

se 2

- An algorithm is a way of defining computation
for solving some problem
- Generate output from an input
 - To do this
 - Is the most analog of a computer program

We write them in pseudocode
in a way that makes sense

It runs in model of computation
that says what your computer
is allowed to do

→ specifies what
operations can
be done and
the cost (time)
of them

other model

like RAM and other machine
models (reference)

- 1 "listy" = array $\rightarrow O(1)$
- 2 objects with $O(n)$ attributes (with reasonable number of attributes)

If we do `L.append(x)`, instead
of copying all elements in
a new array, Python uses
table shifting with $O(1)$

`L.append` $\rightarrow O(n \log n)$

`dict[key]` $\rightarrow O(1)$

Document distance problem

$d(D_1, D_2)$

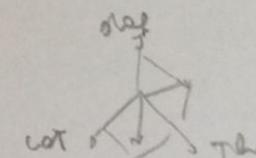
↓ ↓
document document

whether a document
is similar to
another and how much

document \rightarrow sequence of words,
each word is a string
of alphanumeric chars.

Ideas: shared words, to define distance
A document can be a vector: $D_{[w]}$

$$d(D_1, D_2) \approx D_1 \cdot D_2$$



→ coordinates
of word w
in A

find the angle
between the
vectors

Algorithm analysis - meant sort theory of every

input

- Input size, varies based on the problem, could be array, two integers, or graphs. In each case we count differently

- Running time, given an input, n the number of steps needed. We want to define steps as machine independent as possible, so we define each line as 1 step.
The notation is good from our RAI and Pendleton viewpoint.
But it is still messy and we don't need all of that info.

Worst-case

- We always want to consider the worst case,
- that it gives us an upper bound, we know that it will never take any longer
 - The worst case happens pretty often, for example especially when we search an element not present in a list.
 - The average case is usually roughly as bad as worst-case

Order of growth

We now consider the leading terms since the lower order terms are insignificant as n grows (the running time). We also ignore the constant coefficient on leading term

An algorithm is more efficient if its the worst case running time has lower order of growth

Due to constant factors and lower order terms, an algorithm with higher order of growth can take less time on algorithms with lower order on small inputs.

In large inputs the worst case of the fastest still run quickly than the slower ones
 $O(n)$ $O(n^2)$ $O(n^3)$

Divide and conquer
When an algorithm contains recursive calls to itself we can describe its running time by using a recurrence equation. If the problem is small enough it could take $\Theta(1)$ by straightforward, but suppose our problem is big enough and we can divide on a subproblem, each being $1/b$ the size of n (the initial size). It takes $T(n/b)$ to solve one of them, so takes $bT(n/b)$ to solve all of them plus the time to divide it again

divide and conquer

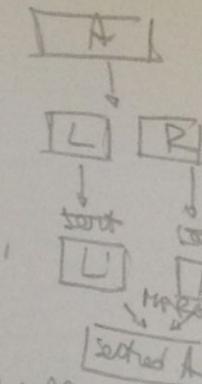
divide and conquer is typically a recursive approach.

It has three steps

- Divide into subproblems that are instances of the same problem

- Conquer the subproblems recursively, in the case the subproblem is small enough, resolve it straightforwardly

- Combine the solution to the subproblems into the solution for the original problem



Merge sort is the most common example.

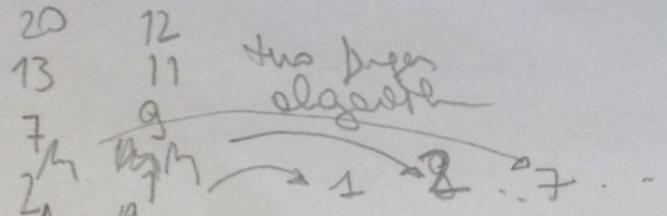
Here it goes; Divide an array into two subsequences of size $n/2$, where n is the length. Lengthen the two subsequences using merge sort. Merge the sorted sequences. Ends note n equals to 1

flat

the
interval
(L and R are
sorted)

Merge ($A[p, q, r]$) takes a subarray in A going from index p to q , and another from q to r , assuming that they are sorted, and it merges them into the subarray $p-q$. It takes $O(n)$ where $n = q-p+1$.

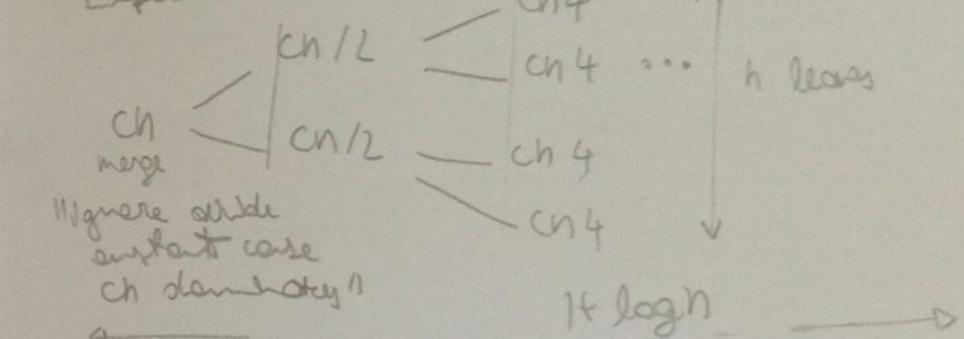
The comparison in merging is really effective since you only compare one element from each list



for splitting or
 $n/2$ complexity $\Theta(n)$

$$\text{Complexity } T(n) = \underbrace{C_1}_{\text{work done}} + \underbrace{2T(n/2)}_{\text{divide}} + \underbrace{C_2 n}_{\text{conquer + merge}}$$

Recall recurrence



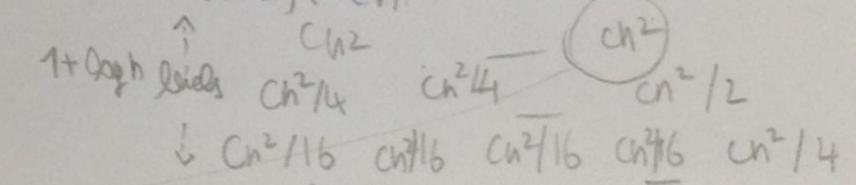
$\rightarrow h$ for each level and $\log n$ levels. $\Rightarrow n \Theta(n \log n)$

In-place sort does
in-place sorting will
be better, in merge
sort you need extra
space

in-place sort $\rightarrow O(1)$, the
temp variable
needed for merge

Recursive Solution

$$T(n) = 2T(n/2) + Cn^2$$



$\Theta(n^2)$ \rightarrow
Most work
done at the root

$n \approx 1$, we need m leaves $\frac{Cn^2}{K}$

Asymptotic notation

used to describe running time

It applies to functions for example insertion sort
 $\rightarrow cn^2 + bn + c$, we subtracted some details of
 this function to have $\Theta(n^2)$.

Asym. notation can be used to characterize other aspects of
 an algorithm, even with running time
 however, sometimes we're interested in
 worst case running time, others we
 want the running time no matter
 the input

Doubt of
 natural n
 numbers

for instance $\Theta(g(n)) = \{ f(n) : \text{running time } T(n) \text{ is running the}$
 $\{ f(n) \text{ there exists positive } n_0 \text{ such that we}$
 constants c_1, c_2 and n_0 such
 that $0 \leq c_1 g(n) \leq f(n) \leq c_2 g(n)$ for
 all $n \geq n_0\}$

+ function $f(n)$ belongs to the set $\Theta(g(n))$ if
 it exist constants c_1 and c_2 such that it can
 be sandwiched between $c_1 g(n)$ and $c_2 g(n)$
 for sufficiently sufficiently large n

by $R(n) = \Theta(g(n))$ we mean $R(n) \in \Theta(g(n))$

$R(n) \in \Theta(g(n))$ members must be
 asymptotically positive. This
 must hold for every asymptotic
 notation

For asymptotic notation we throw away constant
 leading coefficient and lower order terms
 to justify it more we have

next hand below. $C_1 n^2 \leq \frac{1}{2} n^2 - 3n \leq C_2 n^2$, we divide by n^2
 $\leq \frac{1}{2}$, left hand $\leq C_1 \leq \frac{1}{2} - \frac{3}{n} \leq c$ $\frac{1}{2} n^2 - 3n = \Theta(n^2)$
 We see, result holds so

lower order terms of an asymptotically
 positive function can be ignored
 because they are interpreted for
 large n . Some for a highest order.

