 $\Theta(1)$ time push if capacity exists,

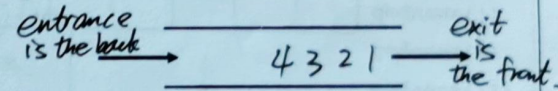
Cost over $O(n)$ pushes is $\Theta(n)$ for an AVERAGE of $\Theta(1)$ per push. *changes*

Why consider an array? *Faster in practice (array is a block of memory) tracing down the pointer for the linkedlist takes time. if array takes 1min, the linkedlist takes 10 min*

Queues

Queue ADT:

- enqueue
- dequeue
- isEmpty

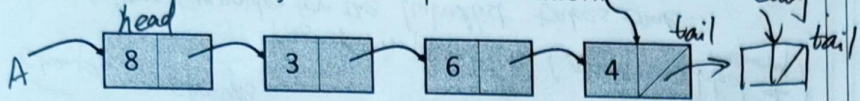


Applications:

- Hold jobs for a printer
- Store packets on network routers
- Make waitlists fair
- Breadth first search

- enqueue (1)
- enqueue (2)
- enqueue (3)
- dequeue () $\rightarrow 1$
- dequeue () $\rightarrow 2$
- enqueue (4)

Queue — linked memory based implementation: B



```

1 template<class SIT>
2 class Queue {
3 public:
4 // ctors dtor
5 bool empty() const;
6 void enqueue(const SIT & e);
7 SIT dequeue();
8 private:
9     struct queueNode {
10         SIT data;
11         queueNode * next;
12     };
13     queueNode * entry;
14     queueNode * exit;
15     int size;
16 };
    
```

Which ptr is "entry", which is "exit"?

B
if B is exit, *dequeue* need to set a pointer with $O(1)$ operation

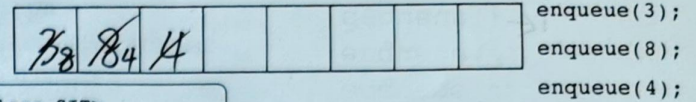
What is running time of enqueue?

$O(1)$

What is running time of dequeue?

$O(1)$

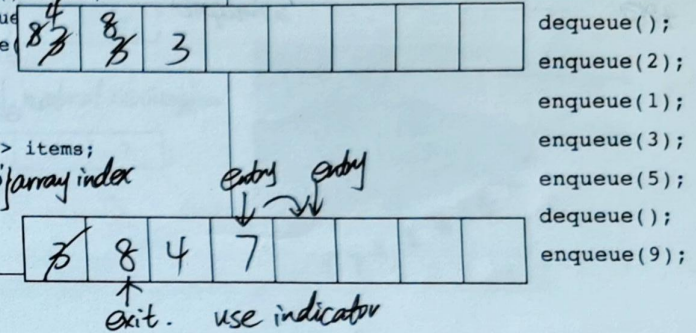
Queue array based implementation:



```

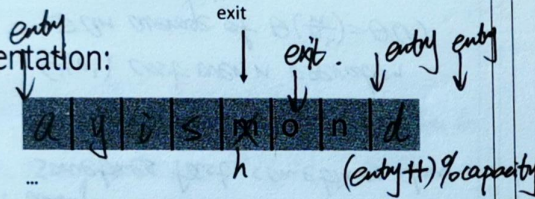
1 template<class SIT>
2 class Queue {
3 public:
4 // ctors dtor
5 bool empty() const;
6 void enqueue
7 SIT dequeue
8 private:
9
10
11 vector<SIT> items;
12 int entry; {array index
13 int exit;
14
15 };
    
```

shifting data is bad



exit. use indicator

Queue array based implementation:



```

... enqueue(m); enqueue(y);
enqueue(o); enqueue(i);
enqueue(n); enqueue(s);
... dequeue();
enqueue(d); enqueue(h);
enqueue(a); enqueue(a);
    
```

```

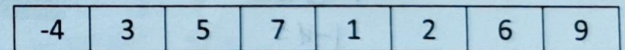
1 template<class SIT>
2 class Queue {
3 public:
4 // ctors dtor
5 bool empty() const;
6 void enqueue(const SIT & e);
7 SIT dequeue();
8 private:
9     int size;
10     vector<SIT> items;
11     int entry;
12     int exit;
13
14
15 };
    
```

double the array start with the exit.

Circular array

Something new...

Given: an array _____
Output: a sorted list

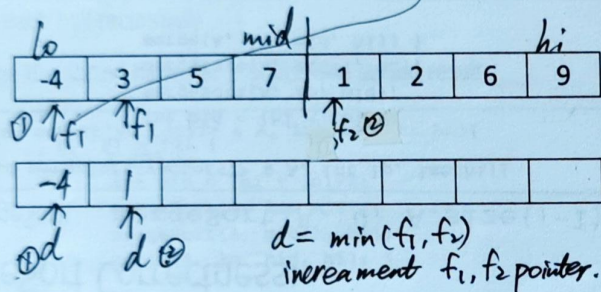


Announcements: PA1, 2/3. Exam1, 2/5.

Something New:

Given: an array that we can divide into 2 sorted parts

Output: a sorted list
merge (vector A, int lo, int mid, int hi)



Running Time? n : # of data elements we want to sort
 $O(n)$

mergeSort

A "divide and conquer" algorithm. ^{1) solve smaller subproblems}
^{2) reassemble solutions into the solution to original problem.}

- If the array has 0 or 1 elements, it's sorted. Stop.
- Split the array into two approximately equal-sized halves.
- Sort each half recursively
- Merge the sorted halves to produce one sorted result.

```

1 void mergeSort(vector<T> & A, int lo, int hi){
2     if (hi > lo) {
3         int mid = (hi + lo)/2;
4         mergeSort(A, lo, mid);
5         mergeSort(A, mid+1, hi);
6         merge(A, lo, mid, hi);
7     }

```

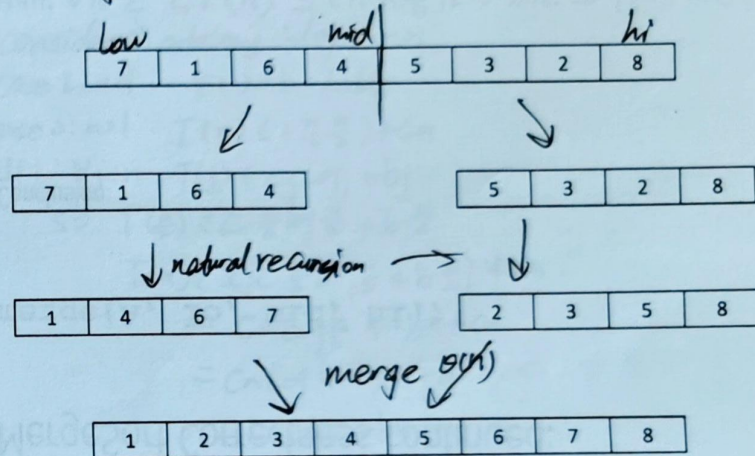
Base hi ≤ lo
0, 1 element

RT: $T(n)$: running time of merge sort on data of size $n = hi - lo + 1$
 $T(0) = T(1) = b$

$$T(n) \leq 2T\left(\frac{n}{2}\right) + cn, n > 1$$

two merge sort merge

OK, then sort this array:
merge sort



no recursive term and no sums
Finding a closed form (two approaches, there are others):

$$\begin{aligned}
 1) \text{ Expand and generalize: } T(0) = T(1) = b \\
 T(n) &\leq 2T\left(\frac{n}{2}\right) + cn & T\left(\frac{n}{2}\right) &\leq 2T\left(\frac{n}{4}\right) + \frac{cn}{2} \\
 &\leq 2\left(2T\left(\frac{n}{4}\right) + \frac{cn}{2}\right) + cn & T\left(\frac{n}{4}\right) &\leq 2T\left(\frac{n}{8}\right) + \frac{cn}{4} \\
 &= 4T\left(\frac{n}{4}\right) + 2cn \\
 &\leq 4\left(2T\left(\frac{n}{8}\right) + \frac{cn}{4}\right) + 2cn \\
 &= 8T\left(\frac{n}{8}\right) + 3cn.
 \end{aligned}$$

pattern $= 2^k T\left(\frac{n}{2^k}\right) + kcn$. *express the pattern in terms of recursive variable.*

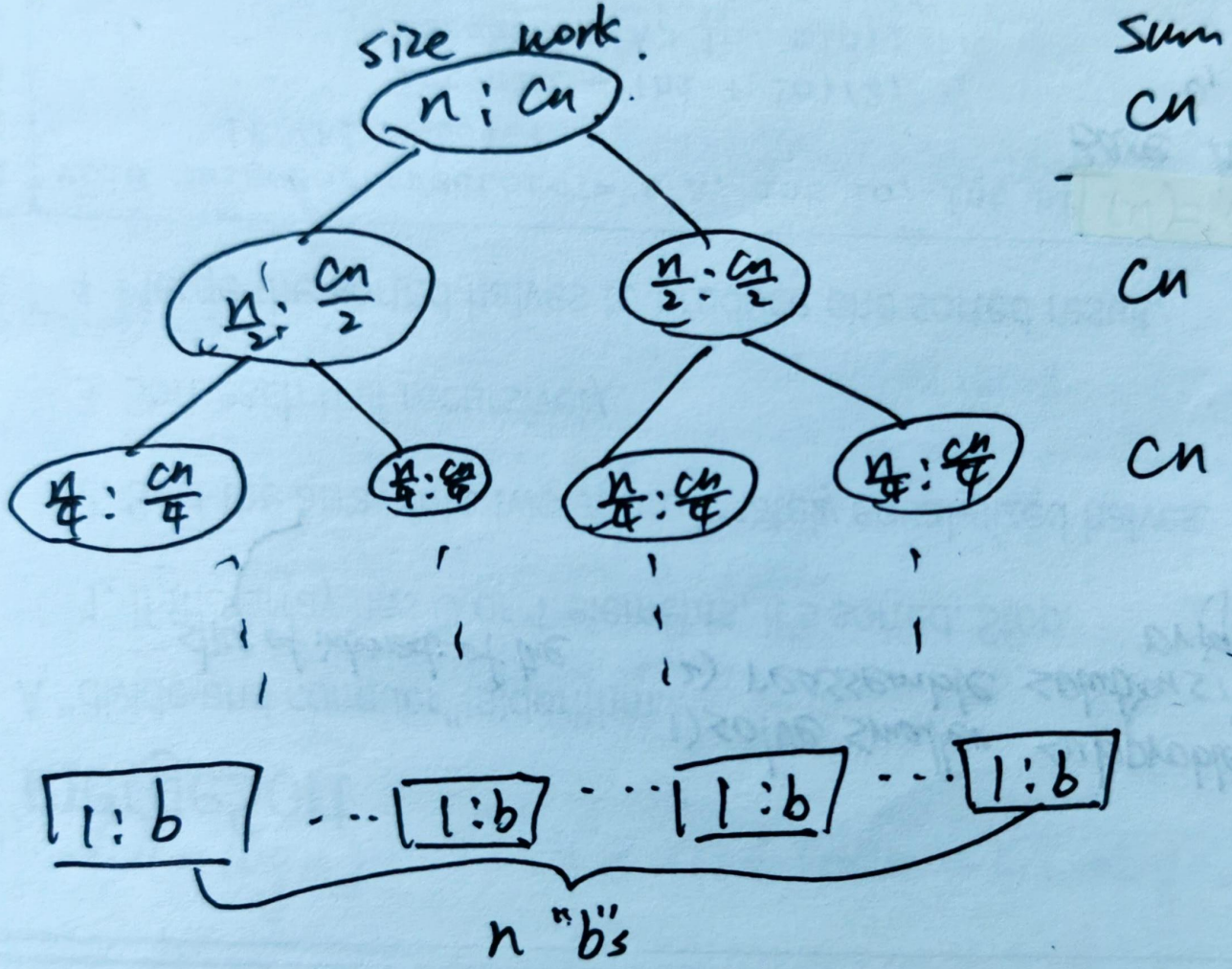
force the base case $\left(\frac{n}{2^k} = 1\right)$
 $n = 2^k \cdot k = \log_2 n$

$$T(n) \leq nb + \log_2 n \cdot cn$$

$$T(n) = O(n \log n)$$

Finding a closed form (two approaches, there are others):

2) Recursion Tree: $T(n) \leq 2T(\frac{n}{2}) + cn$.



Announcements: PA1, today! Exam1, 2/5.

mergeSort:

A "divide and conquer" algorithm.

1. If the array has 0 or 1 elements, it's sorted. Stop.
2. Split the array into two approximately equal-sized halves.
3. Sort each half recursively
4. Merge the sorted halves to produce one sorted result.

```

1 void mergeSort(vector<T> & A, int lo, int hi){
2     if (hi > lo) {
3         int mid = (hi + lo)/2;
4         mergeSort(A, lo, mid);
5         mergeSort(A, mid+1, hi);
6         merge(A, lo, mid, hi); }
7 }
    
```

RT: $O(n \log n)$ $T(n) \leq 2T(\frac{n}{2}) + cn = O(n \log n)$

$$T(n) \leq 2T(\frac{n}{2}) + cn$$

Proving the closed form is correct: $T(0) = T(1) = b$. given

Thm: $\forall n \geq 1, T(n) \leq cn \log n + bn. \Rightarrow T(n) = O(n \log n)$

consider an arbitrary integer $n \geq 1$

Case 1: $n=1$ $T(1) = b$ holds

Case 2: $n > 1$ $T(n) \leq 2T(\frac{n}{2}) + cn$

IH: $\forall j < n, T(j) \leq cj \log j + bj$

$$\text{so } T(\frac{n}{2}) \leq c \cdot \frac{n}{2} \log \frac{n}{2} + b \cdot \frac{n}{2}$$

$$T(n) \leq 2 \left[c \cdot \frac{n}{2} \log \frac{n}{2} + b \cdot \frac{n}{2} \right] + cn$$

$$= cn \log \frac{n}{2} + bn + cn$$

$$= cn \log n - cn + bn + cn$$

$$= cn \log n + bn$$

MergeSort Correctness:

Call MergeSort: `mergeSort(A, 0, A.size()-1);`

```

1 void mergeSort(vector<T> & A, int lo, int hi){
2     if (hi > lo) {
3         int mid = (hi + lo)/2;
4         mergeSort(A, lo, mid);
5         mergeSort(A, mid+1, hi);
6         merge(A, lo, mid, hi); }
7 }
    
```

Claim: mergeSort correctly sorts $A[lo..hi]$ for any $(hi-lo+1) \in \mathbb{Z}$

Let $n = (hi - lo + 1)$, arb

$n \leq 1$: (Base Case) $hi - lo + 1 \leq 1, hi \leq lo$, if $hi = lo$ no elements, if $hi = lo - 1$ 1 element.

$n > 1$: IH: Assume mergeSort correctly sorts a slice of size j , for any $j < n$. *nothing is done.*

`int mid = (hi + lo)/2;` *mid $\geq lo$ and mid $\leq hi$.*

`mergeSort(A, lo, mid);` sorts $A[lo, \dots, mid]$. because sizes smaller.
`mergeSort(A, mid+1, hi);` sorts $A[mid+1, \dots, hi]$. (by IH)

MergeSort Correctness continued:

`merge(A, lo, mid, hi);`
 merge demands $A[lo, \dots, mid]$ and $A[mid+1, \dots, hi]$ sorted
 and returns $A[lo, \dots, hi]$ sorted.

Conclusion:

`mergeSort(A, lo, hi)` works for any size.

\Rightarrow `mergeSort(A, 0, A.size()-1)` sorts A.

To contemplate:

does the value of mid matter in the correctness proof? *yes, but it only needs to be somewhere between mid and hi.*

does the value of mid matter in the analysis of the runtime?

yes.

Where are we in the sorting picture?

	Best Case	Average Case	Worst Case
Insertion	$\Theta(n)$	$\Theta(n^2)$	$\Theta(n^2)$
Selection	$\Theta(n^2)$	$\Theta(n^2)$	$\Theta(n^2)$
Merge	$\Theta(n)$	$\Theta(n \log n)$	$\Theta(n \log n)$

no merge needed.

<https://www.toptal.com/developers/sorting-algorithms>

Complexity of the Sorting Problem:

The *complexity* of a problem is the runtime of the *fastest* algorithm for that problem.

We'll only consider comparison-based algorithms. They can compare two array elements in constant time.

They cannot manipulate array elements in any other way.

For example, they cannot assume that the elements are numbers and perform arithmetic operations (like division) on them.

Insertion, Merge, Quick, ~~Radix~~, Selection

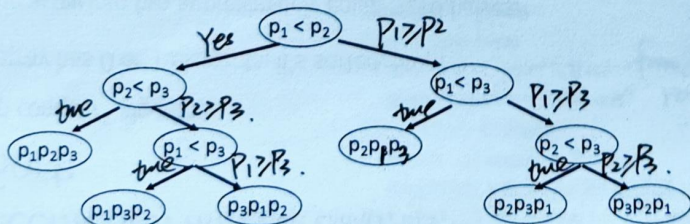
Lower bound on sorting:

An example – different encoding of a sorting procedure

Suppose you have an unsorted array of length 3

Need to produce a result for every possible arrangement ($p_1 p_2 p_3$)

Fundamental operation: Compare 2 elements using $<$



need $n!$ leaves for ends to a sequence of comparison.

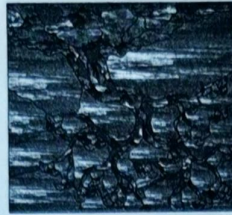
Lower bounds on ($<$ based) sorting algorithms (n items):

Consider a decision tree that arises from a comparison based sorting algorithm on n items:

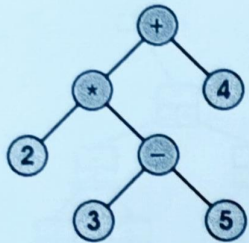
- Every permutation must be differentiated in the algorithm
- There are _____ permutations.
- Comparisons induce a *binary tree* whose height determines the running time of the algorithm.
- Observe that a binary tree of height k has at most _____ leaves.

What is the minimum height of a comparison tree for sorting input of size n ?

Announcements: HW2 due 2/14.



Branching: d-ary trees (*binary* if $d = 2$)



A d-ary tree T is either

- empty $T = \{\}$
- OR
- not empty. $T = \{\text{root}, T_L, T_R, \dots, T_d\}$.
 T_k are d-ary trees

Full d-ary tree:

every node has either 0 or d children.

Perfect d-ary tree of height h: max number of nodes in a tree of height h
 $P(h) = dP(h-1) + 1$

Complete d-ary tree: Like a perfect tree except deepest level may be missing nodes from the right.

Complete trees: not necessarily full.
 Perfect trees are complete trees.

Practice

- 1) Draw a FULL tree with 9 nodes. Is it the same as your neighbor's?
- 2) Draw a COMPLETE tree with 12 nodes. Is it the same as your neighbor's?
- 3) Give an upper bound for the number of nodes in a tree of height h:

$$O(c^h) \quad N(h) = 2^{h+1} - 1$$

$$n \leq N(h) = 2^{h+1} - 1$$

- 4) Give a lower bound for the height of a tree containing n nodes:

$$\Omega(\log n) \quad h(n) = \log_2(n+1) - 1$$

$= \lfloor \log_2 n \rfloor$ based on fact that heights are ints.

