

COL 351 : ANALYSIS & DESIGN OF ALGORITHMS

## LECTURE 18

# GREEDY ALGORITHMS II : HUFFMAN'S ALGORITHM

SEPT 04, 2024

|

ROHIT VAISH

# OPTIMAL PREFIX-FREE CODES

# OPTIMAL PREFIX-FREE CODES

input : a non negative frequency  $p_i$  for each symbol  $i$  of  
alphabet  $\Sigma$  of size  $|\Sigma| = n \geq 2$

# OPTIMAL PREFIX-FREE CODES

**input:** a non negative frequency  $p_i$  for each symbol  $i$  of alphabet  $\Sigma$  of size  $|\Sigma| = n \geq 2$

**output:** A prefix-free binary code with minimum possible average encoding length

$$\sum_{i \in \Sigma} p_i \cdot \text{number of bits used to encode } i$$

# BINARY CODES $\longleftrightarrow$ LABELED BINARY TREES

$$\Sigma = \{A, B, C, D\}$$

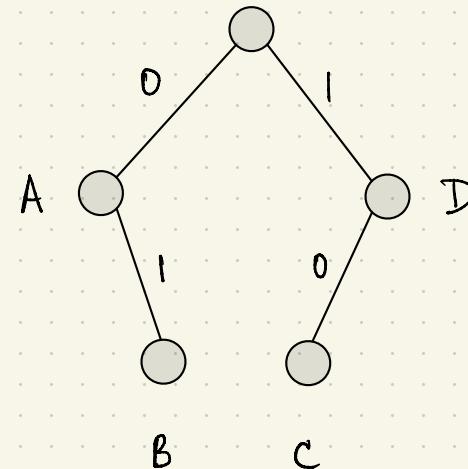
Symbol      variable length

A            0

B            01

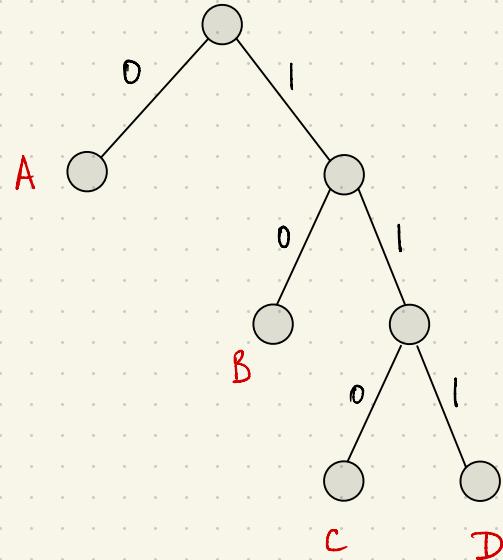
C            10

D            1



# PREFIX-FREE CODES AS TREES

\* prefix-free code  $\Leftrightarrow$  only leaves are labeled

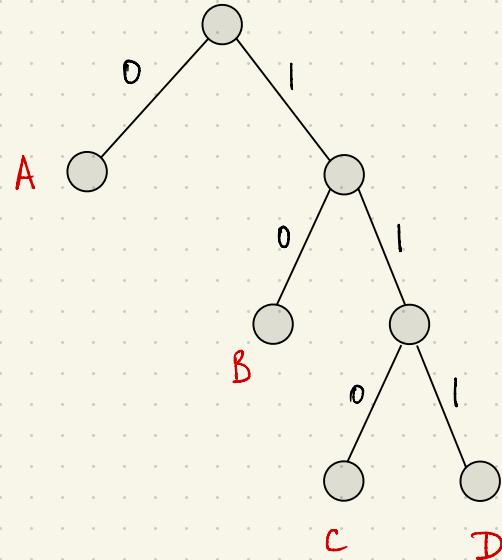


# PREFIX-FREE CODES AS TREES

\* prefix-free code  $\Leftrightarrow$  only leaves are labeled

\* Encoding length of symbol  $i \in \Sigma$

= depth of leaf labeled  $i$  in the tree



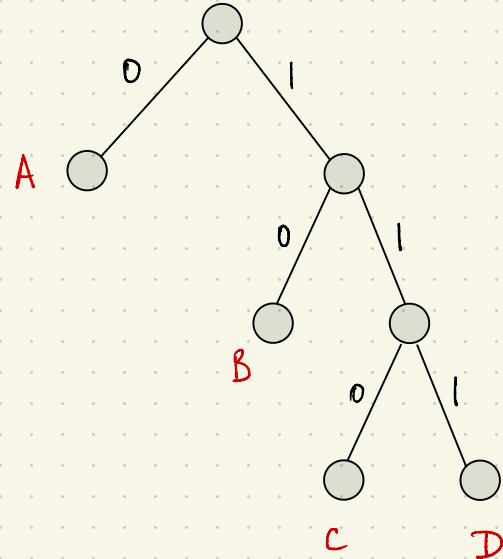
# PREFIX-FREE CODES AS TREES

\* prefix-free code  $\Leftrightarrow$  only leaves are labeled

\* Encoding length of symbol  $i \in \Sigma$

= depth of leaf labeled  $i$  in the tree

\*  $\Sigma$  tree: a binary tree with leaves labeled in one-to-one correspondence with symbols of  $\Sigma$ .



# PREFIX-FREE CODES AS TREES

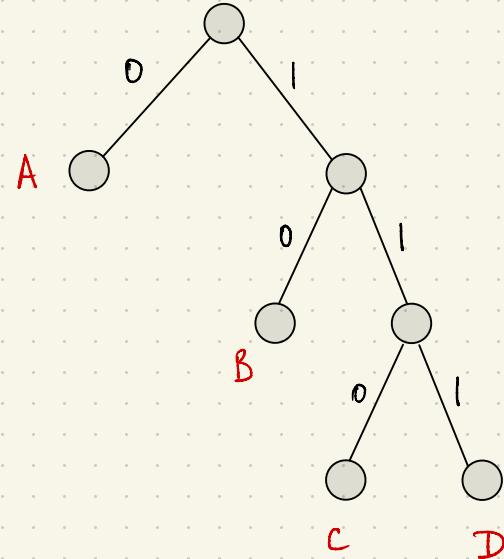
