

Distributed shortest path

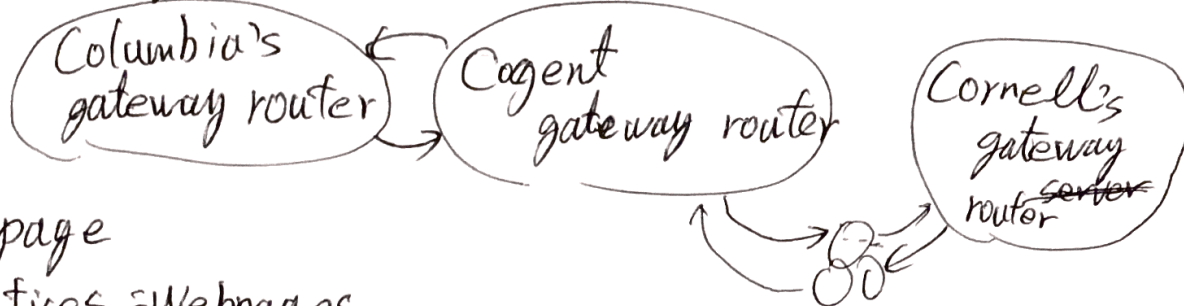
⇒ Why it is important?

⇒ Internet traffic.

Internet is a graph

- Vertices: end hosts + router
- Directed edges: physical or wireless connections

e.g.



- Webpage

- Vertices: Webpages
- Directed edges: Hyperlinks

- Social networks

- vertices: People
- Edges: Friend/Follow relationship.

⇒ How to get the shortest path?

⇒ Dijkstra?

⇒ Problem: In order to apply Dijkstra, Columbia gateway router need to know entire Internet graph.

⇒ Need a shortest-path algorithm that uses only local computation.

⇒ Bellman-Ford algorithm

Sequence Alignment

(2)

⇒ Operation to find how "similar" the two strings are.
(Used in computational genomics) DNA

- ⇒ Ex. ① Extrapolate function of genomic substring.
② Check the proximity in evolutionary tree.

What does "similar" mean?

- A G G (G) C T
- A G G (-) C A

space mismatch

⇒ can be nicely aligned.

⇒ Get the cost according to the cost and penalties of space and mismatches.

⇒ Needs the algorithm to compute fast enough to be useful. Brute-force search does not work.

⇒ Dynamic programming.

Algorithm ~~Design~~ Design Paradigms

⇒ No silver bullet.

Some paradigms

- Divide and conquer
- Randomized algorithms
- Greedy algorithms
- Dynamic programming

Greedy Algorithm

Def: Iteratively make "myopic" decisions, hope everything works out at the end.

Example: Dijkstra's shortest path

⇒ Processed each destination once, irrevocably.

Pros Cons Contrast with Divide & Conquer (No revisit)

- ① easy to propose multiple greedy algorithm for many probs.
- ② easy running time analysis. (Contrast with master method)
- ③ Hard to establish correctness.

Danger: Most greedy algorithms are not correct all the time
(~~unless~~ Regardless what intuition says)
⇒ e.g. Dijkstra don't work when it has a negative number for the distances.

How I can prove it is correction?

Method 1: Induction

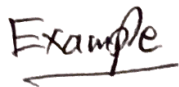
e.g.: Correctness proof for Dijkstra algorithm.

Method 2: "exchange argument"

method 3: Whatever it works.

④

- Microprocessor
CPU
Unit
Register



Request sequence: $c \rightarrow d \rightarrow e \rightarrow f \rightarrow a \rightarrow \text{oh, no!} \rightarrow b \rightarrow \text{Fuck!}$

$\downarrow \quad \downarrow \quad \downarrow \quad \downarrow$
OK OK not in cache and insert 'f'
evict 'a' and insert 'e'

-2 were inevitable (e, f)

- ⇒ Greedy algorithm is solution. (Furthest-in-future algo.)
- ⇒ Become a benchmark of algorithm.

Scheduling Problem

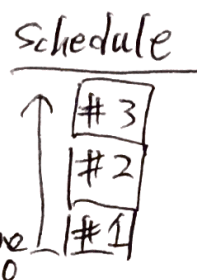
Setup: - One shared resource (e.g. a processor)
 - Many "Jobs" to do (e.g. processes)

Q: In what order should we sequence the jobs?

Assume: each job j has ① weight w_j ("Priority") ② length l_j

Def: The completion time C_j of job j = sum of job lengths up to and including j .

Example: 3 job, $l_1=1$, $l_2=2$, $l_3=3$
 $C_1=1$ $C_2=3$ $C_3=6$



Goal: Minimize the weighted sum of completion times: $\min \sum_{j=1}^n w_j C_j$

e.g. if $w_1=3$, $w_2=2$, $w_3=1$ total C time is
 $3 \cdot 1 + 2 \cdot 3 + 1 \cdot 6 = 15$

Intuition

Purpose: Want to minimize $\sum_{j=1}^n w_j \cdot C_j$

Goal: Devise correct greedy algorithm

⇒ How to figure out the algorithm?

⇒ First, think about the special case.

If all the length (C_j) are same, prioritize the ^{bigger} w_j .

If all the weight (w_j) are same, prioritize the smaller C_j .

How to prioritize the task?

(b)

- Larger the weight \Rightarrow High prior.
- Smaller the length \Rightarrow High prior.

Two way of getting cost.

