

# PDSA

Week -1

CLASSTIME Page No.

Date / /

- L1.1 Intro. to Jupyter NB & Google Colab
- manual - text editor to write code
    - run at the command line
  - IDE - single app. to write and run code
    - quick update - run cycle
    - debugging, testing, e.g. repl.it
  - share your code → collaborative dev.,  
export your results
  - documentation - interleave with the code
  - switching betw. diff. versions of code, export  
and import your proj., preserve your output
  - Jupyter NB
    - a seq. of cells, 1-D spreadsheet
    - cells hold on text, markdown nota" for  
formatting ([www.markdownguide.org](#))
    - edit and re-run indi. cells to update  
environment
    - supports diff. kernels → Julia, Python, R
    - widely used to document and disseminate  
ML projects (kaggle)
    - ACM software sys. award 2017
  - google colab → free to use  
[colab.research.google.com](#)
    - customized jupyter nb
    - all std. packages of ML are preloaded -  
scikit-learn, tensorflow
    - access to GPU hardware

L1.2 Implementation of Python Codes (Part-1)

- # bigger heading, ## smaller heading,
- list → gives' list ⇒ markdown lang.  
for text documents

Eg. def fac(n):  
 if n == 1:  
 return 1  
 else:  
 fl = [ ]  
 for i in range(1, n+1):  
 if n % i == 0:  
 fl.append(i)  
 if prime(i):  
 pt.append(i)  
 return pt

L1.3 Python Recap - 2: ~~fib~~ - word problem -  
 - gcd(m, n):  
 to largest k that divides both m & n  
 $\text{gcd}(8, 12) = 4$ ,  $\text{gcd}(18, 25) = 1$   
 - it always exists, if nothing then 1  
 -  $\text{gcd}(m, n) \leq \min(m, n)$   
 - def gcd(m, n):  
 indicates  $c = [ ]$  otherwise  
 list for i in range(1, min(m, n)+1):  
 if  $(m \% i) == 0$  and  $(n \% i) == 0$ : taken  
 c.append(i)  
 return (c[-1])

→ after if loop -  
 max = max of all factors -

- return (max) -  
 max = max of all factors -

L1.4 Python Recap - 3: ~~fib~~ or count -  
 - checking primality  $\Rightarrow$  Counting primes  
 - list the first n primes

def fp(m):  
 (count, i, pt) = (0, 1, [ ])  
 while (count < m):  
 if prime(i):  
 pt.append(i)  
 count = count + 1  
 i = i + 1  
 return (pt)

- directly check if n has a factor b/w.  
 def fp(n):  
 if True:  
 for i in range(2, n):  
 if  $(n \% i) == 0$ :  
 a = False  
 break # about loop  
 return a  
 - Properties of prime -  
 (i)  $\infty$  primes  
 (ii) twin primes:  $p, p+2 \Rightarrow \infty$  twin pairs  
 (iii) prime diff.  
 def pd(n):  
 up = 2  
 pdd = { }  
 for i in range(3, n+1):  
 if prime(i):  
 d = i - up  
 up = i  
 if d in pdd.keys():  
 pdd[d] = pdd[d] + 1  
 else:  
 pdd[d] = 1  
 return (pdd)

L1.5 Python Recap - 3

- d divides m & n  $\Rightarrow m = ad, n = bd$   
 $(m - n) = (a - b)d$   
 $\Rightarrow d$  also divides  $(m - n)$

def gcd(m, n):  
 $(a, b) = (\max(m, n), \min(m, n))$   
 if  $a \% b == 0$ :  
 return b # base case  
 else:  
 return (gcd(b, a - b))

- Euclid's algorithm

same as before

use:

`return gcd(b, a%b)`

- It is one of the first non-trivial algorithm.

### L1.6 Exception handling

- Common errors
  - (i)  $y = x/2, z = 0$
  - (ii)  $y = \text{int}(3)$  → not as int
  - recover gracefully, (iii)  $y = 5^x, x$  not defined using Exception
  - (iv) read a file, file doesn't exist
- syntax error → u can't connect lines
- correct  $\Rightarrow$  possible in names, zero division, index error
- unhandled exceptions aborts executions

try:

→ code

except IndexError:

→ handle index error

except:

→ handle other exceptions

else:

→ execute if try block w/o errors

- we can use exceptions freely, as in a dictionary

if b in scores.keys():
 scores[b].append(s)

else:

`scores[b] = [s]`

try:

`scores[b].append(s)`

except KeyError:

`scores[b] = [s]`

### L1.7 Classes & Objects

- Class - template for a data type, how data is stored, have public functions to manipulate data
- obj. - concrete instance of template

class Point:

~~def \_\_init\_\_(self, a=0, b=0)~~

~~self.x = a~~

~~self.y = b~~

~~initializes int. val.~~

~~self~~

~~a=0, b=0~~

~~parameter is self~~

Example:

class

from

on, dy

def translate(self, delx, dely)

~~self.x += delx~~

~~self.y += dely~~

def ad(self):

import math

d = math.sqrt(self.x\*\*2 + self.y\*\*2)

return d

- Polar coordinates

(a, b) instead of (x, y)

import math

class point

def \_\_init\_\_(self)

~~self.a = math.sqrt~~

~~(a\*\*2 + b\*\*2)~~

~~self.a = math.sqrt~~

~~(a\*\*2 + b\*\*2)~~

~~self.b = math.atan~~

~~(b/a)~~

~~if a == 0:~~

~~self.theta = math.pi~~

~~else:~~

~~self.theta = math.pi~~

~~b = math.atan(b/a)~~

~~but translate becomes more~~

~~elongated.~~

- \_\_init\_\_( ) - constructor

- str( ) - convert obj. to string

- add( ) - preceded by +

- mult( )

- div( )

- ge( )

### L1.8 Implementation of Python Codes (Part - 2)

- all codes in google colab -
- $\rightarrow p = \text{Point}(3, 4), q = \text{Point}(7, 10)$
- $\rightarrow p.\text{sd}(), q.\text{sd}()$

### L1.9 Timing Our Code

- library time  $\rightarrow \text{perf-time}()$
- 2 consecutive executions gives the int
- default unit is seconds.

import time

```
start = time.perf_counter()
```

# execute code

```
end = time.perf_counter()
```

```
elapsed = start - time
```

import time

```
class Timer:
```

```
def __init__(self):
```

```
    self.start_time = 0
    self.elapsed_time = 0
```

```
def start(self):
```

```
    self.start_time = time.perf_counter()
```

```
def stop(self):
```

```
    self.elapsed_time = time.perf_counter()
    self.end_time = self.start_time
```

```
def elapsed(self):
```

```
    return (self.end_time - self.start_time)
```

### L1.10 Implementation of Python Codes (Part - 3)

- refer to code
- create a timer - Time series
- raise 'an exception'

- trivial loop to see Python performs  $10^7$  operations in 1 second

### L1.11 Why Efficiency Matters?

- SIM cards linked to Aadhar Card
- Validation of Aadhar details to SIM card
- simple nested loop for matching the 2
  - $\hookrightarrow$  no. of sim cards  $\times$  no. of aadhar card
  - $10^{10}$  times of nested loop

$10^7$  in 1s

$$1 \rightarrow \frac{1}{10^7} \rightarrow 10^{11} \text{ seconds}$$

$= 3200 \text{ years}$

- propose a date (bdy)  $\rightarrow$  interval of (yes, earlier, later) possibilities  $\Rightarrow$  query the mp.  $\rightarrow$  halves the interval

$$365 \rightarrow 182 \rightarrow 91 \rightarrow 45 \rightarrow 22 \rightarrow 11 \rightarrow 5 \rightarrow 2 \rightarrow 1$$

- assume Aadhar details are sorted out as Aadhar number  $\rightarrow$  use the halving strategy to check each SIM card
  - $\hookrightarrow 30 \text{ queries} \times 10^9 = 3000 \text{ sec.} = 50 \text{ min}$
- arranging the data results in a much more efficient solution.
- both algorithms and data structures matter

### L1.12 Implementation of Python Codes (Part - 4)

- naive algo's stop work
- recursion viewer  $\rightarrow$  limit of remembering answers

## Week-2

### L2.1 Analysis of Algorithms

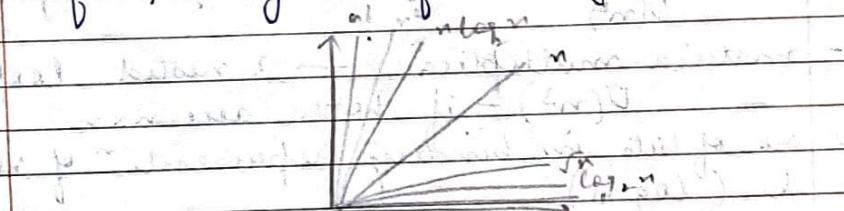
- naive approach takes 1000s of yrs, smart sol' takes a few minutes
- running time and space depends on storage is itd. by processing-time available memory (width of hardware)
- time depends on input size.  
input size  $\rightarrow n$   
run time  $\rightarrow t(n)$
- diff. inputs of size  $n$  take diff. amt. of time  
 $10^4 \Rightarrow n$  Naive  $\Rightarrow t(n) \approx n^{1/2}$   
(sum card q.) clever  $\Rightarrow t(n) \approx n \log_2 n$
- find closest pair of obj's.  
naive:  $\rightarrow n^2$   
clever:  $\rightarrow n \log_2 n$