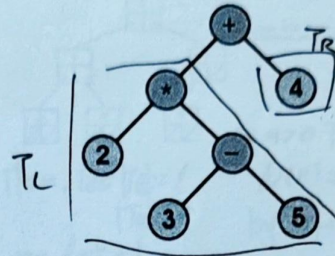
- 5) Give upper and lower bounds for the number of nodes in a complete tree of height h:

upper $c^{h+1} - 1$ $2^{h+1} - 1$

lower $c^{h+1} - c^{h+1}$ $2^{h+1} - 2^h + 1$

Binary Tree Height

height(r) -- length of longest path from root r to a leaf



Given a binary tree T, write a recursive defn of the height of T, height(T):

- $T = \{\}$ empty $h(T) = -1$
- $T = \{r, T_L, T_R\}$
 $h(T) = \max\{h(T_L), h(T_R)\} + 1$

Number of nodes in a perfect tree of height h, N(h):

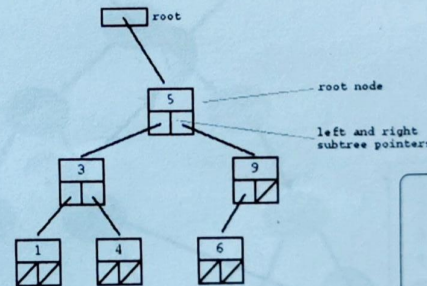
$$N(h) = dN(h-1) + 1$$

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

$$N(1) = dN(0) + 1$$

$$N(h) = 1 + d + d^2 + \dots + d^{h-1} = \frac{d^{h+1} - 1}{d - 1}$$

Rooted, directed, ordered, binary trees



Tree ADT:

- insert
- remove
- traverse

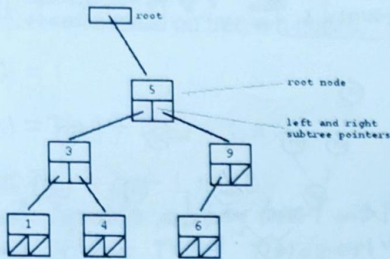
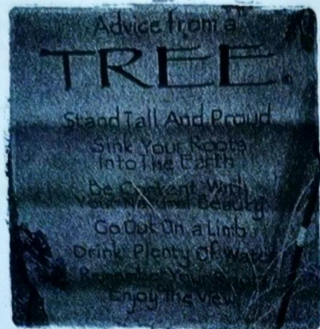
```

1  template<class T>
2  class tree {
3  public:
4      ...
5  private:
6      struct Node {
7          T data;
8          Node * left;
9          Node * right;
10     };
11     Node * root;
12     ...
13 };
    
```

Announcements:

HW2 due 2/14.

Rooted, directed, ordered, binary trees



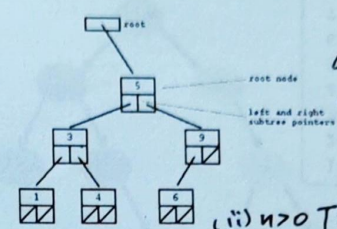
```

1  template<class T>
2  class tree {
3  public: // discuss later
4  private:
5      struct Node {
6          T data;
7          Node * left;
8          Node * right;
9      };
10     Node * root; // + more?
11 };
    
```

Tree ADT:

- insert
- remove
- traverse

$N(n)$: # of nulls in our implementation of a tree of size n
 Theorem: Our implementation of an n item binary tree has $n+1$ null pointers. $\forall n \geq 0$.

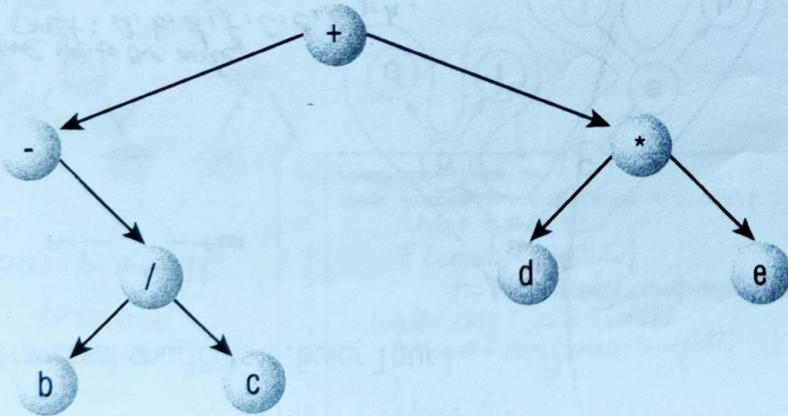


Consider an arbitrary binary tree implemented as on previous page. Let n denote the amount of data elements it holds.

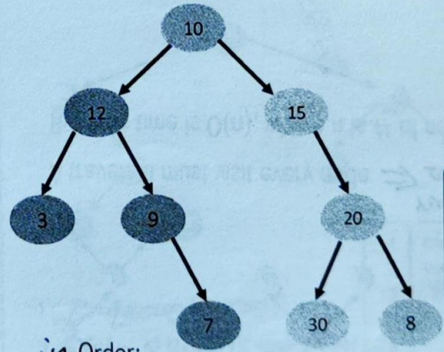
(i) $n=0$ $T=\{\}$, we implement with one NULL pointer which is $0+1 - N(0) = 1$ ✓.

(ii) $n>0$ $T = \{\text{root}, T_L, T_R\}$
 $N(n) = N(L) + N(R)$
 by IH $N(L) = l+1$ $N(R) = r+1$
 $N(n) = N(L) + N(R) = l+1+r+1 = n+1$.

Traversal... a scheme for visiting every node.



Traversal... a scheme for visiting every node.

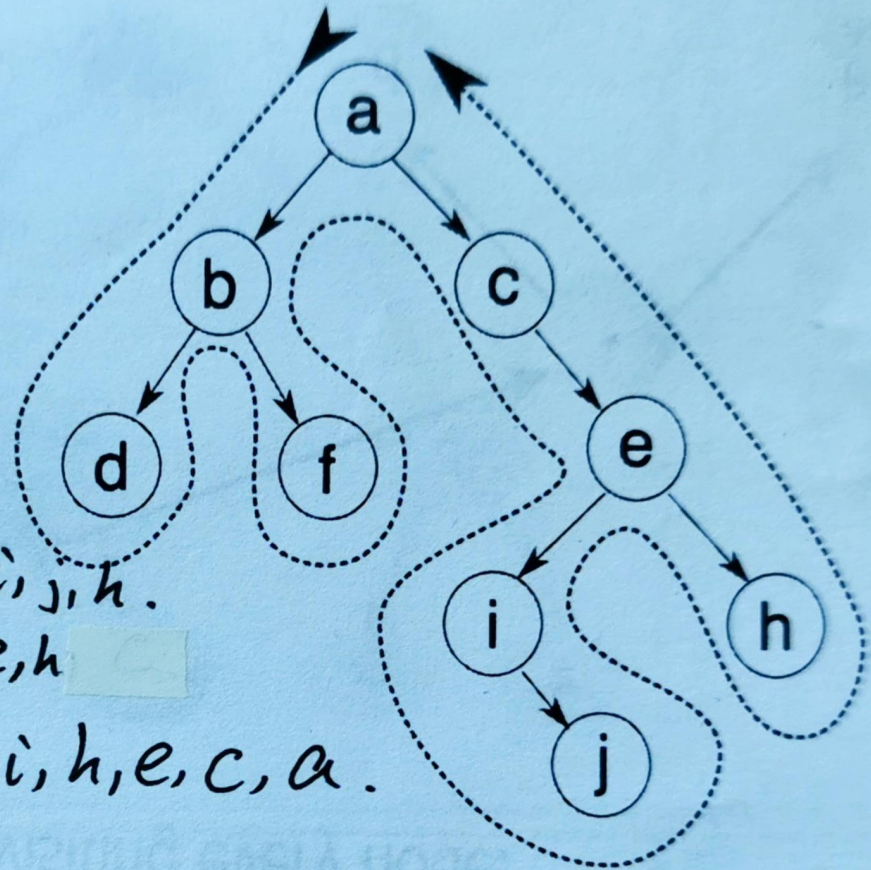
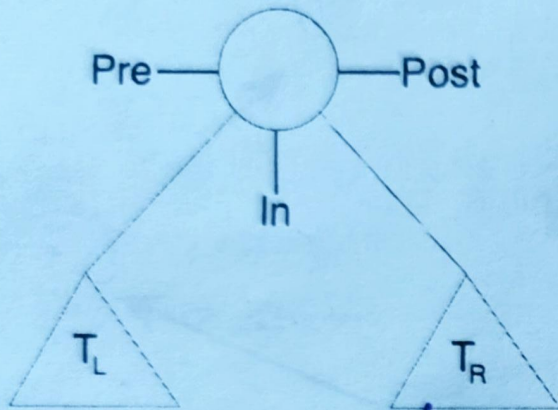


in Order: 3, 12, 9, 7, 10, 15, 30, 20, 8
 preOrder: 10, 12, 3, 9, 7, 15, 20, 30, 8
 postOrder: 3, 7, 9, 12, 30, 8, 20, 15, 10

```

1  template<class T>
2  void tree<T>::_Order
3      (Node * croot){
4      if (croot != NULL){
5          cout << croot->data;
6          _Order(croot->left);
7          cout << croot->data << " ";
8          _Order(croot->right);
9          cout << croot->data;
10     }
11 }
12 }
13 }
14 }
    
```

Traversal shortcuts... Euler Tour



first time got to the node

pre: a, b, d, f, c, e, i, j, h.

second time in: d, b, f, a, c, i, j, e, h

third time post: d, f, b, j, i, h, e, c, a.

Announcements:

HW2 due 2/14.

Today: traversals, ADT dict, BST

evaluation of if's with
 $T(n)$: # of traversals on tree with n nodes.

$T(0) = 1$

$T(n) = T(n_L) + T(n_R) + 1, n > 0.$

Thm: $T(n) = 2n + 1, n \geq 0.$

*Proof. Consider arbitrary tree T with $|T|=n$
 Base case: if $n=0$ $T(0)=1$ $T(0)=2 \cdot 0 + 1$*

if $n > 0$, $T(n) = T(n_L) + T(n_R) + 1$

$n = n_L + n_R + 1$

IH $\Rightarrow T(n_L) = 2n_L + 1$ $T(n_R) = 2n_R + 1.$

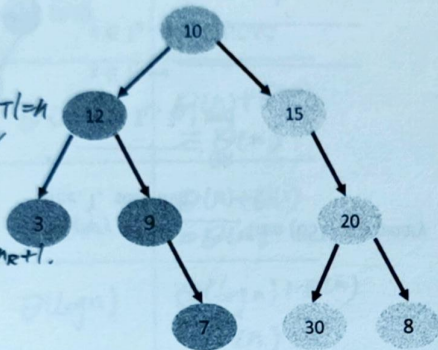
so $T(n) = 2n_L + 1 + 2n_R + 1 + 1$

$= 2(n_L + n_R + 1) + 1$

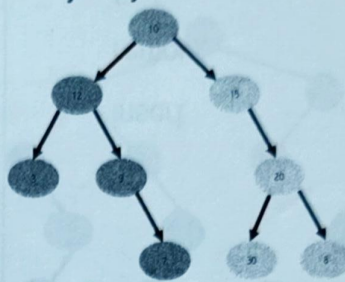
$= 2n + 1 \Rightarrow T(n) \in \Theta(n).$

```

1  template<class T>
2  void tree<T>::preOrder
3      (Node * croot){
4      if (croot != NULL){
5          cout << croot->data;
6          preOrder(croot->left);
7          preOrder(croot->right);
8      }
9  }
    
```



Mystery function: What's a good name?



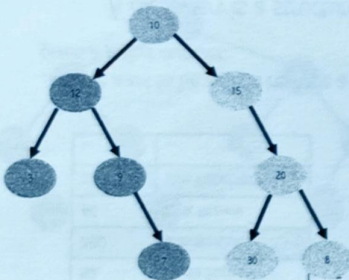
RT:

$\Theta(n)$ by analogy to post order traversal

```

1  template<class T>
2  void tree<T>::cleanUp (Node * croot){
3      if (croot != NULL) {
4          cleanUp (croot->left);
5          cleanUp (croot->right);
6          delete croot;
7          croot = NULL;
8      }
9  }
    
```

Traversals: a broader view...



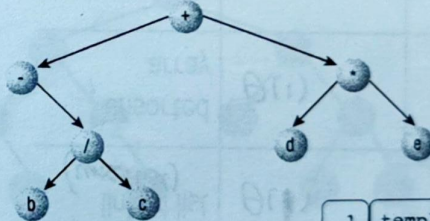
RT:

$\Theta(n)$ by analogy to traversal

```

1  template<class T>
2  Node * tree<T>::copy(Node * croot){
3      Node * t = NULL;
4      if (croot != NULL){
5          t = new Node(croot->data);
6          t->left = copy(croot->left);
7          t->right = copy(croot->right);
8      }
9      return t;
10 }
11 }
    
```

Traversals: something totally different...



~~abcde~~
 $+ - x / d e b c.$

$T(n) = \Theta(n)$

1. enqueue the root -
2. while queue is not empty
 - a) dequeue a node.
 - i. print its data
 - ii. enqueue its children

```

1  template<class T>
2  void tree<T>::levelOrder(Node * croot){
3      queue<Node> Q;
4      Q.enqueue(croot);
5      while (!Q.empty())
6          Node * t = Q.dequeue();
7          if (t != NULL)
8              cout << t->data;
9              Q.enqueue(t->left);
10             Q.enqueue(t->right);
11 }
12 }
    
```

2n+1 iterations

$\Theta(n)$ $\Theta(n)$

data <key, value>. Key: unique ID.

ADT Dictionary:

Suppose we have the following data...

ID #	Name
103	Jay Hathaway
92	Linda Stencel
330	Bonnie Cook
46	Rick Brown
124	Kim Petersen
...	...

...and we want to be able to retrieve a name, given a locker number.

More examples of key/value pairs:

CWL -> Advising Record

Course Number -> Schedule info

Color -> PNG

Vertex -> Set of incident edges

Flight number -> arrival information

URL -> html page

A *dictionary* is a structure supporting the following:

```
void insert(kType & k, dType & d)
```

```
void remove(kType & k)
```

```
dType find(kType & k)
```

Announcements:

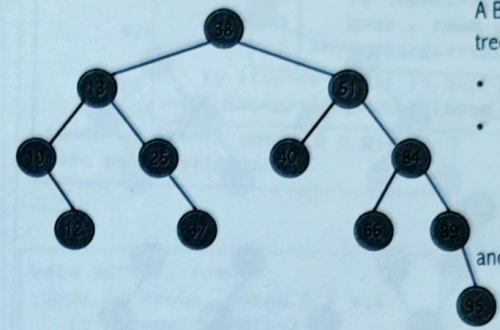
HW2 due 2/14.

Today: BST

ADT Dictionary: Implementation options...

	insert	find	(find +) remove
linked list	$\Theta(l)$	$\Theta(n)$	$\Theta(n) + \Theta(l)$ $= \Theta(n)$
unsorted array	$\Theta(l)$	$\Theta(n)$	$\Theta(n) + \Theta(l)$ $= \Theta(n)$
sorted array	$\Theta(\log n) + \Theta(n)$ <i>find open spot</i>	$\Theta(\log n)$	$\Theta(\log n) + \Theta(n)$ $= \Theta(n)$

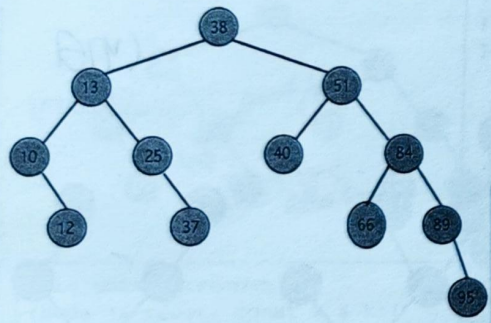
Binary Search Tree in-order traversal.



A Binary Search Tree (BST) is a binary tree, T , such that:

- $T = \{ \}$ empty
 - $T = \{ r, T_L, T_R \}$ and
 - $x \in T_L \rightarrow x.key < r.key$
 - $x \in T_R \rightarrow x.key > r.key$
 and $T_L + T_R$ are BST.
- keys are unique.

Dictionary ADT: (BST implementation)

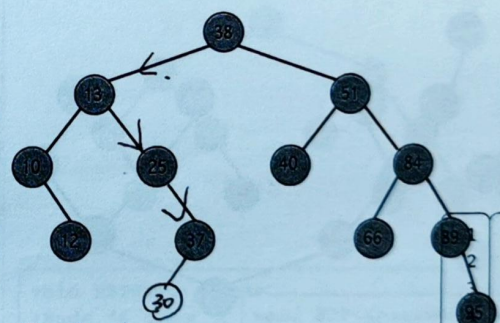


public: {
 insert
 remove
 find
 traverse $\Theta(n)$.

```

1 template<class K, class D>
2 class Dictionary{
3 public: //ctors etc
4 ...
5 private:
6     struct Node {
7         D data;
8         K key;
9         Node * left;
10        Node * right;
11    };
12    Node * root;
13    ...
14 };
    
```

Binary Search Tree - Insert insert(30);



Running time:
 Not like traversal
 $\Theta(h)$ where h is height of tree.

```

1 template<class K>
2 void BST::insert
3 (Node * & croot, const K & key){
4 ...
5     if (croot == NULL)
6         croot = new Node(key);
7 ...
8     else if (key < croot.key)
9         insert(croot->left);
10 ...
11     else if (key > croot.key)
12         insert(croot->right);
13 }
    
```


Announcements:

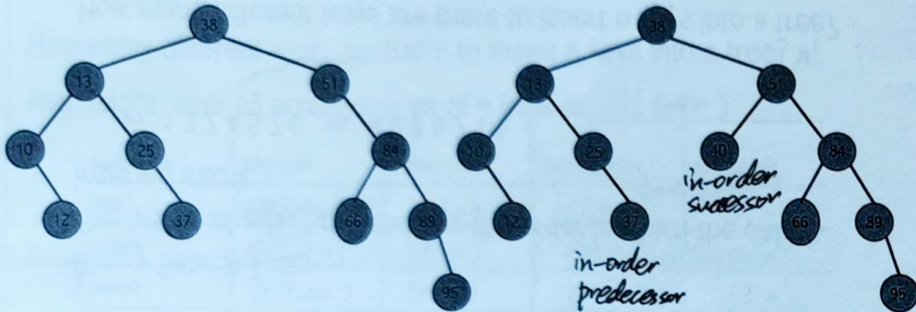
HW2 due 2/14.