\* prefix-free code  $\Leftrightarrow$  only leaves are labeled

\* Encoding length of symbol  $i \in \Sigma$

= depth of leaf labeled  $i$  in the tree

\*  $\Sigma$  tree: a binary tree with leaves labeled in one-to-one correspondence with symbols of  $\Sigma$ .

\* For any  $\Sigma$ , some  $\Sigma$ -tree always exists  
e.g.,  $\{0, 10, 110, 1110, \dots\}$



# OPTIMAL PREFIX-FREE CODES

**input :** a non negative frequency  $p_i$  for each symbol  $i$  of alphabet  $\Sigma$  of size  $|\Sigma| = n \geq 2$

**output :** A prefix-free binary code with minimum possible average encoding length

$$\sum_{i \in \Sigma} p_i \cdot \text{number of bits used to encode } i$$

# OPTIMAL PREFIX-FREE CODES

**input:** a non negative frequency  $p_i$  for each symbol  $i$  of alphabet  $\Sigma$  of size  $|\Sigma| = n \geq 2$

**output:** A  $\Sigma$ -tree  $T$  with minimum possible average leaf depth

$$L(T) := \sum_{i \in \Sigma} p_i \cdot \text{depth of leaf labeled } i \text{ in } T$$

# BUILDING A TREE

What's a principled approach for building a  $\Sigma$ -tree?



Top down

(divide and conquer)

Bottom up

(greedy)

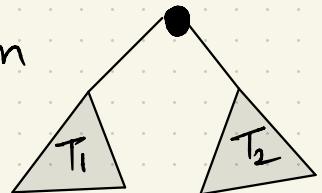
# TOP DOWN APPROACH

[Shannon - Fano encoding]

# TOP DOWN APPROACH

[Shannon - Fano encoding]

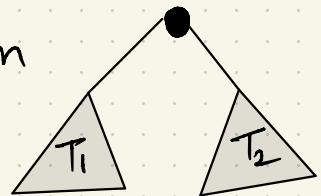
- \* Divide and conquer
- \* partition  $\Sigma$  into  $\Sigma_1, \Sigma_2$  each with  $\approx 50\%$  of total frequency
- \* recursively compute  $T_1$  for  $\Sigma_1$ ,  $T_2$  for  $\Sigma_2$ , return



# TOP DOWN APPROACH

[Shannon - Fano encoding]