O notation

gives an asymptotic upper bound. Give
 $g(n)$, we define $\Theta(g(n))$ as

$\Theta(g(n)) = \{ f(n) : \text{there exist positive } c$
 and n_0 such that

$$0 < f(n) \leq cg(n) \quad \forall n \geq n_0$$

Note that $f(n) = \Theta(g(n))$ implies
 that $f(n) = O(g(n))$

When $c > 0$, any given function with the
 form $cn + b$ is $O(n^2)$.

Usually O notation can cover every input
 being in upper bounded, differently for
 Θ that depends on the type of input

Ω pro notation

It provides an asymptotically lower bound
 for a given $g(n)$, $\Omega(g(n)) \approx$

$\Omega(g(n)) = \{ f(n) : \text{there exist } c \text{ and } n_0$
 such that $0 \leq g(n) \leq f(n)$ for
 $n \geq n_0$

6 for any $f(n)$ and $g(n)$, we have $f(n) = \Theta(g(n))$ iff
 $f(n) = O(g(n))$ and $\Omega(g(n))$

so a running time of $\Omega(g(n))$ means that no matter the
 input size n there, the running time is at least $g(n)$

Heaps.

Doubly queue

Implements a set of elements, each element is associated with a key

operations

$\text{insert}(S, x)$: insert element x into set S

$\text{max}(S)$: return the element with largest key

extract_max : // and remove it from S

$\text{increase_key}(S, x, k)$: increase key x 's key to k

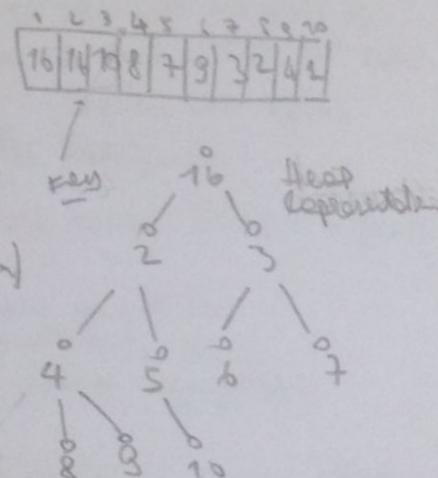
Heap is an implementation of the priority queue, it is a tree structure visualized as a nearly complete bin tree (should have a multiple of two to be complete, like 4, 8, 16, etc.)

Heap as a tree

root: first element ($i=1$)

parent(i) = $i/2$

$\text{left}(i) = 2i$ $\text{right}(i) = 2i+1$



Max-Heap property: the key of a node is \geq the keys of its children

Min-Heap ... \leq to

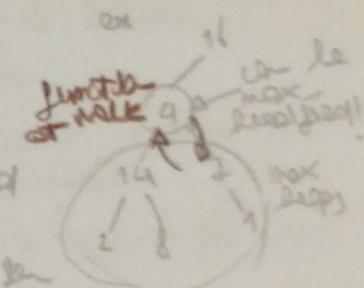
Heap operators

build_max_heap : produces a max heap from an unordered array (A)

max_heapsify : corrects a single violation of the heap property in a subtree's root

heapsort (A)

What of time you're gonna take up, because everything you add must be largest. Max-heapsify and by working from bottom up to less work



It repeatedly goes to max-heapsify down to the leaves when needed. $O(\log n)$

Heapsort

We take an unsorted array $A[1:n]$ to a max heap

for $i = n/2$ to 1
 $\text{max_heapsify}(A, i)$ leafs are $A[n/2+1:n]$

you are calling max-heapsify multiple times, into every child, you verify the properties, to start from $n/2$ to jump the leaves and then keep you from work your way up

Max-heapsify takes constant $O(1)$ time for nodes one level above the leaves and $O(1)$ for nodes 1 levels above and we have $n/4$ at level 1, $n/8$ at 2 ... 1 at $\log n$ level constant for the heap $O(n)$

A binary search tree is the tree representation of binary search.

Binary insertion problem

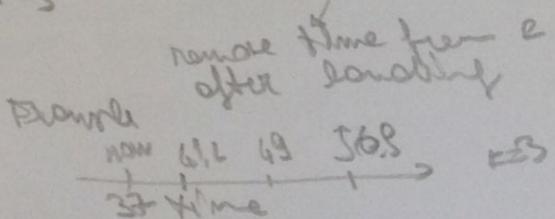
insert into a single binary

- reservations for future landings

- Reserve requests until leaving the T $|R|=n$

Add T to the set F, if no other landings are scheduled in K minutes

insert position will result in T to be cleared & ready to be able to insert



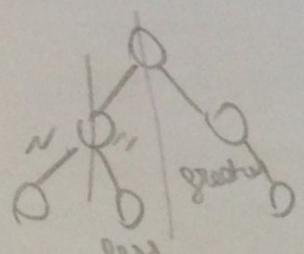
Visited array, BST:

Does not meet requirement as scheduling requests may take $O(n^2)$, max

Sorted array

Does not meet, word is well, insertion takes O(n) where you have to shift

Sorted list solves the insertion problem but word takes O(n)



Insert for all nodes x, y is the left subtree of x, $\text{key}(y) \leq \text{key}(x)$. y is right ... $\text{key}(y) \geq \text{key}(x)$.

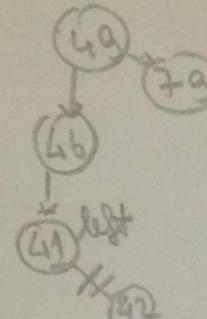
BST
insert in sorted array
example

- node x $\text{key}(x)$
- all pointers
like a tree
parent (y)
leftchild(x)
right(x)

-

inserting

insert 48
insert 79
insert 46
insert 41



If h is left,
insertion will
never be done in
 $O(h)$

(h because the tree
isn't necessarily balanced)

find_min()

go to the full list
leaf $O(h)$

next_larger(x)
(next larger of
e given value)

insert and
delete mostly
use number
that correspond
to subtree size
to compute rank(x)

insert

clear the
tree done
and insert
by 1 and
the new node
at 1

running problem

root 42
(k=3)

You proceed like
in binary search

42 with 49 ✓
42 with 46 ✓
42 with 41 ✗
(not k=3)

Segmented BSTs

For example we want
to be able to compute
rank(T):

How many places
are scheduled
to land at
time $\leq t$

Dont forget the
complexity!

Example

using the
root and
the leaves

subtrees
on 1 or leafs

rank
stuff

In game
first run
the regular
insert if it
has been done
immediately

Rec 3

In BSTs

Algorithm type
queries
updates

Not, bin
next level
result

If up rep.
invariant does not
hold, guides give many
answers

If you want to speed up
your code, you have
to make more tree
invariants still holds

+ Good way to keep
it to to reuse
dele_node(). Then delete it

red
invariant

DATA TYPES \dagger

to check
downward
in the code

invariant
all nodes left smaller, all
nodes right bigger of a
given element

- Fast search!

Python code

find()

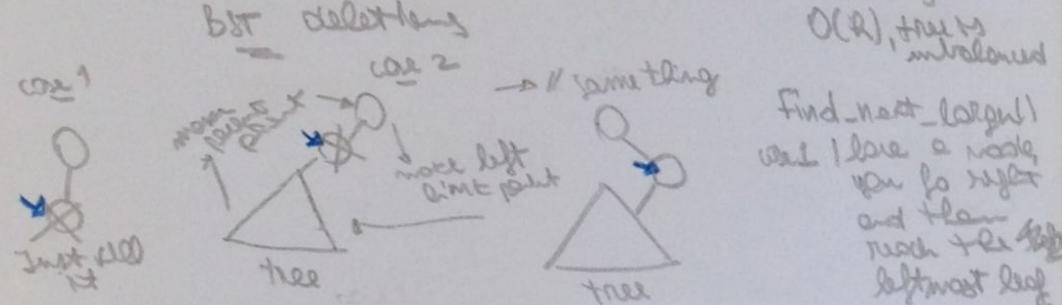
\rightarrow binary search
 $O(\log n)$, tree is
imbalanced

find_next_left()

well I lose a node,
you go right
and then
reach the left
leftmost leaf.

case 2
if it's a leaf
or has only
left, you go
up and then
do the thing

case 3
go up until
else to go
left



case 3
find_greatest
 \rightarrow
otherwise next greatest and
make it to the node
you're trying to
delete (x), delete
by replace

$O(h)$:

- Search the key

$O(1)$

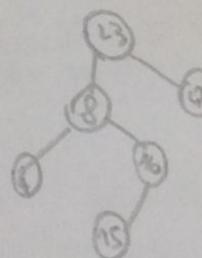
- case 3 (worst) takes $O(h)$ to find successor

BST example

nodes (py objects)

left
right
parent
key

each node
does it



I can augment a tree
by move a node further
by storing more info like

RSTs are trees,
some nodes have
only left/right,
others two children

Inverted if flat

$x \geq \text{left}$ and $x \leq \text{right}$.