Today: BST removal + AVL intro <https://cs221viz.netlify.com/src/bst>

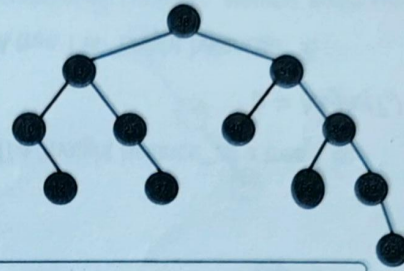
```
T.remove(37);
```

```
T.remove(51);
```

```
T.remove(38);
```



Binary Search Tree - Remove

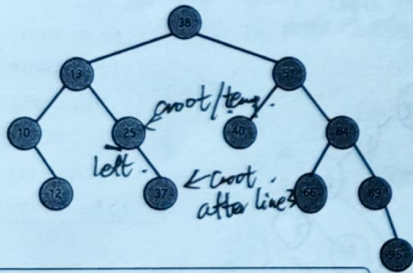


do Removal(croot);
if (croot->key < key)
 remove(croot->right);
else
 remove(croot->left);

```
void BST<K>::remove
(Node *& croot, const K & k){
    if (croot != NULL){
        if (croot->key == key)
            doRemoval(croot);
        else if (croot->key < key)
            remove(croot->right);
        else
            remove(croot->left);
    }
}
```

```
1 void BST<K>::doRemoval(Node * & cRoot) {
2     if ((croot->left != NULL) && ((croot->right != NULL)){
3         two ChildRemove(cRoot);
4     }
5     else
6         zeroOneChildRemove(cRoot);
}
```

Binary Search Tree - Remove



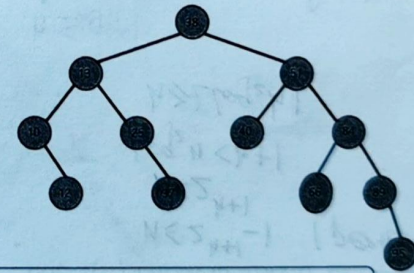
```
void BST<K>::remove
(Node *& croot, const K & k){
    if (croot != NULL)
        doRemoval(croot);
}

void BST<K>::doRemoval(Node * & cRoot) {
    if ((croot->left != NULL) && ((croot->right != NULL)){
        two ChildRemove(cRoot);
    }
    else if (croot->left != NULL)
        zeroOneChildRemove(cRoot);
    else if (croot->right != NULL)
        zeroOneChildRemove(cRoot);
    else
        delete cRoot;
}
```

```
1 void BST<K>::zeroOneChildRemove(Node * & cRoot) {
2     Node * temp = cRoot;
3     if (cRoot->left == NULL) cRoot = cRoot->right;
4     else cRoot = cRoot->left;
5     delete temp;
6 }
```

Binary Search Tree - Remove

running time: $O(h)$
proportional to the height.



```
void BST<K>::remove
(Node *& croot, const K & k){
    if (croot != NULL)
        doRemoval(croot);
}

void BST<K>::doRemoval(Node * & cRoot) {
    if ((croot->left != NULL) && ((croot->right != NULL)){
        two ChildRemove(cRoot);
    }
    else if (croot->left != NULL)
        zeroOneChildRemove(cRoot);
    else if (croot->right != NULL)
        zeroOneChildRemove(cRoot);
    else
        delete cRoot;
}
```

```
1 void BST<K>::twoChildRemove(Node * & cRoot) {
2     Node * iop = rightMostChild(cRoot->left);
3     cRoot->key = iop->key;
4     zeroOneChildRemove(iop);
5 }
```

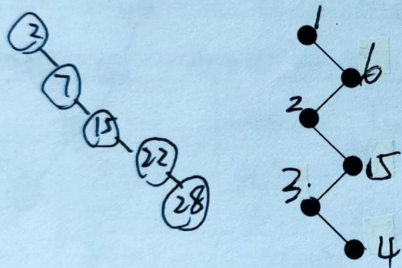
```
1 Node * & BST<K>::rightMostChild(Node * & cRoot) {
2     if (cRoot->right == NULL) return cRoot;
3     else return rightMostChild(cRoot->right);
4 }
```

ADT *dictionary* supports the following:

```
void insert(kType & k, dType & d)
dType find(kType & k)
void remove(kType & k)
```

Binary Search Tree implementation of a dictionary incurs running time $\theta(h)$ for all of these.

```
dict<int> myT;
myT.insert(2);
myT.insert(7);
myT.insert(15);
myT.insert(22);
myT.insert(28);
...
```



How many "bad" n-item trees are there?
 $n-1$ R or L choices, 2^{n-1}

The *algorithms* on BST depend on the height (h) of the tree.

The *analysis* should be in terms of the amt of data (n) the tree contains.

So we need relationships between h and n :

$h \geq f(n)$ (what is max # of nodes in tree of height h ?)

$$n \leq 2^{h+1} - 1 \quad (\text{perfect tree}).$$

$$n < 2^{h+1}$$

$$\log_2 n < h+1$$

$$h \geq \lceil \log_2 n \rceil$$

$$h \leq g(n)$$

$h \leq n-1$ ← room for improvement.

Good news, or bad?

Announcements:

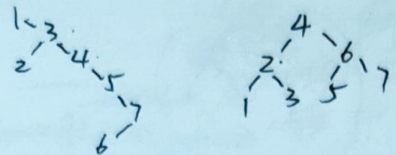
PA2 due 03/09. MT2 03/04.

Today: AVL <https://cs221viz.netlify.com/src/bst>

The height of a BST depends on the order in which the data is inserted.

ex. 1 3 2 4 5 7 6 vs. 4 2 3 6 7 1 5

$h=5$ vs. $h=2$

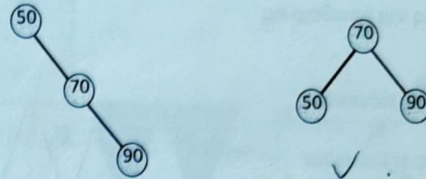


How many different ways are there to insert n keys into a tree? $n!$

Avg height, over all arrangements of n keys is $\Theta(\log n)$

operation	BST		sorted array	sorted list
	avg case	worst case		
find	$\Theta(\log n)$	$\Theta(n)$	$\Theta(\log n)$	$\Theta(n)$
insert	$\Theta(\log n)$	$\Theta(n)$	$\Theta(n)$	$\Theta(n)$
delete	$\Theta(\log n)$	$\Theta(n)$	$\Theta(n)$	$\Theta(n)$
traverse	$\Theta(n)$	$\Theta(n)$	$\Theta(n)$	$\Theta(n)$

something new... which tree makes you happiest?



The "height balance" of a tree T is:

$$b = \text{height}(T_R) - \text{height}(T_L)$$

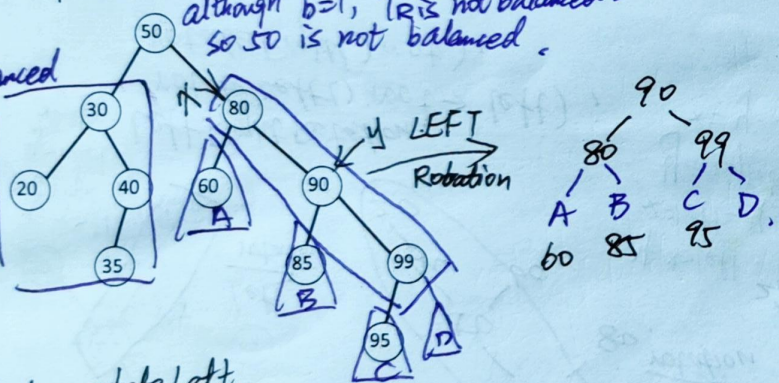
A tree T is "height balanced" if:

- $T = \{\}$ empty tree
- $T = \{r, T_L, T_R\}$ and $|b| \leq 1$ and T_L & T_R are height balanced. } if $b=0$ perfect tree.

operations on BST - rotations

although $b=1$, T is not balanced. so 50 is not balanced.

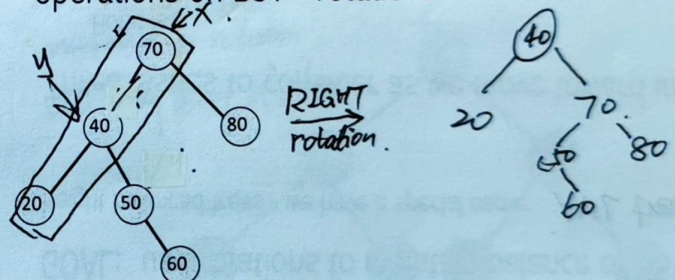
balanced



function rotateLeft
 $y = x \rightarrow \text{right};$
 $x \rightarrow \text{right} = y \rightarrow \text{left};$
 $y \rightarrow \text{left} = x;$
 $x = y;$

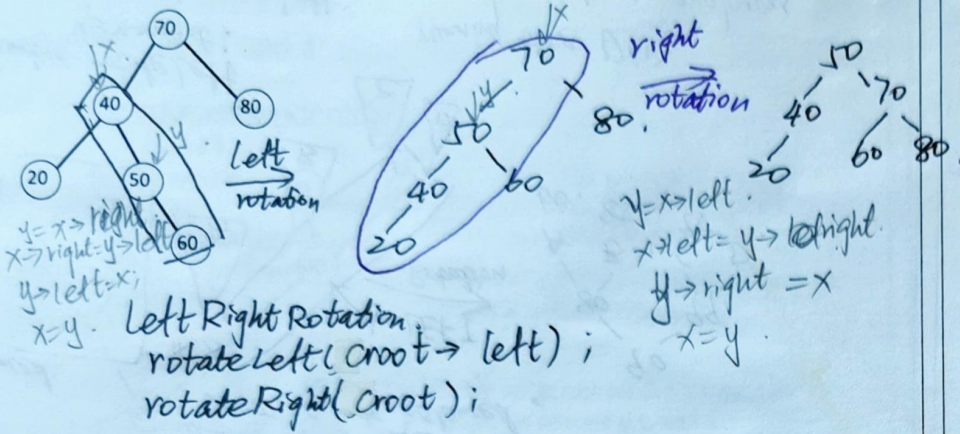
running time $\Theta(1)$
 right rotation is analogous
 $y = x \rightarrow \text{left};$
 $x \rightarrow \text{left} = y \rightarrow \text{right};$
 $y \rightarrow \text{right} = x;$
 $x = y;$

operations on BST - rotations



function rotateRight
 $y = x \rightarrow \text{left};$
 $x \rightarrow \text{left} = y \rightarrow \text{right};$
 $y \rightarrow \text{right} = x;$
 $x = y;$

operations on BST - rotations



balanced trees - rotations summary:

- there are 4 kinds: left, right, left-right, right-left (symmetric!)
- local operations (subtrees not affected)
- constant time operations
- BST characteristic maintained

GOAL: use rotations to maintain balance of BSTs.

height balanced trees - we have a special name: *AVL trees*.

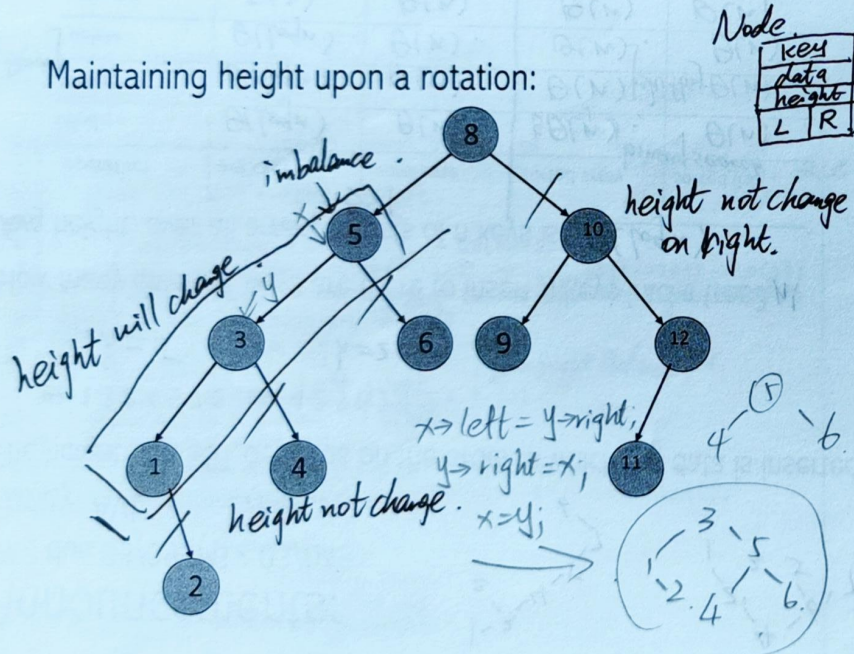
Three issues to consider as we move toward implementation:

Rotating

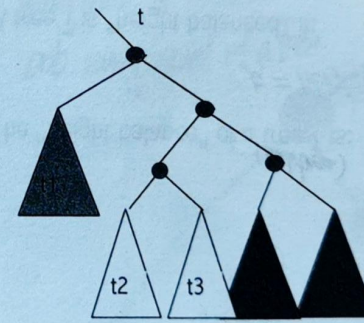
Maintaining height

Detecting imbalance *diagnosing rotations.*

Maintaining height upon a rotation:



AVL trees: rotations (identifying the need)



if an insertion was in subtrees t_4 or t_5 , and if an imbalance is detected at t , then a _____ rotation about t rebalances the tree.

We diagnose this by noting that the balance factor at $t \rightarrow \text{right}$ is _____

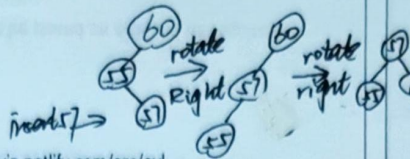
If an insertion was in subtrees t_2 or t_3 , and if an imbalance is detected at t , then a _____ rotation about t rebalances the tree.

We diagnose this by noting that the balance factor at $t \rightarrow \text{right}$ is _____

Announcements:

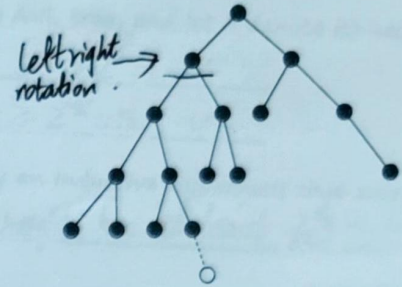
PA2 due 03/09. MT2 03/04.

Today: AVL insert and analysis <https://cs221viz.netlify.com/src/avl>



AVL tree insert:

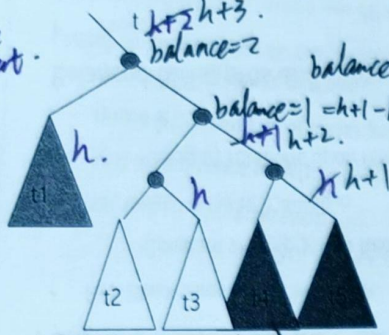
```
struct treeNode {
    T key;
    int height;
    treeNode * left;
    treeNode * right;
};
```



Insert:

- insert at proper place
- check for imbalance
- rotate if necessary
- update height

AVL trees: rotations (identifying the need)



if an insertion was in subtrees t4 or t5, and if an imbalance is detected at t, then a left rotation about t rebalances the tree.

We diagnose this by noting that the balance factor at $t \rightarrow$ right is 1. $(h+1-h)$.

If an insertion was in subtrees t2 or t3, and if an imbalance is detected at t, then a right left rotation about t rebalances the tree.

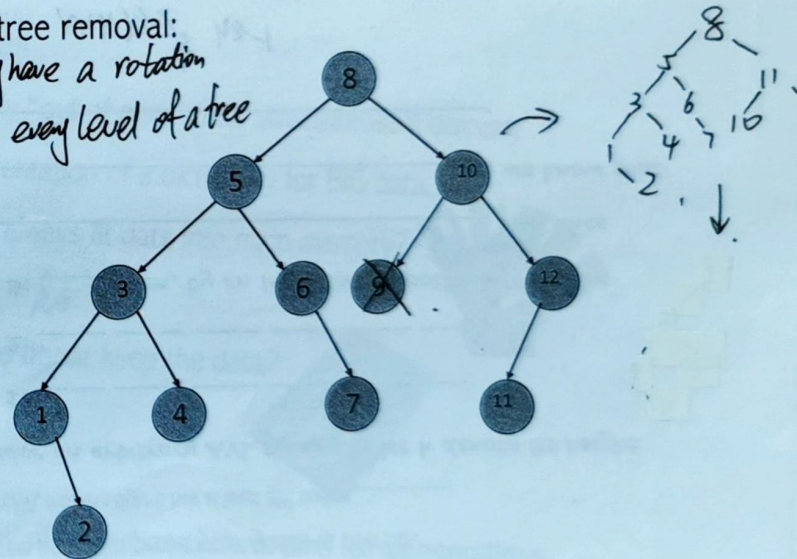
We diagnose this by noting that the balance factor at $t \rightarrow$ right is -1. $(h-(h+1))$.

if height(t_i) = $h+1$ then adding one node to the right (making $h+1$ to $h+2$) will not change the structure \rightarrow right is -1 $(h-(h+1))$

AVL tree insertions: time = $\Theta(n)$

```
1 template <class T>
2 void AVLTree<T>::insert(const T & x, treeNode<T> * & t){
3     if( t == NULL ) t = new treeNode<T>( x, 0, NULL, NULL); // Base case
4     else if( x < t->key ){ // left
5         insert( x, t->left );
6         int balance = ht(t->rt) - ht(t->left);
7         int leftBalance = ht(t->left->rt) - ht(t->left->left);
8         if( balance == -2 ) // long on the left -> right rotation.
9             if( leftBalance == -1 )
10                rotate right ( t );
11            else
12                rotate left right ( t );
13     }
14     else if( x > t->key ){ // right
15         insert( x, t->rt );
16         int balance = ht(t->rt) - ht(t->left);
17         int rtBalance = ht(t->rt->rt) - ht(t->rt->left);
18         if( balance == 2 )
19             if( rtBalance == 1 )
20                rotate left ( t );
21            else
22                rotate right left ( t );
23     }
24     t->ht = max(ht(t->left), ht(t->rt)) + 1;
25 }
```

AVL tree removal: may have a rotation at every level of a tree



Announcements:

PA2 due 03/09. MT2 03/04.

Today: AVL analysis, B-Trees <https://cs221viz.netlify.com/st/avl> *we already know h is $\sim \log n$*

AVL tree analysis:

Since running times for Insert, Remove and Find are $O(h)$, we'll argue that $h = O(\log n)$.

Putting an upper bound on the height for a tree of n nodes is the same as putting a lower bound on the number of nodes in a tree of height h .

• Define $N(h)$: least # of nodes in AVL tree of ht h

• Find a recurrence for $N(h)$: $N(h) = 1 + N(h-1) + N(h-2)$

• We simplify the recurrence:

$$N(h) \geq 2N(h-2)$$

• Solve the recurrence: (guess a closed form)

$$N(h) = 2^{\frac{h}{2}} \quad h \geq 0$$

$$\begin{aligned} N(-1) &= 0 \\ N(0) &= 1 \\ N(1) &= 2 \end{aligned}$$

AVL tree analysis: prove your guess is correct.

Thm: An AVL tree of height h has at least $2^{h/2}$ nodes, $h \geq 0$.

Consider an arbitrary AVL tree, and let h denote its height.

Case 1: $h=0$ $N(0)=1 \geq 2^{0/2}=1$ ✓

Case 2: $h=1$ $N(1)=2 \geq 2^{1/2}=\sqrt{2}$ ✓

Case 3: $h > 1$ then, by an Inductive Hypothesis that says $\forall j < h$, an AVL tree of height j has at least $2^{j/2}$ nodes, and since

$N(h) \geq 2N(h-2)$, we know that $N(h) \geq 2 \cdot 2^{(h-2)/2} = 2^{h/2}$

Let n denote the actual # of nodes of an AVL tree of height h .

Punchline:

$$n \geq N(h) \geq 2^{h/2} \quad h \geq 0$$

$$h \leq 2 \log_2 n, \quad n \geq 1 \quad h = O(\log n)$$

Classic balanced BST structures:

• Red-Black trees – max ht $2 \log_2 n$.

Constant # of rotations for insert, remove, find.

