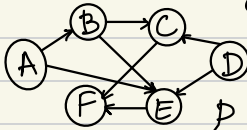


0.1 Topological Orderings

Def. Directed Acyclic Graph (DAG) - directed graph with no directed cycles

Precedence constraints: edge (v_i, v_j) means v_i must happen before v_j

Def. Topological ordering of directed graph $G = (V, E)$ - ordering of nodes (V) v_1, v_2, \dots, v_n s.t. for every edge (v_i, v_j) $i < j$



constraints: before outgoing edges, after incoming edges

$A < B, E$ $B, D < C < F$ $A, B, D < E < F$

$A < B < C$ $D < C, E$ $C, E < F$

$\Rightarrow A \rightarrow B \rightarrow E \rightarrow C \rightarrow F \Rightarrow \{A, B, D, E, C, F\}, \{A, D, B, E, C, F\}, \{D, A, B, E, C, F\}$

function TopologicalOrdering(G)

if $|V| = 0$ then

return empty list

find a node v with no incoming edges

remainder \leftarrow TopologicalOrdering($G - \{v\}$)

return v + remainder

$O(n+m)$, $|V|=n$, $|E|=m$

Lemma $G - \text{DAG} \rightarrow G$ has a node with no incoming edges

Proof by contradiction

$G - \text{DAG}$

Suppose \nexists node $\rightarrow \exists$ incoming edge to node.

Pick a node and follow $(|V|)$ edges backwards.

Must have visited a node twice. $\Rightarrow G$ has a cycle. \Rightarrow Contradiction.

Lemma $G - \text{DAG} \rightarrow G$ has topological ordering

Proof by induction

Base case: $|V|=1$, holds because G is a topological ordering.

Induction Hypothesis: $G - \text{DAG}$, $|V| \leq k \rightarrow G$ has topological ordering.

Induction step:

Given DAG G , $|V|=k+1$.

Find node v with no incoming edges ($\exists v$ by previous lemma).

$G - \{v\}$ is DAG, removing v cannot create cycles.

By IH, $G - \{v\}$ has $|V|=k \rightarrow$ has topological ordering T .

$\Rightarrow \{v, T\}$ is valid topological order of G , v has no incoming edges.

0.2 Time Complexity

Def. $T(n)$ is $O(f(n))$ iff $\exists c > 0 \in \mathbb{R}, n > 0 \in \mathbb{N}$ s.t. $\forall n \geq n_0: T(n) \leq c \cdot f(n)$.

Def. $T(n)$ is $\Omega(f(n))$ iff $\exists c > 0 \in \mathbb{R}, n > 0 \in \mathbb{N}$ s.t. $\forall n \geq n_0: T(n) \geq c \cdot f(n)$.

Def. $T(n)$ is $\Theta(f(n))$ iff $T(n)$ is $O(f(n))$ and $T(n)$ is $\Omega(f(n))$.

$\forall d > 0 \in \mathbb{R}, \log n$ is $O(n^d)$

$O(\log_a(n)) = O(\log_b(n))$ for any constant $a, b > 0$ because $\log_a(n) = \frac{\log_b(n)}{\log_b(a)}$ (constant / ratio)

1.1 Cycling Trip

Problem: Selecting breakpoints

route of length L , camp at points with distances b_1, b_2, \dots, b_n from start point at most C kilometers per day, reach end in as few days as possible

Solution: Camp at last possible campsite each day.

sort breakpoints in descending order ($0 = b_1 \leq b_2 \leq \dots \leq b_n = L$)

$S \leftarrow \{0\}$

$x \leftarrow 0$

while $x \neq L$ do

p - largest integer s.t. $b_p \leq x + C$

 if $b_p = x$ then

 return No solution

$x \leftarrow b_p$

$S \leftarrow S \cup \{p\}$

return S

$O(n \log n)$

1.2 Proof by Contradiction

Th. Greedy is optimal

Assume greedy is not optimal.

Let $0 = g_1 < g_2 < \dots < g_r = L$ denote campsites chosen by greedy.

Let $0 = h_1 < h_2 < \dots < h_r = L$ denote campsites in optimal solution where $h_1 = g_1, h_2 = g_2, \dots, h_r = g_r$

for the largest possible value of r , but $g_{r+1} \neq h_{r+1}$

Note that $g_{r+1} \geq h_{r+1}$ by greedy choice of algorithm, so $g_{r+1} > h_{r+1}$

Since $g_{r+1} - g_r \leq C$ and $h_{r+1} - g_{r+1} < h_{r+1} - h_r \leq C \Rightarrow$ replacing h_{r+1} with g_{r+1} maintains optimality.

1.3. Proof by Induction

Lemma $\forall r, f_r \leq g_r$, f -Greedy campsites, f -optimal solution (Greedy stays ahead)

Base case ($r=1$): $f_1 \leq g_1$

I.H: $f_k \leq g_k$

Induction step:

To prove: $f_{k+1} \leq g_{k+1}$

By I.H, $f_k \leq g_k$

f_{k+1} - within reach ($\leq c$) but as far as possible from f_k .

g_{k+1} - not further than the farthest reachable from g_k , but by I.H, definitely reachable from f_k .