- \* Divide and conquer
- \* partition  $\Sigma$  into  $\Sigma_1, \Sigma_2$  each with  $\approx 50\%$  of total frequency
- \* recursively compute  $T_1$  for  $\Sigma_1$ ,  $T_2$  for  $\Sigma_2$ , return
- \* Suboptimal 😞 [Tutorial sheet 6]



HUFFMAN'S GREEDY ALGORITHM

## HUFFMAN'S GREEDY ALGORITHM



leveraging partial knowledge of optimal solution

Recall : prefix-free code  $\longleftrightarrow$  labeled binary tree  
(only at leaves)

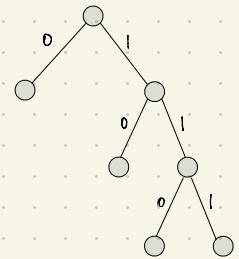
Recall : prefix-free code  $\longleftrightarrow$  labeled binary tree  
(only at leaves)

an optimal  $\Sigma$  tree  $T^*$

Want

a labeling of  $T^*$

Recall : prefix-free code  $\longleftrightarrow$  labeled binary tree  
(only at leaves)



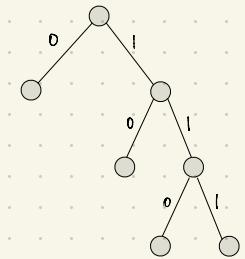
an optimal  $\Sigma$  tree  $T^*$

Want

if we knew this

a labeling of  $T^*$

Recall : prefix-free code  $\longleftrightarrow$  labeled binary tree  
(only at leaves)



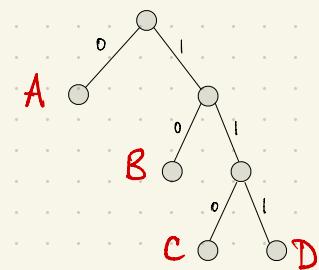
an optimal  $\Sigma$  tree  $T^*$

Want

if we knew this

Can we determine this?

a labeling of  $T^*$



## LABELING AN OPTIMAL SOLUTION

**Lemma:** Let  $u$  and  $v$  be leaves of  $T^*$  such that  $\text{depth}(u) < \text{depth}(v)$ .

Further, suppose in the labeling of  $T^*$  corresponding to an optimal prefix code,  $u$  is labeled  $a \in \Sigma$  and  $v$  is labeled  $b \in \Sigma$ . Then,  $p_a \geq p_b$ .

## LABELING AN OPTIMAL SOLUTION

**Lemma:** Let  $u$  and  $v$  be leaves of  $T^*$  such that  $\text{depth}(u) < \text{depth}(v)$ .

Further, suppose in the labeling of  $T^*$  corresponding to an optimal prefix code,  $u$  is labeled  $a \in \Sigma$  and  $v$  is labeled  $b \in \Sigma$ . Then,  $p_a \geq p_b$ .

**Proof :** Exercise

**Hint :** Apply exchange argument to  $L(T^*) := \sum_{i \in \Sigma} p_i \cdot \text{depth of leaf labeled } i \text{ in } T^*$ .