You can lose the sorted
order by doing in-order
travel.

Operations are usually $O(h)$

AVL Trees

Requires heights of left and right children of every node to differ by at most ± 1 .

AVL trees are balanced.

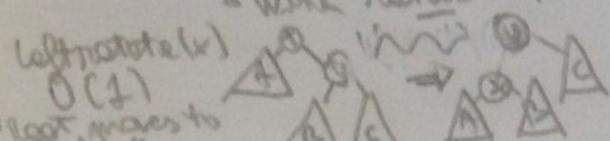
In the worst case is
when the right subtree
has height one more
than left for every node



AVL insert

- Simple BST insertion
- Fix AVL property

With rotations



The height of a tree is given
by the length of the longest
path from root to leaf

A tree is balanced if
the height is $O(\lg n)$

We use AVL trees to
maintain balance

The height of a node is
the longest path from
node to the leaf.

The root has height 0, i.e.,
max height among left
and right, + 1.

Icon we data structure objn.
to store height value

"will store null pointers
after leaves or -1 to
make it work"

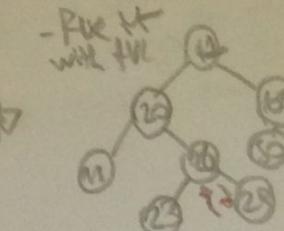
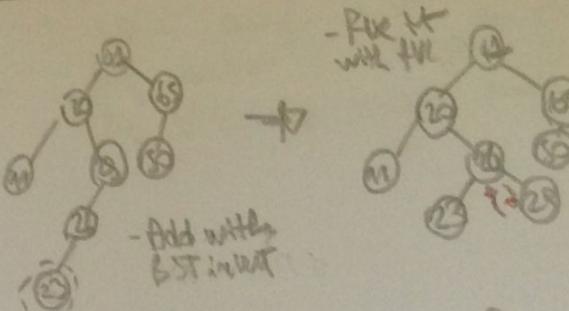
$N_h = \text{no. of nodes}
possible in an
AVL tree of height
h$

$$N_h = 1 + N_{h-1} + N_{h-1}$$

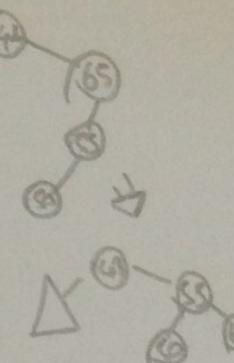
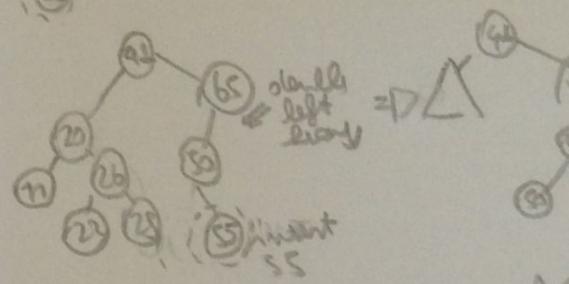
$$h \leq 2 \lg n$$

If we can ensure
this, the tree is
balanced!

restores the original
BST property



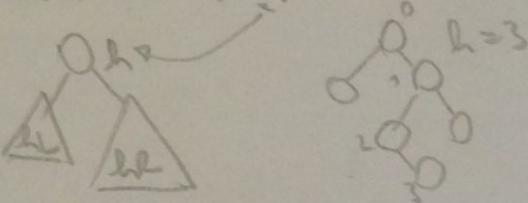
with a right rotation
of node "29", that
is double left
heavy



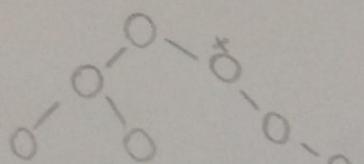
With a right rotation
I get the symmetric,
without rotation.
I make a double
rotation; the first
to eliminate the
L pattern to S,
then another one
since we are to
the previous cell

AVL trees get every
operation seen to
be in $O(\lg n)$

$$h = \max(h_L, h_R) + 1$$



This is how we update the height for a given node. We should do this for insertions and deletions.



Not AVL because not balanced, the node is less height bigger in right than left of mentioning.



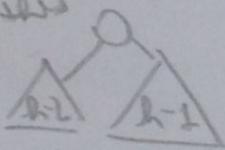
Invariant

$$\text{height difference} |h_L - h_R| \leq 1$$

is a BST

height $\sim O(\log n)$

If you want to build an AVL tree with the greatest nodes possible given a height, "do this"



For this problem only
 $N(h) = N(l-1) + N(k+1)$

$$N(h) \approx 1.6^h$$

$\approx \phi(\Phi)$

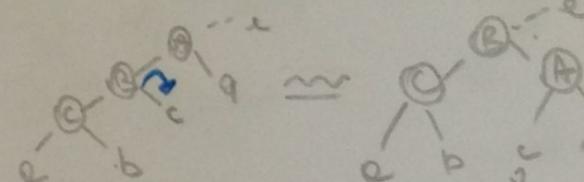
Insertion

- usual BST insertion going down path
- update the height of every node in the path

Deletion \sim similar

Rotation

You can left-rotate or right-rotate



periods

$$B.p = p$$

$$B.p = 2, B.p = 4 \\ l.p = B$$

$$l.l = B$$

c got moved in such a way in a BST, others are in same order

$$B.r = A$$

$$A.p = B$$

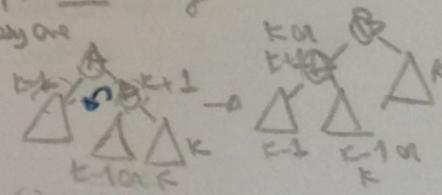
$$A.l = C$$

$$C.p = A$$

See a left rotation I did and $l.l$ and $B.p$.

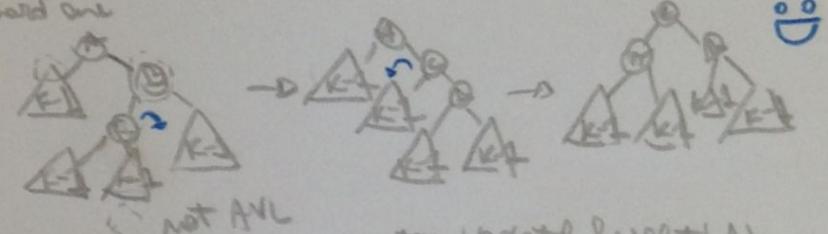
Rebalancing

= easy one



() difference $\geq 1 \Rightarrow \tilde{n}$

= hard one



NOT AVL

insert, update height(A)
 left, update height(B)

(the formula looks like the fibonaccii formula, with ± 1 , so they are higher)

Linear time sorting

comparison model, usually used for lower bound

- In that model, all input stems are black boxes (ADT)
- The data structure supports very comparison

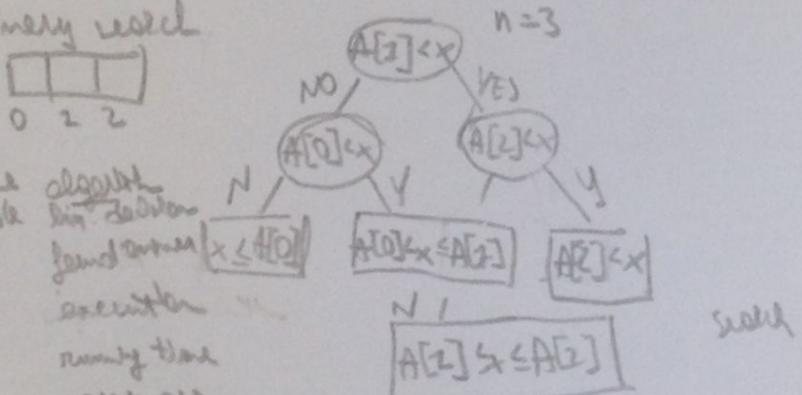
every sorting algorithm seen so far about comparison

