

Huffman Codes : are example of greedy algorithms.

Encoding data on arbitrary alphabets in binary.

Setup:

$\Sigma$  - alphabet

- Think of the english alphabet or a subset of it.
- OR  $\{0, 1\}$  - - {0, 1, 10, 11}

Binary code for  $\Sigma$  is a way of writing each symbol of  $\Sigma$  as a distinct string over  $\{0,1\}$

- A binary code for  $\Sigma$  gives us a natural way of encoding strings over  $\Sigma$  as binary strings.  
Necessary for storing inherently non-binary data
- Desirable from such an encoding
  - no collision → should be able to 'decode'
  - efficiency (encoding, decoding complexity length of encoding)

E.g.:

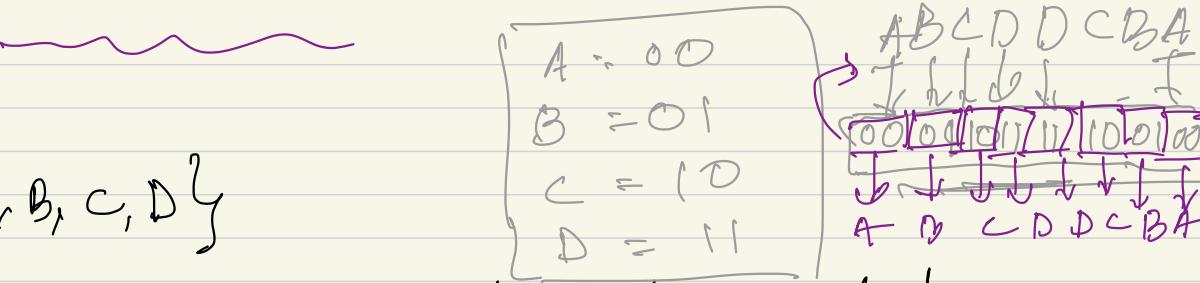
$$1) \sum_{i=1}^4 \{A, B, C, D\}$$

Associate a 2-bit string to each symbol.

$$2) \sum \rightarrow 26 \text{ symbols. } \leftarrow \{\textcircled{A}, \textcircled{B}, \dots, \textcircled{Z}\} \quad |\sum| = 26$$

$\underline{2^5 \geq 26}$

Q. How do we encode and decode decode → natural requirement.



Note: - the length of encoding of each symbol is the same.

- very natural solution

But: can also think of variable length binary code.

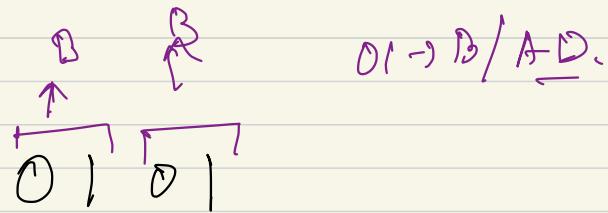
Why?

Might help with 'efficiency'.

	$C_1$	$C_2$
A	00	0
B	01	01
C	10	10
D	11	1

① Efficiency of  $C_2$  over  $C_1$

A A A D A D D B C ✓



② what about decoding.

$\underline{C_1}:$

$\underline{C_2}:$

Want such codes that are uniquely decodable,

- can always do fixed length encoding

- but can we do better.

Prefix free codes.

$\forall \alpha, \beta \in \Sigma$ ,  $\text{Enc}(\alpha)$  is not a prefix  
of  $\text{Enc}(\beta)$  and vice versa.

- $C_2$  - not prefix free
  - All fixed length codes are prefix free.
  - How does prefix freeness help?
- In many cases, not all alphabet symbols in  $\Sigma$  are equally likely.
- Makes sense to encode more frequent symbols by shorter codewords
- ↳ variable length coding appears more efficient.

Ex:

Advantages of variable length prefix free coding.