• AVL trees – max ht $1.44 \log_2 n$.

$O(\log n)$ rotations upon remove.

*no rotation on a find.
1 rotation on insert.*

$O(\log n)$ rotation on remove.

Balanced BSTs, pros and cons:

• Pros:

- Insert, Remove, and Find are always $O(\log n)$

- An improvement over: *BST, linked list, arrays*

- Range finding & nearest neighbor

• Cons:

- Possible to search for single keys faster *via hashing*

- If data is so big that it doesn't fit in memory it must be stored on disk and we require a different structure.

B-Trees

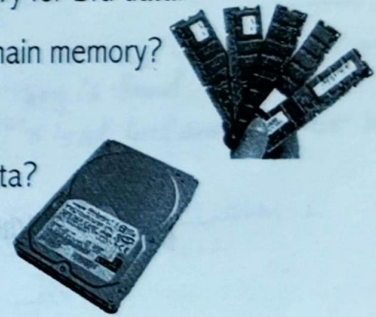
The only "out of core" data structure we'll discuss.

Implementation of a dictionary for BIG data.

Can we always fit data into main memory?

No.

So where do we keep the data?



Big-O analysis assumes uniform time for all operations.

But... *remote data access is slower.*

The Story on Disks

2GHz machine gives around 2m instructions per millisecond.

Seek time around 5ms for a current hard disk.

Imagine an AVL tree storing North American driving records.

How many records? 200,000,000

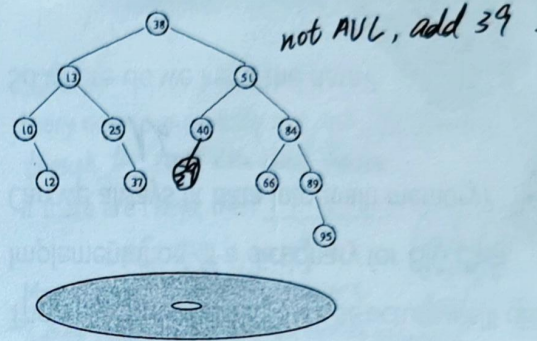
How much data, in total? 200m x 1MB = 2TB.

How deep is the AVL tree? $\log_2 200m \approx 56$.

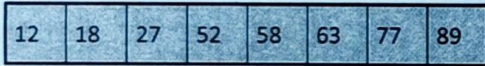
How many disk seeks to find a record?

The Story on Disks, continued.

Suppose we weren't careful...

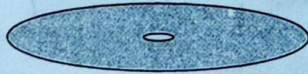


B Tree of order m



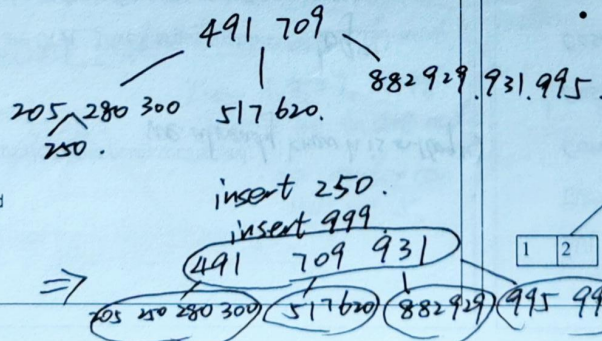
has capacity $m-1$.

B tree node.
contains $\lceil \frac{m}{2}, m-1 \rceil$ keys.
except the root.
(which may have only 1 key)



Goal: Minimize the number of reads from disk

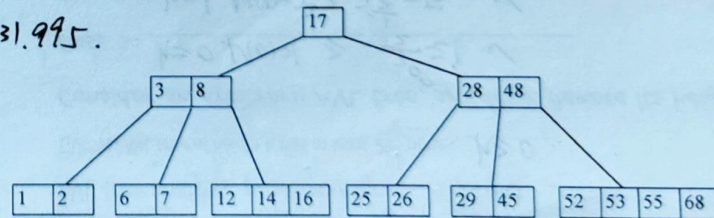
- Build a tree that uses 1 disk block per node
 - Disk block is the fundamental unit of transfer
- Nodes will have more than 1 key
- Tree should be balanced and shallow
 - In practice branching factors over 1000 often used



Definition of a B-tree

B-tree of order m is an m -way tree

- For an internal node, # keys = #children - 1
- All leaves are on the same level
- All leaves hold no more than $m-1$ keys
- All non-root internal nodes have between $\lceil m/2 \rceil$ and m children
- Root can be a leaf or have between 2 and m children.
- Keys in a node are ordered.



Announcements:

PA2 due 03/09. MT2 03/04.

Today: Btrees

BTree of order m

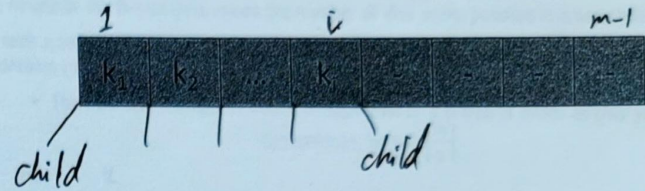
12	18	27	52	58	63	77	89
----	----	----	----	----	----	----	----

Goal: Minimize the number of reads from disk *or remote sources*.



- Build a tree that uses 1 disk block per node
 - Disk block is the fundamental unit of transfer
- Nodes will have more than 1 key
- Tree should be balanced and shallow
 - In practice, m over 1000 often used!!

B Tree of order m : Nodes



Every node has capacity $m-1$, but every node contains at least $\lceil \frac{m}{2} \rceil - 1$ keys.
though the root can hold fewer.

If there are i keys, then $i+1$ child pointers.

$i \in [\lceil \frac{m}{2} \rceil - 1, m-1]$ *least # of children for a non-leaf node*

May also maintain sibling pointers.

greatest #: m

Level order pointers, expedites removal.

Note: *root & leaves have slightly different constraints.*

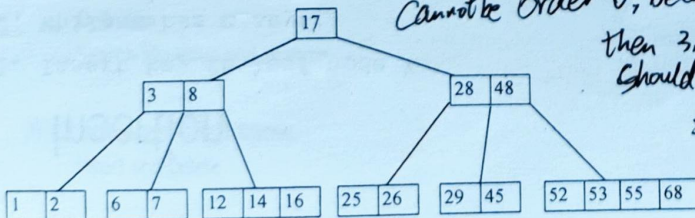
Definition

B-tree of order m is an m -way tree

- For an *non-leaf* internal node, # keys = #children - 1
- All leaves are on the same level
- All leaves hold no more than $m-1$ keys
- All non-root internal nodes have between $\lceil m/2 \rceil$ and m children
- Root can be a leaf or have between 2 and m children.
- Keys in a node are ordered.

Using binary search here. → change to $\Theta(\log m)$ from $\Theta(m)$. but does not change too much, since disk-read is dominant factor.

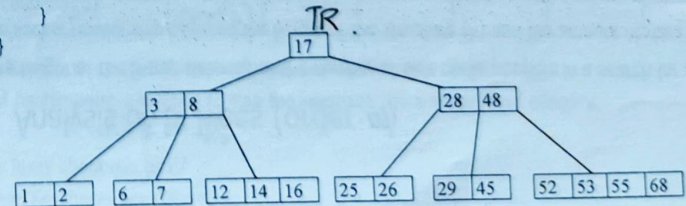
Cannot be order 6, because if $m=6$, then 3, 8, 17, 28, 48 should be in one node so order can only be 5.



Search

```
bool B-TREE-SEARCH(BTreeNode & x, T key) {
    int i = 0;
    while ((i < x.numkeys) && (key > x.key[i])) } find the correct child i
        i++;
    if ((i < x.numkeys) && (key == x.key[i])) key is found.
        return true;
    if (x.leaf == true) x is at leaf and does not have the key.
        return false;
    else {
        BTreeNode b = DISK-READ(x.child[i]);
        return B-TREE-SEARCH(b, key);
    }
}
```

Search(TR, 5)



Insertion

1. Insert key in leaf node X
2. While X has m keys:

Split X into two nodes:

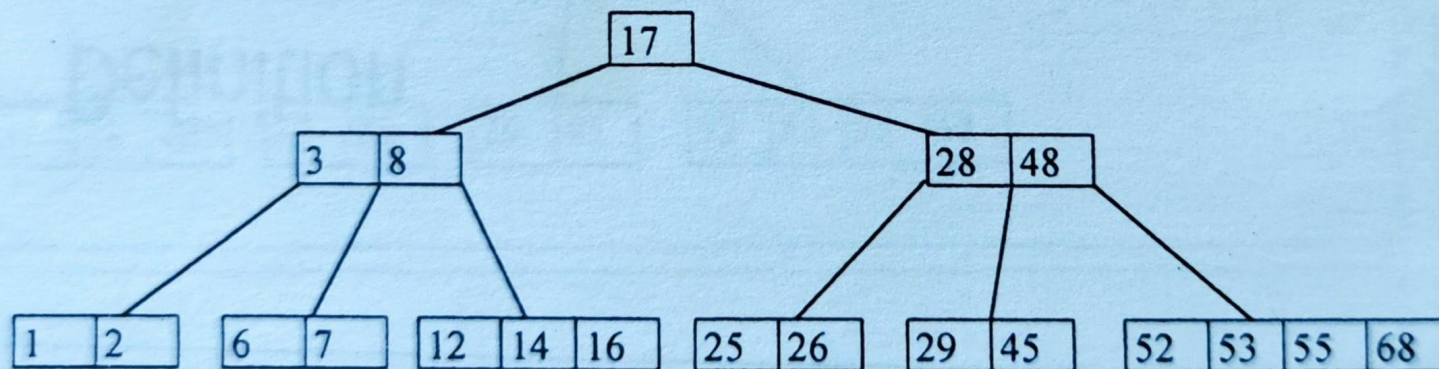
orig has $m/2$ smallest keys

new has $m/2$ largest keys

Insert middle key into parent & attach new node

(If X is root, create new node and attach both)

Set X to be the parent.



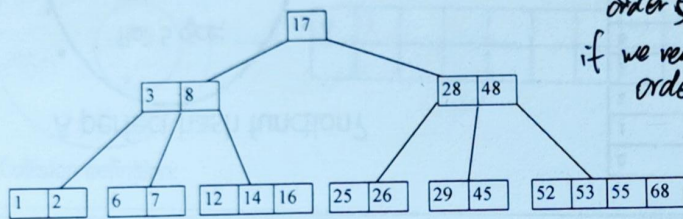
Announcements:

PA2 due 03/09.

Today: finish Btrees, start something new!!

B-tree of order m is an m -way tree

- For an internal node, # keys = #children - 1
- All leaves are on the same level
- All leaves hold no more than $m-1$ keys
- All non-root internal nodes have between $\lceil m/2 \rceil$ and m children
- Root can be a leaf or have between 2 and m children.
- Keys in a node are ordered.



height of B-tree: $\log_m \frac{n+1}{2}$
 $= O(\log_m n)$
 disk seek = height.

if only insert, it can only be order 5,
 if we remove, it maybe order 6.

Analysis of B-Trees (order m)

The height of the B-tree determines the number of disk seeks possible in a search for a key. We seek a relationship between the height of the structure (h) and the amount of data it contains (n).

- The minimum number of nodes in each level of a B-tree of order m . (For your convenience, let $t = \lceil \frac{m}{2} \rceil$)

root	1
level 1	2
level 2	$2t$
...	
level h	$2t^{h-1}$

- The total number of nodes is the sum of these:

$$1 + 2 \sum_{i=0}^{h-1} t^i = 1 + 2 \cdot \frac{t^h - 1}{t - 1}$$

- So, the least total number of keys is:

$$1 + 2 \cdot \left(\frac{t^h - 1}{t - 1} \right) \cdot (t - 1) = 2t^h - 1$$

$n \geq 2t^h - 1$
 $h \approx \log_{\frac{m}{2}}(n+1)$
 $h = O(\log n)$

$= O(\log n)$
 since $m(t)$ is constant.

Summary

B-Tree search:

$O(\log m)$ instead with binary search.
 $O(m)$ time per node to find key or appropriate child.
 $O(\log_m n)$ height implies $O(m \log_m n)$ total time
 BUT: m is constant, search is $O(\log n)$.

Insert and Delete have similar stories.

What you should know:

- Motivation
- Definition
- Search algorithm and analysis

What you should not know:

- Insert and Delete

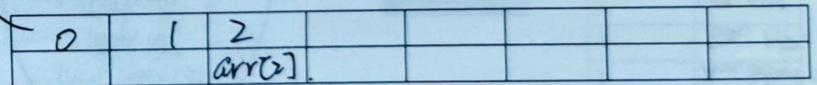
Hashing - using "hash tables" to implement dictionaries.

Structure of a dictionary:

- Key -> Value
- Locker # -> student
- Course Number -> Schedule info
- Vertex -> Set of incident edges
- URL -> html page
- dice roll -> payoff amt

Associative Array:

Dictionary with a particular interface
 Overloads operator [] for insert and find
`myDict["Miguel"] = 22;`
`int d = myDict["Miguel"];`



(defn) Keyspace - a (math) description of the keys for a set of data. K

Goal of hashing: use a function to map the keyspace into a small set of integers.

$$h: K \rightarrow \mathbb{Z}_0$$

What's fuzzy about this goal?

Problem: Keyspaces are often large...

Hashing - using "hash tables" to implement dictionaries.

$= O(\log n)$
since $n(t)$ is
constant.

Structure of a dictionary:

Key -> Value

Locker # -> student

Course Number -> Schedule info

Vertex -> Set of incident edges

URL -> html page

dice roll -> payoff amt

Associative Array:

Dictionary with a particular interface

Overloads operator[] for insert and find

```
myDict["Miguel"] = 22;
```

```
int d = myDict["Miguel"];
```

0	1	2					
		arr[2]					

(defn) Keyspace - a (math) description of the keys for a set of data. K

Goal of hashing: use a function to map the keyspace into a small set of integers.

$$h: K \rightarrow \mathbb{Z}_0$$

What's fuzzy about this goal?

Problem: Keyspaces are often large...

Overview:

client code

declares an object of ADT dictionary

```
dict<ktype, vtype> d;
```

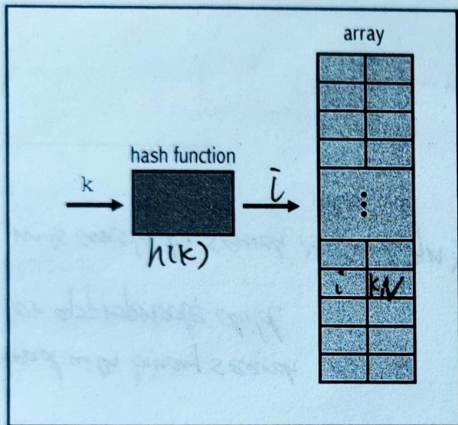
```
d.operator[](k);
```

```
ex: insert is d[k] = v;
```

A Hash Table consists of:

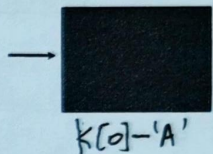
- hash function .
- An array
- CRS

class dict (implemented with a hash table)



A perfect hash function:

(Ann, black cat)
(Ben, HP)
(Cory, spy)
(David, bball player)
(Ellen, butterfly)
(Finn, cereal killer)
(Gus, ghost)
(Harmony, bee)



$k[0] - 'A'$

0	Ann, black cat.
1	Ben, HP.
2	Cory, SPY
3	⋮
4	⋮
5	
6	
7	

A contrived example: perfect hash function - a bijection
these keys have a fabulous hash fn.

- a. each key hashes to a different int injective.
- b. collection of keys hash to a seq of ints surjectivity.

Our purpose is to investigate
general purpose hash function.

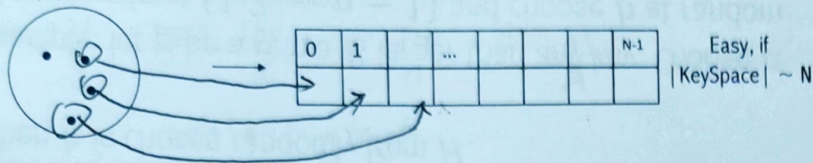
Announcements:

PA2 due today! PA3 out soon!
 Today: Hash functions, continued
 General Purpose Hash Functions

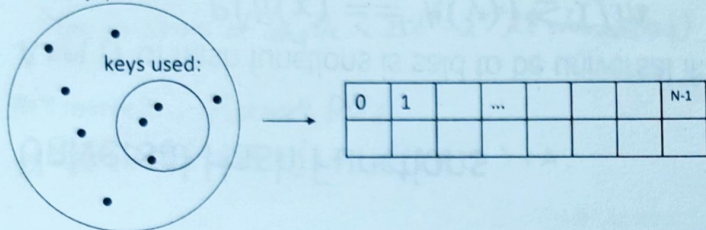
What characterizes a good hash function?

- $h(k)$ computes in $\Theta(1)$
- $k_1 = k_2 \Rightarrow h(k_1) = h(k_2)$
- $k_1 \neq k_2 \Rightarrow P(h(k_1) = h(k_2)) = \frac{1}{N}$
 where N is the table size.

KeySpace



KeySpace



Collision definition:

Expected Value *discrete, finite.* *weighted average of possible values.*

How many pips do we expect to see on a die?

Definitions that will help us:

X : # of pips. $X \in \{1, 2, 3, 4, 5, 6\}$

$$E[X] = \frac{1}{6} \cdot 1 + \frac{1}{6} \cdot 2 + \dots + \frac{1}{6} \cdot 6 = 3.5$$

How many pips do we expect to see on 2 dice?

X, Y : independent dice.

$$E[X + Y] = \frac{1}{36} \cdot 2 + \frac{2}{36} \cdot 3 + \dots + \frac{1}{36} \cdot 12 = 7$$

Linearity of Expectation:

$$E[X + Y] = E[X] + E[Y]$$

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

Probability

$$P(\text{one random student's birthday is not today}) = \frac{364}{365} = 1 - \frac{1}{365}$$

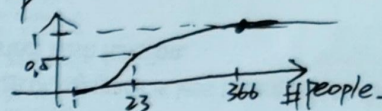
What's the probability that someone in this room has a birthday today?

$$1 - P(\text{no one has a birthday today}) = 1 - \left(\frac{364}{365}\right)^{\# \text{ of people}}$$

What's the probability that there are 2 people with the same birthday?

$$1 - P(\text{all on different days}) = 1 - \left(\frac{365}{365} \cdot \frac{364}{365} \cdot \frac{363}{365} \dots\right)^K$$

$$= 1 - \frac{P_{365}^K}{365^K}$$



Devise and analyze an algorithm to find the people with the same bday...

index on days and place people on table.

Expected # of collisions

What is expected # of people who share a bday w someone else?

Definitions that will help us:

X_{ij} : 1 if student i has the same birthday as student j , 0 otherwise.

$E[X_{ij}]$: $P(x_i \text{ and } x_j \text{ has same birthday}) \cdot 1 + P(\dots) \cdot 0 = \text{first term}$

Sum over all X_{ij} are the number of shared bdays:

$$E[X] = \sum_{i=0}^{K-1} \sum_{j=i+1}^{K-1} E[X_{ij}] = \sum_{i=0}^{K-1} \sum_{j=i+1}^{K-1} P(i \text{ and } j \text{ have the same birthday})$$

$$= \frac{1}{365} \sum_{i=0}^{K-1} \sum_{j=i+1}^{K-1} 1$$

$$= \frac{1}{365} \cdot \frac{K(K-1)}{2}$$

Collisions: if we randomly put k items into m bins, we expect $\frac{1}{N} \cdot \frac{K(K-1)}{2}$

pairs to collide. Implication: $\frac{1}{N} \cdot \frac{K(K-1)}{2} \geq 1$

so if # keys $K \geq \sqrt{2N}$ there will be a collision.