- time cost = # comparisons

Decision tree

Any comparison algorithm can be viewed as a tree of all possible comparisons and their outcomes, and resulting answer. For any particular n :

e.g. binary search



decision tree algorithm	internal node	leaf
internal node	bin decision	fundamental
leaf	$x \leq A[0]$	$x \leq A[1]$
root-leaf	execution	
path length	running time	
height of tree	cost of sorting	running time

Sorting lower bound

to represent n things, picking a place them among the computer memory requires $\Omega(n \lg n)$ in the worst case

takes $\Omega(n \lg n)$ time

→ proof decision tree is binary and must have at least n leaves, one for each answer (can be more in other algorithms)

Sorting lower bound

Decision tree

YES
NO

$A[5] \leq A[3]$

leaf
the number of comparisons must be high

RAM model, memory is an array, i.e. access everything in constant time

Linear time sorting

- assume the keys were sorting are integers and they're positive. Each bits in a word can be manipulated in constant time

- can do more than just comparisons

→ then the height must be at least $\lg n$, where n is the running time

Sorting

Decision tree is binary and # leaves = # possible answers ($n!$)

$\Rightarrow \text{height} \geq \lg(n!)$

... not true

$\Rightarrow n \lg n - O(n) = \Omega(n \lg n)$

Counting sort

3 5 5 7 5 3 6

3 3 5 5 5 6 7

- Allocate an array of integer size k , sum + length in order and insert the key corresponding in the array when met

3 5 5 7 5 3 6

If you want to sort those trees, we use lists instead of counters

$O(k)$ $L = \text{array of } k \text{ empty lists}$

$O(n)$ for j in range(n):

$O(1)$ $L[\text{key}(A[j])].append(A[j])$

$O(1)$ $\text{output} = L$

$O(n)$ for i in range(k):

$O(1)$ $\text{output}.extend(L[i])$

$O(n+k)$

Sort by $S(n)$

insert $O(n^2)$
merge $O(n\lg n)$
Heap $O(n\lg n)$
counting $O(n+k)$
radix

Counting sort

4 1 2 3 3

Makes assumption that
the values goes from
0 to $K-1$ (K values)

0	1	2	3	4	5
0	1	2	4	0	3

1 2 3 3 7 - print into output

However we were
about stability: If you have two
key that look equal to the algorithm (but they)
might be different objects, the one
that shows up first in the input should
also show up first in the output

4 1 2 3^a 3^b 7 different objects

To solve stability:
1. Can implement
2. Insert list at each array element

0	1	2	3	4	5
1	2	3 ^a	4	3 ^b	7

for i=n-1 to 0 -
 pos[i] = sum[pos[i]] -
 num = pos[i]

 suppose we
 are doing this in c

0	0	1	2	4	5	5
0	1	2	1	0		

- We see how many keys there are
smaller than the searched
value

(not fast enough)

key (object):
2 x def -- cmp -- (self, other)
 $\begin{cases} < 0 \\ = 0 \\ > 0 \end{cases}$
(is not the
last time)

3 x def -- lt (self, other)
 $\begin{cases} < & < \\ = & = \\ > & > \end{cases}$
 $\begin{cases} < & < \\ = & = \\ > & > \end{cases}$
 $\begin{cases} < & < \\ = & = \\ > & > \end{cases}$
 $\begin{cases} < & < \\ = & = \\ > & > \end{cases}$
 $\begin{cases} < & < \\ = & = \\ > & > \end{cases}$
 $\begin{cases} < & < \\ = & = \\ > & > \end{cases}$

Radix sort

A key is a sequence of digits

2 3 4 1

1 4 3 2

2 4 1 3

1 2 4 3

2 1 4 3

1 2 3, 4 → already sorted

2 4 1 3

1 4 3 2 → stable

2 2 3 4 → unstable

2 3 4 1 → unclear

1 2 4 3 → top

2 1 4 3

As I go to the left
I will see the
sorted keys

Have last stable because
three unstable sorting, they
are also sorted according
to last the others

For each digit I run counting sort, $O(n+k)$,
so total $\rightarrow O(nk)$

$$\begin{aligned} d_n &= O(n) & \rightarrow O(n) \\ d &= O(1) \\ k &= O(n) \end{aligned}$$

HASHING !!

Dictionary ADT

mention a set of items, each with a key

insert (item) - (overwrites existing key)

delete (item)

search (key) - return item with given key or report doesn't exist

"In DST you can look up for next larger/bigger in dictionaries you are not allowed to do this"

We can have a $O(1)$ wif diff probability

Python: dict

$D[key]$ - search

$D[key] = val$ - insert

$\text{del } D[key]$ - delete

item = (key, value)

Motivation

- database
- database
- compilers and interpreters
- network router server
- indexing search
- strong communities
- file/disk synchronization

Simple approach

Direct access Table

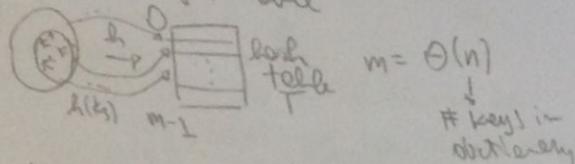
- store items in array indexed by a key
- associate with only an integer

Bad:

- keys may not be integers
- gigantic memory lag

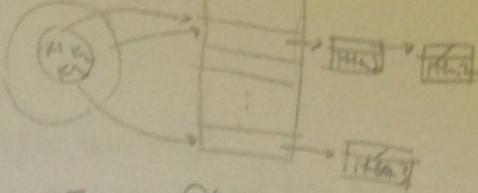
solution, hashing (but into pre-arranged bins)

- Reduce universe of all possible keys (integers) to a reasonable size m for table



If more keys map to slot by pigeonhole principle

↳ Clashing - if you have multiple items in the same slot, store them as a list



Unless you're really unlucky, the item will be at the top

of the list, with file assumption

Simple Uniform Hashing

This is false, but useful for analysis

Example: each key is equally likely to be hashed to any slot of the table independently of other key hashing.

Hash functions

Division method: $h(k) = k \mod m$

- feed only of m prime and n not very close to a power of 2 or 10

Analysis

expected length of slot

for N keys, m slots

$$\frac{N}{m} = d \text{ is the } (\frac{1}{m} \text{ prob}) \text{ load factor}$$

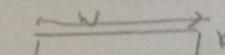
$O(1)$ if $m = O(n)$

running time = $O(1+d)$

(deter)

(intuition why lessening never)

Multidigit method: $[h(k) = ((k \cdot r) \mod 2^n) \gg (n-r)]$



$w = \text{word} \gg n$

/ shift right

right word

(take the right word and shift right by $n-r$)

$$m = 2^n, n \text{ is between } 0 \text{ and } m-1$$

r must be odd and not close to a power of two

UNIVERSAL HASHING

$$h(k) = [(ak+b) \mod p] \text{ mod } m$$

file want one key $k_1 \neq k_2$

$$\Pr[h(k_1) = h(k_2)] = \frac{1}{m}$$

some number p

random between 0 and $p-1$

Analysis

total expected slot length

$$n/m = d$$

$O(1+d)$ for

In python $\text{hash}(x) \gg$
The pre-hashed function,
given an object, produces
a string

- sometimes different strings are mapped to equal integers
- $\text{hash('186')} = \text{hash('D10C')}$

- hash - for user objects; if you don't want it in memory,

How to choose m ? $\rightarrow O(n)$ to make space smaller
 - want $m = \Theta(n)$
 $\rightarrow \Theta(n)$ in order for d to be constant

To do this, start small, $m=10$
 grow and shrink as necessary
 $\rightarrow n > m$: grow table $\rightarrow m \rightarrow m^*$
 Table doubling
 - insert takes $\Theta(k)$ time, $\rightarrow \Theta(1)$ amortized
 Also: \leftarrow deletions take $\Theta(k)$ time

Deletion
 - if $m = \frac{n}{2}$, then shrink $\rightarrow m$
 But if I delete 2 and insert $(2 \times 2 \rightarrow 2+1)$, we can get $\Theta(n)$ parity error

- if $m = \frac{n}{4}$, then the amortized time is $\Theta(1)$

Python lists: are resizable arrays, each operation is amortized constant (append, pop), they use table doubling

Simple algorithm
 any($s = target[1:i+1]$)
 for i in range($len(t)-len(s)$)

$\Theta(|s| \cdot (1 + |s|))$

Instead of clearing the array every time we could keep the last pointer of the string. Having k to 2^k reasonable size and put it in a vector

Rolling hash
 - given a list r , N like to be able to append a character and delete the first character by
 character $\rightarrow r.append(c)$
 character $\rightarrow r.pop(0)$ it means the string x
 $\rightarrow r()$ last value of $x = h(x)$