## LABELING AN OPTIMAL SOLUTION

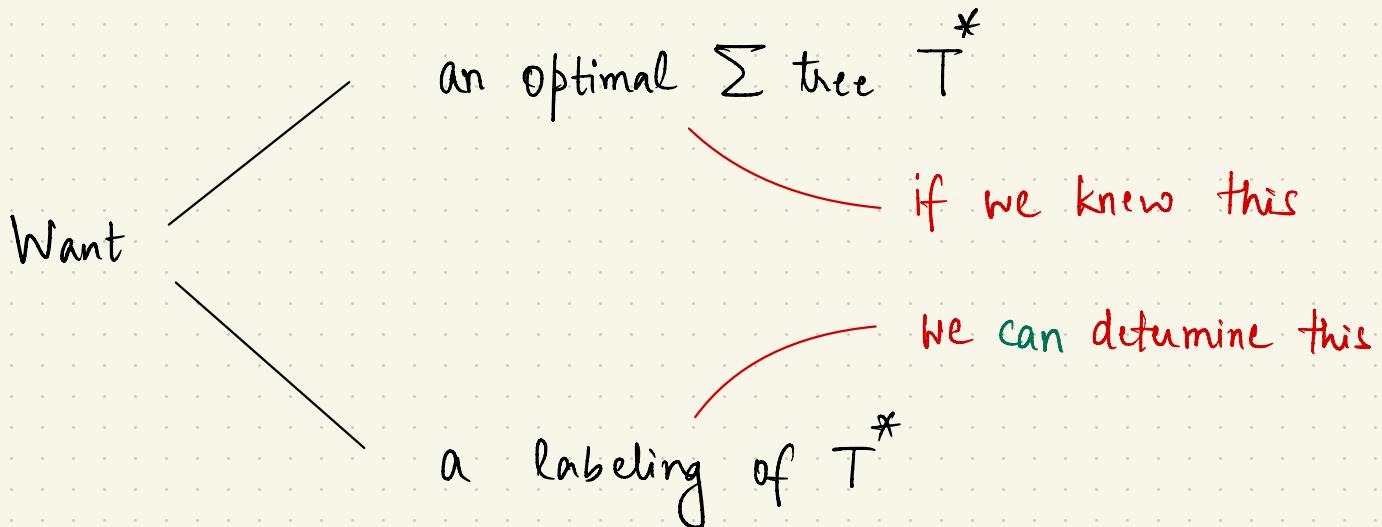
**Lemma:** Let  $u$  and  $v$  be leaves of  $T^*$  such that  $\text{depth}(u) < \text{depth}(v)$ .

Further, suppose in the labeling of  $T^*$  corresponding to an optimal prefix code,  $u$  is labeled  $a \in \Sigma$  and  $v$  is labeled  $b \in \Sigma$ . Then,  $p_a \geq p_b$ .

Labeling algorithm :

- Label depth 1 leaves (if any) with the highest-frequency symbols
- Label depth 2 leaves (if any) with the next-highest-frequency symbols
- and so on ..

Recall : prefix-free code  $\longleftrightarrow$  labeled binary tree  
(only at leaves)



Recall : prefix-free code  $\longleftrightarrow$  labeled binary tree  
(only at leaves)

Want

The diagram shows two black lines originating from the word "Want" on the left. One line points upwards and to the right towards the text "an optimal  $\Sigma$  tree  $T^*$ ". The other line points downwards and to the right towards the text "a labeling of  $T^*$ ". A red curved arrow originates from the text "how to determine this?" and points towards the text "a labeling of  $T^*$ ".

an optimal  $\Sigma$  tree  $T^*$

how to determine this?

a labeling of  $T^*$

Recall : prefix-free code  $\longleftrightarrow$  labeled binary tree  
(only at leaves)

Want

an optimal  $\Sigma$  tree  $T^*$

how to determine this?

brute force : exponentially-many  
binary trees

a labeling of  $T^*$



Recall labeling algorithm (given  $T^*$ ):

- Label depth 1 leaves (if any) with the highest-frequency symbols
- Label depth 2 leaves (if any) with the next-highest-frequency symbols
- And so on . . .

Recall labeling algorithm (given  $T^*$ ) :

- Label depth 1 leaves (if any) with the highest-frequency symbols
- Label depth 2 leaves (if any) with the next-highest-frequency symbols
- And so on . . .



Lowest frequency symbols are assigned to the deepest leaves !

Recall labeling algorithm (given  $T^*$ ) :

- Label depth 1 leaves (if any) with the highest-frequency symbols
- Label depth 2 leaves (if any) with the next-highest-frequency symbols
- And so on . . .



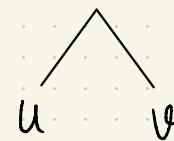
Lowest frequency symbols are assigned to the deepest leaves!

Recall labeling algorithm (given  $T^*$ ):