- when comparing  $t(n)$ , focus on orders of magnitude, ignore constant factors
- asymptotic complexity - what happens in the limit, as  $n$  becomes large
- typical growth func<sup>ns</sup> = log, poly., expn.
- analysis should be independent of the underlying hardware  $\rightarrow$  don't use actual time, measure in basic op<sup>ns</sup> (comparisons, assignment)
- typically a natural parameter  $\rightarrow$  size of a list, no. of objects, no. of vertices and edges in the graph
- numeric problems? is  $n$  a prime  $\Rightarrow$  no. of digits is a natural measure of input size

- performance varies across input instances, take the 'avg.' behaviour  $\Rightarrow$  but this isn't practical.  $\therefore$  we take 'worst case' input  $\Rightarrow$  upper bound for worst case guarantees good performance

### L2.2 Comparing Orders of Magnitude

- $f(n) = n^2$  grows faster than  $g(n) = 5000n$
- $f(n)$  is said to be  $O(g(x))$ , if we can find constants  $c$  and  $x_0$  such that  $c \cdot g(x)$  is an upper bound for  $f(x)$  for  $x$  beyond  $x_0$ .  
 $f(x) \leq c \cdot g(x)$  for every  $x \geq x_0$



Ex.  $100n + 5$  is  $O(n^2)$

$$\rightarrow 100n + 5 \leq 100n + n = 101n^1, \text{ for } n \geq 5$$

$$\rightarrow 101n \leq 101n^2 \Rightarrow$$

$$\Rightarrow n_0 = 5, c = 101 \text{ on } n_0 = 1, c = 105$$

Ex.  $100n^2 + 20n + 5$  is  $O(n^2)$  but  $n^3$  is not  $O(n^2)$

If  $f_1(n)$  is  $O(g_1(n))$  and  $f_2(n)$  is  $O(g_2(n))$  then  $[f_1(n) + f_2(n)]$  is  $O(\max(g_1(n), g_2(n)))$   
 $\hookrightarrow$  least efficient phase is the upper bound for the whole algorithm

- $f(x)$  is said to be  $\Omega(g(x))$ , such that  $cg(x)$  is a lower bound for  $f(x)$  for  $x$  beyond  $x_0$   
 $f(x) \geq cg(x) \text{ if } x \geq x_0$  Ex.  $n^3$  is  $\Omega(n^2)$

- gen.: we establish lower bounds for a problem rather than an indi. algo.
- $f(x)$  is said to  $\Theta(g(x))$  if it is both  $O(g(x))$  and  $\Omega(g(x)) \Rightarrow c_1 g(x) \leq f(x) \leq c_2 g(x) \quad x \geq x_0$
- ex.  $\frac{n(n-1)}{2}$  is  $\Theta(n^2)$

### L2.3 Calculating Complexity - Exgs:

- iterative and recursive progs.
- max. element in a list  $\rightarrow O(n)$
- check whether a list contains duplicates  
 $\hookrightarrow$  worst case  $\rightarrow \text{len}(l)^2$  no duplicates  $O(n^2)$
- matrix multiplication  $\rightarrow$  3 nested loops  $O(n^3)$  - if both are  $n \times n$
- no. of bits in binary representation of  $n$   $\rightarrow O(\log n)$
- Towers of Hanoi  $\rightarrow$  3 pegs A, B, C  
 $\hookrightarrow$  recurrence

$$\begin{aligned} M(n) &= 2M(n-1) + 1 \\ &= 2^2 M(n-2) + (2+1) \dots \\ &= 2^k M(n-k) + (2^k - 1) \\ \dots &= 2^n - 1 \quad \Rightarrow \text{exponential} \end{aligned}$$

- iterative  $\rightarrow$  focus on loops
- recursive  $\rightarrow$  write and solve a recurrence

### L2.4 Searching in List

- is value  $v$  present in list  $c$ ?
- naive sol<sup>n</sup>  $\rightarrow$  scans the list, input size

$n$ , length of list  $\Rightarrow$  worst case complexity is  $O(n)$

- what if  $l$  is sorted in asc. order, comparison of  $v$  with the mp. of  $c$ , concept of binary search  $\rightarrow O(\log n)$
- $T(n)$ : the time to search a list of length  $n$

$$\left. \begin{array}{l} n = 0 \quad T(n) = 1 \\ n > 0 \quad T(n) = T(n/2) + 1 \end{array} \right\} \text{Recurrence}$$

$\Rightarrow$  solve by unwinding

$$T(n) = T(n/2) + 1$$

$$\begin{aligned} &= (T(n/4) + 1) + 1 = T(n/2^2) + 2 \\ &= T(n/2^k) + \underbrace{1 + \dots + 1}_k \end{aligned}$$

$$\begin{aligned} &= T(1) + k \quad \text{for } k = \log n = 2 + \log n \\ &\Rightarrow O(\log n) \end{aligned}$$

### L2.5 Selection Sort

- sorting list helps in  $\Rightarrow$  binary search, finding median, checking duplicates, building freq. table of values
- ex. TA for course
- Select the next element in sorted order, but try to avoid using the second list  $\Rightarrow$  eventually the list is arranged in place in ascending order
- correctness follows from the invariant efficiency  $O(T(n)) \rightarrow O(n^2)$
- this is intuitive algo. to sort a list, every case is worst case complexity itself with same time every time.

### L2.6 Insertion Sort

- same eg.
- move the first paper, second paper check whether higher or lower, and so on, insert each paper
- start building a new sorted list  $\Rightarrow$  pick next element and insert it into the sorted list  $\Rightarrow$  an iterative formulation
  - $\hookrightarrow$  assume  $L[:i]$  is sorted
  - insert  $L[i]$  in  $L[:i]$
- recursive way  $\rightarrow$  inductively sort  $L[:i]$ 
  - insert  $L[i]$  in  $L[:i]$
- correctness follows from the invariant
- efficiency  $\rightarrow T(n) \in O(n^2)$
- in case of recursive method,  
 $T(n) \rightarrow \text{insert } TS(n) \rightarrow T\text{sort}$   
 $T(n) \in O(n^2)$
- card induction  $\rightarrow$  we have to create a sorted list
- worst case is  $O(n^2)$  but the gen. term taken is  $O(n)$

### L2.7 Merge Sort

- to bring complexity below  $O(n^2)$ .
- divide the list into 2 halves
- combine 2 sorted lists A & B into a single sorted list C  $\Rightarrow$  merging A & B
- let 'n' be the length of L  $\Rightarrow$  now do the same halving thing recursively till only 1 element left
- divide and conquer algorithm  
 break - solve - combine the soln

- if  $\text{len}(A) \leq 1$ , nothing to do otherwise  $\Rightarrow$  sort and merge
- this was done to reduce the complexity of  $O(n^2)$  of selection and insertion sort

### L2.8 Analysis of Merge Sort

- merge A of length m, B of length n
  - $\hookrightarrow$  output list C has length  $(m+n)$
  - $\hookrightarrow$  in each iteration we add at least 1 element to C
  - $\rightarrow$  merge  $\rightarrow O(m+n) \xrightarrow{(m=n)} O(n)$
- in recurrence,
 
$$T(n) = 2^n T(n/2^n) + kn$$
- $O(n \log n) \rightarrow$  effective  $\Rightarrow$  best time possible in sorting
- recursions on merge are possible
- merge always creates a new list to hold the merged elements, this and recursive calls and returns are expensive

### L2.9 Implementation of Searching & Sorting Algo.

- set up Timer class to time executions
- naive search & binary search
 

10.6 s	0.83 s
--------	--------

$\hookrightarrow$  for looking for a odd no. in  $10^5$  even nos.
- selection sort  $\rightarrow 1.26$  sec.
- insertion sort  $\rightarrow$  on already sorted, performance is good ( $0.008$  sec) but when not sorted even bad than selection sort ( $2$  sec.)

Week - 3L 3.1 Quick Sort

- new list and costly new storage and recursive calls and cutovers are expensive in merge sort
- divide and conquer w/o merging  
 $T(n) = 2T(n/2) + n \Rightarrow T(n) = O(n)$   
 $\therefore T(n) = O(n \log n)$
- how to find the median? Instead <sup>pick</sup> some element in L — Pivot by (C.A.R. Hoare)  
 ↳ partition L into lower and upper parts w.r.t to pivot.
- partitioning is imp. — Pivot, Lower, Upper, Unclassified
- in this we avoid a merge step  $\Rightarrow$  allows ~~can~~ in place sort

L 3.2 Analysis of Quick Sort

- partitioning w.r.t. pivot  $\Rightarrow O(n)$   
 If the pivot is median  $\Rightarrow O(n \log n)$
- Worst case? Pivot is max or min  $\rightarrow O(n^2)$   
 ↳ already sorted array : worst case!
- However, the avg. case is  $O(n \log n)$
- any fixed choice of pivot allows us to construct worst case input, instead, choose pivot position randomly at each step  
 ↳ expected running time  $\Rightarrow O(n \log n)$
- recursive calls, in practice, quick sort is very fast  $\rightarrow$  excel, python

L 3.3 Implementation of Quick Sort algo.

