Huffman Encoding Learning Goals · Describe "Binary code," "Prefix free, "Average letter length" · Explain connection between binary codes + trees · Describe Huffman's alg. · Analyze runtime of Huffman's alg / · Describe impact of data structures alg runtime / · Prove correctness of Huffman's alg ~~~~ ~

Do we use induction to greedy algorithms? prove correctness of A) Never B) Sometimes C) Always What is true about strong + regular induction A) They are equivalent, the choice is stylistic B) You can always use strong induction C) Regular induction can be used to prove Strong induction. K+1 F FK/2

Binary Codes ex; Z= Za, b, c, d, ..., ZZ def: Given an alphabet Z_1 , a binary code is a function $f: Z \rightarrow Z_{0,1}$? ex: Braille [] [OIDIOI], ASCI, Morse Code Suppose you have a message where the letter "a" OCCURS 50% of the time, "b" 30%, and "c" 20%. best binary encoding of Z= Za, b, c]? Which is the A): f(a) = 00 B) f(a) = 0 C) f(a) = 0f(b) = O(b)f(b) = 1f(b) = 10 $f_{\alpha}(c) = 1$ f(c) = 10f(c) = 01

.

C) f(a) = OA): f(a) = OOB) f(a) = Df(p) = 0f(b) = 10f(b) = 1f(c) = 1f(c) = 10f(c) = ODoesn't take Ambrquous Not ambiguous, for Jecoding. average bits advantage of per letter is differing occurance rates 61 → c? small . Ol Jab? # of bits in f(i) Average letter length: $L(f) = \sum_{i \in \mathbb{Z}} |f(i)| \cdot p(i)$ ex. C] $|f(a)|p(a) + |f(b)| \cdot p(b) + |f(c)| \cdot p(c)$ = 1 · 0.5 + 2 · 0.3 + 2 · 0.2 = 1.5

Binary Trees & Binary Codes root 0/1 f(a) = 0f(a) = Of(b) = 0 $f(b) = 10 \longleftarrow$ (b) $f(L) = \langle \rangle$ f(c)=1There is ambiguity for decoding if multiple letters lie on same path O | O | | O | ...a b m'a def: A code is "prefix free" if all letters are at leaves in corresponding binary tree.

Merge Trees to Create Prefix Free Codes é α 6 С 9 6 9 J Ù 6 6 ġ g

Optimal Binary Encoding Problem
Input:
$$\Sigma$$
 (alphabets of symbols)
 $P: \Sigma \rightarrow \mathbb{R}$ (probabilities/frequency for each symbol)
Output: $f: \Sigma \rightarrow \Xi_{0,1}Z^{\dagger}$ s.t.
• f is prefix free
• minimize average letter length
 $L = \Sigma P(i) | f(i) |$
 $i \in \mathbb{Z}$
Objective function

Huffman's Algorithm For each it z: 1. Create a tree with one node, label i . Give tree weight p(i) While there is more than one tree: . Merge two trees with smallest weight

. Set weight of merged tree to be sum of weights.

i p(i)	· Use Huffman's algorithm to create a binary
a .3	Code
b.25	· What is the average letter length of your
C . 2	code.
d 1.15	· What is the runtime of Huffman's in terms of
e).[[Z]=n? Ideas to improve?
	· Why greedy?

.25 .[5 .2 .3 Ē (\mathcal{A}) 6 .25 .25 .2 .3 E \hat{a} $\overline{\mathbb{O}}$ e d • 45 .25 .3 e d С **,** 55, •45 α C J. e d Ċ α (E d $L = .3 \cdot 2 + .25 \cdot 2 + .2 \cdot 2 + .15 \cdot 3 + .13$ Ô = 2.25 Ο $\langle \overline{a} \rangle$ 010 0 è

Z=N $\rightarrow (\mathcal{M}^2)$ O(n) For each i e Z: ·Create a tree with one node, label ·Give tree weight p(i) O(i) $O(1)^{1}$ O(n) While there is more than one tree: . Merge two trees with smallest weights O(n) . Set weight of merged tree to be sum of weights O(1) Why greedy:

 $O(n \log n)$ Huffman's Algorithm · Create a tree with one node, label i o(n) · Give tree mpichil - ~ For each ie 2: Initiatize heap with Jour In Trees ~ O(nlogn) While there is more than one tree: O(n) Merge two trees with smallest weights O(logn)
Set weight of merged tree to be sum of weights
Peinsert merged tree into neap O(logn) $\mathcal{O}(\mathbf{1})$ Improve runtime? · At each iteration of while loop, find minimum value tree. · Helpful data structure?? Min Heap O(nlogn) · Initialize n items in heap

O(logn) Remove minimum we item

O(logn) · Insert new item

Go Program! (Lots of details to figure out!) (see programming assignment) Ethical Matrix tree merges Thm: Huffman's Algorithm produces a prefix free code That minimizes average letter length. PF: We will prove correctness by induction on N = |Z|. Base Case: If N=2, there are 2 characters: a, b. Huffman does a a b a (b) This has optimal average letter length. (1)

Inductive Step: Assume for induction that Huffmans alg produces à prefix free code with optimal average letter length for any alphabet with k characters. Consider an input Z s.t. |Z| = k+1. Let a, b be the characters of Z with the Smallest p-values. Define Z=Z-Za, bj UZa/bj where a/b is a new character with p(a/b) = p(a) + p(b). ex: Z = Ze, f, g, h, f = p(e) = . l = p(f) = .7 p(g) = .15 p(h) = .05 $\Sigma^{-} = \{e/h, f, g\} p(e/h) = .15 p(f) = .7 p(g) = .15$

ex: $\Sigma = \Xi e, f, g, h \Xi p(e) = . | p(f) = .7 p(g) = .15 p(h) = .05$ They 5= Huffman 2 .15 J °.D5 .7 (\mathcal{O}) (E .15 .7 .15 .15 .7 .15 F (9) g én (Ì) (7) à elv (q 9

In general: Huffman By inductive assumption, T-gives an optimal binary code (b/c has K characters, created using Huffman) Lemma: There is an optimal tree for Z where a, b are siblings. (will prove later)

Suppose for contradiction that T is not optimal. Let T++T be an optimal tree with a, b siblings (by Lemma) on 2 Define T*-9+1 (fill in with d, p(a), p(b)) φ Then $L(T^*) = \sum_{i \neq q, b \in S} P(i)$ $P(\alpha)(d+1) + P(b)(d+1)$ + $L(T^{*-}) = \sum_{i \neq a/b} P(i) d(i)$ + P(a|b) d $(p(\alpha) + p(b)) d$

So $L(T^*) - L(T^{*-}) = P(a) + P(b)$ Similarly $L(T) - L(T^{-1}) = P(\alpha) + P(b)$ $\frac{1}{1} = L(T) - L(T^{*}) = L(T) - L(T^{*})$ Thus Rearranging: $L(T) - L(T^*) = L(T) - L(T^*)$ This is a contradiction because Left side is positive because T* by our assumption, is a tree with optimal average letter length, while T, our Huffman alg, does not. Right hand side is less than or equal to 0, because T by inductive assumption has minimal av. letter length, and Tt-

is any other free on Z. => Contradiction. Thus Huffmans alg works correctly for any alphabet with Krl letters, so we have proved the inductive step. Lemma: There is an optimal tree for Z with a, b (characters with smallest p-values) siblings. Let T' be any optimal prefix free tree where a b are not siblings. Let x, y be siblings in T' at Max depth. da) I do I like exchange x = a, and y = b, the new tree Tex will have the same or smaller average letter length as T' because... 1 dmax