If we can do all of this in constant time, we can do binary matching with the following:

KARP-RABIN ALGORITHM

HOW TO BUILD THIS:
 divisor method
 $h(k) = k \bmod p$
 random prime

Treat x as multi digit number y in base c

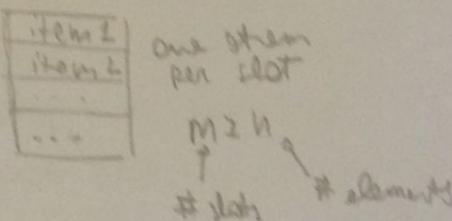
alphabet size
 $h = h \cdot c + ord(c)$
 extend
 $y \rightarrow y - c \cdot e$

STRING MATCHING
 Given two strings s, t :
 what google does is occur as a substring of t

$s = '6006' + [inset]$

Open addressing

Another approach for collision



Probing

We try to insert something, if you fail we compute a slightly different hash for that

Therefore, we continue to probe until we find an empty slot

$$h: V \times \{0 \dots m-1\} \rightarrow \{0 \dots m-1\}$$

key \mapsto index Total count produces a number from 0 to $m-1$

You ensure that you are using open addressing properly you want

$$h(k, 1) \ h(k, 2) \dots \ h(k, m-1)$$

relatively prime to m \rightarrow must be a permutation of $0, 1 \dots m-1$

Whether there is some free slot, eventually the sequence of probes will allow you to insert into one

Insert (k, v)

keep probing until an empty slot is found

New flag

Search (k)

A long as the slot encounters are not k , keep probing until other encounter may or find an empty slot

Delete (k)

comes used to find empty slot, so we replace delete stream with 'delete me' flag, so that next can still work. (Insert will now overwrite delete me flag)

Hashing strategies

Linear probing: $h(k, i) = (h'(k) + i) \bmod m$
cluster, concentrate groups of permutations ✓
clumped slots, which grows \rightarrow allows an insert encounters it

Double hashing

$$h(k, i) = (h_1(k) + i \cdot h_2(k)) \bmod m$$

\uparrow $h_2(k)$ is relatively prime to $m \rightarrow$ permutation

Uniform Hashing Assumption

Each key is equally likely to have any one of the $m!$ permutations and its probe sequence

$$\alpha = \frac{\eta}{m} \text{ lot of operations} \leq \frac{1}{1-\alpha}$$

Graph search

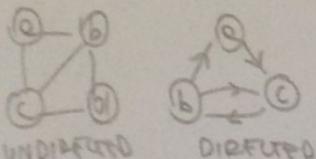
"explore" a graph

graph $G(V, E)$

V = set of vertices
 E = set of edges

UNDIR ↳ unordered pairs

DIR ↳ ordered pairs
 $e(v, w)$



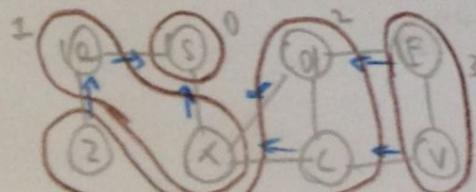
$$E = \{\{a, b\}, \{b, c\}, \{a, c\}\}$$

breadth first search

- Visit all nodes reachable from given $s \in V$
- $\Theta(V+E)$ time

- look at nodes reachable in $O(1, 2, 3, \dots)$ many steps required
 $\Theta(V+E)$

- Configure to avoid visiting vertices



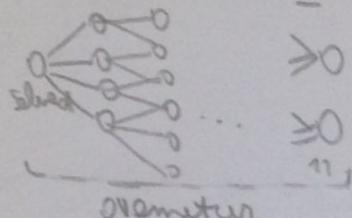
Applications

- web crawling
- social networking
- netflix recommendation

→ Array representation
is really terrible and expensive

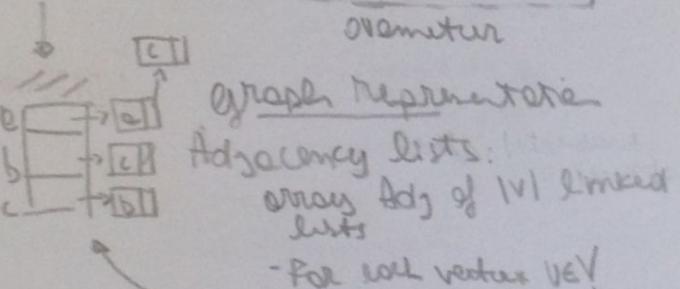
Poofset table ($2 \times 2 \times 2$)

- configuration graph
- vertex for each possible state
↳ neither $8!$ or 3^8
- edge for each possible move
0-0 undirected



Graph representation

Adjacency lists:
array Adj of $|V|$ linked lists



(could be arrays, hash table, objects, etc.)

shortest paths

$v \leftarrow parent[v]$
 $\leftarrow parent[parent[v]]$

$\leftarrow \dots$
 \leftarrow the shortest path from s to v .
its length is $length[s]$

BFS (s, Adj):

level = $\{s : 0\}$

parent = $\{s : \text{None}\}$

$i = 1$

frontier = $[s]$

while frontier:

next = []

for v in frontier:

for u in $Adj[v]$:

if u not in level:

level[u] = i

parent[u] = v

next.append(u)

frontier = next

$i + 1$

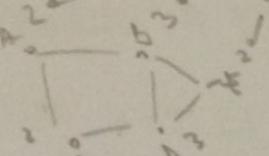
$$\sum_{v \in V} |Adj[v]| = 2|E| \rightarrow \text{Handshaking lemma}$$

$|E|$ directed

$O(V+E)$

you touch every vertex once when added to level

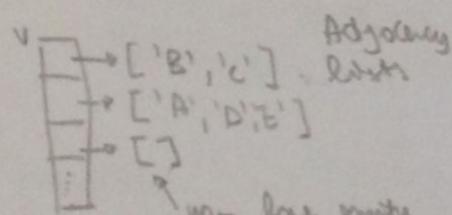
graph
degree, number of edges.



connected, there's a path from any node to any other node

If you add of all edges, you get $2E$. Handshaking Lemma

Python graph

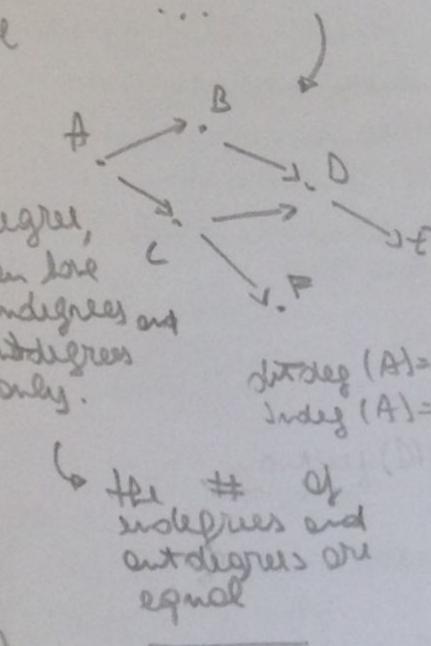
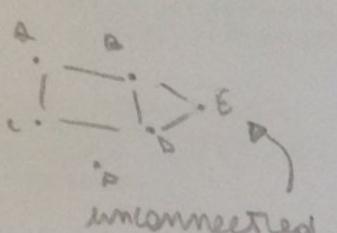


$O(V+E)$
running empty list

But neighbors
 $O(\deg(v))$

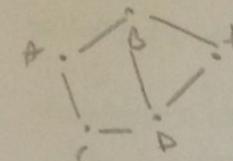
	A	B	C	D	E
A	1	1	1	0	0
B	1	1	0	1	1
C	0	1	1	0	
D	0	1	1	1	1
E	0	1	0	1	1

$2E + V$

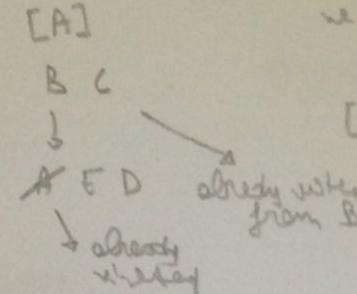


BFS

we have a graph



cannot be reached



For visited vertices we use checkboxes
set by $v \rightarrow \text{True}$

$[A, B, C, D, E]$ \rightsquigarrow a queue

$O(V + V)$
visit each node once

$O(V+E)$
=
already visited
already visited
 $O(E)$

code in python
you don't need to initialize when visiting a already

We want to find the shortest path from starting node (A in example)

Adjacency matrix

first neighbor
 $O(V)$

DFS used to explore the whole graph (not only the part reachable from s , like BFS)

- recursively explore the graph, backtracking as necessary
- careful not repeat vertices

parent $\{s : \text{None}\}$

DFS-Visit (V, Adj, s):

For v in $\text{Adj}[s]$:

"we have a vertex s , look at all outgoing edges if v not in parent, parent [v] = s

For each edge we check whether I already visited it (by setting parent)!"