- merge sort  $\rightarrow 10^6$  input  $\Rightarrow 9$  sec. } random
- quick sort  $\rightarrow 10^6$  input  $\Rightarrow 7$  sec. }

### L3.4 Concluding remarks on sorting algos.

- stable sorting — often list values are tuples, sometimes we have to sort out things which are already sorted → this is crucial in many app.
- quick sort is not stable (swapping values disturbs the order)
- merge sort is stable if we merge carefully, minimizing the data movement
- merge sort is used for 'external' sorting
- heapsort —  $O(n \log n)$
- can use hybrid structure, use divide and conquer till  $n$  is very large, later change the strategy

### L3.5 Diff. b/w lists & arrays

- storing a seq. of values — list, array.
- list
- flexible length, easy to modify the structure, values are scattered in memory
- array
- fixed size, allocate a block of memory, supports random access
- list → seq. of nodes, 'linked' list, easy to modify, insertion and del" is easy via local 'plumbing', flexible size, takes time  $O(i)$
- arrays → fixed size, contiguous block of memory, random access, accessing any  $i$  takes same time, inserting and deleting elements is expensive

- exchanging  $A[i]$  &  $A[j]$  → constant - array
- $O(n)$  - lists
- delete  $A[i]$ , insert  $v$  → constant - list
- $O(n)$  - arrays
- need to keep implementation in mind while analysing data struc.
- binary search work well with array, but insertion is problem with array
- either way —  $O(n)$
- how does it work in Python?

### L3.6 Designing a Flexible List and Ops on the same

- Python class Node → list is a seq. of nodes
- self. value is the stored value
- self. next pts. to next node
- l = Node() → empty list
- node(5) → singleton list
- appending to a list
- insert at the start of the list
- delete a value  $v$
- use a linked list of nodes to implement a flexible list

### L3.7 Implementation of Lists in Python

- accessing an element in lists →  $O(n)$
- python lists are not implemented as flexible linked lists → it maps to the memory and assign a fixed block for the list
- l.append() and l.pop() —  $O(1)$
- insertion & deletion —  $O(n)$

- effectively, Python lists behave more like arrays
- arrays are useful representing matrices  

$$\begin{pmatrix} 0 & 1 \\ 1 & 0 \end{pmatrix} \quad [[0, 1], [1, 0]]$$

`zlist = [0, 0, 0]`  
`zmatrix = [zlist, zlist, zlist]`  
 (mutability aliases different values)  
`zmatrix = [[0 for i in range(3)] for j in range(3)]`

on

`import numpy as np`

- `zmatrix = np.zeros(shape=(3, 3))`
- `arrange` is eg. of range of lists
- numpy arrays are easier to use

- L 3.8 Implementa<sup>n</sup> of dict. in Python
- array / list allows access thru positional indices
  - dict. allows access thru arbitrary keys, access time is same for all keys
  - underlying storage is an array
  - Hash func<sup>n</sup>  $\rightarrow h(s) \rightarrow X$ , maps a set of values  $S$  to small range of integers  
 $\rightarrow |X| \ll |S|$ , collisions  $h(s) = h(s')$ ,  $s \neq s'$   
 $\rightarrow$  a good "hash func" will minimize collisions
  - $\rightarrow$  SHA-256  $\rightarrow$  range is 256 bits  
 $\hookrightarrow$  avoid uploading duplicate for cloud storage
  - Hash-table - An array of size  $n$  combined with  $h$

- deal with collisions
  1. open addressing (closed hashing)
    - $\hookrightarrow$  probe a seq. of alternative slots in the same array
  2. open hashing
    - $\hookrightarrow$  each slot in the array pts. to a list of values, insert into the list for the given slot
- dict. key must be immutable  $\Rightarrow$  if value  $\Delta s$ , hash also  $\Delta s$
- many heuristics / optimizations possible for dealings

- L 3.9 Implementa<sup>n</sup> of diff. betw. list & arrays
- create an <sup>empty</sup> list of  $10^5$  appends  $\rightarrow 1.549$  sec.
  - create an empty list,  $10^5$  inserts  $(0, i)$   $\rightarrow 2.4$
  - array searching  $\rightarrow$  in numpy naive search  $\rightarrow$  takes too long time  
 binary search  $\rightarrow$  works much better
  - but other sorting algo on arrays takes much longer time on indexing

## Week 4

### L 4.1 Introduction to Graphs

- visualizing relations as graphs  $A \subseteq T \times C$
- $A = \{(t, c) | (t, c) \in T \times C, t \text{ teaches } c\}$
- Graphs:  $G = (V, E)$  → set of edges is a set of vertices
- Directed graph  $\Rightarrow (v, v') \in E \Rightarrow (v', v) \notin E$
- Undirected graph  $\Rightarrow (v, v') \in E \Rightarrow (v', v) \in E$
- A path is a seq. of vertices  $v_1, v_2, \dots, v_k$  connected by edges. It does not visit a vertex twice. A seq. that visits a vertex is called a walk.
- Paths in directed graphs. Vertex is reachable from vertex  $u$ , if there is a path  $u \rightarrow v$ . Shortest? Connected graph?
- Map colouring - states neighbours should have a different colour, how many colours needed?  $(v, u) \in E \Rightarrow c(v) \neq c(u)$
- 4 colour theorem - for planar graphs from geographical maps, 4 colours suffice
- vertex cover - min. no. of cameras needed so that all coverings are covered.  $V \rightarrow \text{vertices}$ ,  $E \rightarrow \text{coverages}$  segments
- matching; class project with funds
- A parent B (directed),  $A \xleftarrow{\text{fund}} B$  (undirected)

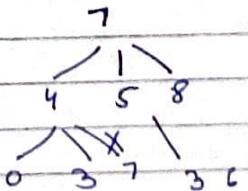
### L 4.2 Representing Graphs

- adjacency matrix  $A[i][j] = 1, \text{ if } (i, j) \in E$
- edges = [list]
- import numpy as np
- $A = \text{np. zeros}(\text{shape}=(10, 10))$
- $\text{for}(i, j) \text{ in edges}$
- $A[i][j] = 1$

- undirected graph is symmetric across main diagonal.
- neighbours of  $i \rightarrow$  column  $j$  with entry 1, scan row  $i$  to identify neighbours of  $i$
- degree of vertex  $i$ , for directed graphs, outdegree and indegree
- adjacency matrix  $\rightarrow$  size is  $n^2$   
undirected:  $k \leq n(n-1)/2$
- directed:  $|E| \leq n(n-1)$
- adjacency list - list of neighbours for each vertex
- $A \text{ list} = \{\}$       adj. less space
- $\text{for } i \text{ in range}(10):$
- $A \text{ list}[i] = []$
- $\text{for } (i, j) \text{ in edges:}$
- $A \text{ list}[i]. \text{append}(j)$
- choose representation depending on requirements

### L 4.3 BFS

- explore the graph level by level, each visited vertex has to be explored
- initially visited ( $v$ ) = False
- A queue - FIFO (first in, first out)



7 → 4 → 5 → 8 → 0 →

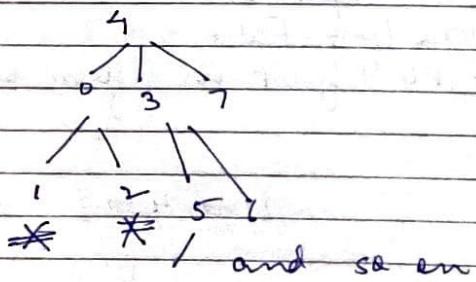
BFS

going like this

- BFS for adj. matrix —  $O(n^2)$
- BFS for adj. list —  $O(n+m)$
- (eg. of amortized analysis)  
if  $m \ll n^2$ , working with adj. lists is much more efficient, for graphs  $O(n+m)$  is best possible complexity
- recover path from  $i$  to  $k$ ? from  $k$ , follow parent links to trace back a path to  $i$ .
- by exceeding the level at which a vertex is visited, we get its distance from the source vertex
- BFS is strategy to explore a graph, lvl. by lvl.
- parent info. + lvl. info. → edges are also labelled with a cost, this is known as weighted graphs.

#### L4.4 DFS

- stack — LIFO (last in, first out)  
(most recently suspended is checked first)



- DFS is most natural to implement recursively, for each unvisited neighbour of  $v$ , call  $\text{DFS}(v)$

- no need to maintain a stack, makes visited and parent global, recursion maintains stack, just separate the initialisation step
- same complexity as of BFS  
 $O(n^2)$  — adj. matrix  
 $O(n+m)$  — adj. list
- paths discovered by DFS are not shortest paths, unlike BFS. It records the order.