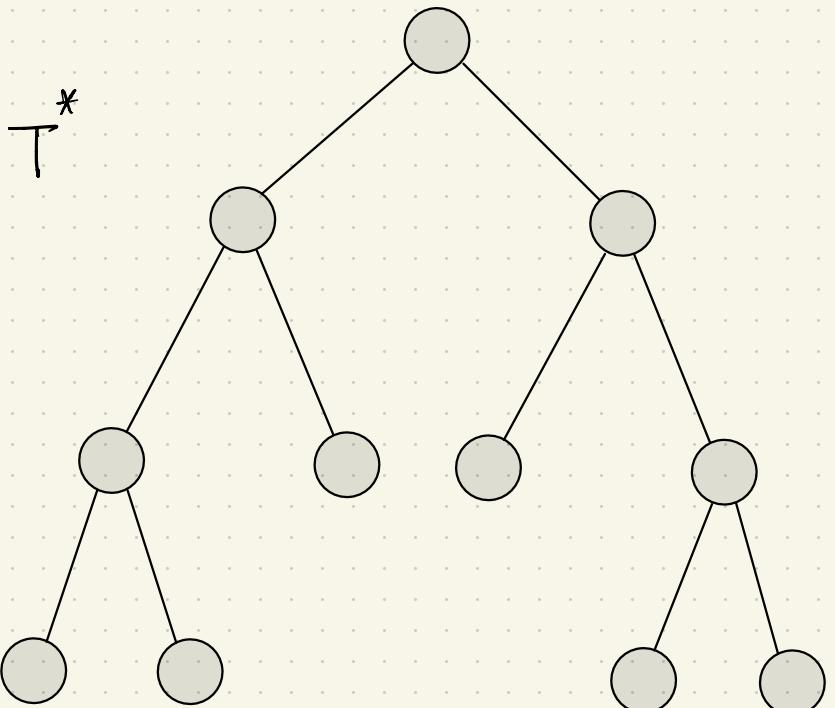
- Label depth 1 leaves (if any) with the highest-frequency symbols
- Label depth 2 leaves (if any) with the next-highest-frequency symbols
- And so on . . .

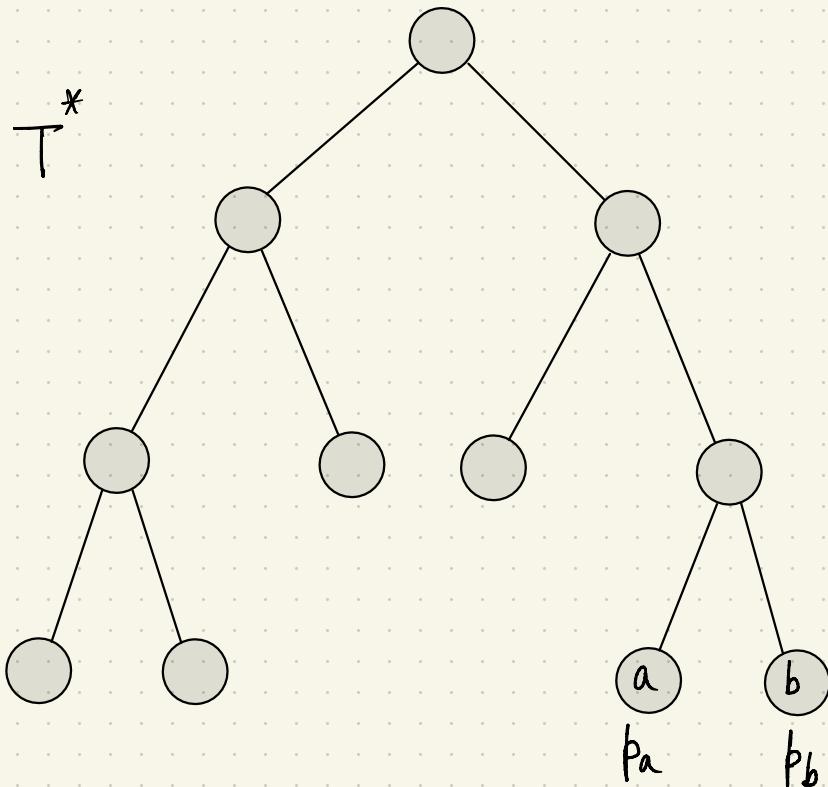
 Lowest frequency symbols are assigned to the deepest leaves!

Check: If  $u$  is a leaf at maximum depth, then  $u$ 's sibling is also a leaf at maximum depth.