<u>Alphabet</u>	<u>Frequency</u>	<u>C<sub>1</sub></u>	<u>C<sub>2</sub></u>	Note
A	60%	00	0	prefix free
B	25%	01	10	
C	10%	10	110	
D	5%	11	111	

$s \in \Sigma^n$ ,  $\left| Enc_{C_1}(s) \right| \leq 2n$ .  
 satisfies freq.  
 requirements

$$\begin{aligned} |Enc_{C_2}(s)| &= \\ &= 1 \times 0.6n + 2 \times 0.25n + \\ &\quad 3 \times 0.1n + 5 \times 0.05n \\ &= 1.55n < 2n. \end{aligned}$$

Suppose in a document, symbols {A, B, C, D} appear with the aforementioned freq.

What is the length of the codeword under C<sub>1</sub>, C<sub>2</sub>?

wtd average of lengths of encoding of symbols in  $\Sigma$ :

Question of interest.

Input: Non negative frequencies  $p_a$  for each symbol  $a$  of an alphabet  $\Sigma$  of size  $n \geq 2$ .

Output: A prefix free binary code with  $\boxed{\min}$  avg encoding length under these frequencies.

$$\text{Avg Enc Length} := \sum_{a \in \Sigma} p_a \cdot (\# \text{ bits in Enc}(a))$$

How do we know these frequencies?

Domain knowledge - linguists  
- Biologists for  
Biological data

Here, we assume that the freq are available to us.

- How do we start?
- What is a trivial solution?
- We need to somehow be able to merge optimality.

→ prefix free condition makes it hard to think about.

- Try to reformulate this problem as a question about labeled trees
- Will make it slightly easier to think about

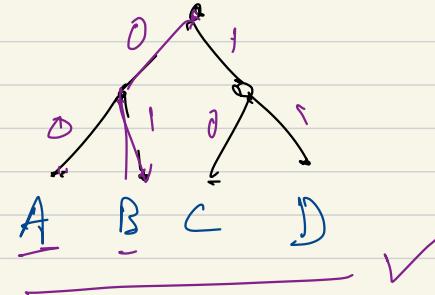
Codes vs Trees.

Ex1 A 0 0 ✓

B 0 1

C 1 0

D 1 1



A binary tree rep. the encoding.

Every internal vertex has a left edge labeled 0  
— Right edge ~ 1

If the label on both from root to a vertex  $v$

equals  $\text{Enc}(\alpha)$  for some  $\alpha \in \Sigma$ , then  $v$  is  
labeled by  $\alpha$ .

Ex 2

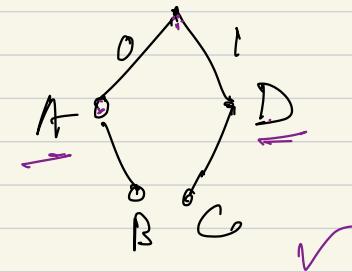
A 0 ✓

B 0 1 ✓

C 1 0

D 1

✓



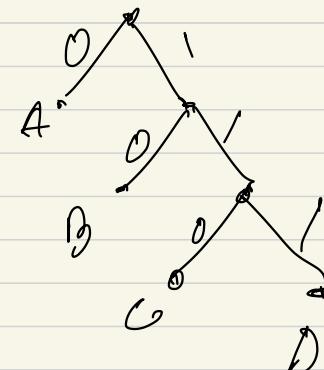
(3)

A 0

B 1 0

C 1 1 0

D 1 1 1



✓

[ Labeled bin trees of this type.]  
vs  
codes.

A Labeled tree from a given code.

- ①  $\ell := \max$  length of  $\text{Enc}(a)$  for any  $a \in \Sigma$
- ② Take a full binary tree of depth  $\ell$ .

③ For every internal vertex, label the edge  
to its left child by 0  
edge to right child by 1.

④ Encoding of each  $a \in \Sigma$  defines a path  
from root to a vertex  $v$ . Label  $v$  by  $\underline{a}$ .

⑤ Delete unlabeled leaves till none remaining.  
(ie all leaves labeled)

⑥ can also delete edge labels. (follows the protocol  $\text{left} \rightarrow 0$   
 $\text{right} \rightarrow 1$ )

A code from a given  $\Sigma$  labelled tree.

- For every  $a \in \Sigma$ , follow the path from the root to the vertex labeled  $a$ .
- The labels of the edges on this path (concatenated) equals  $\text{Enc}(a)$ .

