

9.0

Last week we gave some examples of greedy algorithms.

We still have to prove two lemmas concerning Huffman's alg.

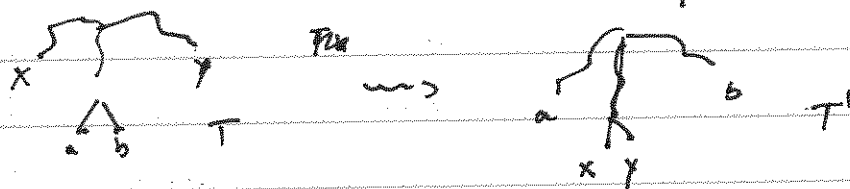
Then we will start discussing graphs and graph algorithms.

Lemma Let x, y be with lowest freq. Then, there exists an optimal prefix-code C in which codewords for x and y have the same length and they differ in only one ~~bit~~ the last bit

Pf We will start with any optimal tree and get another optimal tree in which this property holds.

Let T be an optimal tree and let a, b be siblings at maximal depth. Assume w.l.o.g. $x.\text{freq} \leq y.\text{freq}$, $a.\text{freq} \leq b.\text{freq}$.

We would like to replace x with a and y with b .



what will be the diff between the encodings:

$$\sum_{c \in C} c.\text{freq} \cdot d_T(c) - \sum_{c \in C} c.\text{freq} \cdot d_{T'}(c) =$$

$$x.\text{freq} (d_T(x) - d_{T'}(x)) + a.\text{freq} (d_T(a) - d_{T'}(a)) + y.\text{freq} (d_T(y) - d_{T'}(y)) + b.\text{freq} (d_T(b) - d_{T'}(b))$$

$$\text{Notice } d_T(x) = d_{T'}(a), d_{T'}(x) = d_T(a) \quad \text{same for } y, b. \quad \Rightarrow$$

$$\begin{aligned} &= (a.\text{freq} - x.\text{freq}) (d_T(a) - d_{T'}(a)) + (b.\text{freq} - y.\text{freq}) (d_T(b) - d_{T'}(b)) \\ &= (a.\text{freq} - x.\text{freq}) \underbrace{(d_T(a) - d_T(x))}_{\geq 0} + (b.\text{freq} - y.\text{freq}) \underbrace{(d_T(b) - d_T(y))}_{\geq 0} \geq 0 \end{aligned}$$

By assumption $x.\text{freq} \leq y.\text{freq}$ are smallest. So it ~~(a,b) < (x,y)~~ then ~~at least one of them is different~~. Thus, the entire sum is non-negative. Thus T' at least as good as T .

The proof also gives that x, y are siblings at maximal depth!

9.2

2.4

we are almost done. We saw that in ^{some} the optimal code, x, y will be siblings at ~~max~~ depth. But ~~is then~~ the question is, do we get an optimal code by our process?

Lemma: Let C' be the new alphabet after "merging" x, y . Let T' be any optimal ~~prefix~~ prefix-code / Tree for C' . Let T be obtained by adding two ~~leaves~~ for children to z in the obvious way. Then T is optimal for C .

Proof: ~~The~~ Let $B(T'), B(T)$ be the costs of T', T .

Then $B(T) = B(T') + x \cdot \text{freq} + y \cdot \text{freq}$ (everything stayed the same except z got two new children)

$$\Rightarrow B(T') = B(T) - x \cdot \text{freq} - y \cdot \text{freq}.$$

Assume for a contradiction that T is not optimal. Let S be a better tree, wlog S has x, y as siblings. ~~Let~~ Let S' be same as S with z instead of x, y .

$$B(S') = B(S) - x \cdot \text{freq} - y \cdot \text{freq} < B(T) - x \cdot \text{freq} - y \cdot \text{freq} = B(T') \quad \text{our assumption}$$

contradicting T' 's optimality.

The two lemmas imply that Huffman is optimal.

9.3

We now move to graphs and related algorithms.

A graph has a set of vertices V (in our case V will always be finite) and a set of edges $E \subseteq V^2$. If $(u, v) \in E$ we imagine an edge going from u to v .

~~We say the graph~~ A graph as we described is called a directed graph.

Notice that we can have both $(u, v) \in E$ and $(v, u) \in E$.

We can also have a self-loop - $(u, u) \in E$.

A graph is undirected if edges are unordered pairs $\{v_i, v_j\}$.

The degree of a vertex is the number of neighbors in an undir. graph.

In a directed graph we have in-degree and out-degree.

A path is a sequence v_1, v_2, \dots, v_k where $v_i (v_i, v_{i+1}) \in E$.

v is reachable from u if there is a path from u to v .

A path is simple if all vertices are distinct.

A cycle is a ~~non~~ path from u to u with

~~A graph is~~ An undirected graph is connected if there is a path between any two vertices.

A dir. graph is strongly connected if there is a dir. path between any two vertices. (both directions).

The connected component of u in an undir. graph is all vertices reachable from u .

The strongly connected component of u in a dir. graph is the set of all vertices that have a dir. path to and from u .

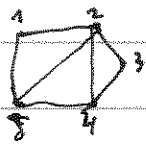
9.4

We would like to work with graph so we have to find a way to represent them that is efficient and easy to work with.

One representation is by the adjacency matrix and the other by adjacency list.

Adj-list: for every u we have a list of all v ~~that~~ st $(u,v) \in E$

E.g. for an undir graph.