N is table size.

Announcements:

PA3 out soon!

Today: Hash functions, continued

Mapping strings to integers

How can we map 'Andy' to a number?

$$256^3 \frac{121}{y} + 256^2 \frac{160}{d} + 256^1 \frac{110}{n} + 256^0 \frac{65}{A}$$

$256^4 \times \text{string length}$

Is this a good mapping scheme?

Yes for strings of length $< 256^8 = 2^{64}$ (8 characters).

We'll rewrite it... Horner's Rule.

$$256(256(256(y) + d) + n) + A.$$

Hashing Strings (an example)

Given: 8 character strings are easy to hash

The idea: Select 8 random positions from long strings and hash that substring.

A bunch of strings:

Lookyhere, Huck, being rich ain't going No! Oh, good-licks; are you in real dead Just as dead earnest as I'm sitting here nto the gang if you ain't respectable, y Can't let me in, Tom? Didn't you let me Yes, but that's different. A robber is m irate is -- as a general thing. In most Now, Tom, hain't you always ben friendly ut, would you, Tom? You wouldn't do that Huck, I wouldn't want to, and I DON'T wa ay? Why, they'd say, 'Mph! Tom Sawyer's t!' They'd mean you, Huck. You wouldn't uck was silent for some time, engaged in Well, I'll go back to the widder for a m can come to stand it, if you'll let me All right, Huck, it's a whiz! Come along Will you, Tom -- now will you? That's go he roughest things, I'll smoke private a hrough or bust. When you going to start Oh, right off. We'll get the boys together

Hashing Strings (an example)

Given: 8 character strings are easy to hash

The idea: Select 8 random positions from long strings and hash that substring.

A bunch of strings:

http://en.wikipedia.org/wiki/Le%C5%9Bna_Grobla
http://en.wikipedia.org/wiki/Blow_the_Man_Down
http://en.wikipedia.org/wiki/Swen_K%C3%B6nig
[http://en.wikipedia.org/wiki/2/7th_Cavalry_Commando_Regiment_\(Australia\)](http://en.wikipedia.org/wiki/2/7th_Cavalry_Commando_Regiment_(Australia))
http://en.wikipedia.org/wiki/Salman_Ebrahim_Mohamed_Ali_Al_Khalifa
http://en.wikipedia.org/wiki/Alice_High_School
http://en.wikipedia.org/wiki/Beautiful_Dirty_Rich
[http://en.wikipedia.org/wiki/RFA_Sir_Bedivere_\(L3004\)](http://en.wikipedia.org/wiki/RFA_Sir_Bedivere_(L3004))
[http://en.wikipedia.org/wiki/BirThright_\(band\)](http://en.wikipedia.org/wiki/BirThright_(band))
http://en.wikipedia.org/wiki/Jacky_Vimond
<http://en.wikipedia.org/wiki/Vachon>
http://en.wikipedia.org/wiki/McCarthy_426_Stone
http://en.wikipedia.org/wiki/Salisbury,_New_Hampshire
http://en.wikipedia.org/wiki/A_Line_of_Deathless_Kings
http://en.wikipedia.org/wiki/Newfoundland_Irish
http://en.wikipedia.org/wiki/Beatrice_Politi
http://en.wikipedia.org/wiki/Bona_Sijabat
http://en.wikipedia.org/wiki/Sour_sanding
http://en.wikipedia.org/wiki/Dr_Manmohan_Singh_Scholarship
http://en.wikipedia.org/wiki/Religion_in_Jordan

Hashing strings

```
1 int hash ( string s, int p) {  
2     int h = 0;  
3     for ( i = s.length-1; i >= 0; i-- )  
4         h = (256 * h + s[i]) % p;  
5     return h;
```

Horner's Rule.

Running time?

Problem: Suppose $\text{ascii}(A) = 61$

$$\text{hash}('A', 61) = 0.$$

$$\text{hash}('AA', 61) = (256 \times 61 + 61) \% 61 = 0.$$

$$\text{hash}('AAA', 61) = (256^2 \times 61 + 256 \times 61 + 61) \% 61 = 0.$$

Solution:

Keep hash function secret by creating a random hash function.

Universal Hash Functions

A set H of hash functions is said to be universal if:

$$P(h(x) == h(y)) \leq \frac{1}{m}, \quad m \text{ is the table size.}$$

When h is chosen randomly from H . (satisfies SUHA)

Example: let p be a prime # larger than any key.

- Choose a at random from $\{1, 2, \dots, p-1\}$ and
- choose b at random from $\{0, 1, \dots, p-1\}$ then let

$$h(x) = ((ax + b) \bmod p) \bmod m$$

Find the largest set of keys that collide:

$$\textcircled{1} h(x) = (3x + 2) \bmod 9$$

keys: {2, 5, 8}

$$\textcircled{2} h(x) = (3x + 2) \bmod 11$$

keys: x=0:2, x=1:5, x=2:8, x=3:0, x=4:3, x=5:6, x=6:9, x=7:1, x=8:4, x=9:7, x=10:10, x=1:2.

Observations:
 ① spreads out the data.
 ② can still map infinitely many keys to the same int.

explores the entire space.

Collision resolution -

SUHA: simple universal hashing assumption

Separate Chaining: (an example of open hashing) allow collisions.

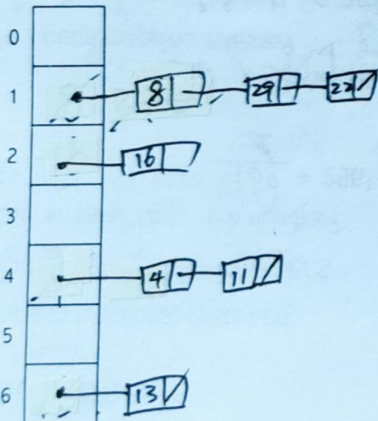
$$S = \{16, 8, 4, 13, 29, 11, 22\}$$

$$|S| = n \quad h(k) = k \% 7$$

of keys

Load factor:

$$\alpha = \frac{n}{m} \left(\frac{\# \text{ keys}}{\text{table size}} \right)$$



	Worst case	Under SUHA
Insert	$\Theta(n)$	α
Remove/find	$\Theta(n)$	α

$\alpha > 0$. unconstrained. under separate chaining.

Hash function summary

We can only avoid collisions if the size of keyspace is less than equal to the size of our hash table and we have a perfect hash function.

m is the table size.

We have to deal with collisions: even with only \sqrt{m} keys we will expect one.

if $m=10,000$, we expect a collision with only 100 keys.

We need a collision resolution strategy

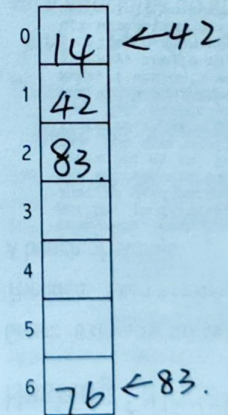
- Separate chaining
- Open addressing
- Other strategies we won't talk about

Collision Handling - Probe based hashing: (example of closed hashing)

$$S = \{76, 14, 42, 83, 11, 22\}$$

$$|S| = n \quad h(k) = k \% 7$$

of failed attempts



Try $H(k,0) = (h(k) + 0) \% 7$. If full...

try $H(k,1) = (h(k) + 1) \% 7$. If full...

try $H(k,2) = (h(k) + 2) \% 7$. If full...

try...

find: keep looking until an empty block. if remove (42) then find (83), it cannot find it. - use flag { never occupied 0, once occupied, now empty } find. - [now occupied . 2]

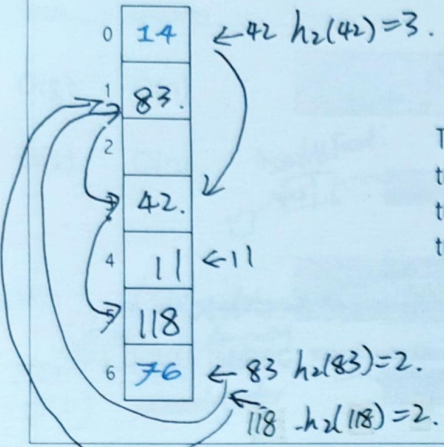
Announcements:

PA3 out, due 03/30!

Today: Hash functions, fin. Something new!! <https://cs221viz.netlify.com/src/linearhashing.html>

Collision Handling - Probe based hashing: (another example of closed hashing)

Double hashing (two hash functions)
 $S = \{76, 14, 42, 83, 11, 22\}$ $|S| = n$



$$H(k,i) = h_1(k) + i \cdot h_2(k) \cdot P$$

$h_1(k) = k \% 7$ initial hash
 $h_2(k) = (5 - k \% 5) \% 7$ step 5

Try $H(k,0) = (h(k) + 0 \cdot h_2(k)) \% 7$. If full...
 try $H(k,1) = (h(k) + 1 \cdot h_2(k)) \% 7$. If full...
 try $H(k,2) = (h(k) + 2 \cdot h_2(k)) \% 7$. If full...
 try...

P, Q are primes.
 P is the table size.
 $2 < Q < P$.

Hash table performance: expected # of probes for Find(key) under SUHA

Linear probing -

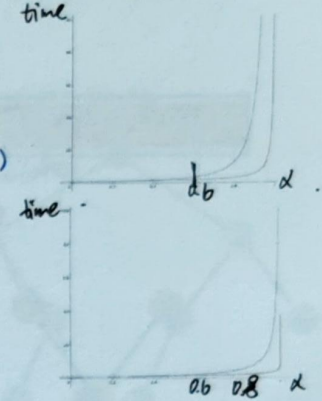
successful: $\frac{1}{2} (1 + 1/(1-\alpha))$
 unsuccessful: $\frac{1}{2} (1 + 1/(1-\alpha)^2)$

Double hashing -

successful: $1/\alpha \ln 1/(1-\alpha)$
 unsuccessful: $1/(1-\alpha)$

Separate chaining -

successful: $1 + \alpha/2$
 unsuccessful: $1 + \alpha$



Do not memorize!

Observe: $\alpha = \frac{n}{m}$.

- As α increases, running times increase!
- If α is held constant then running times are constant!

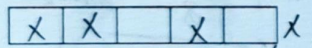
α can be greater than 1.

ReHashing:

hold α constant if we control table size m .

What if the array fills? Double the array and copy the data.

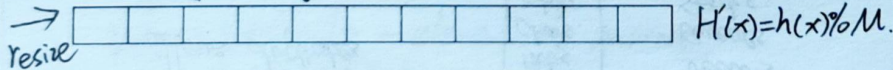
what if the load factor exceeds 0.6?



Resize to first prime bigger than double the array and rehash the data.

$$H(x) = h(x) \% m$$

$$H'(x)$$



Hashing Miscellaneous Discussion -

Which collision resolution strategy is better?

$\alpha > 1$: separate chaining.

clustering: probe based strategies

What structures do hash tables replace for us? Dictionaries

AVL: $\Theta(\log n)$ for worst and avg time. Hash Table: $\Theta(n)$ worst case $\Theta(1)$ on average

There is a constraint on Keyspaces for BST that does not affect hashing... amortized.

Keys for AVL/BST must be comparable

Why do we talk about balanced BST if hashing is so great?

More resources: Comparison based structures allow for search when exact keys are n't known - max, nearest neighbor.

<http://jeffe.cs.illinois.edu/teaching/algorithms/notes/05-hashing.pdf>

http://en.wikipedia.org/wiki/Hash_function

Applications of hashing?

Area of active research in mathematics to develop general purpose hash functions.

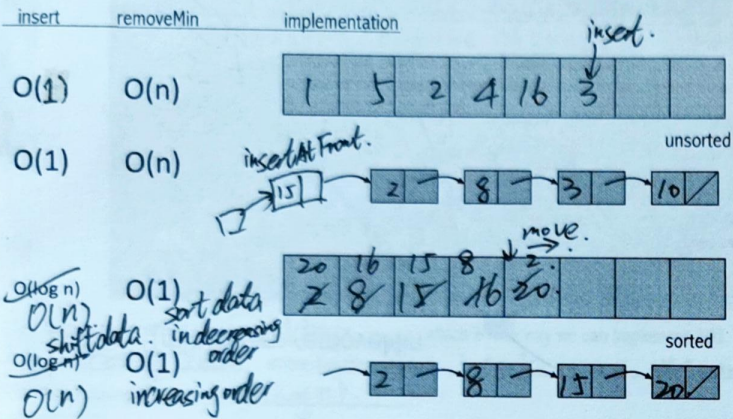
all collision resolution algorithms give the same worst case time complexity for find

Announcements:

PA3 out, due 03/30!

Priority Queue ADT:

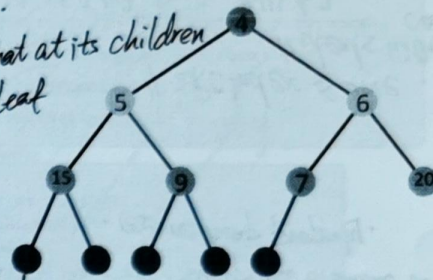
insert \rightarrow comparable keys
 remove min \rightarrow key
 key - priorities



Priority Queue: another implementation option

Tell me everything you can about this structure:

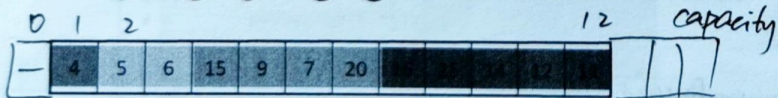
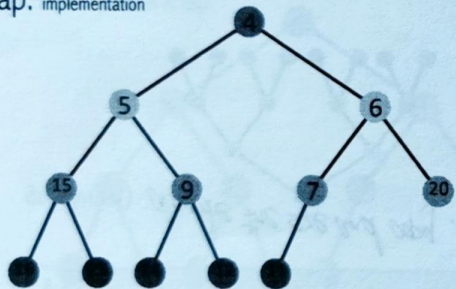
- binary tree - complete.
- each node's key is \leq that at its children
- every path from root to leaf is non-decreasing



Heap: min element is at the top.

(min)Heap: implementation

Any complete tree-based structure leads itself to an array implementation.

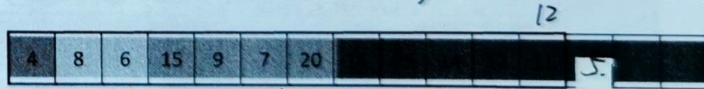
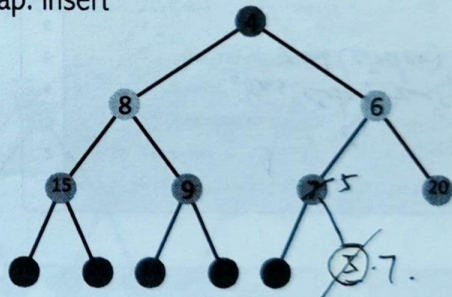


keys in level order - elements of an array.

$LeftChild(i) = 2i + 2$
 $RightChild(i) = 2i + 1$
 $parent(i) = \lfloor \frac{i}{2} \rfloor$

How tall is a heap of size n ?
 $h \leq \lfloor \log_2 n \rfloor$

(min)Heap: insert



insert at end.
 bubble it up ("heapify up")

<https://cs221viz.netlify.com/src/heap.html>

Running time: $O(\log n)$.
 constant time to insert at end.
 heapify up is proportional to the height.

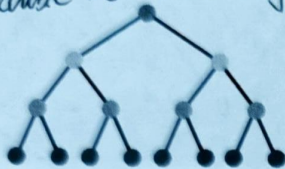
Code:

```

1 void heap<t>::insert(const T & key){
2     if (size == capacity) growArray(); // array is full
3     size++;
4     items[size] = key;
5     heapifyUp(size);
6 }

```

growArray() double the size and copy. \Leftrightarrow add a new level to the tree.



Code:

```

1 void heap<t>::insert(const T & key){
2     if (size == capacity) growArray();
3     size++;
4     items[size] = key;
5     heapifyUp(size);
6 }

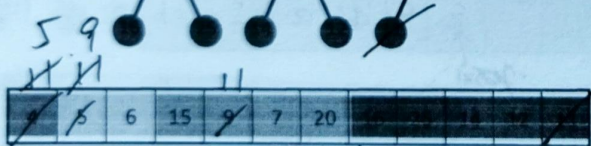
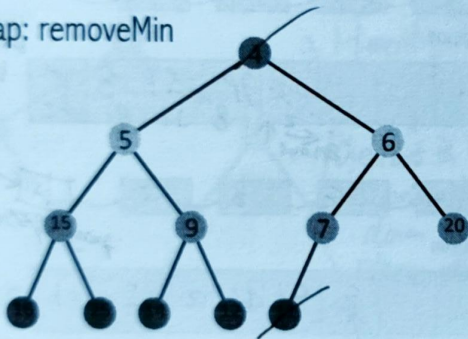
```

```

1 void heap<t>::heapifyUp(int cIndex){
2     if (cIndex > 1)
3         if (items[cIndex] < items[parent(cIndex)]){
4             swap(items[cIndex], items[parent(cIndex)]);
5             heapifyUp(parent(cIndex));
6         }
7 }

```

(min)Heap: removeMin



Swap last key to root
restore order by swapping downward.

Code:

```

1 T heap<t>::removeMin(){
2     T minVal = items[1]; // grab the min
3     items[1] = items[size]; // swap last element to front.
4     size--;
5     heapifyDown(1); // restore heap properly.
6     return minVal;
7 }

```

```

1 void heap<t>::heapifyDown(int cIndex){
2     if (hasAChild(cIndex)) // 2 * cIndex <= size
3         minChildIndex = minChild(cIndex); // deals with one child case.
4     if (items[cIndex] > items[minChildIndex]){
5         swap(items[cIndex], items[minChildIndex]);
6         heapifyDown(minChildIndex);
7     }
8 }
9 }

```

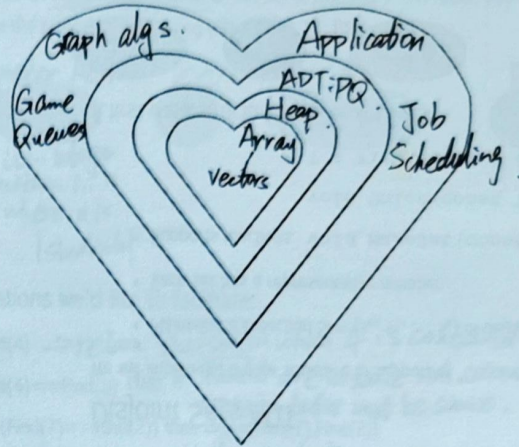
$O(\log n)$.

Announcements: PA3 due 03/30, 11:59p.

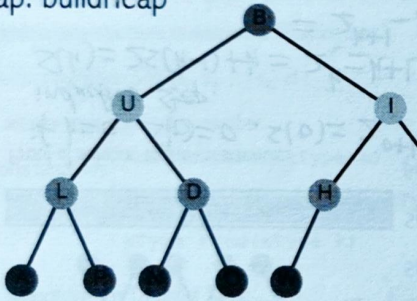


This image reminds us of a Heap, which is one way we can implement ADT Priority Queue, whose functions include insert and remove min/max with running times $\Theta(\log n)$.