Analysis

$$\Theta(V+E)$$

- visit each vertex once in DFS alone $\Theta(v)$

- DFS-Visit called once per vertex v

↳ $\Theta(|\text{Adj}[v]|)$

↳ $\Theta(E)$ excluding lemma for directed graphs

"Tree edges and back edges can't start in another unvisited graph!"

"an edge can only be of one type"

- forward edges: goes forward along the tree (like $a \rightarrow b$)
- backward edges: goes to an ancestor in the tree (like $b \rightarrow a$)
- cross edges: all other edges

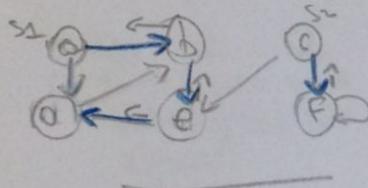
Edge classification

- tree edges (parent pointers): visit new vertex via that edge. They form a tree, like shown on ex.

- forward edges: goes forward along the tree (like $a \rightarrow b$)
- backward edges: goes to an ancestor in the tree (like $b \rightarrow a$)
- cross edges: all other edges

DFS ($V, \text{Adj}:$
parent = $\{\}$
for $s \in V$:
IF s not in parent
parent[s] = None
DFS-Visit (V, Adj, s)

↓
iterating over outgoing edges from s

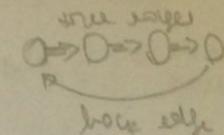


edge classifier are handy for cycle detection and topological sort

cycle detection

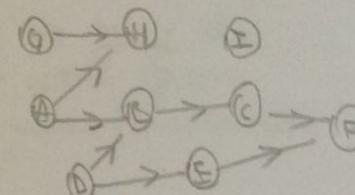
G has a cycle iff DFS has a back edge

Proof



In scheduling

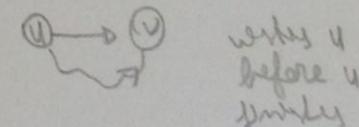
we are given a directed acyclic graph, i want to order vertices such that all edges point from lower order to higher order



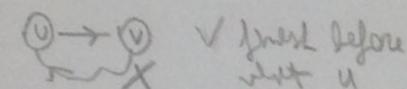
topological sort: run DFS and output the reverse of finishing times of vertices

for any edge $e = (u, v)$, v finishes before u finishes.

case 1: u starts before v

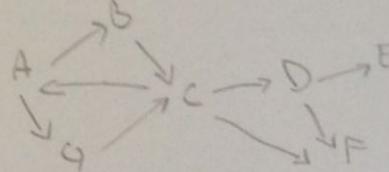


case 2: v starts before u



dependency list

```
{
    'A': ['B', 'C'],
    'B': ['C'],
    'C': [],
    'D': [],
    'E': []
}
```



class G

adj = { ... }

vertices():
 return adj.keys()

neighbors(v):
 return adj[v]

DFS(g, s):
 graph
 start point
 r = DFS result()
 r.parent[v] = None
 DFS_vset(g, s, r)

DFS_VISIT(g, v, r); r.visited[v] = True
 for n in g.neighbors(v)
 if n not in r.parents
 r.parent[n] = v
 DFS_VISIT(g, n, r)

DFS('A')

DFS_VISIT(A) B G

DFS_VISIT(B) C

DFS_VISIT(C) A D F

DFS_VISIT(D) E

DFS_VISIT(E) F

DFS_VISIT(F)

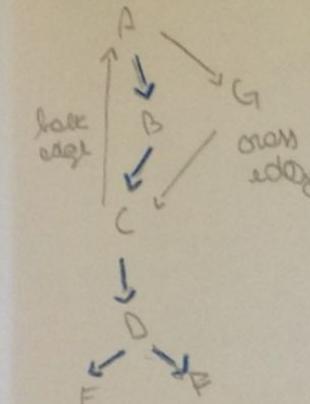
DFS_VISIT(G) C

not in parents, we must visit it

DFS result

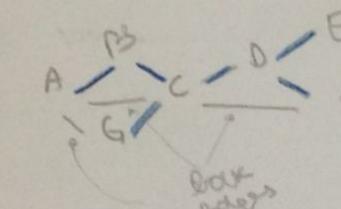
$v \rightarrow \text{parent}[v] = \text{None}$

$v \rightarrow \text{parent}[v] = \text{some other node}$



- if DF follows it; tree edge
- if v is descendant v - forward edge
- else v is ancestor v -> back edge
- else cross edge

For an undirected graph



1 (out + late cross edges and forward edges)

(when you're at an edge while PPT, in a directed graph you always see the forward edge later.)

In directed graphs see that immediately and you therefore mark it as backward edge)

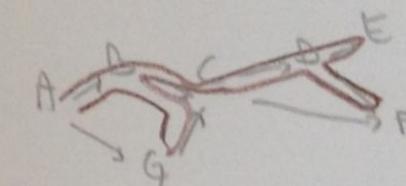
Topological Sorting!

graph of class prerequisites
ex: 6042 for 6.006

You can't have a cycle because you wouldn't have a dependency there
(In the graph in the previous page, you cannot start from A because you need C first)

must be directed

DAG (Directed, acyclic, graph)

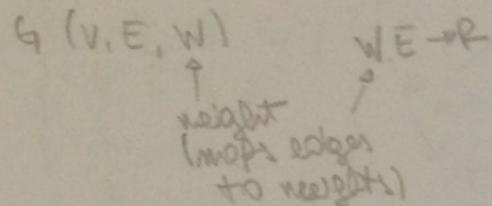


E F D C B G A

parents {
 'A': None
 'B': 'A'
 'C': 'B'
 'D': 'C'
 'E': 'D'
 'F': 'D'
}

tree edge
(comes
in of other
types)

Shortest path



Dijkstra's non-negative weight edges

Most of the time
E dominates. $O(V \log V + E)$

Dijkstra is a
nice algorithm
because it linear
on E $E = O(V^2)$ size

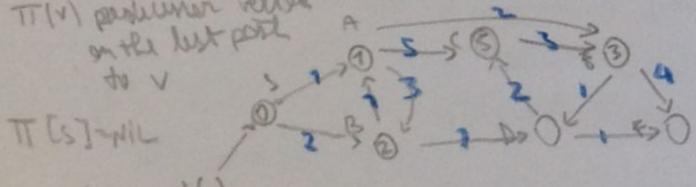
"W doesn't exist" \rightarrow the
complexities!

Weighted graphs

$$V_0 \xrightarrow{p} V_F$$

path from V_0
to V_F

$\pi(v)$ predecessor vertex
on the best path
to v



$d(v)$
weight
(shortest weight
to get here)

becomes delta
and you're
done

then start from the
final point to find
shortest path

Bellman-Ford
 $O(VE)$ works on positive
and negative weight
edges

path $p = \langle v_0, v_1, \dots, v_k \rangle$

$(v_i, v_{i+1}) \in E$ for $0 \leq i < k$

"a path is a sequence
of edges, each of them
must be in the graph"

$$w(p) = \sum_{i=0}^{k-1} w(v_i, v_{i+1})$$

"sum of the weight of edges"

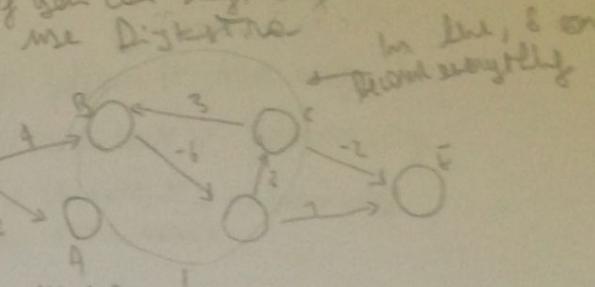
$$d(u, v) = \begin{cases} \min\{w(p) : u \xrightarrow{p} v \exists \text{ any such path}\} \\ \infty \text{ otherwise} \end{cases}$$

"Initially everything is infinity,
starting point is 0, then try to
reduce these infinities for all
reachable vertices