Note: - can also just drop the edge labels.  
- treat left as 0, right as 1

$\underbrace{\text{Properties of Codes}}_{C}$        $\underbrace{\text{as properties of Trees}}_T$

① Length of  $\text{Enc}(a)$  = depth of the vertex labelled  $a$   
 ( # edges from the root to  $a$  )  
 ( on the path )

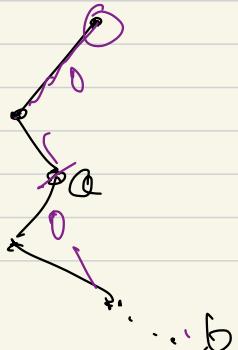
② Given frequencies  $p_a$  for  $a \in \Sigma$ ,  $\left( \sum_{a \in \Sigma} p_a \cdot \text{Enc}(a) \rightarrow \right)$   
 avg length of encoding under the code  $C$   
 $= \sum_{a \in \Sigma} p_a \cdot (\text{depth of vertex labelled } a)$ .

③ Prefix freeness:

$\text{Enc}(a)$  is a prefix of  $\text{Enc}(b)$

iff

vertex with label  $a$  in  $T$  is  
an ancestor of the vertex with label  $b$



Lemma. Code  $C$  is prefix free  
iff

for every  $a \in \Sigma$  the root of  $T$  labeled  $a$   
is a leaf. ( $T$  is prefix - no unlabeled leaves)

Pf.  $\Rightarrow$  Easy



Exercise

15

Restating our question in terms of trees,

Defn:

For alphabet  $\Sigma$ ,

$\Sigma$ -tree : binary tree with leaves labeled  
in 1-1 correspondence with symbols in  $\Sigma$

(can treat edge labels  
implicitly)

Prefix free codes for  $\Sigma \leftrightarrow$   
 $\Sigma$ -trees.

## Average depth of a $\Sigma$ -tree

$T$ :  $\Sigma$ -tree

$\{p_a\}_{a \in \Sigma}$ : frequencies

$L(T, p)$   $\equiv$  average leaf depth in  $T$

$$= \sum_{a \in \Sigma} p_a \cdot (\text{depth of the leaf labeled } a \text{ in } T)$$

$\hookrightarrow$  Avg encoding length of the corresponding code.

Task:

Input:  $\{p_a\}_{a \in \Sigma}$  for some alphabet  $\Sigma$

Output:  $\Sigma$ -free with min  $L(T, p)$ .

Huffman: - efficient algo. for this problem(1951)

- term paper for a course

- corresponding codes are called Huffman codes

## Huffman's algo.

① Start with  $n$ -isolated vertices each labelled by one alphabet symbol in  $\Sigma$ .  $\rightarrow$  view this as a forest

② Iteratively, merge the trees in this forest two at a time based on Huffman's criterion.

③ Stop when one tree remains.