① ~~$W - L$~~ (diff)

② $\frac{W}{L}$ (ratio)

only one can always work.

e.g. $L_1 = 5$ $L_2 = 2$
 $W_1 = 3$ $W_2 = 1$

ratio: $\frac{5}{3} > \frac{1}{2}$
 diff: $-2 < -1$



$5 \times 3 + 2 \times 1 = 17 \Rightarrow$ Get smaller.
 Gave optimal

$2 \times 1 + 3 \times 5 = 17 \Rightarrow$ Not always
 end in optimal.

Prim's minimum spanning tree algorithm (MST)

What? \Rightarrow Algorithm to connect a bunch of points together as cheaply as possible.

Application \Rightarrow clustering, networking

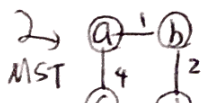
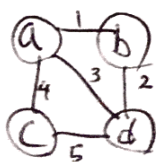
Input: Undirected graph $G=(V, E)$

- assume adjacency list representation.
- OK if edges costs are negative.

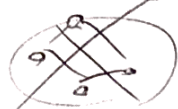
Output: Minimum cost tree $T \subseteq E$ that spans all vertices.
 i.e., sum of edge costs.

① T has no cycles

② the subgraph (V, T) is connected
 (i.e., contains path between each pair of vertices)



cost = 7



disallowed

(1)

Assumption #1: Input graph G is connected

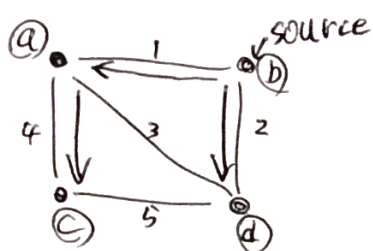
- Else no spanning trees.

- Easy to check in preprocessing (DFS, BFS)

Assumption #2: Edges costs are distinct

- Prim + Kruskal remain correct with ties.

Algorithm



1. Pick source node randomly. (a) (b)

2. Choose the cheaper edge (Greedy) and expand to a. Now tree is a-b

3. From a-b, 3 options to expand which are (c, 4), (d, 3), (d, 2).

4. Choose the cheapest one which is (d, 2). Now the tree is a-b-d.

5. From a-b-d, 2 options to expand (c, 4), (c, 5).

6. Choose the cheaper one which is (c, 4)

7. all nodes connected! Done!

Pseudo code

- initialize $X = [s]$ [$s \in V$ chosen arbitrarily]
 - $T = \text{None}$ [invariant: $X = \text{vertices spanned by tree-so-far } T$]
 - while $X \neq V$:
 - Let $e = (u, v)$ be the cheapest edge of G with $u \in X, v \notin X$
 - Add e to T
 - Add v to X
- i.e. increase # of spanned vertices in cheapest way possible.

Implementation of MST Prim's algorithm.

⑧

⇒ Running time of Prim's

linear {
- initialize $X = \{s\}$
- $T = \emptyset$
loop for $n-1$ times {
- while $X \neq V$:
- let $e = (u, w)$ the cheapest } Use "heap" to get min in $O(\log n)$.
- add e to T , add v to X
} $\Rightarrow O(n)$

How to implement the heap?

- Natural idea: Use heap to store edges, with keys = edge cost

⇒ Better idea?

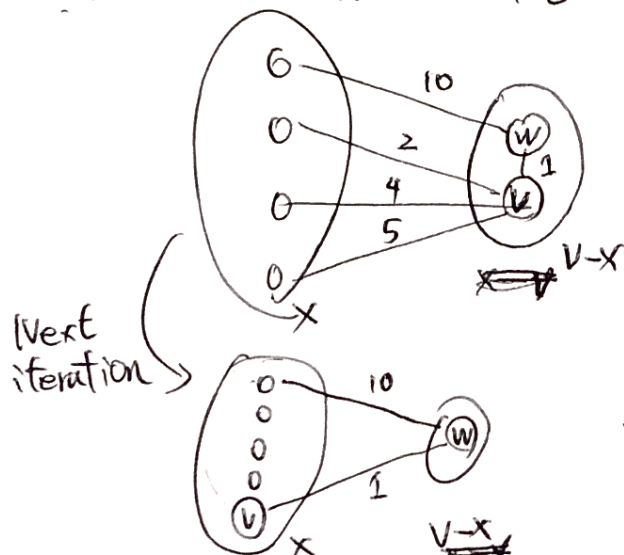
⇒ Store the nodes!

How to store the nodes in heap?

Invariant #1: Elements in heap = vertices of $V - X$

Invariant #2: For $v \in V - X$, $\text{key}[v]$ = cheapest edge (u, v) with $u \in X$.

Check: Initialize the heap with $O(m \log n)$ preprocessing.



Current $\text{key}[v] = 2$ (cheapest in $\{2, 4, 5\}$)
current $\text{key}[w] = 10$ (Only edge)
Next iteration $\text{key}[w] = 1$ (Cheapest in $\{1, 10\}$)

⇒ Have to recompute the key every iteration.

How to update the key value?

(9)

when v added to X :

- for each edge $(v, w) \in E$:

- if $w \in V - X$ \longrightarrow the only vertices whose key might have changed

- Delete w from heap

- Recompute $\text{key}[w] := \min[\text{key}[w], \text{Cost}_{vw}]$ update key if needed.

- Re-Insert w into heap

$\min[\text{key}[w], \text{Cost}_{vw}]$

$\Rightarrow O(m \log n)$