What have we done? *Abstraction layers.*

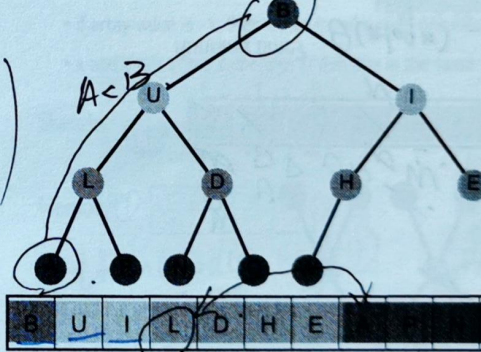


(min)Heap: buildHeap

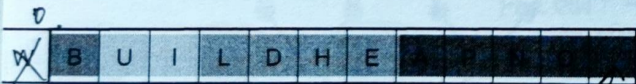


3. Heapify Down from tail
Running time: $O(n)$
The root is the only one that needs to go all the way down to the leaf others don't have to

(min)Heap: buildHeap - 3 alternatives



1. Sort the array:
 $\Theta(n \log n)$
merge sort.



first cell should be empty
1. Heapify Up from End: Not good, from beginning: Repeated insert. $\Theta(n \log n)$.
2. Sort the data: will satisfy heap property if the array is sorted, then any subsequence (trip from root to leaf) is also sorted. $\Theta(n \log n)$

```

1 void heap<t>::buildHeap(){
2   for (int i=2; i<=size; i++)
3     heapifyUp(i);
4 }
    
```

$\Theta(n \log n)$

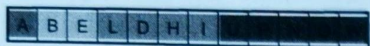
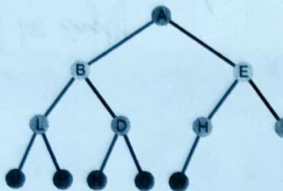
```

1 void heap<t>::buildHeap()
2   for (int i=parent(size); i>0; i--)
3     heapifyDown(i);
4 }
    
```

$\Theta(n)$

$\frac{n}{2}$ calls avoided by starting with parent.
Leaves are already heaps, no need to call heapifyDown on leaves.

(min)Heap: buildHeap



Thm: The running time of buildHeap on an array of size n is $\underline{O(n)}$. $O(n)$ then argue $\sum(h)$.

Instead of focusing specifically on running time, we observe that the time is proportional to the sum of the heights of all of the nodes, which we denote by $S(h)$.

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

$$S(0) = 0.$$

$$\text{Soln } S(h) = 2^{h+1} - h - 2.$$

Proof of solution to the recurrence:

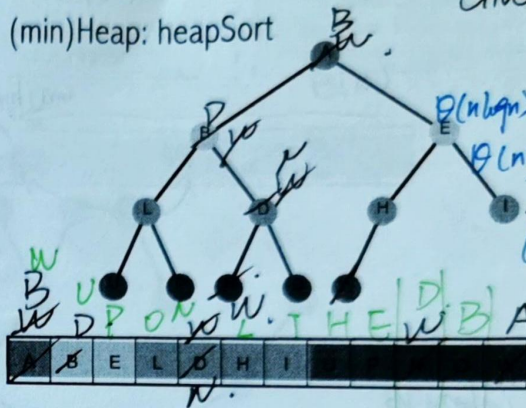
if $h=0$ $S(0)=0$. $S(0) = 2^{0+1} - 0 - 2 = 0$.
inductive step.

$$S(h) = 2S(h-1) + h = 2(2^h - (h-1) - 2) + h = 2^{h+1} - h - 2.$$

But running times are reported in terms of n , the number of nodes...

$$T(n) = S(h) = 2^{h+1} - h - 2, \text{ but we know } h \leq \log_2 n \\ \leq 2^{\log_2 n + 1} - \log_2 n - 2 \leq 2n - \log_2 n - 2 \leq 2n = \Theta(n).$$

(min)Heap: heapSort



Running time? $\Theta(n \log n)$.

Why do we need another sorting algorithm?

Same as merge Sort.

Given an unsorted array

1. build a heap. $\Theta(n)$.

2. Remove min n times placing min in the last spot position \rightarrow size.

3. Adjust the sort to satisfy expectation of caller.

(Reverse array) (Place data at 0)

Depends on the

A "semi-sort" that the

heap gives us,

Otherwise, it's much like selection sort.

in-place algorithm.

not stable — relative

locations of tied values are not maintained.

Announcements

PA3 available, due 03/30, 11:59p.



If this structure were a heap, we could build it in time $\Theta(n)$, and use it to support heapSort, which runs in time $\Theta(n \log n)$. $\sqrt{n \log n}$ lower bound in sorting algorithm.

Disjoint Sets ADT

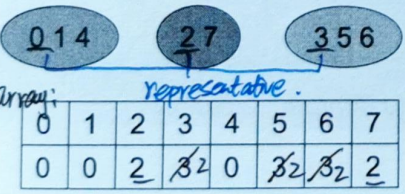
all intersections of two sets will be empty (no intersection operation)

We will implement a data structure in support of "Disjoint Sets":

- Maintains a collection $S = \{s_0, s_1, \dots, s_k\}$ of disjoint sets.
- Each set has a representative member.
- Supports functions:


```
void MakeSet(kType k)  constructor
void Union(kType k1, kType k2)
kType Find(kType k)
```

A first data structure for Disjoint Sets:

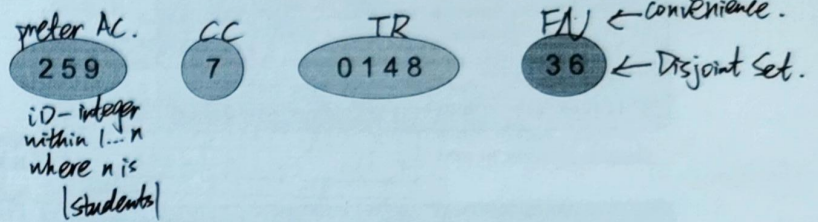


Find(7) = 2.

Find: Returns $arr[k] \rightarrow \Theta(1)$
 Union: $Find(2) \cup Find(3) \rightarrow \Theta(n)$.

A Disjoint Sets example:

partitions a set into non-overlapping subsets.
 Let R be an equivalence relation on the set of students in this room, where $(s, t) \in R$ if s and t have the same favorite among {FN, FB, TR, CC, PMC, AG}.

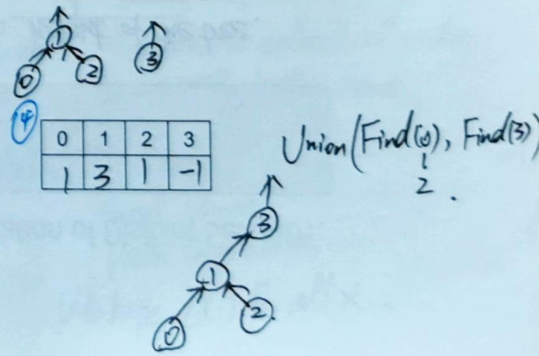
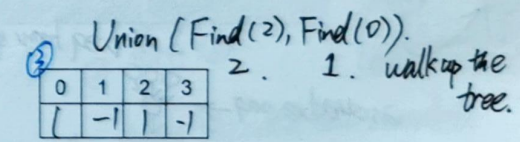
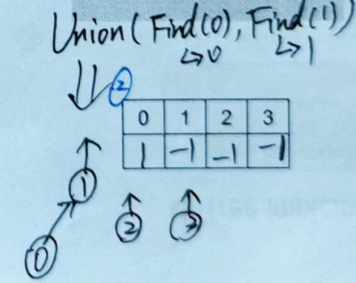
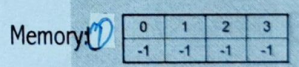
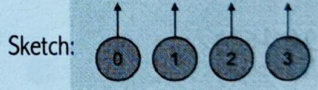


Operations we'd like to facilitate:

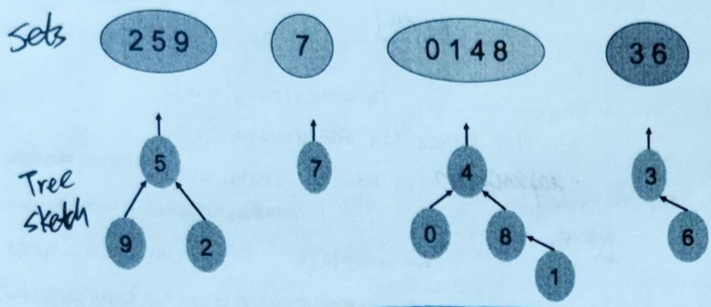
- Find(4) → return the set in which 4 is contained (returns some 'representative' of the set, 0, 1, 4, 8, any is fine)
- Find(4) == Find(8) True if 4 and 8 are in same set. representative must be same.
- If $!(Find(7) == Find(2))$ then $Union(Find(7), Find(2))$
 not in the same set. put them in the same set by taking union.
 After union, all elements of the resulting set have the same representative.

A better data structure for Disjoint Sets: UpTrees

- if array value is -1, then we've found a root, o/w value is index of parent which is in the same set as current index.
- x and y are in the same tree iff they are in the same set.



Example: partition of elements into disjoint sets



array

0	1	2	3	4	5	6	7	8	9
4	8	5	3	1	-1	3	-1	4	5

Find(1)

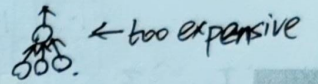
UpTree implementation of Disjoint Sets ADT:

```

1 int DS::Find(int i) {
2     if (s[i] < 0) return i;
3     else return Find(s[i]);
4 }
    
```

Running time depends on height of the tree.

Worst case? $\Theta(n)$.



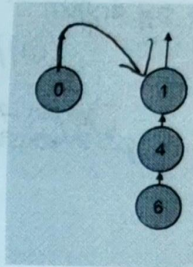
What's an ideal tree?

constant height.

```

1 int DS::Union(int root1, int root2) {
2     s[root2] = root1;
3 }
    
```

Running time $\Theta(1)$.

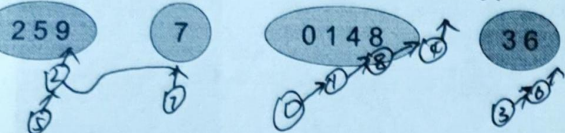


Announcements

available, due 03/30, 11:59p. GS submission this evening.

call, Disjoint Sets example:

use the members of each set below have the same favorite among {AC, FN, FB, TR, CC, CR}.



0	1	2	3	4	5	6	7	8	9
1	8	-1	6	-1	2	-1	-1	4	5

up trees in an array.

are the results of the following statements?

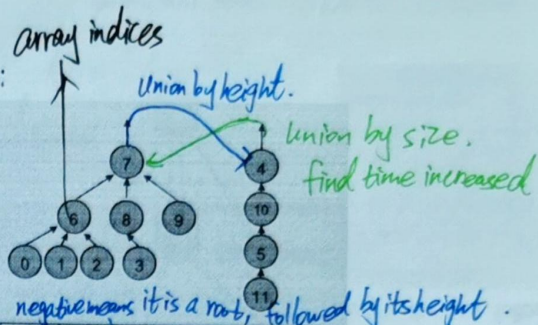
$\text{Find}(0) = 4$ (return the set in which 0 dwells)

$\text{Find}(0) == \text{Find}(8)$ True

$\text{Find}(7) == \text{Find}(9)$ then $\text{Union}(\text{Find}(7), \text{Find}(9))$ pointing 7 to 2 keeps tree short.

ADT - Find
Union

Smart unions:



Union by height:
store $-(\text{height} + 1)$ in root value.

0	1	2	3	4	5	6	7	8	9	10	11
6	6	6	8	-4	10	7	-3	7	4	5	

Keeps overall height of tree as small as possible.

Union by size:
store $-(\# \text{ of nodes})$ in root value.

0	1	2	3	4	5	6	7	8	9	10	11
6	8	6	8	-4	10	7	-8	7	4	5	

Increases distance to root for fewest nodes.

Both of these schemes for Union guarantee the height of the tree is $\Theta(\log n)$.

Smart unions:

```
int DS::Find(int i) {
    if (s[i] < 0) return i;
    else return s[i] = Find(s[i]);
}
```

Keeps the same. $\Theta(\log n)$, improved from $\Theta(n)$.

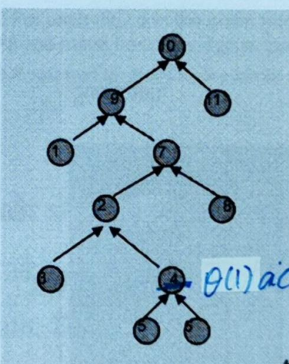
```
int DS::UnionBySize(int root1, int root2) {
    int newSize = s[root1] + s[root2];
    if (isBigger(root1, root2) {
        s[root2] = root1;
        s[root1] = newSize;
    } else {
        s[root1] = root2;
        s[root2] = newSize;
    }
}
```

by making Union more complex, we reduce runtime of find. \rightarrow store negative magnitudes. resulting root.

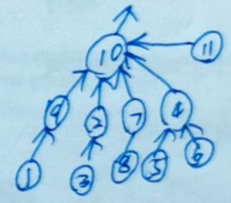
Runtime $\Theta(1)$.

Path Compression:

Find(4)



example execution. $\Theta(1)$ access from array.



along the path from a node to the root, it redirects all nodes on the path to the root.

Up trees
Smart Union
Path Compression
We hope $\Theta(1)$ for all x.

Analysis:

Something we'll need -- Iterated log $\log^* n := \begin{cases} 0 & \text{if } n \leq 1; \\ 1 + \log^*(\log n) & \text{if } n > 1 \end{cases}$

Counts the number of times

Example. we can take log of a number before we hit 1.

$2^{65536} \xrightarrow{\textcircled{1}} 65536 \xrightarrow{\textcircled{2}} 16 \xrightarrow{\textcircled{3}} 4 \xrightarrow{\textcircled{4}} 2 \xrightarrow{\textcircled{5}} 1$. $\log^* 2^{65536} = 5$.

Relevant result:

In an upTree implementation of Disjoint Sets using smart union and find with path compression...

any sequence of m union and find operations results in worst case running time of $O(m \log^* n)$, where n is the number of items.

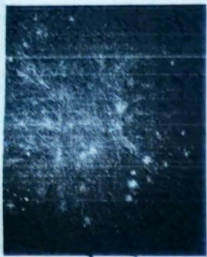
\searrow *nearly constant.*

actually a loose bound.

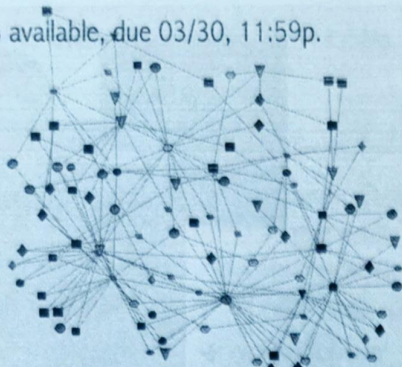
<http://research.cs.vt.edu/AVresearch/UF/>

ouncements PA3 available, due 03/30, 11:59p.

phs!! $G = (V, E)$



2003 internet.
ph: vertex: router
edge: network



Adjacent Vertices, $N(v) = \{u \mid (v, u) \in E\}$.

Incident Edges, $I(v) = \{(u, v) \mid (u, v) \in E\}$.
 V has $|I(v)|$ # of incident edges.

Aside: Vertices are labelled, no pair of adjacent vertices is the same.

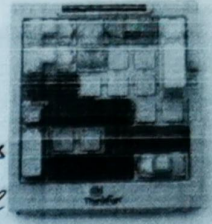
subgraph: $G' = (V', E')$
 is a subgraph of a graph G
 iff $\forall v \in V', E' \subseteq E,$
 $e = (v, u) \in E',$ then $v \in V', u \in V'$.

degree of vertex:
 (# of neighbors that)

winning state

Path: sequence of vertices connected by edges.

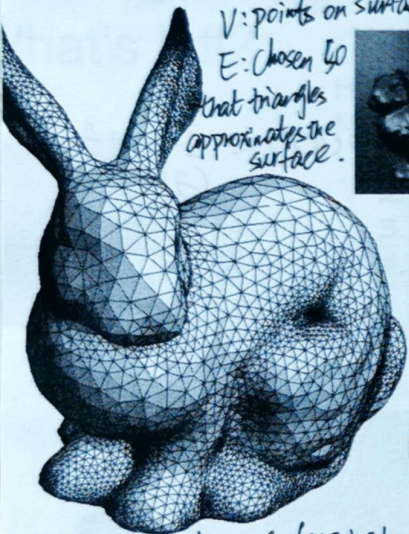
V : game states
 E : single move



← path is the solution to the puzzle

Aside: there is a great alg for solving this one - BFS.

V : points on surface of an object.
 E : chosen Δ that triangles approximate the surface.



Planar surface mesh.
 - Can be embedded in a plane drawn with no crossing edges.

This graph is used to calculate whether a given number is divisible by 7.

1. Start at the circle node at the top.
2. For each digit d in the given number, follow d blue (solid) edges in succession. As you move from one digit to the next, follow 1 red (dashed) edge.
3. If you end up back at the circle node, your number is divisible by 7.

directed graph:
 $(u, v) \neq (v, u).$
 $(a, b) \neq (b, a).$

3703

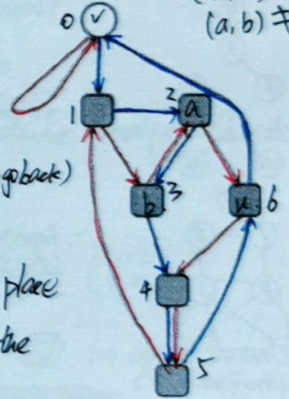
Walk: sequence of vertices between which there are edges.

Trail: walk with no repeated edges (can't go back)

simple Path: trail with no repeated vertices

Circuit: trail that begins and ends at same place

Cycle: path that begins and ends at the same place.
 (path that allows one repeat).



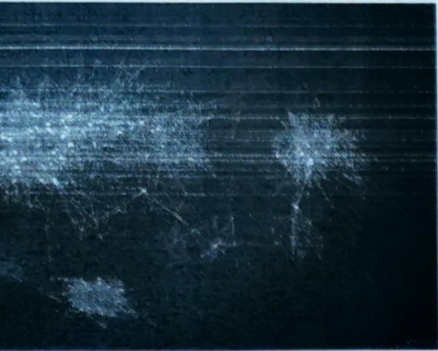
$V: k \text{ is } k \text{ mod } 7.$



not a path but it is a trail

ouncements PA3 available, due 03/30, 11:59p.

Graphs!! $G = (V, E)$



Tree: *connected, acyclic graph.*

\equiv between every pair of vertices, there is one path.

Simple (sub)Graph:

graph with no self loops or multiedges.

Spanning subgraph: (u, v) is unique $u \rightarrow v$

Let G' be a subgraph of G .

$G' = (V', E')$, $G = (V, E)$

G' is a spanning subgraph if $V' = V$

Spanning tree:

G' is a spanning tree of G if it is a tree and $V' = V$.

connected (sub)Graph: A graph is connected if between any pair of vertices there is a path.

connected component: there is a maximal connected subgraph.