- a) What is the 'merge' operation
- b) What is Huffman's criterion

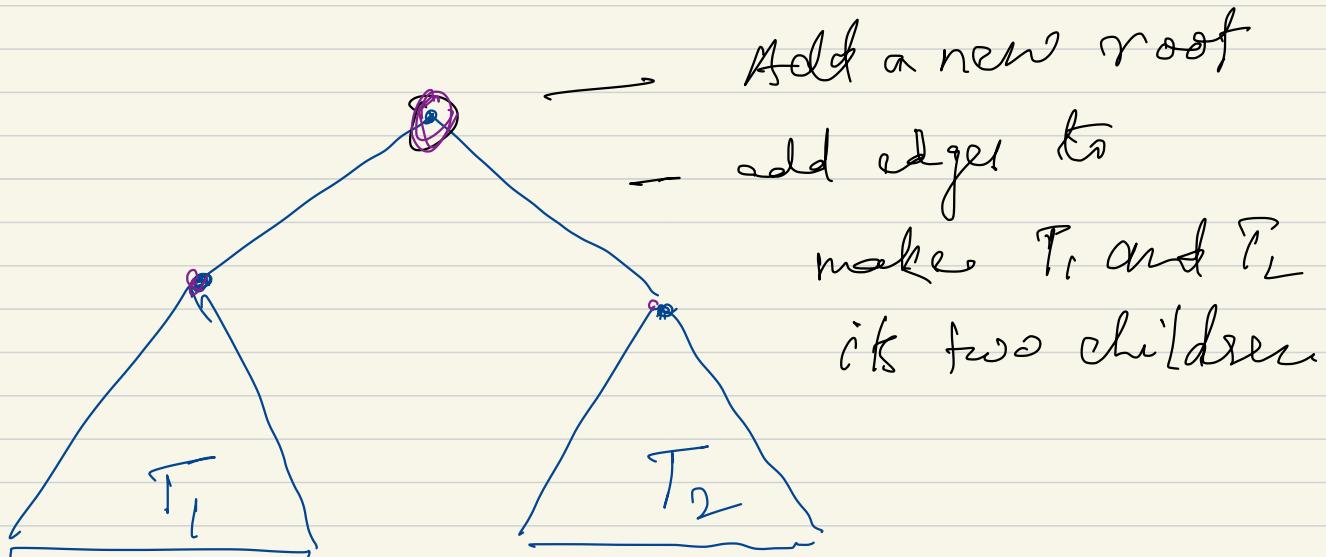
# Trees from Merge operations

- ① Start with  $n$ -vertices, isolated, labelled by symbols in  $\Sigma$

a      b      c      d  
A      B      C      D

- ② Merge operations

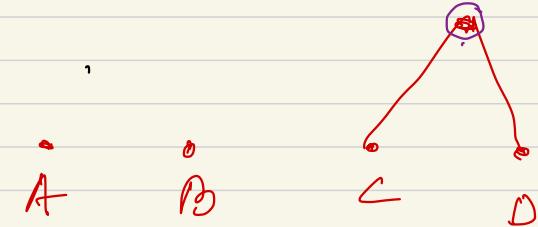
-Two merge any two trees  $T_1, T_2$   
(all trees in this discussion are rooted)



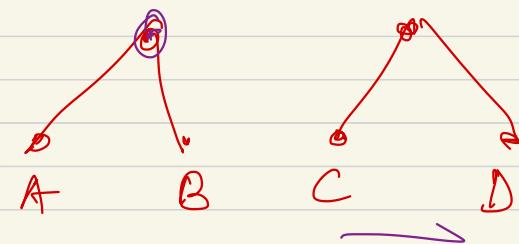
- ③ Stop when a single tree remains

Example:

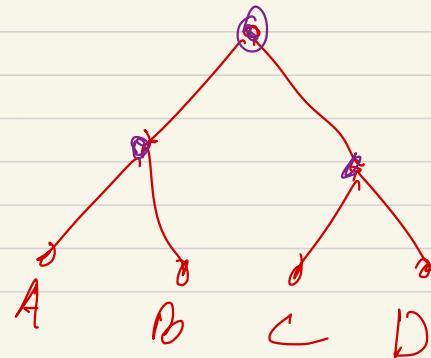
① Merge C, D



② Merge A, B

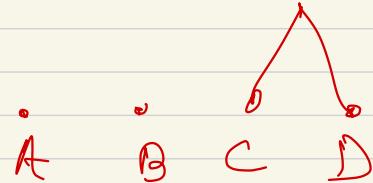


③ Merge AB, CD

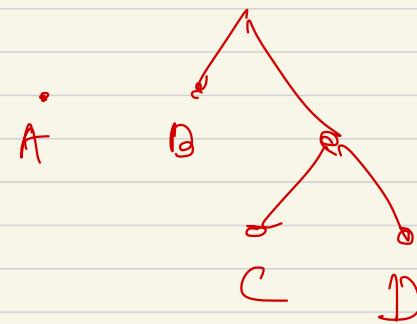


~~2~~

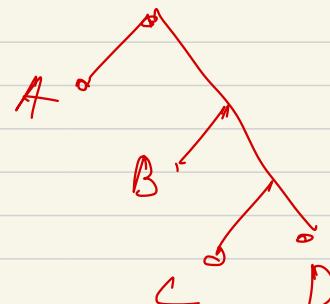
① Merge C, D



② Merge B with CD



③ Merge A with BCD



Huffman's algo.

- ① Start with  $n$ -isolated vertices each labelled by one alphabet symbol in  $\Sigma$ .
- ② Sequence of merges based on some criterion.
- ③ Stop when one tree remains.