#### L4.5 App. of BFS and DFS

- undirected graph is connected if every vertex is reachable from any other vertex
- in a disconnected graph, we can identify the connected components → use the component number strategy
- a cycle is a path, that starts and ends at the same vertex, it shouldn't repeat edges.
- a graph is acyclic, if it has no cycles
- edges explored by BFS form a tree (tree per component, collection of trees is forest)
- in a directed graph, a cycle must follow the same direction
- diff. type of non-tree edges — Forward, Back, Cross edges  
only back corresponds to cycles
- Tree edge | forward edge  $(u, v)$   
Internal [pre( $u$ ), post( $u$ )] contains [pre( $v$ ), post( $v$ )]
- Back edge  $(u, v)$
- Cross edge  $(u, v)$  → intervals are disjoint

- L4.6 Intro. to DAGs
- tasks and dependencies e.g. startup moving into new office spaces
- 
- ```

graph TD
    A[conducts (E)] --> B[Tiling]
    C[Conduit (T)] --> D[Plastering]
    C --> E[Painting]
    B --> F[Wiring (E)]
    B --> G[Cabling (T)]
    E --> H[Fitting (E)]
    style A fill:none,stroke:none
    style C fill:none,stroke:none
    style D fill:none,stroke:none
    style E fill:none,stroke:none
    style F fill:none,stroke:none
    style G fill:none,stroke:none
    style H fill:none,stroke:none
    
```
- we have to find a schedule.  $\rightarrow$  Topological Sorting, longest paths

- overall,  $O(n^2)$
- using 'adj. list' -  $O(m+n)$
- always, more than 1 topological sort is possible, it depends on choice of which next vertex

#### L4.8 Longest paths in DAGs

- Topological sort gives feasible schedule
- e.g. DAGs of processing of queues
- longest path
- if indegree ( $i$ ) = 0, longest-path-to ( $i$ ) = 0
- longest-path-to ( $i$ ) =  $1 + \max_{(j,i) \in E} \{ \text{longest-path-to}(j) \}$
- longest path algo.
  1. compute indegrees of each vertex
  2. initialize to 0 (longest path)
  3. list a vertex with indegree 0 and remove it
  4. update indegrees, longest path
  5. repeat till all vertices are visited
- we use adj. lists for the same  $\rightarrow O(m+n)$
- longest path makes sense for cyclic graphs (path has at most  $(n-1)$  edges)
- computing longest paths in arbitrary graphs is much harder than for DAGs - no good strategy is present

- L4.7 Topological Sorting
- enumerate  $V = \{0, \dots, i\}$
- claim: Path  $i \rightarrow j$ ,  $i$  must be listed before  $j$ . Every DAG can be topologically sorted.
- start with vertices and no dependencies, then list vertices whose dependency has already been listed
  - every DAG has a vertex with indegree 0.
  - topological sort algo.
    1. compute indegrees of each vertex
    2. list a vertex (0 indegree) and remove it from DAG
    3. update indegrees
    4. find another vertex with indegree 0
    5. repeat till all vertices are listed
  - maintaining digrees is  $O(n^2)$ , keep n times

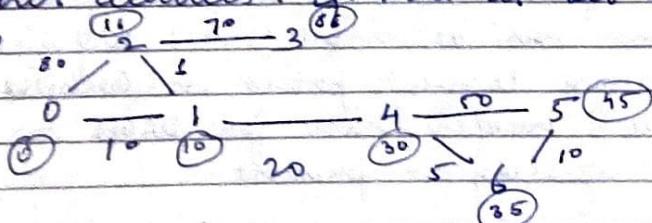
Week-5

### L5.1 Shortest Paths in weighted graphs

- BFS computes shortest path in terms of no. of edges
- may assign value to edges  $\rightarrow$  weighted graph
- adj. matrix  $\rightarrow$  weight is always 0 if no edge
- adj. list  $\rightarrow$  records weight with edge info.
- weighted shortest path need not have least no. of steps
- single source shortest paths — find shortest paths from a vertex to every other vertex
- all pairs shortest paths — find shortest paths b/w. every pair of vertices
- -ve edge weights. Eg. Taxi driver going home at the end of day
- -ve cycles  $\rightarrow$  the concept of shortest path does not exist

### L5.2 Single Source Shortest Path (Dijkstra's Algo)

- computes shortest paths from  $v_0$  to all other vertices. Eg. Fire at oil depot  $v_0$



- compute exp. burn-time for each vertex  
 $\hookrightarrow$  Edsor w/ Dijkstra's Algo.
- each new shortest path we discover extends an earlier one — Greedy strategy

- by induction, assume we have found shortest path to all vertices already burnt
- we can not use Dijkstra's algo for -ve edge weights
- maintain 2 dictionaries with vertices as keys
  - $\rightarrow$  visited, initially false for all  $v$
  - $\rightarrow$  distance, initially  $\infty$  for all  $v$
- $\infty$  takes  $O(n^2)$  time, main loop runs  $n$  times  $\rightarrow$  Overall  $O(n^3) \rightarrow$  adj. matrix  $\hookrightarrow$  adj. list too

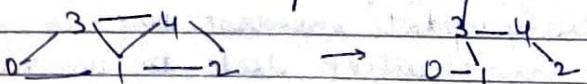
### L5.3 Single Source Shortest Path with -ve wts. (Bellman Ford Algo.)

- in prev. algo., correctness req. non--ve edge weights only
- shortest route to every vertex is a path, no loops
- Min. weight path from  $v_0 \rightarrow v_k$  is  $v_0 \xrightarrow{w_1} j_1 \xrightarrow{w_2} j_2 \xrightarrow{w_3} \dots \xrightarrow{w_{k-1}} j_{k-1} \xrightarrow{w_k} v_k$
- repeatedly update shortest dist. to each vertex based on shortest dist. to its neighbours
- Initialize (source vertex  $v_0$ )  
 $D(j)$ : min dist. till vertex  $j$   
$$D(j) = \begin{cases} 0, & \text{if } j = v_0 \\ \infty, & \text{otherwise} \end{cases}$$
- works for directed and undirected graphs
- if there was a -ve cycle, dist. would continue to decrease
- adj. matrix takes  $O(n^3)$
- adj. list takes  $O(mn)$
- if this algo. does not converge after  $(n-1)$  iterations, there is a -ve cycle.

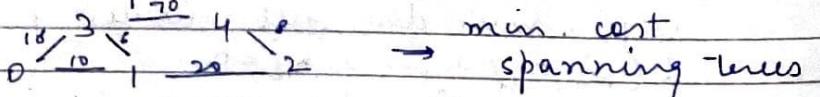
- L5.4 All pairs shortest paths (Floyd-Warshall Algo)
- transitive closure,  $A^+ = A + A^+ + \dots + A^{n-1}$
  - $B^K[i, j] = 1$  if there is path from  $i \rightarrow j$  via vertices  $\{0, 1, \dots, k-1\}$
  - Warshall's algo. — also computes transitive closure
  - $SP^0$  &  $SP^1$  are same because 0 has no indegree
  - shortest path matrix  $SP$  is  $n \times n \times (n \times 1)$
  - time complexity is  $O(n^3)$

L5.5 Min. Cost Spanning Trees

- roads are damaged, needs restoration acc. to the penisity on fibre optic cables example
- certain minimal set of edges remains connected. Minimal connected graph is a tree. Want a tree that connects all vertices — spanning-tree



- Spanning trees with cost



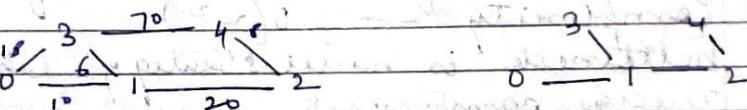
- A tree is connected acyclic graph. A tree of  $n$  vertices has exactly  $(n-1)$  edges. Adding an edge to a tree must create a cycle. In a tree, every pair of vertices is connected by a unique path.

$G$  is connected  
 $G$  is acyclic  
 $G$  has  $n-1$  edges

- Prim's algo. — start with smallest  $\Delta$  genes
- Kruskal's algo. — scan the edges in asc. order to connect comp. w/o forming cycle

L5.6 Min. cost spanning trees (Prim's Algo.)

- incrementally growing from 1 edge (which is a tree) → extend the tree



- smallest edge  $(1,3) \rightarrow (1,0)$  → can't add  $(0,2) \rightarrow (1,2) \rightarrow (2,4)$   
 ↗ this would form a cycle

- correctness — minimum spanning lemma
- in the implementa", keep track of -
  - visited ( $v$ ) — is  $v$  in the spanning tree
  - distance ( $v$ ) — shortest dist. from  $v$  to the tree
  - tree Edges — edges in the current spanning tree

- we are looking into list  $\rightarrow O(mn) \rightarrow O(n^2)$

- unvisited
- mba( $v$ ) — nearest neighbour of  $v$  in tree
- very similar to Dijkstra's algo
- complexity is  $O(n^2) \Rightarrow$  kind of greedy strategy

- bottleneck is identifying unvisited vertex with min. distance.

L5.7 MCST (Kruskal's Algo.)

- start with  $n$  comp., process edges in asc.