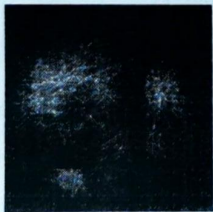
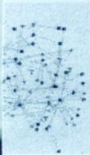
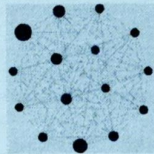
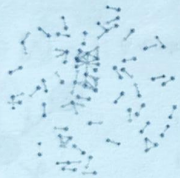
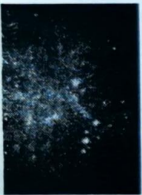
if we add anything to subgraph, then it is no longer connected.

What's left?

How do we get from here to there?

Need:

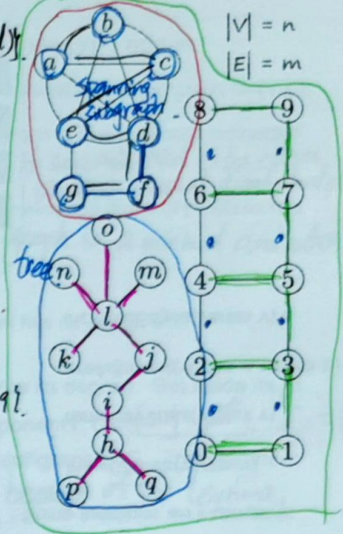
1. Common Vocabulary
2. Graph implementation
3. Traversal
4. Algorithms.



Graph Vocabulary: Use the graph G to answer these questions.

1. List the edges incident on vertex b : $\{(b, a), (b, c), (b, e), (b, d)\}$.
2. What is the degree of vertex h ? 3 .
3. List all the vertices adjacent to vertex i : $\{h, j\}$.
4. Describe a path from p to o : p, h, i, j, l, o .
5. Describe a path from 6 to g : DNE.
6. List the vertices in the largest complete subgraph in G : $\{a, b, c, d, e\} = K_5$. *every node is connected.*
7. Describe the largest connected subgraph in G : *Tie between the $\{p, o\}$ and $\{6-9\}$.*
8. Describe the connected components in G : 3 connected components.
9. How many edges in a spanning forest of G ? $9 + 9 + 6 = 24$.
10. How many paths from 0 to 9 ? $2^4 = 16$ (encoding of binary #).
11. Can you draw G with no crossing edges? *not possible for red component. Cannot draw K_5 in the plane. if a graph has a K_5 minor then it is not planar.*

Connected subgraph $G = (V, E)$

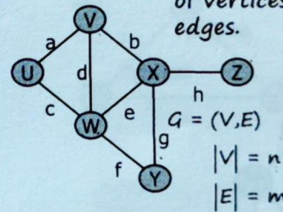


graphs

Connected subgraph is not necessarily a connected component. e.g. $\{4-5\}$ is connected subgraph, but not a connected component.

Graphs: theory that will help us in analysis

Running times often reported in terms of n , the number of vertices, but they often depend on m , the number of edges.



How many edges?

At least:

connected - $n-1$

we can remove an edge in a cycle to make it a tree.

not connected - 0

At most:

Complete graph. simple - $\frac{n(n-1)}{2}$. (C_n)

not simple - \diamond

no upper bound can have multiple edges, self-loops.

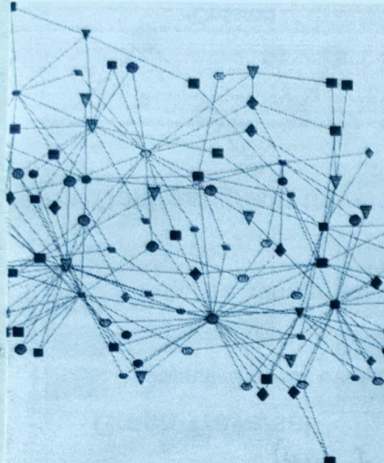
Relationship to degree sum:

$$\sum_{v \in V} \deg(v) = 2m$$

n vertices, each with $(n-1)$ incident edges. $n(n-1)$ endpoints. $\Rightarrow \frac{n(n-1)}{2}$ edges.

ouncements PA3 available, due 04/09, 11:59p.

phs!! $G = (V, E)$



Suppose $n = 100$,
 Give a lower bd for the number of edges: 99
 $|E| = |V| - 1$ for a minimally connected graph (tree).

Give an upper bd for the number of edges:
 $m \leq 4950 = \frac{n(n-1)}{2}$

Suppose the average degree is 22.
 How many edges? $\frac{\# \text{ of incident edges of a vertex}}{2}$

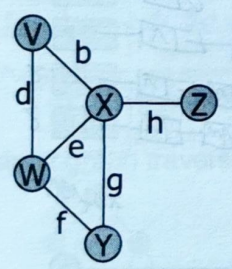
$$2|E| = \sum_{v \in V} \text{deg}(v)$$

$$\frac{\sum_{v \in V} \text{deg}(v)}{n} = 22$$

$$\Rightarrow \sum_{v \in V} \text{deg}(v) = 2200 \Rightarrow 1100 \text{ edges.}$$

every edge is counted twice.

raphs: Toward implementation...(ADT)



Functions: (merely a smattering...)

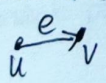
- insertVertex(pair keyData)
- insertEdge(vertex v1, vertex v2, pair keyData)
- removeEdge(edge e);
- removeVertex(vertex v);

Vertices - hashtable.
 Edges - includes endpoint info
 + some structure that reflects the connectivity of the graph
 2 alternatives.

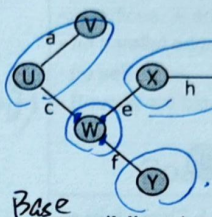
incidentEdges(vertex v); \rightarrow set of neighbors

areAdjacent(vertex v1, vertex v2);
 existence of an edge.

origin(edge e); $\rightarrow u$
 dest(e) $\rightarrow v$
 $(u,v) \neq (v,u)$



Thm: Every minimal connected graph $G=(V,E)$ has $|V|-1$ edges.



Proof: Consider an arbitrary minimal connected graph $G=(V,E)$.

Lemma: Every connected subgraph of G is minimally connected.
 (easy proof by contradiction) if subgraph had extra edges, then the original graph would too.
 IH: For any $j < |V|$, any minimal connected graph of j vertices has $j-1$ edges. Any smaller graph is a minimal connected graph.

Base
 Suppose $|V| = 1$: A minimal connected graph of 1 vertex has no edges, and $0 = 1-1$.

Inductive.
 Suppose $|V| > 1$: Choose any vertex w and let d denote its degree. Set aside its d incident edges, partitioning the graph into $d+1$ components, $C_0 = (V_0, E_0), \dots, C_i = (V_i, E_i)$.
 $C_0 = (V_0, E_0)$, each of which is a minimal connected subgraph of G . This means that $|E_k| = |V_k| - 1$ by IH, which we can apply because of the lemma, by observing that C_k is smaller than G .

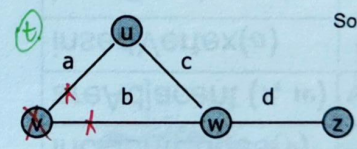
Now we'll just add up edges in the original graph:

$$|E| = d + |E_0| + |E_1| + \dots + |E_d|$$

$$= d + \sum_{k=1}^d (|V_k| - 1)$$

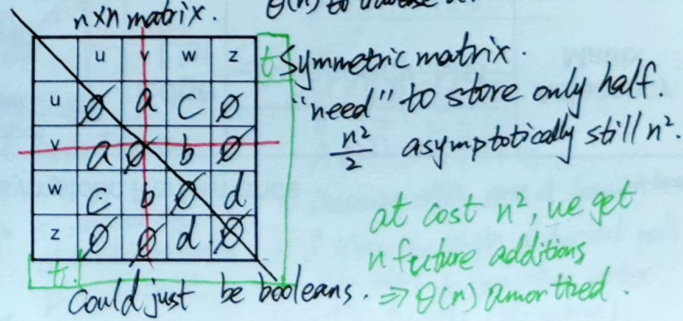
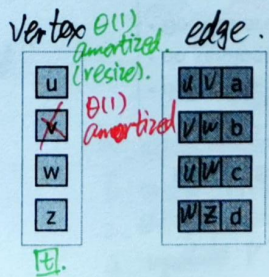
$= \sum_{k=1}^d |V_k| = |V| - 1 = n - 1$ all vertices except w .

Graphs: Adjacency Matrix $G=(V,E)$ $|V|=n$ $|E|=m$



Some functions we'll compare:

- insertVertex(vertex v) $\theta(n^2)$ if we double, $\theta(n)$ if we double and copy.
- removeVertex(vertex v) $\theta(n)$ to remove all incident edges.
- areAdjacent(vertex v, vertex u) $\theta(1)$ justifies choice of matrix
- incidentEdges(vertex v) $\theta(n)$ to traverse a row of the matrix

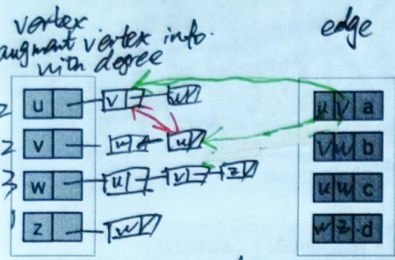
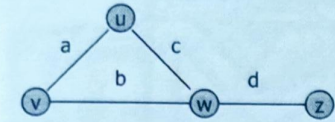


Announcements PA3/HW3 av, due 04/09, 11:59p.

Graphs: Adjacency List

Some functions we'll compare:

insertVertex(vertex v) $\Theta(1)$ amortized.
 removeVertex(vertex v) $\Theta(\text{deg}(v))$
 areAdjacent(vertex v, vertex u) $\Theta(\min(\text{deg}(v), \text{deg}(u)))$
 incidentEdges(vertex v) $\Theta(\text{deg}(v))$



vertex augment vertex info. with degree
 edge
 $2m$ total nodes
 storage is $\Theta(n+m) = O(n^2)$ (max of $n+m$)

$m = O(n^2)$
 or $\Theta(n)$
 $\Theta(1)$ } Anything $\leq n^2$.

Graphs: Asymptotic Performance

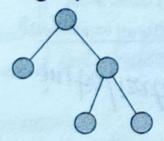
	Edge List <small>if no additional structure $n+m$</small>	Adjacency List	Adjacency Matrix
n vertices, m edges	simple	$n+m$	n^2
no parallel edges			
no self-loops			
Bounds are big- Θ			
Space		$n+m$	n^2
incidentEdges(v)	m	$\text{deg}(v)$	n
areAdjacent(v, w)	m	$\min(\text{deg}(v), \text{deg}(w))$	1
insertVertex(o)	1	1	$\Theta(n)$ amortized
insertEdge(v, w, o)	1	1	1
removeVertex(v)	m	$\text{deg}(v)$	$\Theta(n)$ amortized
removeEdge(e)	1	1	1

Graph Traversal

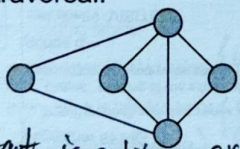
Objective: Visit every vertex and every edge in the graph, while honoring the connectivity of the graph, walk around on edges of graph.

Purpose: We can search for interesting substructures in the graph. Like shortest path from current location to another vertex. List of low degree.

Contrast graph traversal to BST traversal:

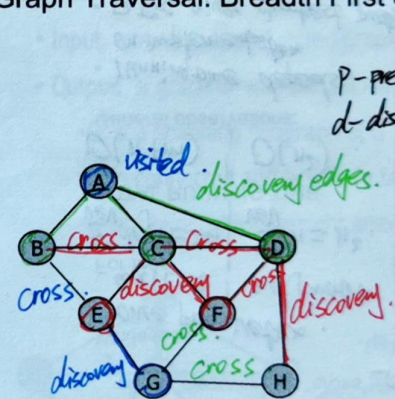


- Ordered
- Obvious start
- good way of knowing when done



- start is arbitrary or application dependent.
- iterate over incident edges
- mark progress

Graph Traversal: Breadth First Search (by example)



vertex structure / adj list data

v	d	p	edges
A	0	-	B, C, D
B	1	A	A, C, E
C	1	A	A, B, D, E, F
D	1	A	A, C, F, H
E	2	C	C, B, G
F	2	C	C, D, G
G	3	E	E, F, H
H	2	D	D, G

Queue: ~~G H E F D B C A~~

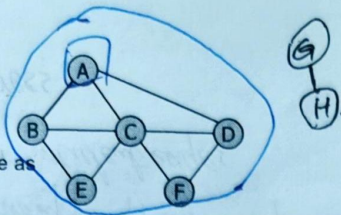
$\begin{matrix} 0 & 1 & 1 & 2 & 2 & 3 \\ A & C & B & D & E & F & H & G \end{matrix}$

- For every neighbor of current vertex.
1. mark visited.
 2. edge between is discovery.
 3. update p & d, $d = 1 + \text{current } d$.
 4. put on queue. $p = \text{current}$.

- ① Discovery edges are a spanning tree
- ② d gives the length of shortest path from start to every other vertex.
- ③ p gives the paths (like up tree).

Graphs: Traversal – BFS

Visits every vertex and classifies each edge as either "discovery" or "cross"



initialization

Algorithm BFS(G)

Input: graph G

Output: labeling of the edges of G as discovery edges and back edges

For all u in G.vertices()

setLabel(u, UNEXPLORED)

For all e in G.edges()

setLabel(e, UNEXPLORED)

For all v in G.vertices()

if getLabel(v) = UNEXPLORED

BFS(G,v)

grabs all components of graph.

Algorithm BFS(G,v)

Input: graph G and start vertex v

Output: labeling of the edges of G in the connected component of v as discovery edges and cross edges

queue q;

setLabel(v, VISITED)

q.enqueue(v);

While !q.isEmpty()

q.dequeue(v)

For all w in G.adjacentVertices(v)

if getLabel(w) = UNEXPLORED

setLabel((v,w), DISCOVERY)

setLabel(w, VISITED)

q.enqueue(w)

else if getLabel((v,w)) = UNEXPLORED

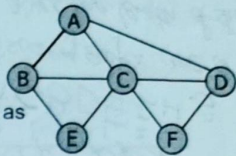
setLabel((v,w), CROSS)

pseudo code for the algorithm.

update pld.

Graphs: Traversal – BFS

Visits every vertex and classifies each edge as either "discovery" or "cross" $\Omega(n+m)$



While loop once per vertex.

adj list | adj matrix
For loop

$$\sum_{v \in V} \text{deg}(v) = 2m \quad \sum_{v \in V} n = n^2$$

TOTAL RUNNING TIME:

$$O(n+2m) \quad O(n^2)$$

General observations:

- Running time depends on implementation
- BFS gives shortest path with pld book keeping
- easy to count components & detect cycles

Algorithm BFS(G,v)

Input: graph G and start vertex v

Output: labeling of the edges of G in the connected component of v as discovery edges and cross edges

queue q;

setLabel(v, VISITED) $O(1)$

q.enqueue(v); $O(1)$

While !q.isEmpty()

q.dequeue(v) $O(1)$

For all w in G.adjacentVertices(v) $\text{deg}(v)$

if getLabel(w) = UNEXPLORED

setLabel((v,w), DISCOVERY)

setLabel(w, VISITED)

q.enqueue(w)

else if getLabel((v,w)) = UNEXPLORED

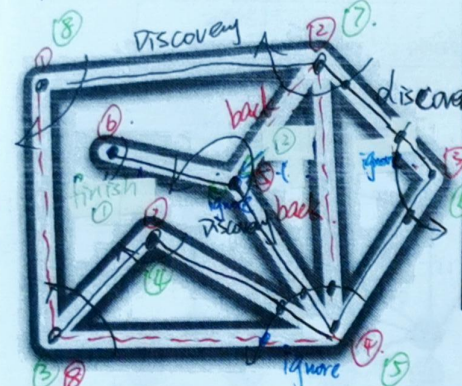
setLabel((v,w), CROSS)

{O(1)}

BFS runtime $\Theta(n+m)$.

Announcements PA3/HW3 av, due 04/09, 11:59p.

Graphs: Traversal - DFS



Ariadne, Theseus, and the Minotaur

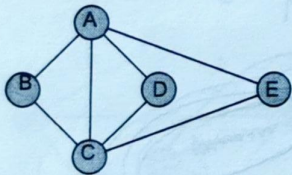
DFS(v)
 1. visits a vertex v.
 2. for each neighbor u if u unvisited, marks (u,v) as discovery DFS(u).
 else marks (u,v) as back.

- <http://www.cs.duke.edu/csed/jawaa2/examples/DFS.html>
- <http://www.student.seas.gwu.edu/~idsv/idsv.html>
- <http://www.youtube.com/watch?v=8qrZ1clEp-Y>

in graph: Discovery
 Back
 Ignore.
 Finish.

Discovery edges spanning tree.
 not shortest path.

Graphs: DFS example



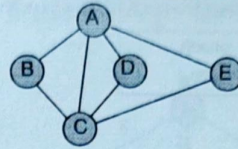
A	B	C	D	E
B	A	C		
C	B	A	D	E
D	A	C		
E	A	C		

Application: Topological sort, Eulerian paths

- Discovery edges form spanning tree.
- 2 convenient visit times.
- can count components.
- detect edges.

setting/getting labels:
 every vertex labeled twice $\Theta(n)$
 every edge is labeled twice $\Theta(m)$
 querying vertices: (for loop) each vertex $\deg(v)$ (adj list)
 total over algorithm $\sum_{v \in V} \deg(v) = 2m$
 querying edges: $\Theta(1)$ per edge.
TOTAL RUNNING TIME:
 $\Theta(n+m)$
 adjacency list.
 $\Theta(nm) = \Theta(n^2)$
 adj matrix.

DFS: "visits" each vertex, classifies each edge as either "discovery" or "back"



Algorithm DFS(G)
 Input: graph G
 Output: labeling of the edges of G as discovery edges and back edges

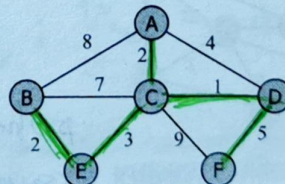
For all u in G.vertices() setLabel(u, UNVISITED)
 For all e in G.edges() setLabel(e, UNEXPLORED) } Initialization
 For all v in G.vertices() if getLabel(v) = UNVISITED DFS(G,v) } pickup components.

Algorithm DFS(G,v)
 Input: graph G and start vertex v
 Output: labeling of the edges of G in the connected component of v as discovery edges and back edges
 setLabel(v, VISITED)
 For all w in G.adjacentVertices(v) if getLabel(w) = UNVISITED setLabel((v,w), DISCOVERY) DFS(G,w) else if getLabel((v,w)) = UNEXPLORED setLabel(e, BACK)

Minimum Spanning Tree Algorithms:

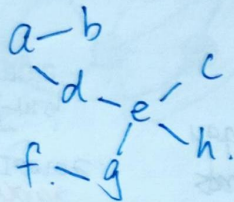
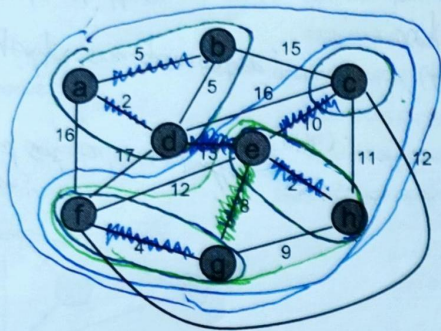
- Input: connected, undirected graph G with unconstrained edge weights *positive/negative.*
- Output: a graph G' with the following characteristics -
 - G' is a spanning subgraph of G
 - G' is connected and acyclic (a tree)
 - G' has minimal total weight among all such spanning trees -

Observation: if we partition the vertices into 2 sets, and look at the edges between, then least weight edge is part of some MST. We don't care about total path length between any specific pair.