Crucial missing detail: what is the merging criterion?

Merging increases the depth of all leaves in the two participating trees.

Huffman's criterion:

Merge the two trees that cause the minimum increase in the average leaf depth  
→ greedy choice at each stage

Increase in average tree depth of the forest.

$$T_1, T_2, T_3, T_4 \leftarrow \text{Tree } \Sigma_1, \Sigma_2, \dots, \Sigma_m \subseteq \Sigma$$

$\Sigma_i$  = set of leaf labels of  $T_i$

Sum of average tree depths

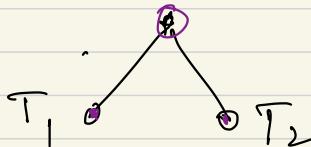
$$= \sum_{i=1}^m \left( \sum_{\alpha \in \Sigma_i} p_\alpha \cdot \frac{\text{depth of leaf } \alpha \text{ in } T_i}{\text{labelled } \alpha \text{ in } T_i} \right)$$



Suppose  $T_1, T_2$  are merged.

So, we have  $T'_1, \underline{T_3}, \underline{T_4}, \dots$

$$T'_1 =$$



Sum of avg free depths

$$= \sum_{i=3}^m \left( \sum_{a \in \Sigma_i} \text{pa. depth of leaf (lab: } a \text{ is } T'_i) \right)$$

$$+ \sum_{a \in \Sigma_1 \cup \Sigma_2} \text{pa. (depth of leaf labelled } a \text{ in } T'_1)$$

$$\begin{aligned}
 &= \textcircled{1} + \sum_{a \in \Sigma_1} p_a \cdot (1 + \text{depth } a \text{ in } T_1) \\
 &\quad + \sum_{a \in \Sigma_2} p_a \cdot (1 + \text{depth } a \text{ in } T_2)
 \end{aligned}$$

( \* \* )

Increase due to merge of  $T_1, T_2$

$$= \sum_{a \in \Sigma_1} p_a + \sum_{a \in \Sigma_2} p_a$$

Intuition:

for every Tree  $T$ , assign a weight  $p(T) = \sum_{a \in \Sigma(T)} p_a$ .

Huffman criterion: Merge trees with the two smallest weights in each step.

Why does this work?

Needs a proof. Not totally trivial. -

## Huffman's Algorithm

① for each  $a \in \Sigma$

$T_a :=$  tree with one node  
with label  $a$

$$p(T_a) := p_a$$

$$F := \{ T_a \}_{a \in \Sigma}$$

② while at least two trees in  $F$

$T_1, T_2 :=$  trees with smallest  $p(T_a)$   
values in  $F$

$\tilde{T} :=$  Merge  $T_1, T_2$

$$p(\tilde{T}) = p(T_1) + p(T_2)$$

$$P \leftarrow F \setminus \{T_1, T_2\} \cup \{T\}$$

③

Return the unique tree

Example:

0.6

A

0.25

B

0.10

C

0.05

D

Running time?

### Correctness

Theorem: For every alphabet  $\Sigma$  and non negative frequencies  $\{p_a\}_{a \in \Sigma}$ , Huffman's algo outputs a prefix free code with minimum possible average encoding length.

Equivalently

Huffman's algorithm outputs a 2-tree  
with min possible average tree depth.

Proof:

## Induction on $|\Sigma|$

Two main ideas: First some notation.

- ① Let  $a, b \in \Sigma$  be symbols s.t  $p_a, p_b$  are the two smallest frequencies.
- ② By Huffman's Algo:  $T_a, T_b$  are merged in the first iteration.
- ③ Reg of the other freq. (as long as they are not smaller than  $p_a, p_b$ ), output of algo

will have a, b as siblings.

Two steps:

- ① Prove that there exist optimal  $\Sigma$ -trees where a, b are siblings.
- ② Prove that the output of Huffman's Algo minimizes avg tree depth for all  $\Sigma$ -trees where a, b are siblings -

Is this sufficient for correctness?

Step 1:

$\text{Tab}_b$  := set of all  $\Sigma$ -trees where  $a, b$  are siblings.