$\begin{matrix} & 3 & 4 \\ x & \diagdown & \diagup \\ 0 & - 1 & - 2 \end{matrix}$

- start with  $(1, 3) \rightarrow (2, 4) \rightarrow (0, 1)$ ,  
→ then try for  $(0, 3)$  skip → then  $(1, 2)$
- use disjoint path and their later form the tree
- correctness follows the same from Kruskal's
- complexity —  $O(n^2)$
- bottleneck is naive strategy of label and merge component
- efficient union-find brings the complexity  $O(m \log n)$
- it has bottom up approach
- if edge weights repeat, MCST is not unique

## Week - 6

### 16.1 Union - Find Data Structures

- keep track of components and merge them efficiently
- data struc. to maintain collec<sup>n</sup> of disjoint sets -
  - find ( $u$ ) - return set containing  $u$
  - Union ( $u, v$ ) - merge sets of  $u, v$
- a set  $S$  partitioned into components  $(C_1, C_2 \dots C_j)$ , each  $s \in S$  belongs to exactly one  $C_j$ .
- always merge smaller component into the larger one
- individual merge operations can still take time  $O(n)$  → more careful accounting
- Union() complexity →  $O(\log m)$  (amortized)
- Kruskal's Time →  $O((m+n) \log n)$

### 16.2 Priority Queues

- job scheduler - maintains a list of pending jobs with their priorities
- when the processor is free, scheduler picks up the job with max. priority and schedules it
  - also new jobs keep coming up
- unsorted list → insert (easy), delete (hard)
  - sorted list → vice-versa
- keep the track of size of each queue → insert into the 1<sup>st</sup> queue that has space
  - takes time →  $O(\sqrt{N})$  → insert()
- max in each queue is the last element → identify the max → delete it →  $O(N)$

- maintain a spl. binary tree - heap  
→ height  $O(\log N)$

### L6.3 Heaps

- binary trees - values are stored as nodes in a rooted tree, each node has upto 2 children (L child & R child) → order is comp't.
- other than root, each node has a unique parent
- leaf node - no children
- size - no. of nodes
- height - no. of levels
- binary tree filled level-by-level, L → R
- non negs.
- no holes allowed
- can't keep a full incomplete problem
- heap prop. is violated
- complexity → insert()  $\propto$  height of tree
- no. of nodes at level  $j$  is  $2^j$
- if we fill  $k$  levels,  $2^0 + 2^1 + \dots + 2^{k-1} = 2^k - 1$  Nodes
- if we have  $N$  nodes, at most  $(1 + \log N)$  levels → insert() is  $O(\log N)$
- delete max → max. value is at the top of the tree → after 1 value, tree shrinks → restore the heap property downwards
- Complexity →  $O(\log N)$
- building a heap - heapify()  $\rightarrow O(N)$

### L6.4 Using Heaps in Algos

- heaps can be represented as list / array
- updating values in a min-heap
- updates are  $O(\log N)$
- maintaining 2 additional dictionaries
- Dijkstra's algo - using min-heap  
→ heap sort - starts with an unsorted list  $\rightarrow O(n \log n)$

### L6.5 Search Trees

- sorting is useful for efficient searching, the data is dynamically changing
- insertion takes time  $O(n)$
- if move to tree structure like heaps for efficiency quick search
- Binary search tree - no duplicates, each node has a value and pointers to its children that are also binary search trees
- worst case complexity - height of the tree
- time  $= O(\log N)$

Week 7L7.1 Balanced search trees

- search trees - find(), insert & delete() all work on a single path
  - worst case - height of the tree
  - balanced trees have height  $O(\log n)$
  - balance - left & right subtrees should be equal. 2 measures - size & height  $\Rightarrow$  this is possible for comp. binary trees
  - self.height() - no. of nodes on longest path from root to self
  - AVL - Adelson-Velskii, Landis (<sup>height</sup><sub>balance</sub>)
  - min. size height-balanced trees
- |       |          |       |          |       |          |       |
|-------|----------|-------|----------|-------|----------|-------|
| $n=1$ | $\vdots$ | $n=2$ | $\vdots$ | $n=3$ | $\vdots$ | $n=4$ |
|-------|----------|-------|----------|-------|----------|-------|

- $S(h)$ , size of smallest height-balanced tree of height  $h$
- Fibonacci  $\rightarrow F(0)=0, F(1)=1, F(n) = F(n-1) + F(n-2)$
- $S(n)$  grows exponentially with  $n$
- balanced tree - slope is  $\{-1, 0, 1\}$
- t.insert(v), t.delete(v)  $\rightarrow$  slope from +2 to -2
- left rotate  $\rightarrow$  converts slope  $\{-2\}$  to  $\{0, 1, 2\}$
- rebalance each time the tree is modified
- automatically balances bottom up

L7.2 Greedy Algos - Interval Scheduling

- need to make a seq. of choices to achieve a global optimum
- always moves forward
- at each stage, make the next choice based on some local criterion
- drastically reduces space to search for sol<sup>n</sup>

- Ex.
- Dijkstra's algo., Prim algo.
  - IITM has sp. video classroom for delivering online lectures, slots may overlap, so not all bookings honoured, so how to maximize the no. of lectures
  - take the 1<sup>st</sup> slot with any local strategy
  - greedy strategy is optimal
  - greedy algo takes local optimum choices to get a global optimum
  - it never considers its choice worth
  - just pushes that local strategy is globally optimum at that moment
  - interval scheduling - many natural greedy strategy
  - correct strategy needs a pass of  $O(n^2)$
  - greedy sol<sup>n</sup>: step ahead! step-by-step of any optimal sol<sup>n</sup>

L7.3 Greedy algs - Minimizing lateness

- Eq. of IITM has single 3D printer
- each user will get access to the printer, but may not finish before deadline
  - minimize the maximum lateness
  - try to finish first shortest jobs
  - or give priority to request with smaller slack time
  - renumber the jobs so that they are sorted by deadline
  - There is an optimal schedule with no idle time
  - eliminate idle time by shifting jobs earlier

- this can reduce lateness  
exchange argument & inversions  
any 2 schedules with no inversions and  
no idle time have same max. lateness
- no inversions, idle-time - only diff  
can be the order of sig. with same  
deadline
- scheduling jobs with same deadline  
produces same lateness
- there is an optimal schedule with no  
inversions & no idle time
- schedule sig. with start times  $t(i)$  and  
deadlines  $d(i)$ , to minimize max.  
lateness
- simple greedy algo. →  $O(n \log n)$ 
  - sorts the sig. →  $O(n \log n)$
  - makes off schedule →  $O(n)$
- correctness by exchange argument.
- consider an O optimus argument
- transform it → greedy soln.

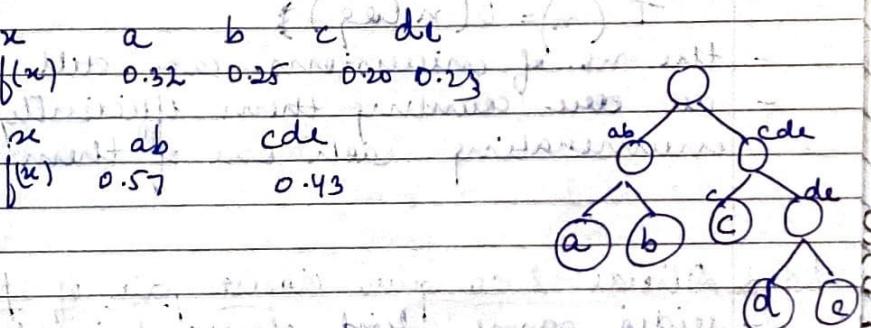
L7.4 Greedy algo → Huffman coding  
 - efficient communication  
 - variable length encoding → most efficient  
 Example is 0, except is 1, a principle

|      |    |    |     |    |     |
|------|----|----|-----|----|-----|
| x    | a  | b  | c   | d  | e   |
| E(x) | 11 | 01 | 001 | 10 | 000 |

totalizing length - length of a string  
 $\sum_{x \in A} n \cdot f(x) \cdot |E(x)|$

- avg. no. of bits / letter in encoding
- any optimal prefix code produces a full-tree full tree - each node has 0 or 2 children
- in an optimal tree, for any leaf at max. depth, its sibling is also a leaf
- in an optimal tree, if leaf labelled x is at smaller depth than leaf labelled y.  $f(x) \geq f(y)$  Huffman's algo.

|      |      |      |      |      |      |
|------|------|------|------|------|------|
| x    | a    | b    | c    | d    | e    |
| f(x) | 0.32 | 0.25 | 0.20 | 0.18 | 0.05 |



- Huffman coding is greedy
  - recursively combine letters with lowest freq.
  - locally optimal choice at next step
  - never go back & do other pairing
- Claude Shannon → info. theory