$\Rightarrow f_{k+1} \leq g_{k+1}$

Th Greedy is optimal

Proof by contradiction

$f \rightarrow k$ campsites, $f \rightarrow m$ campsites.

Suppose greedy is not optimal $\Rightarrow k > m$.

By Lemma, $f_k \leq g_k$.

$\Rightarrow m \geq k$

1.4. Interval Scheduling

Job (activity) j : starts at s_j and finishes at f_j .

Two jobs are compatible if no overlap.

Find maximum number of compatible intervals (jobs).

Solution: Consider jobs in ascending order of finishing time f_i .

sort jobs by ascending finish time ($f_1 \leq f_2 \leq \dots \leq f_n$)

$A \leftarrow \{\}$

for $k \leftarrow 1$ to n do

if job k compatible with A then

$A \leftarrow A \cup \{k\}$

return A

$O(n \log n)$

1.5. Greedy Stays Ahead

Lemma $\forall r \leq k: f(i_r) \leq f(j_r)$, i - Greedy schedule, j - optimal one (Greedy stays ahead)

Proof by induction

Base case ($r=1$): by order of jobs $f(i_1) \leq f(j_1)$

I.H: $f(i_r) \leq f(j_r)$

Induction step:

To prove: $f(i_{r+1}) \leq f(j_{r+1})$

$f(j_r) \leq s(j_{r+1})$ no overlap.

$f(i_r) \leq s(j_{r+1})$ by I.H

Greedy can take j_{r+1} , but chooses smallest end time $f(i_{r+1}) \leq f(j_{r+1})$

Th Greedy is optimal

Proof by contradiction

k - #jobs in i (Greedy), m - #jobs in j (optimal)

Suppose Greedy is not optimal $\Rightarrow k < m$.

By Lemma, $f(i_k) \leq f(j_k)$

$\Rightarrow j_{k+1}$ that starts after j_k thus after i_k .

$\Rightarrow j_{k+1}$ that is compatible with i . Contradiction!

2.1. Student Scheduling

Problem: Minimizing Lateness

Single resource processes one job at a time. Job j requires t_j time and is due d_j .

If j starts at s_j , then it finishes at $f_j = s_j + t_j$. Lateness of j is $L_j = \max(0, f_j - d_j)$.

Minimize the maximum lateness $\max(L_j)$

Solution: Consider jobs in ascending order of deadline d_j
sort jobs by ascending deadline ($d_1 \leq d_2 \leq \dots \leq d_n$)

$t \leftarrow 0$

$A \leftarrow \{0\}$

for $j \leftarrow 1$ to n do

$t \leftarrow t + t_j$

$A \leftarrow A \cup \{t\}$

return A

$O(n \log n)$

2.2 Exchange Argument

Def: Inversion in schedule - a pair of jobs s.t. $d_i < d_j$ but j is scheduled before i .

Claim: Swapping two adjacent inverted jobs reduce the number of inversions by 1 and does not increase the maximum lateness.

L - lateness before swap, L' - after.

$$L'_k = L_k \quad \forall k \neq i, j$$

$$L'_i \leq L_i \quad (i \text{ done earlier})$$

$$L'_j = f_j - d_j = f_i - d_j \leq f_i - d_i = L_i$$

Claim: If a schedule (with no idle time) has an inversion, then it has one with a pair of inverted jobs scheduled consecutively (adjacent jobs).

Suppose there is an inversion.

There is a pair of jobs i and j s.t. $d_i < d_j$, but j scheduled before i .

Walk through the schedule from j to i .

Increasing deadlines = no inversions, at some point deadline decreases

Th: Greedy schedule S is optimal.

Suppose S is not optimal.

Let S' be an optimal schedule with the least inversions and no idle time $S \neq S'$.

If S' has no inversions, then the lateness is same as S . Contradiction!

If S has an inversion, let $i-j$ be an adjacent inversion.

Swapping i and j does not increase the maximum lateness and decreases the number of inversions \Rightarrow contradiction with definition of $S' \Rightarrow S$ is optimal.

3.1 Connecting all machines

Problem: Creating the internet

n machines and m connections, each with cost $c(i)$

Find cheapest way to connect all the machines.

Def: Spanning tree T of graph $G=(V, E)$ - subset of edges $T \subseteq E$

s.t. graph $G'=(V, T)$ is strongly connected.

Minimum Spanning Tree (MST) of graph $G=(V, E)$ - spanning tree T

s.t. $c(T) = \sum_{e \in T} c(e)$ is minimal, $\forall T'$ s.t. $c(T') < c(T)$

Solutions: Prim (Jarník), Kruskal, Reverse Delete, Borůvka

3.2. The Cut Property

S -subset of nodes V , e - min cost edge with exactly 1 endpoint in $S \Rightarrow e \in \text{MST}$.