Negative weight
what?
never tell
that network

Negative weights
are a problem
because of negative
cycles that makes
the nodes in the
cycle indeterminate

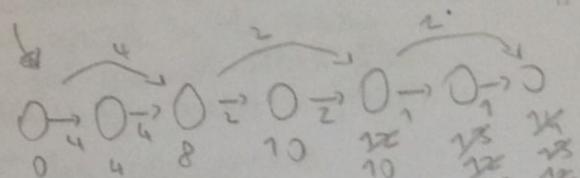
Bellman Ford computes δ
for all positive vertices
and marks the others as
indeterminate



Negative
cycle

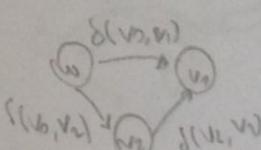
You can get to
B with weight -1,
there is no more
a lower bound of 0

Initialize for $\forall v$ $\delta[v] = \infty$ $\pi[v] = \text{NIL}$
 $d[v] = 0$ $\forall e$ $d[e] = \infty$ $w[e] = \text{NIL}$
repeat select edge (u, v) \rightarrow relax edge (u, v) \rightarrow means ($\delta[u] > d[u] + w(u, v)$)
 $d[v] = d[u] + w(u, v)$
 $\pi[v] = u$



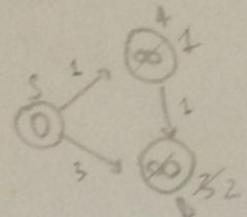
Optimal substructure

- subproblem of the shortest path
are shortest paths



expensive
relaxation
the Bellman
Ford
algorithm is
O(n^2)

Dijkstra



$d[v]$ = length of the current shortest path from s

$\delta(s, v)$ = length of a shortest path

$\pi(v)$ = predecessor of v in shortest path from s to v

Relax (u, v, w)

$$\begin{aligned} \text{if } d[v] > d[u] + w(u, v) \\ d[v] &= d[u] + w(u, v) \\ \pi[v] &= u \end{aligned}$$

Lemmas

The relaxation algorithm maintains the invariant that $d[v] \geq \delta(s, v)$ for all vertices

(you'll never get something that is less than your $\delta(s, v)$)

Proof. by induction on the # steps. Assume

$$d[v] \geq \delta(s, v) \text{ (base case)}$$

By Δ inequality $\delta(s, v) \leq \delta(s, u) + \delta(u, v)$

$$\delta(s, u) \leq d[w] + w(u, v)$$

DAGs

can't have (u, v)

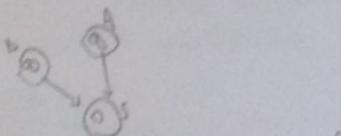
cycles

(can have negative edges)

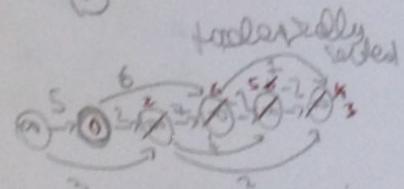
1) Topologically sorted DAGs

path from u to v
implies v is later
 v in ordering

$O(V+E)$ time



2) One pass over vertices in topologically sorted order relaxing each edge that leaves each vertex



pass a vertex everything on left is going to be marked as infinity

By taking the dems and pull it up (for Dijkstra).
Vertices on top are the predecessors.

If we take another vertex, the order will still be right

just piles

the min ->
from Q
and processes
it

Q is a priority queue, pushes
in diff values

that makes it
greedy

while
 $Q \neq \emptyset$:

• extract-min(Q)
Delete v
from Q

$S \leftarrow S \cup \{v\}$
for $v \in \text{Adj}[w]$
Relax(v, w)

I'm done, I
know the
shortest path

for $v \in \text{Adj}[w]$

Relax(v, w)

Dijkstra (Q, W, S)

InitHelds (Q, S)

$S = \emptyset$

$V[G]$

some nodes
marked from
as more
than no
options

Speeding up Dijkstra

"Not overall complexity
but average case"

We're solving single source
to any/ell destination
(single source single)
→ If you have
e target
→ meet

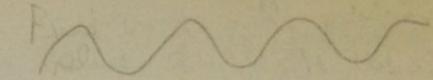
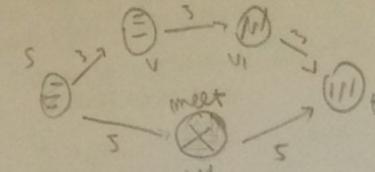
Dijk.
pseudo
Initialise ()
 $d[s] = 0$ $d[v \neq s] = \infty$

while $Q \neq \emptyset$

do $v \leftarrow \text{extract_min}(Q)$
for each vertex $v \in \text{Adj}[v]$
do Relax (v, v, w)
update priority
monitors parent pointers

(stop if $v = t$)

w may not be in Δ



Find a \sqrt{w} (out from w) with
 $d[v_1] + d[v_2] + d[v_3]$

After finishing scan, check
other two vertices

$$d[v_1] + d[v_2] = 3 + 6 = 9$$
$$d[v_1] + d[v_3] = 6 + 3 = 9$$

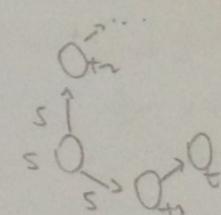
are less than
value found

Directed Level

Modify edge weight with potential
functions

$$w(v, v) = w(v, v) - \lambda(v) + \gamma(v)$$

↑
new weight



Correctness:

$$w(p) = w(p) - \lambda_t(t) + \gamma_t(t)$$

hold max level
preserves $\delta(u, l)$

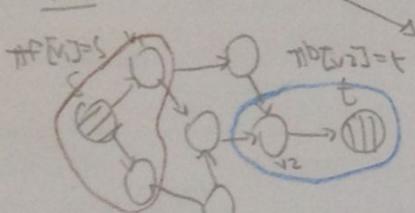
HLLV → for any
vertex, final δ
to l

any path from
 $u \rightarrow t$ will
be lifted if
this amount.
You will shorten
it faster

$$\gamma_t(u) = \delta(u, l) - \delta(t, l)$$

With the correct
choice of u , Dijkstra
will run faster

Bi-directional search



Edges must be traversable
backward also (directed
graph)

1 step standard Dijkstra then
1 step backward Dijkstra from

We have to make shorter
distance queue

$$d_f(s) = 0 \quad d_b(t) = 0$$

↓ everything else
in Δ into

distane in
forward
search

distance
backward
search

Q_F Q_B , priority queue
backward

Π_f → Π_b
priorities
forward

Π_b → Π_f
priorities
backward

Combine record when frontiers
meet, specifically,

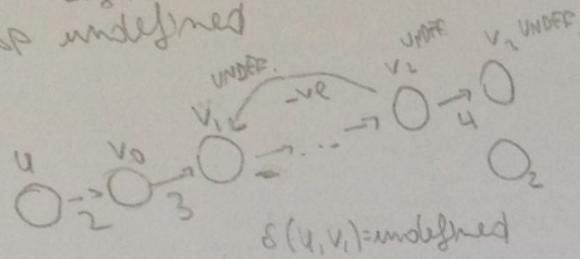
When you pull out s made from
 O_1 and O_2 , finish the record

loop: if w was produced
first from both O_1 , O_2 .

Find s in Π_f Π_b from Δ .
Find s in Π_b Π_f from Δ .