## Week - 8

|           |          |
|-----------|----------|
| CLASSTIME | PAGE NO. |
| Date      | / /      |

|           |          |
|-----------|----------|
| CLASSTIME | PAGE NO. |
| Date      | / /      |

- L8.1 Divide & Conquer Counting inversions -
- break up a problem into disjoint subproblems
  - combine the subs efficiently
  - merge sort, quick sort
  - recommendation systems - online services recommend items to you, compare your profile with others
  - so how to compare profiles?
  - no. of inversions changes from 0 to  $n(n-1)/2$  — measure of dissimilarity
  - recurrence is similar to merge sort
- $$T(0) = T(1) = 1$$
- $$T(n) = 2T(n/2) + n$$
- $$T(n) = O(n \log n)$$
- the no. of inversions can still be  $O(n^2)$
  - we are counting them efficiently without enumerating each one of them

- L8.2 Divide & conquer closest pair of pts.
- video game, find closest pair of objects (n objects)  $\rightarrow O(n^2)$
  - points p1, p2 in 2D  $\rightarrow p = (x, y)$
  - $d(p_1, p_2) = \sqrt{(y_2-y_1)^2 + (x_2-x_1)^2}$
  - we have to find the closest pair...
  - the brute force  $\rightarrow O(n^2)$
  - $P_x$  &  $P_y$  sorted on x and y coordinate
  - splitting helps to take time  $O(n)$
  - now, how to combine the recursive solutions  $\rightarrow$  we can't have 2 pts inside the same box
  - k-dimension scan

- L8.3 Integer multiplication : Divide & Conquer
- form partial products  $\rightarrow$  add them up
  - works the same in any base — e.g. binary
  - to multiply 2 n-bit nos  $\rightarrow O(n^2)$ 
    - $\rightarrow$  n partial prod.
    - $\rightarrow$  add to get cumulative sum
  - each partial prod. seems 'necessary'
  - Karatsuba's algo.  $\rightarrow$  fast-multiply algo.
  - $T(1) = 1$
  - $T(n) = 3T(n/2) + n$ 
    - $\hookrightarrow$  this forms the geometric series
    - $D(n^2) \rightarrow O(n^{1.59})$
  - divide & conquer reduces the time complexity of integer multiplication
  - 1959, Andrei Kolmogorov  $\rightarrow$  multiplication could <sup>not</sup> be done in sub-quadr. time
  - Karatsuba  $\rightarrow$  came with divide & conquer algorithm
  - analysis simplification by Donald Knuth

- L8.4 Divide & Conquer: Recursion Trees
- complexity  $T(n)$  is expressed as recurrence
  - Binary Search:  $T(n) = T(n/2) + 1 \rightarrow O(\log n)$
  - Merge Sort:  $T(n) = 2T(n/2) + n \rightarrow O(n \log n)$
  - recursion tree — nested tree (with 1 node  $+ ((1 + nv) - n)$  for each recursive subproblem)

- value of each node is spent on the subproblem excluding recursive calls
- on an input size  $n$ ,
  - (i)  $f(n)$  is the time spent on non-recursive work
  - (ii)  $m$  is the no. of recursive calls
  - (iii) each recursive call works on a subproblem of size  $n/c$

$$T(n) = mT(n/c) + f(n)$$

- cost of recursive tree for  $T(n)$  has value  $f(n)$
- each node at lev.  $d$  has value  $f(n/c^d)$
- leaves correspond to the base case  $T(1)$
- level  $i$  has  $c^i$  nodes, each with value  $f(n/c^i)$
- tree has  $L$  levels,  $L \geq \log_c n$
- Total cost is  $T(n) = \sum_{i=0}^{L-1} c^i \cdot f(n/c^i)$

- L8.5 Divide & Conquer: Quick Select
- find the  $k^{\text{th}}$  largest value in a seq. of length  $n$
  - sort in desc. order & look at pos<sup>n</sup>  $k \rightarrow O(n \log n)$
  - if we can find median in  $O(n)$ , quick select becomes  $O(n \log n)$
  - recurrence is similar to quicksort
- $$T(1) = \text{constant}$$
- $$T(n) = \max(T(m), T(n-(m+1))) + n,$$

- worst case  $\rightarrow T(n) \text{ is } O(n^2)$
- if pivot is within a fixed fract, then quicksort is  $O(n \log n)$
- other method  $\rightarrow$  median of medians uses median of block medians to locate pivot for quickselect at  $O(n)$  pivot
- medians of block medians helps to find a good pivot in  $O(n)$
- select becomes  $O(n)$ , quicksort becomes  $O(n \log n)$
- C.A.R + same described quickselect in the same paper as is quicksort (1973) — median of medians

- L8.6 Implementation of Quick Select & Fast Select Algo
- median is not exact but still it's better
  - Quick select  $\rightarrow$  9999, 7 sec. iteration
  - Fast select  $\rightarrow$  9999, 0.005 sec.

## Week 9

- L9.1 Dynamic Prog. (n)T → how to draw -
- inductive def. o. nivels is tried i -
  - Factorial  $\text{fact}(n) = n \times \text{fact}(n-1)$
  - induction nivels ← bottom up
  - sol<sup>n</sup> to original problem can be derived by combining sol<sup>n</sup>s to subproblems
  - interval scheduling problems
  - each subset of bookings is 1 or subproblems
    - ↳ guess in bookings, and subproblems
  - greedy strategy looks at only a small fraction of subproblems
  - weighted interval scheduling → scenarios as number of visitors before, but each req. comes with a wt.
  - inductive sol<sup>n</sup> generates same subproblem at diff. stages
  - recursive implementation evaluates each instance of subproblem from scratch
  - memoization and dynamic prog.

## L9.2 Memoization

- build a table of values already computed → memory table
- memoization - check if the value to be computed was already seen before
- store each newly computed value in a table
- look up the table before making a recursive call
- computation tree becomes linear

- Dynamic prog. - anticipate the structure of answers to subproblems -
- derive from inductive def
  - dependencies form a DAG
  - solve subproblems in topological order
    - never need to make a recursive call
- L9.2
- L9.3
- how to find tiles
- Grid Pathways visit each cell in manhattan grid + upper + right - one way: visit only up & right from  $(0,0)$  to  $(m,n)$ . total tiles -
- in gen.  $(m+n)$  segments from  $(0,0)$  to  $(m,n)$  at tile length 200
  - out of 15, exactly 5 are right moves, → 10 are up moves
  - inclusion-exclusion (→ counting is messy)
  - inductive formula

- L9.4
- Common subsequences & subseq.
- guess 2 strings, find lengths of the longest common subsequence
  - subproblem dependency
  - subsequence (→ can drop some letters in between) : 2A
- (n) → analyzing genes via string
- diff command in Unix/Linux
  - algorithm (mn), using DP or xitance
- (2A) A = 2(2A) = 2A
- Two overlapping threads produced

8

- L9.5 Edit distance
- used in merges -
  - document similarity  $\rightarrow$  insert a charac.
  - ~ it's insertion and deletion charac.
  - $A \rightarrow A$  is need  $\rightarrow$  substituted 1 char by other inserted in between another
  - misreadings of edit operat. needed
  - eg. 24 charac. inserted, 18 deleted & 2 subs.  
so edit dist. at most 44.
  - also called Levenshtein distance
  - applica  $\rightarrow$  suggestions for spelling correc.
  - helps in genetics similarity of species
  - edit dist. = LCS (longest common)
  - at  $(0, 0)$  need insertion (subsequence)  $\rightarrow$
  - LCS equivalent to edit distance up to  
deletions & insertions  $\neq$  edits,  $\rightarrow$  true
  - Edit distance  $\rightarrow$  transforms  $w \rightarrow v$
  - $O(mn)$  using dynamic prog.  
see memoiz.

memoized recursion -

### L9.6 Matrix multiplication

- multiply matrices  $A, B$   $\rightarrow$  result is  $C$
- def  $\text{AB}[i, j] = \sum_{k=1}^{n-1} A[i, k] B[k, j]$   
  second  $k$  = column of  $A$  and row of  $B$
- dimensions must be compatible
- $m \times n, n \times p$   $\rightarrow$   $m \times p$
- $AB$ :  $m \times p$
- computing each entry in  $AB$  is  $O(n)$   
overall  $O(mnp)$
- matrix multiplication is associative
- $ABC = (AB)C = A(BC)$   
 $\hookrightarrow$  breaking doesn't change ans. but

the complexity does changes due to it.

Eq.  $WA = 1 \times 100$ ;  $BA = 10 \times 1$ ;  $BC = 1 \times 100$   $\rightarrow$   
it's seen the diff. b/w  $(ABC)$  &  $(AB)C$   
terminal 20000 steps // 200 steps

- in DP, diagonal entries are base case,
- fill matrix by diagonal from main
- + fill diagonal in next row
- fill the table  $\rightarrow O(n^2) + i : i$   
filling each entry  $\rightarrow O(n)$  time
- Total Overall  $\rightarrow O(n^3)$  time

(mnr)  $\rightarrow$  m rows & n columns

of  $i$  to  $n$  steps of time

of  $j$  to  $m$  steps of time

apart insert - need in diagonal

in next moment

PP1 need speed

"try fast to fast filling entry"  $\Rightarrow$

"fast to fast to"  $\Rightarrow$

for  $i$  to  $n$  - need fast

between  $i$  rows transition, need update

inner part need update, dual update

$i$  times  $[(j \text{ mod } + i) : i]$

(mnr)  $\rightarrow$  m rows & n columns

$i$  to  $n$  part update, parallel op

(121)  $\rightarrow$  pth diagonal, interleaved

spl. init

reorder

need  $m \cdot n$  in parallel for an ans.

written on random arrangement

$\uparrow$  step

## Week -10

- L10.1 String Matching with pattern  $p$  and text  $t$
- searching for a pattern in a fundamental problem? when dealing with text
  - editing a document
  - ans. an internet search query
  - looking for a match in a gene seq
  - find every pos "i" in  $t$  such that  $t[i:i+m] = p$
  - nested loop → bottleneck
  - nested search bails out at first mismatch
  - worst case is  $O(nm)$
  - don't need to check all of  $t$  to search for all occurrences of  $p$ .
  - Formalized in Boyer-Moore Algo.