3.3. Prim-Jarnik

Grow a cloud (S) by repeatedly adding the smallest edge out of the cloud.

function PrimJarnik(G)

$d[v] \leftarrow \text{INF}$ for all v in V

s - arbitrary vertex in V

$d[s] \leftarrow 0$; $c \leftarrow 0$

while not all vertices connected do

$m \leftarrow \arg \min d[v \text{ in } V]$

$c \leftarrow c + d[m]$

Remove m from V

if $d[m] = \text{INF}$ then

return INF

for $e \leftarrow (m, u)$ in E for some u do

if $w(e) < d[u]$ then

$d[u] \leftarrow w(e)$

return c $O(|E| \log |V|)$, even $\Theta(|E| + |V| \log |V|)$ with heap-based PQ

3.4. The Cycle Property

C - any cycle in G , e - max cost edge in $C \Rightarrow e \notin \text{MST}$.

3.5. Kruskal

Sort edges by weight, repeatedly add smallest edge (from cut property), until there is a cycle, then discard it (from cycle property).

function Kruskal(G)

sort E ascendingly ($e_1 \leq e_2 \leq \dots \leq e_m$)

$T \leftarrow \{\}$; $i \leftarrow 0$; $k \leftarrow 0$

while $k < |V| - 1$ do

if $T \cup \{e_i\}$ does not have a cycle then

$k \leftarrow k + 1$

$T \leftarrow T \cup \{e_i\}$

return T

$O(|E| \log |V|)$

3. Union-Find

Union-Find (Disjoint-Set) data structure - track elements partitioned into disjoint subsets; near-constant time merge(union) and contains(find)

MakeSet(n) - n initial sets, with 1 element each

set \leftarrow array of size n , set[i] $\leftarrow i$

rank \leftarrow array of size n , rank[i] $\leftarrow 0$

Find(x) - returns id of set containing x

if set[x] $\neq x$ then

set[x] \leftarrow Find(set[x])

return set[x]

Union(x, y) - merges set of x and y

xSet \leftarrow Find(x)

ySet \leftarrow Find(y)

if xSet = ySet then

return False

else if rank[xSet] < rank[ySet] then

set[xSet] \leftarrow ySet

else

set[ySet] \leftarrow xSet

if rank[xSet] = rank[ySet] then

rank[xSet] \leftarrow rank[xSet] + 1

return True

4. K-clustering

Problem: Maximum Spacing K-Clustering

Divide objects into K non-empty groups.

$d(p_i, p_j) = 0$ iff $p_i = p_j$ otherwise $d(p_i, p_j) \geq 0$; $d(p_i, p_j) = d(p_j, p_i)$

Spacing - minimum distance between any pair of points in different clusters.

Solution: Single-list k-clustering

Form graph (without edges) on vertex set (n initial clusters).

Find closest pair of objects from different clusters and add edge

Repeat $n-k$ times until exactly k clusters left.

Solution: Find MST, delete $k-1$ most expensive edges.

Thm C' - clustering C'_1, \dots, C'_k formed by deleting $k-1$ most expensive edges of a MST by Kruskal $\rightarrow C'$ - maximum spacing clustering

Proof

Spacing of C' - length of d' of the $(k-1)^{th}$ most expensive edge.

C - arbitrary clustering C_1, \dots, C_k .

p, q - in same cluster in C' (C'_i), but different ones in C (C_s and C_t).

Some edge (p', q') on p - q path in C' spans two different clusters in C .

All edges on this path have length $\leq d'$ (Kruskal).

Spacing of $C \leq d'$ since p', q' in different clusters

C - arbitrary \Rightarrow True for $\forall C \neq C'$

5. Huffman Encoding

Problem: Efficient Encoding

non-ambiguous binary encoding of text

Def prefix code for set S - $c: S \rightarrow \{0, 1\}^+$ s.t. $\forall x, y \in S: x \neq y \rightarrow c(x)$ is not prefix of $c(y)$

Average Bits per Letter: $ABL(c) = \sum_{x \in S} f(x) \cdot |c(x)|$

Binary Tree:

\rightarrow Children uniquely identified by edge label (0 or 1).

\rightarrow Nodes labeled x iff path from root labeled $c(x)$

\rightarrow Only leaves can have a label in a prefix code.

Def full binary tree - node has either 0 or 2 children

Lemma u, v - leaves in T' and $\text{depth}_{T'}(u) < \text{depth}_{T'}(v) \rightarrow f_u \geq f_v$

Claim $\forall T, \exists T'$ optimal, where two lowest-frequency - sibling leaves at lowest level.

function Huffman(S)

if $|S| = 2$ then

return tree with root and 2 leaves

else

let y and z be the lowest frequency letters in S

$S' \leftarrow S - \{y, z\} \cup \{w\}$, s.f. $f_w = f_y + f_z$

$T' \leftarrow \text{Huffman}(S')$

$T \leftarrow$ add two children y and z to leaf w in T'