Want to show:

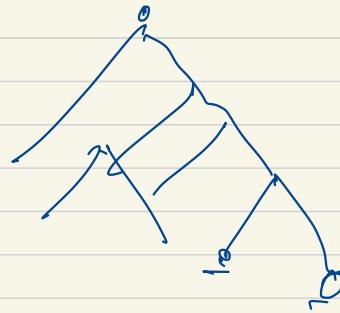
Lemma 1:  $\exists T \in \text{Tab}$  s.t. arg. leaf depth  $L(T, p)$  is  
min over all  $\Sigma$ -trees.

Pf:  $\tilde{T} \leftarrow$  an arbitrary  $\Sigma$ -tree that minimizes  
arg leaf depth.

Want to argue: there is a  $T \in \text{Tab}$  s.t  
 $L(T, p) \leq L(\tilde{T}, p)$ .

W.l.o.g each vertex in  $F$  is either a leaf or has two children. why?

$\Rightarrow \exists$  2 leaves with a common parent at the deepest level.

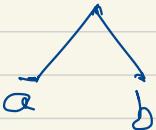


Let  $x, y$  be their labels.

✓

Step 2.

- a, b Merge happens first in Huff. algo

- after flat algo treats  as a single alphabet symbol ab  
with freq  $p_{ab} = p_a + p_b$ .

Formally:

Claim Output of Huff. Algo on  $\Sigma, \{p_\alpha\}_{\alpha \in \Sigma}$   
is the same as its output on  
 $\tilde{\Sigma} = \Sigma \setminus \{ab\} \cup \{ab\}$

$$\{\tilde{p}_\alpha\}_{\alpha \in \Sigma} \sim$$

where,  $\tilde{p}_\alpha = p_\alpha$  for  $\alpha \neq a, b$

$$\tilde{p}_{ab} = p_a + p_b.$$

Proof of claim:

In fact:

$\tilde{\Sigma}$ -trees  $\longleftrightarrow$  Tab  $\xrightarrow{\sim}$   $\begin{cases} \text{I-trees where} \\ a, b \text{ are siblings} \end{cases}$

bijection.  $\beta$

Why?

Lemma: Output of H-algo on  $\tilde{\Sigma}, \beta$  is  $\beta(\tilde{T})$

where  $\tilde{T}'$  is the output of H-algo with input

$\tilde{\Sigma}, \beta'$ .

Further properties:

Lemma: let  $T \rightarrow \Sigma$ -tree in Tab

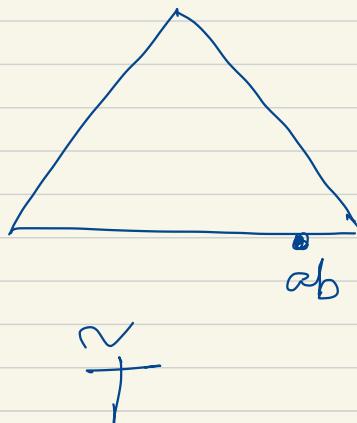
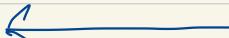
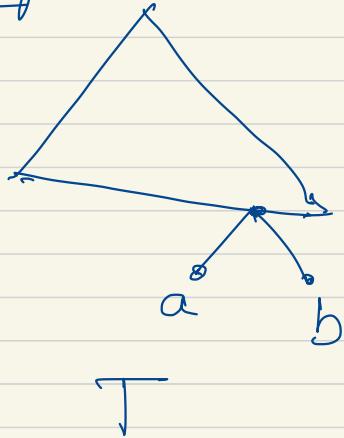
$$\tilde{T} \simeq \beta^{-1}(T).$$

Then,

$$L(T, p) = L(\tilde{T}, \tilde{p}) + (p_a + p_b),$$

- $\overset{1}{\overbrace{\text{Independent of}}}$   
the tree  $T$ .
- just depends on  $a, b$ .

Proof:



Description of  $\beta$ .



Corollary:

$\Sigma^{\checkmark}$ -free  $T^*$  minimizes  $L(\tilde{T}, p)$  over all

$\Sigma^{\tilde{\wedge}}$ -free iff

$T = P(T^*)$  minimizes  $L(T, p)$  over all

$\Sigma$ -trees in Tab.

Are we done?