## L10.2 Boyer-Moore Algorithm

- Boyer, Moore, 1977
- $t$  = "which finally halts. at that point"  
 $p$  = "at that"
- skip over - sped up
- single scan, rightmost value is recorded
- nested loop, compare each segment  $t[i:i+\text{len}(p)]$  with  $p$
- worst case still remains  $O(nm)$
- w/o dictionary, computing last is a bottleneck, complexity is  $O(|\Sigma|)$
- this algo.
  - sublinear
  - avg. no. of comparisons is  $O(2n/\text{char.})$
  - performance improves as pattern length ↑

- in practice, Grouping Unit →  $\Sigma = \{a, b, \dots, z, \#\}$  →  $\Sigma$  is of length  $k$
- L10.3 Rabin-Karp algorithm
- reducing string matching to arithmetic
  - $\Sigma = \{0, 1, \dots, 9\}$  if  $\Sigma$  is a group
  - any string over  $\Sigma$  can be thought of as a no. in base 10
  - pattern  $p$  is an  $m$ -digit no. divided
  - converting a string to a no.
  - can convert a block  $t[i:i+m]$  to an integer  $n_i$  in one scan
  - in computing  $n_i \rightarrow O(mn)$
  - for  $\Sigma = \{a_1, a_2, \dots, a_k\}$  treat each letter as a digit in base  $k$
  - in practice, for realistic  $k$ , the numbers are too large to work with directly
  - instead do all computations and comparisons modulo a suitable prime  $q$
  - false positives
  - preprocessing time is  $O(m)$
  - worst case for general alphabets is  $O(nm)$
  - in practice, no. of spurious matches will be small ( $m$ )
  - then, overall time  $O(n+m)$  or  $O(n)$

## L10.4 String Matching using automata

- traditional string matching →  $O(nm)$
- can we intelligently re-use partial matches
- to reuse, use computed values

- precomputing longest matches as a graph
- pattern  $p$  of length  $m$ , graph with  $(m+1)$  nodes, edges describe how to extend the match at position  $i$
- graph  $\rightarrow$  finite state automaton
  - nodes are states
  - $\rightarrow$  edges are transitions
- build an automaton to keep track of longest matching prefix with a priority queue
  - using this, we can do it in  $O(n)$
  - bottleneck is precomputing the automaton
  - do this in time  $O(m) \rightarrow$  Knuth, Morris, Pratt

### L10.5 Knuth, Morris, Pratt Algorithm

- to compute the automaton for efficiently matching  $p$  against itself
- analysis of fail computation is subtle
- this algo. efficiently computes the automaton describing prefix matches in the pattern  $p$
- complexity of preprocessing a fail function  $\rightarrow O(m)$
- can check matches in  $O(m+n)$
- KMP  $\rightarrow O(m+n)$  matches, need
- the Boyer-Moore algo is faster in practice skipping many funcs.

### L10.6 String Matching: suffix tree

- it often involves search over large fixed body of text
- collected results of all the

- comprehension of set of ref. animals
- genetic data
- make multiple queries on this text
- preprocess this text to make the search easier into a suffix tree
- A tree is a spcl. kind of tree  $\rightarrow$  from information retrieval
- nested tree  $\rightarrow$  children of node have distinct labels
- each maximal path is a word
  - $\rightarrow$  one word shouldn't be a prefix of another
  - $\rightarrow$  add special symbol at the end of word
- build a tree  $T$  from a set of words with  $s$  words and  $k$  total symbols
- basic prop.
  - $\rightarrow$  height of  $T$  is  $\max_{w \in S} \text{len}(w)$
  - $\rightarrow$  a node has at most  $\leq k$  children
  - $\rightarrow$  the no. of leaves in  $T$  is  $S$
  - $\rightarrow$  the no. of nodes in  $T$  is  $n$ , plus  $s$  nodes labelled  $\$$
- can maintain auxiliary info. for each word
- suffix tree  $\rightarrow$  expands to include all suffixes of all words in  $S$
- most ques can be ans. using suffix tree
- what is a substring?
- is a suffix of  $s$ ?
- no. of times  $w$  occurs as a substring of  $s$
- longest repeated substring of  $s$

## Week 11

- how big can a suffix tree be for s of length  $n$ ?  $\rightarrow n^2$
- no. of nodes prep. to  $n$ ?  $\rightarrow n^2$ , an  $b^n$
- main drawback of a tree is size

L10.7 Regular Expressions

- what if we want to look for a pattern?
- $\rightarrow$  unaware of spelling
- $\rightarrow$  check for a subsequence
- anchoring patterns
  - our concern is now that p matches anywhere in the string
  - to match the start  $\wedge p$
  - to match the end  $p \wedge$
- KMP algo built a finite state automaton for the longest prefix matched by a pattern
- this automaton we build had a linear structure, with a single path from start to finish
- the set of words that automata can accept are called regular sets
- the sets we can describe using patterns are exactly the same as those that can be accepted by automata
- our patterns are cd. regular expressions
- python provides a library for matching regular expressions

11.1 Linear Prog.

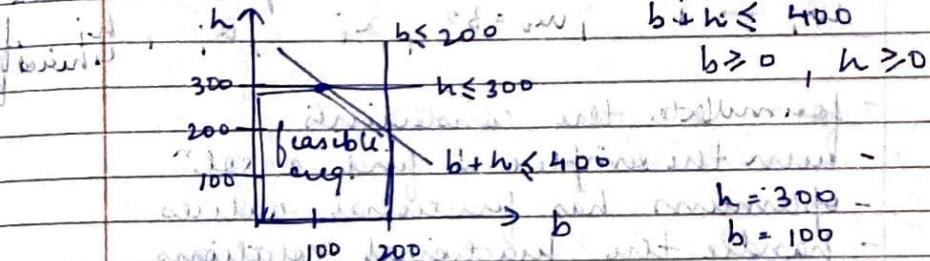
- computational tasks involve "optimizing"
  - shortest path
  - min. cost
  - longest common subseq.
  - LP - constraints & objectives to be optimized are linear functions

e.g. maximize profits

$b \rightarrow$  boxes of books,  $h \rightarrow$  boxes of halwa

$$\text{Profit} = 100b + 600h \quad b \leq 200$$

$$\text{max.} \rightarrow 100b + 600h \quad h \leq 300$$



the ans. lies on the corner (any)

- simplex algo. - can be exp., but efficient in practice

( $\rightarrow$ ) not convex ( $\rightarrow$ ) convex

- Feasible reg. is convex, may be empty or unbounded
- LD Duality - This is always possible. We derived an upper bound on the obj. thru a linear "combination" of constraints

## L11.2 LP Production Planning

can always construct a linear prog.

- Production planning

- 30 employees  $\rightarrow$  20 carpets/month

- salary £20,000, labour £100/carpet

- dep. demand from Jan to Dec.

- overtime  $\rightarrow$  80% extra, cost is 30%/ $\text{worker}$

- hiring £3200/worker, firing £4000/worker

- supplies storage cost £80/carpet

$\rightarrow$  workers  $\rightarrow$  carpets, overtime

$w_i \rightarrow i$ ,  $w_0 = 30$ ,  $x_i \geq 0$ ,  $h_i \leq h_i^*$

- formulate the constraints

- run the simplex & find a soln

- optimality has fractional values

- handle the fractional solutions

$\Rightarrow$  Integer LP  $\rightarrow$  computationally intractable

## L11.3 LP Bandwidth Allocation

- 3 users, A, B, C

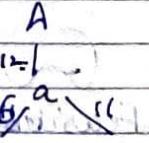
- link have capacity constraints

- each connection should have at least 2 Mbps  $\rightarrow$  C  $\geq 8$

- direct & indirect traffic

connections are allowed

- allocate bandwidth to maximize the revenue



- modelling strategy doesn't scale well
- need to have a better approach to analyze such network flows

## L11.4 Network Flows

- network of pipelines

- ship a map from  $s \rightarrow t$

- Network: graph  $(G) = (V, E)$

- special nodes:  $s$  (source),  $t$  (sink)

- each edge  $e$  has capacity  $c_e$

- flow:  $f_e$  for each edge  $e$

$\Rightarrow$  at each node, sum of incoming flows = sum of outgoing flows

- LP formulation

nat.  $f_e$  for each edge  $e$

- simplex replaces millions of feasible eq. to solve LP

- Ford-Fulkerson algo.