$1 \rightarrow 2 \rightarrow 3$
 $2 \rightarrow 1 \rightarrow 4 \rightarrow 5 \rightarrow 3$
 $3 \rightarrow 2 \rightarrow 4$
 $4 \rightarrow 2 \rightarrow 5$
 $5 \rightarrow 4 \rightarrow 1 \rightarrow 2$

The adjacency matrix is a $|V| \times |V|$ matrix A with $(\forall i,j) A_{ij} = 1 \Leftrightarrow (i,j) \in E$

$$\begin{array}{c}
 \begin{array}{ccccc}
 & 1 & 2 & 3 & 4 & 5 \\
 \begin{array}{c} 1 \\ 2 \\ 3 \\ 4 \\ 5 \end{array} & \begin{pmatrix} 0 & 1 & 0 & 0 & 1 \\ 1 & 0 & 1 & 1 & 1 \\ 0 & 1 & 0 & 1 & 0 \\ 0 & 1 & 1 & 0 & 1 \\ 1 & 1 & 0 & 1 & 0 \end{pmatrix}
 \end{array}
 \end{array}$$

Adj-list is shorter but "harder" to find whether $(u,v) \in E$. (have to go through u 's list). They are of size $|E|$ and $|V|^2$ so if $|E| \approx \mathcal{O}(|V|^2)$ they are of similar sizes.

We now present two canonical alg for searching graphs, which are at the core of many other algorithms.

Breadth First search (BFS)

BFS is an alg. that given a graph $G=(V,E)$ and a source vertex $s \in V$ explores the graph and discovers every vertex reachable from s along with its distance from s .

The distance of u from s is the length of the shortest path from s to u .

The alg. also produces the BFS tree with root s and all reachable vertices.

It works both for directed and undirected graphs.

The alg. works by first discovering s 's neighbors, then their neighbors etc. At each step the frontier advances by between discovered and undiscovered vertices across the breadth of the frontier.

Give a picture.

We shall color vertices white if they were not discovered yet, gray if they are at the frontier and black if we are done with them - i.e. we explored all their nbrs.

The alg. will use a queue - a first-in first-out data structure.

Enqueue puts an element at the end of the queue and Dequeue extracts the top first element from the queue.

9.6

BFS(G, s)

1. For each $u \in G.V \setminus \{s\}$
2. $u.color = \text{white}$
3. $u.d = \infty$ distance
4. $u.\pi = \text{NIL}$ predecessor
5. $s.color = \text{gray}$
6. $s.d = 0$
7. $s.\pi = \text{NIL}$
8. $Q = \emptyset$
9. Enqueue(Q, s)
10. While $Q \neq \emptyset$
11. $u = \text{Dequeue}(Q)$
12. For each $v \in G.Adj[u]$
13. If $v.color == \text{white}$
14. $v.color = \text{gray}$
15. $v.d = u.d + 1$
16. $v.\pi = u$
17. Enqueue(Q, v)
18. $u.color = \text{black}$

Run time: Each vertex enters Q at most once (color change). When we dequeue a vertex we go over all its neighbors. Thus run time is $\mathcal{O}(V+E)$

Claim: BFS discovers every vertex v that is reachable from s , upon termination $v.d = \text{dist}(s, v)$, and a shortest path from s to v is given by whatever path from s to $v.\pi$ and the edge $(v.\pi, v)$.

1.7

Furthermore, let

$$V_\pi = \{v \in V : v.\pi \neq NIL\} \cup \{s\}$$

$$E_\pi = \{(v.\pi, v) : v \in V_\pi \setminus \{s\}\}$$

This is the predecessor subgraph of G and s .

*

Claim: The predecessor graph produced by BFS is a tree consisting of all reachable vertices from s , and for all $v \in V$ it contains a unique simple path of length $d(v)$ from s to v .

Both claims are not hard to prove by induction.

Depth-First-Search (DFS)

This is a different searching alg. that works by going as deep as possible before exploring other vertices.

Basically, DFS explores edges out of the most recently discovered v (so e.g. v 's first grandchild will be discovered before v 's 2nd child: j_2^v). Once all of v 's edges have been explored we back track to explore edges leaving the vertex from which v was discovered.

For each $v \in V$ we will also record the time when it was discovered - and also the time when we finished with it - u.t.

9.3

DFS(G)

1. For each $u \in G.V$
2. $u.color = white$
3. $u.\pi = NIL$
4. $time = 0$
5. For each $u \in G.V$
6. If $u.color == white$
7. DFS-visit(G, u)

If we just want the DFS trees
we should run DFS-visit(G, u)

DFS-visit(G, u)

1. $time = time + 1$
2. $u.d = time$
3. $u.color = gray$
4. For each $v \in G.Adj[u]$
5. If $v.color == white$
6. $v.\pi = u$
7. DFS-visit(G, v)
8. $u.color = black$
9. $time = time + 1$
10. $u.f = time$

u is discovered

we explore along the edge (u, v)

we finished exploring all u 's nbrs.

Run-time of DFS is also $O(V+E)$.

~~Indeed, time is updated whenever~~

Indeed we run DFS-visit only when we reach a white neighbor. Thus, it is called once for each vertex. (a vertex is painted gray immediately at DFS-visit).

Then in DFS-visit the total number of ~~executions~~ commands to execute is $|Adj[v]|$

As each DFS-visit(u) is run only once the total # of operations is $O(\sum_{+O(1)} |Adj[v]|) = O(E+V)$

Some important properties of DFS:

1. (V, E_T) is a forest

2. If u is ancestor of v is the DFS forest iff at time u and there is a $u \rightarrow v$ path consisting of only white vertices (before changing u 's color to gray)

3. If u is ancestor of v then $[v.d, v.f] \subseteq [u.d, u.f]$. If they are unrelated the interval are disjoint. In fact, this is iff statement.