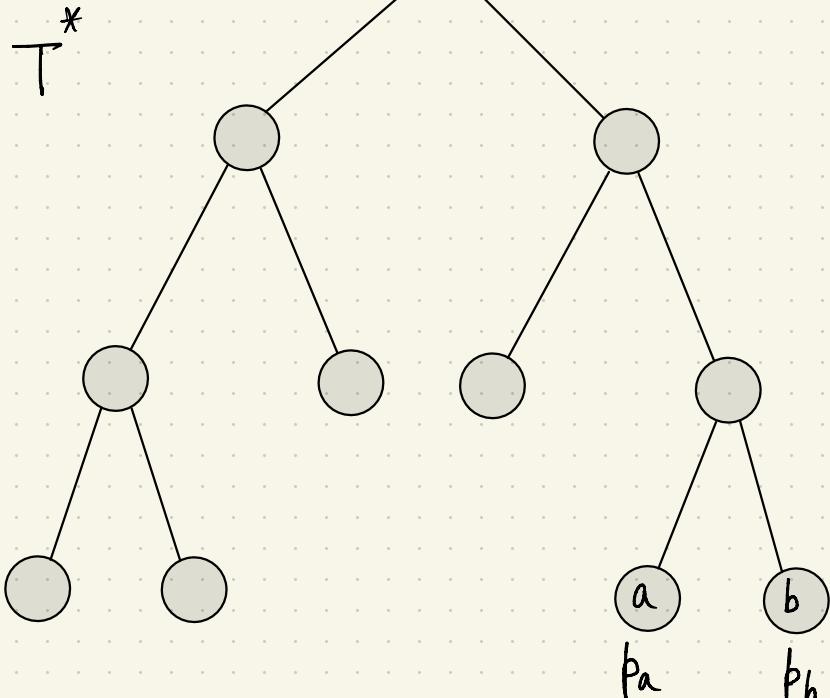
**Corollary:** There is an optimal prefix code for input  $(\Sigma, p)$  with corresponding  $\Sigma$ -tree  $T^*$  in which the two lowest-frequency symbols are assigned to sibling leaves at maximum depth.