- changing BC 7 to 4 does not change the MST.
- change CF 9 to 4 will change the MST (DF to CF).

Kruskal's Algorithm



we are done because we have added $n-1$ edges.

reject (makes a cycle)

(a,d)
(a,h)
(f,g)
(a,b)
(b,d)
(g,e)
(g,h)
(e,c)
(c,h)
(e,f)
(f,c)
(d,e)
(b,c)
(c,d)
(a,f)
(d,f)

edges in increasing cost.

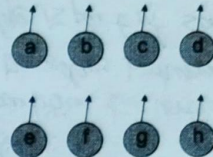
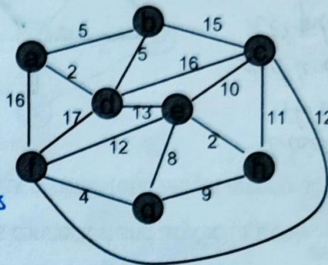
priority queue (in entries)

↓ increasing cost.

(g,h) g & h are in the same component. (g,h) creates a cycle.

maintain the growth of the structure as a collection of disjoint set.

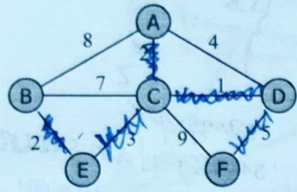
Kruskal's Algorithm (1956)



1. Initialize graph T whose purpose is to be our output. Let it consist of all n vertices and no edges.
2. Initialize a disjoint sets structure where each vertex is represented by a set.
3. Remove Min from PQ . If that edge connects 2 vertices from different sets, add the edge to T and take union of the vertices' two sets, otherwise do nothing. Repeat until $n-1$ edges are added to T .

(a,d)
(e,h)
(f,g)
(a,b)
(b,d)
(g,e)
(g,h)
(e,c)
(c,h)
(e,f)
(f,c)
(d,e)
(b,c)
(c,d)
(a,f)
(d,f)

Announcements PA3/HW3 av, due 04/09, 11:59p.

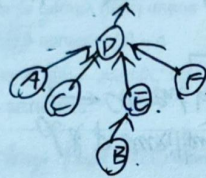
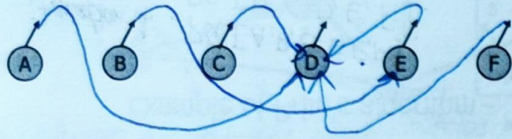


- (a,b)
- (a,c)
- (a,d)
- (a,e)
- (a,f)
- (b,c)
- (b,d)
- (b,e)
- (b,f)
- (c,d)
- (c,e)
- (c,f)

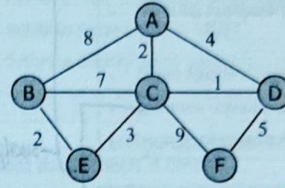
m edges.

$(n-1)$ edges \Rightarrow tree done

disjoint sets.



Kruskal's Algorithm - preanalysis



sorting algorithm $\Theta(n \log n)$.

Algorithm KruskalMST(G)

disjointSets forest;
for each vertex v in V do
forest.makeSet(v);

priorityQueue Q;
Insert edges into Q, keyed by weights

graph $T = (V, E)$ with $E = \emptyset$; $\Theta(n)$.
while T has fewer than $n-1$ edges do

edge $e = Q.removeMin()$
Let u, v be the endpoints of e
if forest.find(v) \neq forest.find(u) then
Add edge e to E $\Theta(1)$ Assume e .
forest.smartUnion
(forest.find(v), forest.find(u))

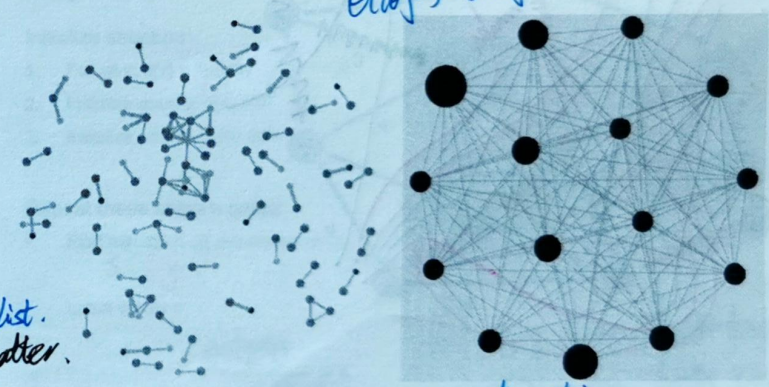
return T

Priority Queue:	Heap	Sorted Array
To build	$\Theta(m)$	$\Theta(m \log n)$
Each removeMin	$\Theta(\log n)$	$\Theta(1)$

add Edge - AdjList add edge to start of each of two lists. $\Theta(1)$.

POPS top restore heap (heapify down)

$m \leq n^2$
 $\log m \leq \log n^2$
 $\log m \leq 2 \log n$
 $\Theta(\log n) = \Theta(\log m)$.



How would you characterize the difference between these two graphs?
connectedness. Dense graph: right graph has many edges $\Theta(n^2)$
completeness. Sparse graph: left has few edges. $\Theta(n)$
adjacency lists.

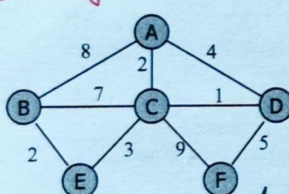
Kruskal's Algorithm - analysis $\Theta(m \log n)$

```

Algorithm KruskalMST(G)
[ disjointSets forest;
  for each vertex v in V do  $\Theta(n)$ .
  forest.makeSet(v);

  priorityQueue Q;  $\Theta(m)$  vs  $\Theta(m \log n)$ .
  Insert edges into Q, keyed by weights

  graph T = (V, E) with E =  $\emptyset$ ;  $\Theta(n)$ .
  while T has fewer than n-1 edges do
    edge e = Q.removeMin()  $\Theta(\log n)$  vs  $\Theta(1)$ .
    Let u, v be the endpoints of e
    if forest.find(v)  $\neq$  forest.find(u) then  $\Theta(1)$ .
      Add edge e to E
      forest.smartUnion
        (forest.find(v), forest.find(u))
    }
  return T
  
```



no iterations over incident edges
 \Rightarrow graph implementation doesn't matter.

Graph Impl = answer - adjlist. but it doesn't matter.

Priority Queue:	Total Running time:
Heap	$\Theta(n + m + n + m \log n + m) = \Theta(n + m + m \log n) = \Theta(m + m \log n)$
Sorted Array	$\Theta(n + m \log n + n + m) = \Theta(n + m + \log n) = \Theta(m + m \log n)$.

we may choose sorted array, since the graph may have already been sorted.
 \rightarrow we can make improvements on $(m \log n)$

Prim's algorithms (1957) is based on the Partition Property:

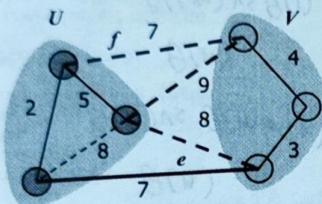
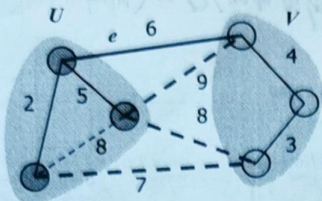
Consider a partition of the vertices of G into subsets U and V .

Let e be an edge of minimum weight across the partition.

Then e is part of some minimum spanning tree.

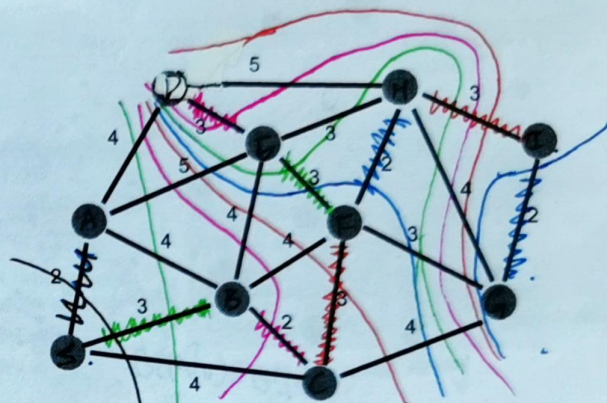
Proof:

See cpsc320?



MST - minimum total weight spanning tree

Theorem suggests an algorithm... (like BFS, with priority queue)



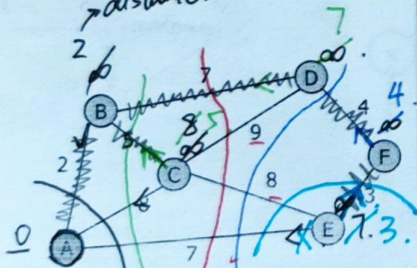
d - distance (cost).
 p - predecessor

Solved/ Done
 "labelled"

Example of Prim's algorithm -

iteration 1. PQ [A, B, C, D, E, F].
 iteration 2. PQ [B, C, D, E, F].
 iteration 3. PQ [C, D, E, F].
 iteration 4. PQ [D, E, F].

arrows are predecessors.
 → distances.



iteration 5 PQ [E, F]

Initialize structure:
 1. For all v , $d[v] = \text{"infinity"}$, $p[v] = \text{null}$
 2. Initialize source: $d[s] = 0$
 3. Initialize priority (min) queue
 4. Initialize set of labeled vertices to \emptyset .

d & p auxiliary structure
Current best cost of attaching to solution
(u, v) is the edge that gives best

Repeat these steps n times:

- Find & remove minimum $d[]$ unlabelled vertex: v PQ operation.
- Label vertex v
- For all unlabelled neighbors w of v ,
 If $\text{cost}(v, w) < d[w]$
 $d[w] = \text{cost}(v, w)$
 $p[w] = v$ update.

Final solution is in the predecessors.
 (no unique solution).

Prim's Algorithm (undirected graph with unconstrained edge weights):

Initialize structure:
 1. For all v , $d[v] = \text{"infinity"}$, $p[v] = \text{null}$
 2. Initialize source: $d[s] = 0$
 3. Initialize priority (min) queue
 4. Initialize set of labeled vertices to \emptyset .

Repeat these steps n times:

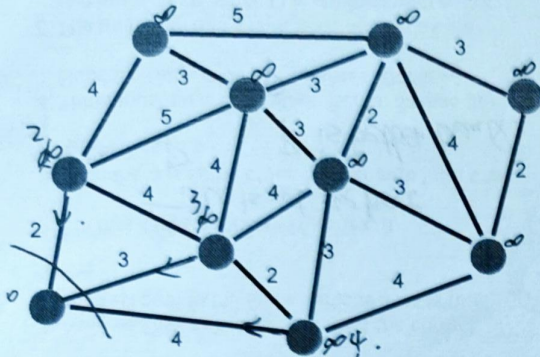
- Remove minimum $d[]$ unlabelled vertex: v
- Label vertex v (set a flag)
- For all unlabelled neighbors w of v ,
 If $\text{cost}(v, w) < d[w]$
 $d[w] = \text{cost}(v, w)$ changes a priority.
 $p[w] = v$

	adj mtx	adj list
heap	$O(n^2 + m \log n)$	$O(n \log n + m \log n)$
Unsorted array	$O(n^2)$	$O(n^2)$

Which is best?
 Depends on density of the graph:
 Sparse
 Dense

Announcements PA3/HW3 av, due 04/09, 11:59p.

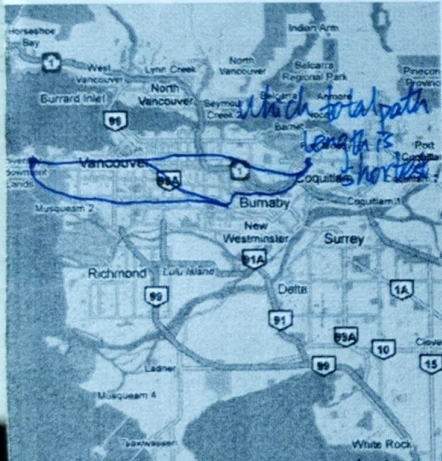
Prim's MST: weighted, undirected edges



For each vertex, we may update $d[v]$ times $\deg(v)$

- choose an arbitrary start. - choice may change final tree but will not change final value of MST (value == total cost)
- init structure pld structures
- proceed algorithm.

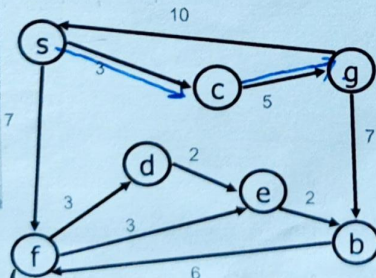
Single source shortest path



Given a start vertex (source) s , find the path of least total cost from s to every vertex in the graph. *dist, time, \$*

Input: directed graph G with non-negative edge weights, and a start vertex s .

Output: A subgraph G' consisting of the shortest (minimum total cost) paths from s to every other vertex in the graph. Dijkstra's Algorithm (1959)



- if a graph has a neg weight cycle, then shortest path doesn't exist.
- if a graph has neg-weight edge and no neg weight cycle, then shortest paths exist, but Dijkstra's won't find them

Prim's Algorithm (undirected graph with unconstrained edge weights):

Initialize structure:

1. For all v , $d[v] = \text{"infinity"}$, $p[v] = \text{null}$ $\Theta(n)$
2. Initialize source: $d[s] = 0$ $\Theta(1)$
3. Initialize priority (min) queue $\Theta(n)$ $\Theta(1)$
4. Initialize set of labeled vertices to \emptyset . $\Theta(1)$

graph impl.	adj mtx	adj list
PR	$O(n^2)$	$n + n \log n + m \log n$ $O(n \log n + m \log n)$
heap	$O(n^2)$	$O(n \log n + m \log n)$
Unsorted array	$O(n^2)$	$O(n^2)$ $O(n^2 + m)$

Repeat these steps n times:

- Remove minimum $d[]$ unlabeled vertex: v $\Theta(1)$
- Label vertex v (set a flag) $\Theta(1)$
- For all unlabeled neighbors w of v $\Theta(\deg(v))$
if $\text{cost}(v,w) < d[w]$ $\Theta(1)$

changes PR. $\begin{cases} d[w] = \text{cost}(v,w) \\ p[w] = v \end{cases}$ $\log n$ $\Theta(1)$ hash table.

Which is best?

Depends on density of the graph:
Sparse (if m is $O(n)$): heap + adj list.
Dense unsorted array.

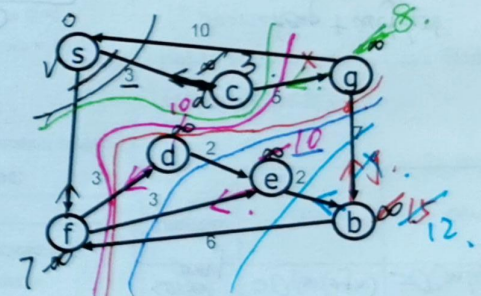
Single source shortest path (directed graph w non-negative edge weights):

Initialize structure:

1. For all v , $d[v] = \text{"infinity"}$, $p[v] = \text{null}$
2. Initialize source: $d[s] = 0$
3. Initialize priority (min) queue

Repeat these steps n times:

- Find minimum $d[]$ unlabeled vertex: v $\Theta(n)$
- Label vertex v
- For all unlabeled neighbors w of v , $\Theta(\deg(v))$
if $d[v] + \text{cost}(v,w) < d[w]$
 $d[w] = d[v] + \text{cost}(v,w)$
 $p[w] = v$



$d[v]$ contain path lengths
 $p[v]$ contain paths themselves in reverse order.

ADT is a description of the functionality of a data structure

Single source shortest path (directed graph w non-negative edge weights):

Initialize structure:

- For all v , $d[v] = \text{"infinity"}$, $p[v] = \text{null}$
- Initialize source: $d[s] = 0$
- Initialize priority (min) queue

Repeat these steps n times:

- Find minimum $d[]$ unlabelled vertex: v
- Label vertex v
- For all unlabelled neighbors w of v ,

If $d[v] + \text{cost}(v,w) < d[w]$
 $d[w] = d[v] + \text{cost}(v,w)$ Same cost.
 $p[w] = v$

analysis is the same as prim's.

Dijkstra's correctness:

1. Assume Dijkstra's algorithm finds the correct shortest path to the first k vertices it visits (the cloud). ...

2. But that it fails on the $k+1$ st vertex, u .

Q is not the best

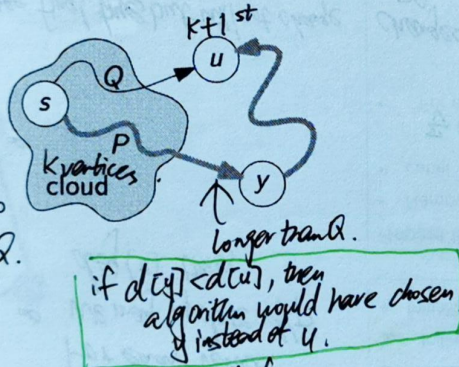
3. Then there is some other, shorter path from s to u . Call it P . P is better than Q .

4. There must be a node other than u , outside the cloud through which P passes. Call it y .

5. The path from s to y is at least as long as the path from s to u , since Q is shortest path out of the cloud. if it weren't, we should have chosen y instead of u .

6. P is even longer!! But that's a contradiction.

7. So our assumption that we failed on the $k+1$ st vertex is incorrect.



Let u be the last vertex ($k+1 = n$).

Dijkstra's

Prim's Algorithm (undirected graph with unconstrained edge weights):

Initialize structure:

- For all v , $d[v] = \text{"infinity"}$, $p[v] = \text{null}$
- Initialize source: $d[s] = 0$
- Initialize priority (min) queue
- Initialize set of labeled vertices to \emptyset .

Repeat these steps n times:

- Remove minimum $d[]$ unlabelled vertex: v
 - Label vertex v (set a flag)
 - For all unlabelled neighbors w of v ,
- If $d[v] + \text{cost}(v,w) < d[w]$
 $d[w] = d[v] + \text{cost}(v,w)$
 $p[w] = v$

	adj mtx	adj list
heap sorted array	$O(n^2 + m \log n)$ $= O(m \log n)$	$O(n \log n + m \log n)$ $= O(m \log n)$
Unsorted array	$O(n^2)$	$O(n^2)$

Which is best?

Depends on density of the graph:

Sparse heap + adj list.

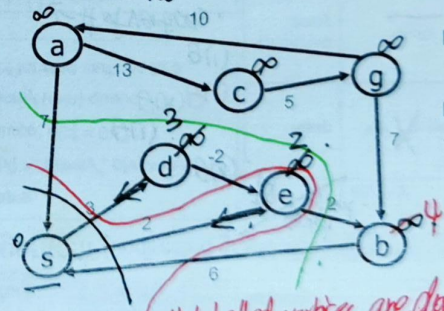
Dense unsorted array.

Single source shortest path (directed graph w non-negative edge weights):

Dijkstra's Algorithm (1959)

Why non-negative edge weights??

neg edge weights may reduce total path length after we call a vertex finished.



Initialize structure:

Repeat these steps:

- Label a new (unlabelled) vertex v , whose shortest distance has been found
- Update v 's neighbors with an improved distance

all labelled vertices are done.
 not $s \rightarrow d \rightarrow e$ is shorter than $s \rightarrow e$.
 and we do not improve it.