start with a flow

- residual graph - for each edge  $e$  with capacity

- certificate of optimality - max flow-min cut theorem

- in this algo, choose augmenting paths wisely, use BFS to find augmenting path

## L11.5 Reductions

- each institution is willing to teach a set of courses

- bipartite matching -
- $V \rightarrow V_0, V_1$  & we've seen
- all edges from  $V_0$  to  $V_1$
- matching: subset of edges so that no 2 of them share an endpoint
- reduce to
  - we want to solve A to B
  - we know how to "B" given A
  - convert input A to B
  - A reduces to B
  - can transfer efficient sol<sup>n</sup> from B → A
  - but pre & post processing must be efficient
  - both should be in polynomial time
  - bipartite reduces to max. flow
  - max. flow reduces to LP
  - no. of vars., constraints is linear in the size of graphs
  - reverse interpretation is also useful
  - big hammers
  - IP & network flows are powerful tools

### L10.6 Interactivity: Checking Algo.

- good algo - have polynomial time algo.
- search space for sol<sup>n</sup> is expo.
- brute force - scan expo. possibilities & choose the best
- large class of 'natural' problems where shortcut isn't available
- generating vs checking
- checking algo. is for problem P
- boolean satisfiability

- $x_i \text{ (neg)} \quad x_i \vee x_j \text{ (or)} \quad x_i \wedge x_j \text{ (and)}$
- conjunct<sup>n</sup> of clauses  $(C_1 \wedge C_2 \wedge \dots \wedge C_n)$
- each clause forces a unique value
- e.g. Travelling Salesman
  - a network of cities with distances
  - b/w each pair, find the shortest tour that visits each city exactly once & also every city has to be visited
  - Union independent set of size k iff  $\forall v \in V$  there is a vertex cover of size  $N-k$
  - indep. set & vertex cover reduce to each other

### L11.7 Interactivity of P & NP

- NP - factorisation, satisfiability, travelling salesman, vertex cover, indep. set able to check
- Non-deterministic Polynomial Time
  - guess a set<sup>n</sup> & check it
  - P vs. NP: polynomial time algo. (most cases)
  - C is included in NP - generate a sol<sup>n</sup> & check it
- $P \neq NP$  NP
- problems are inter-reducible
- Cook-Levin Th. → Every problem in NP can be reduced to SAT
- SAT is said to complete for NP

was 1st algorithmic problem which had a fixed time complexity

- 112.1 Summary (Week 11 → 3 weeks)
- timing our code → time library
  - analysing an algo. → time & memory
    - Time complexity:  $T(n) \rightarrow O(T(n))$  Space complexity:  $O(1)$
    - Worst case time complexity:  $T(n) = f(n)$  Space complexity:  $O(n)$
  - upper bound on worst case gives us an overall guarantee of performance
  - if  $f(n) \in O(g(n))$  means  $g(n)$  is an upper bound for  $f(n)$  (optimal soln)
  - orders of magnitude
  - calculating complexity → type → iterative  
→ recursive
  - searching algs.
    - Linear search
    - Binary search
      - sorted list
      - check every element with target element
      - worst when target is not in list till list becomes empty

|       |        |               |
|-------|--------|---------------|
| Best  | $O(1)$ | $O(1)$        |
| Avg.  | $O(n)$ | $O(n \log n)$ |
| Worst | $O(n)$ | $O(n \log n)$ |

  - selection sort
  - insertion algs.
  - swapping elements to avoid Tandem list
  - $T(n) = n(n+1)/2$
  - time complexity remains same
  - insertion sort
    - insert into the list -  $T(n) = n(n+1)/2$
    - time complexity varies based on seq. of elements

- merge sort
  - divides the list into 2 halves and sort the sublists then combine
  - $T(n) = 2T(n/2) + n$
- Quicksort
  - choose a pivot (1<sup>st</sup> element), partitions into upper & lower, pivot goes in between, sort the partitions recursively. This places an in-place sort

|          | Selection Sort | Merge Sort | Quicksort     |
|----------|----------------|------------|---------------|
| Best     | $O(n^2)$       | $O(n)$     | $O(n \log n)$ |
| Avg.     | $O(n^2)$       | $O(n^2)$   | $O(n \log n)$ |
| Worst    | $O(n^2)$       | $O(n^2)$   | $O(n \log n)$ |
| In-Place | Y              | Y          | N             |
| Stable   | N              | Y          | Y             |

- list
  - flexible length
  - scattered memory vs contiguous block
  - need to follow links randomly
  - inser<sup>n</sup> & del is easy
  - swapping takes c. time
  - list (in python) → maps to array
    - (best and append & pop →  $O(1)$ )
    - inser<sup>n</sup> & del →  $O(n)$
- dict implemented as hashtable (hash func)
  - need a strategy to deal with collisions
  - closed - possible for free space in array
  - open - each slot pts. to a list of key-value pairs

- L12.2 Summary week 4-6
- graph → relationship b/w entities
  - directed or undirected
  - ↳ paths are seq. of connected edges
  - unreachable → from u to v
  - $G = (V, E)$  :  $|V| = n$ ,  $|E| = m$
  - if  $G$  is connected,  $m$  varies from  $(n-1)$  to  $n(n-1)$
  - Adj. matrix →  $n \times n$  matrix
  - " list → dictionary of list
  - BFs - express the graph, built by level queue, adj. matrix ( $O(n^2)$ ), list ( $O(mn)$ )
  - DFS - uses a stack (suspend) same complexity as BFs
  - paths by BFs is shortest paths in terms of no.
  - find strongly connected points
  - DAGs → represent dependencies
  - ① Topological Sorting
  - ② longest paths  $\rightarrow O(m+n)$
  - Single source shortest (Dijkstra's algo)
  - greedy strategy
  - complexity  $(O(n^2))$
  - not work with negative weight
  - min heap  $\rightarrow O((m+n) \log n)$
  - single source + (u,v) wt. (Bellman Ford)
  - adj. matrix  $\rightarrow O(n^3)$
  - adj. list  $\rightarrow O(mn)$
  - all pairs shortest paths (Floyd-Warshall)
  - complexity  $\rightarrow O(n^3)$
  - Trees → connected acyclic path
  - spanning tree → MCST (multiple), but cost'll be unique
    - Prim
    - Kruskal

- Peim's - implementation similar to Dijkstra's (complexity  $\rightarrow O(n^2)$ )
- Kruskal's -  $O(n^2)$  due to naive handling of components,  $O(m \log n)$  using union-find data struc.
- Union-find
  - Union()  $\rightarrow \text{log } m$
  - Find()  $\rightarrow O(1)$
- Persistence Queues
- Heaps - binary tree, no holes
  - insert / delete  $\rightarrow O(\log n)$
  - heapify  $\rightarrow O(n)$
- Search trees
  - CS worst case - an unbalanced tree with  $n$  nodes may have height  $O(n)$
  - find, insert, delete-all  $\rightarrow O(n)$
- Balanced search trees - AVL trees
  - find, insert, delete  $\rightarrow O(\log n)$

- L12.3 Summary week 7 → 9
- greedy algo → global optimum, need go back, needs a proof of optimality
  - eg. Dijkstra, Peim, Kruskal's algo
  - Applications of greedy algo.
    - (i) Interval scheduling - IITM classes
    - (ii) minimizing lateness - 3D printer
    - (iii) Huffman coding  $\rightarrow O(k \log k)$

- Divide & Conquer — eg. merge, quick  
eg. counting inversions — divide & conquer  
 $T(n) = 2T(n/2) + n = O(n \log n)$

- closest pair of pts — divide & conquer  
 $O(n^2) \rightarrow O(n \log n)$

- integer multiplication

Naive —  $O(n^2)$   
Kanatsuba's —  $O(n^{1.5})$  or  $O(n \log^3 n)$

- quick select —  $O(n \log n)$   
 $\hookrightarrow$  median of medians → fast-select

- recursion trees

$$T(n) = aT(n/c) + f(n)$$

- DP  $\rightarrow$  sol<sup>n</sup> of problem can be derived from sol<sup>n</sup> of subproblems

$\hookrightarrow$  solve subproblems in topological order

- Memoization  $\rightarrow$  build a table of values already computed — memory-table

- Grid paths — up & right

$$O(mn) \rightarrow \text{DP}$$

$O(m+n)$  — memoization

- common subword & subseq. & edit dist. complexity —  $O(mn)$

- matrix multiplication —  $O(n^3)$

## 12.4 Summary week 10+11

- string matching — searching for a pattern while dealing with text  
 $\rightarrow O(mn)$

- Boyer-Moore algo  $\rightarrow$  works well in prac. but worst case is  $O(nm)$
- Rabin-Karp algo  $\rightarrow O(n+m)$  or  $O(n)$
- Automata  $\rightarrow O(m^3, l \leq l)$
- Knuth-Morris-Pratt algo  $\rightarrow$  easily calculates this automata  $\rightarrow O(m+n)$

- Tries (info. retrieval)  $\rightarrow$  rooted tree, each maximal path is a word  $\Rightarrow$  suffix tree
- Regular expressions