Bellman-Ford

- graph w/ negtive edges, in the case of negative cycle, this also returns SP undefined



Graph SP algo

Initalise for $v \in V$ $d[v] = \infty$ $\pi[v] = N/A$
 $d[s] = 0$

Main repeat, select edge (iteration)

Relax $(u, v, w) \leftarrow$ sets π pointers
edge

Until you can't
relax more

Bellman-Ford (G, W, S)

Initialise

for $i = 1$ to $|V| - 1$
for each edge $(u, v) \in E$

 Relax (u, v, w)

for each edge $(u, v) \in E$ If you can relax edges
 if $d[v] > d[u] + w(u, v)$ report that -ve cycle exists

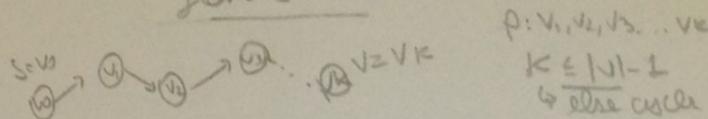
$d[v] = d[u] + w(u, v)$

$O(|E|)$

Relax $(u, v, w):$
if $d[v] > d[u] + w(u, v)$
 $d[v] = d[u] + w(u, v)$
 $\pi[v] = u$

Theorem: if G does not contain -ve weight cycles,
then BF after execution returns $d[v] = \delta(s, v)$
for all $v \in V$

Corollary: if a edge $d[v]$ fails to converge after $|V| - 1$
passes, there exists -ve wt cycle reachable
from s



Theorem
Proof

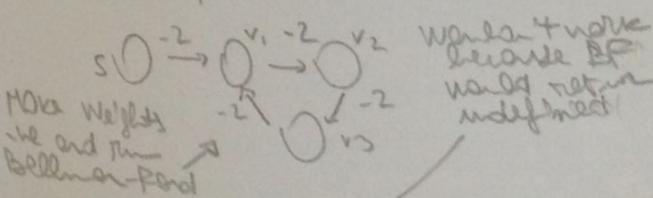
Let $v \in V$ $p = \langle v_0, v_1, v_2, \dots, v_k \rangle$ $v_0 = s$ $v_k = v$.
This path p is a shortest path with min # edges.
No -ve weight cycles $\Rightarrow p$ is simple $\Rightarrow k \leq |V| - 1$

- After 1 pass through E , no lone $d[v_1] = \delta(s, v_1)$, because we will relax (v_0, v_1) .
- After 2 passes, ... $d[v_2] = \delta(s, v_2)$ becomes in 2nd. pass
- After K passes, $d[v_K] = \delta(s, v_K)$
 $|V| - 1$ passes \Rightarrow all reachable vertices have δ value.

After $|V| - 1$ passes we find an edge that can be relaxed
 \Rightarrow current sp from s to some vertex v_N not simple because I have a repeated vertex.
Found a cycle, that is -ve

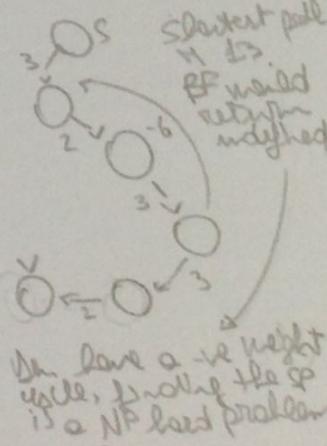
Corollary
Proof

Longest path



It is also
NP-hard

There is no
algorithm
better than
 \exp^{+} time



DP = Careful Whiteforce

Review

DP = Guessing + recursive + memoization

I'll say
at possibilities
for a car
so I can't
test one

In what ways,
usually exponential \rightarrow

Starts off fast enough
becomes polynomial

Time \propto subproblems

$(\text{bottom}(i) + \text{top}(j))$

Time \propto subproblems

DP = shortest paths in some DAG

Dynamic programming of n
the bottom hat hat
problem; given hat
problem is only DP problem

Time = # subproblems - Time / subproblem

5 easy steps to DP (how to apply in start ordered)

- ① Define subproblems \leftarrow subproblems for ~~themselves~~
- ② Guess (part of the solution) \leftarrow choices for just
- ③ Relate subproblem solution \leftarrow Time / subproblem
- ④ Compute & memoize or build table bottom up
- ⑤ Solve original problem

Text justification

Split test to great lines

Test = bit + g

— — — — |

words

— — — — |

bottom(i,j)

How it works

one words(i,j) =

$\begin{cases} (\text{page with note with}) \\ \in \text{one} \end{cases}$

use recursive to
compute for
topological sort

If small my car uses
most of the time,
spare, cables will be
over longer so this is
highly interleaved

① The subproblems are suffixes (words[i:])
 $\# \text{ subprob. } n$
where \propto short line given

② $\# \text{ choices } \leq O(n)$

③ ~~Time recurrence: $\text{DP}(i) \leftarrow \text{value to}$~~
~~subproblem~~
~~min (bottom(i) + top(j))~~
~~for j in range (i+1, n+1)~~
~~Time / subp = $O(n^2)$~~

④ Top order $n, n-1 \dots 0$
total: $\Theta(n^2)$

⑤ ORP
prob

DP(0)

Bottom up

- $\text{DP}[i] = 0, 1 \dots n-1$

- 1 player \rightarrow dealer

- 2 1 let/and

⑥ Interleaved subprob i:
few lines

⑦ Given $\text{bottom}(i)$

$\& \text{bottom}(j)$

many lines should

1 let &

⑧ Recurrence $\text{BJ}(i) = \max($

outcomes $\in \{-1, 0, 1\}$

$\# \text{BJ}(g)$

for j in range ($i+1, n$)

$\# \text{ valid play}$

parent pointers
remember which guess
was best

$\text{parent}[i] = \arg \min(\dots) = j$

$\emptyset \rightarrow \text{parent}[0] \rightarrow \text{parent}[\text{parent}[0]]$

1st line 2nd line 3rd line

DP
Subproblems for strings / sequences

- suffixes $x[i:] \& i$
- prefixes $x[:i] \& i$
- substrings $x[i:j] \& i \leq j$

$\Theta(n^2)$

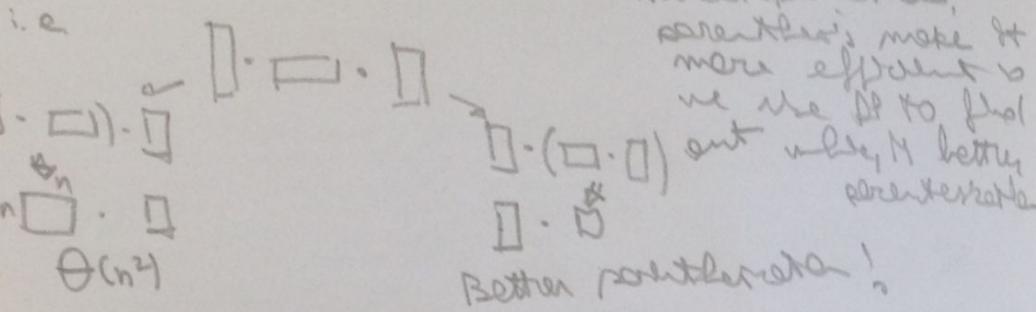
DP subproblems for

- ① in DP procedures.
- Reusing subproblems
- smaller the last part in DP programming

parenthesis matching

optimal evaluation of
bracelet expression

$$(A_0 \cdot A_1 \cdot \dots \cdot A_{n-1})$$



We start with guess

② the ~~optimal~~ multiplication
intermediate

$$(A_0 \cdots A_{k-1}) \cdot (A_k \cdots A_j \cdots A_{j-1}) = O(n)$$

nesting $(x[i:j])$ best multiplication # of pairs

② Recurrence
of $DP(i, j) \Rightarrow$

$$DP(i, j) = \min_{k \in \text{range}(i+1, j)} (DP(i, k) + DP(k, j) + \text{cost of } (A_i \cdot A_k \cdot A_{k+1} \cdots A_{j-1}))$$

For k in range($i+1, j$)

+ time / subprobz $\Theta(n^3)$

BEST WE
CAN DO
FOR KNOTS

not bad
 $y \leq n$ well
pseudo-polynomial
Time = $\Theta(n \cdot S)$

Edit distance
Given two strings x & y calculate the closest way to
convert x to y

three character edits:

- insert c
- delete c
- replace c $\rightarrow c'$

longest common
subsequence

HIERARCHY OF
MICHELANGELO

You can model it as a
edit distance problem

- edit cost / delete = 1
- cost replace of $c \rightarrow c' = \infty$ edit

① subprobz

edit distance or

$$x[i:j] \& y[i:j]$$

subprobz = $\Theta((k-1)k)$

② recurrence

$$DP(i, j) = \min \begin{cases} \text{cost of replace } x[i:j] \rightarrow y[i:j] \\ + DP(i+1, j+1) \end{cases}$$

$$\text{③ knapsack} \quad i \cdots \cdots (\text{nest of } \text{knapsack } x[i:j] + DP(i+1, j))$$

for $i = 1 \cdots 0$
for $j = 1 \cdots 0$

④	$\Theta(n^2)$	time = $\Theta(n^2)$
		best algo for edit distance

max sum
of values for a
subset of items
of $\geq x \leq S$

$$DP(i, S) = \max \begin{cases} DP(i+1, S) \\ DP(i+1, S-x) + v_i \end{cases}$$

② knapsack H : in
subset $\$$ or not

② knapsack = $\max_{i \in S} \#$ of subproblems
satisfying $x \leq S$

KNAPSACK

You can only
bring what you
can fit in your
backpack,
and several
value vi
integers

$$\# = \Theta(n \cdot S)$$