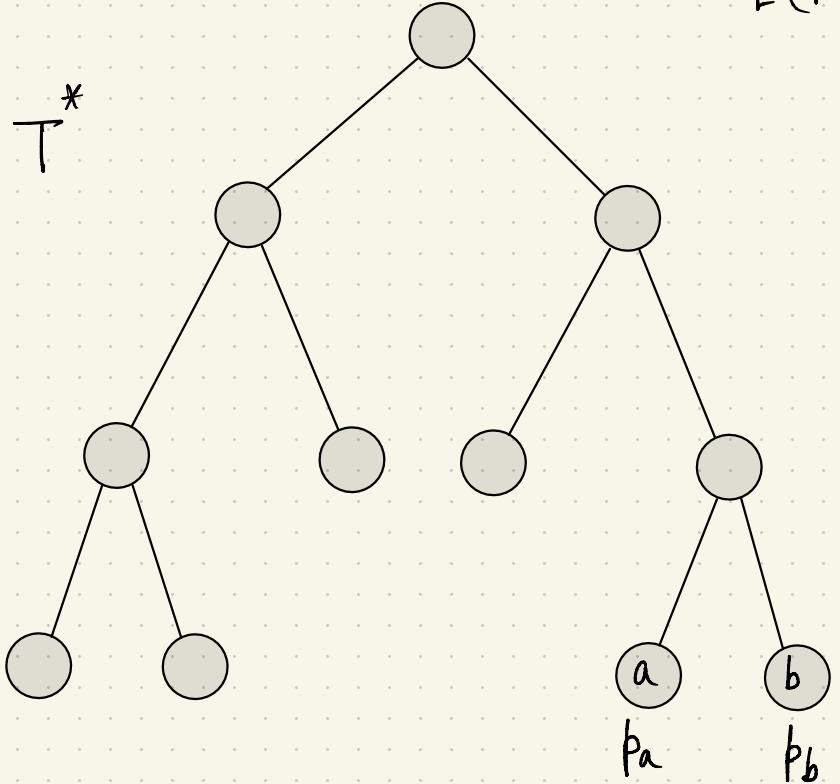


Two lowest-frequency symbols  
*(safely "locked in")*

$$L(T^*) = \sum_{i \in \Sigma} p_i \cdot \text{depth of leaf labeled } i$$

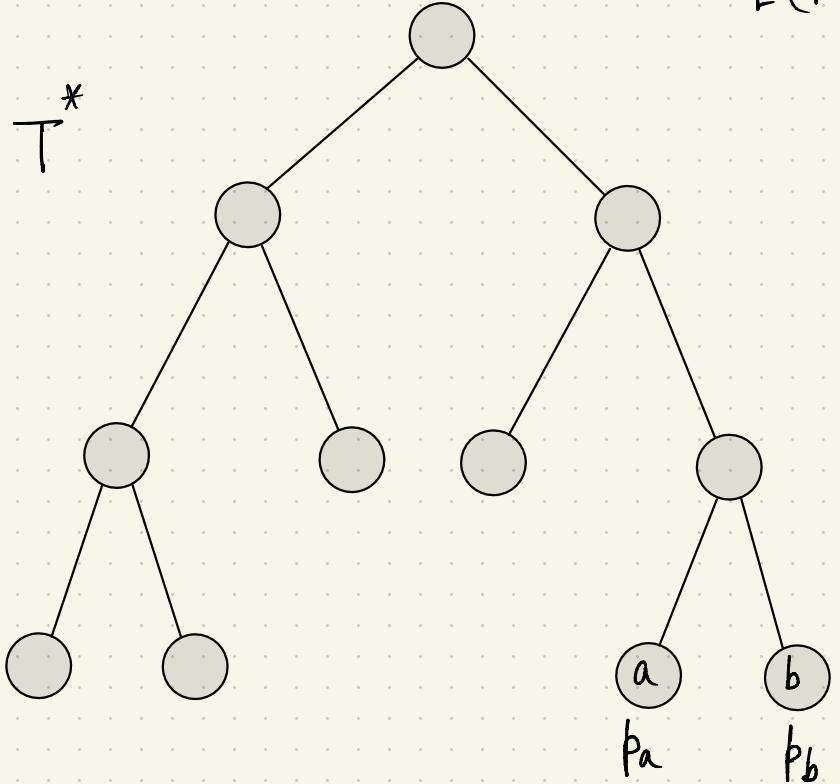


Two lowest-frequency symbols  
*(safely "locked in")*



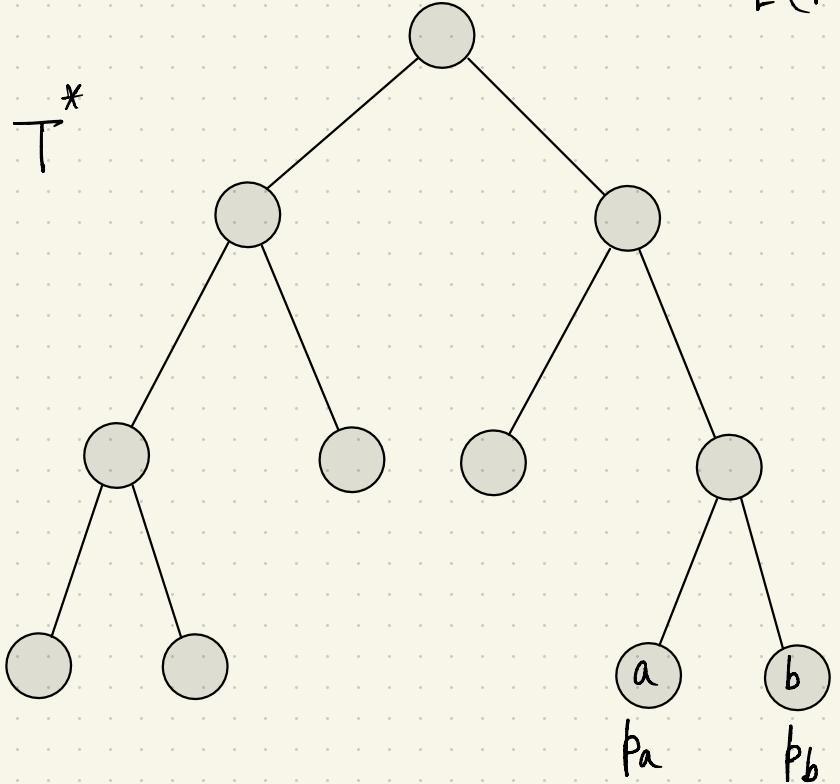
$$\begin{aligned}
 L(T^*) &= \sum_{i \in \Sigma} p_i \cdot \text{depth of leaf labeled } i \\
 &= p_a \cdot \text{depth}(a) + p_b \cdot \text{depth}(b) \\
 &\quad + \text{rest}
 \end{aligned}$$

Two lowest-frequency symbols  
 (safely "locked in")



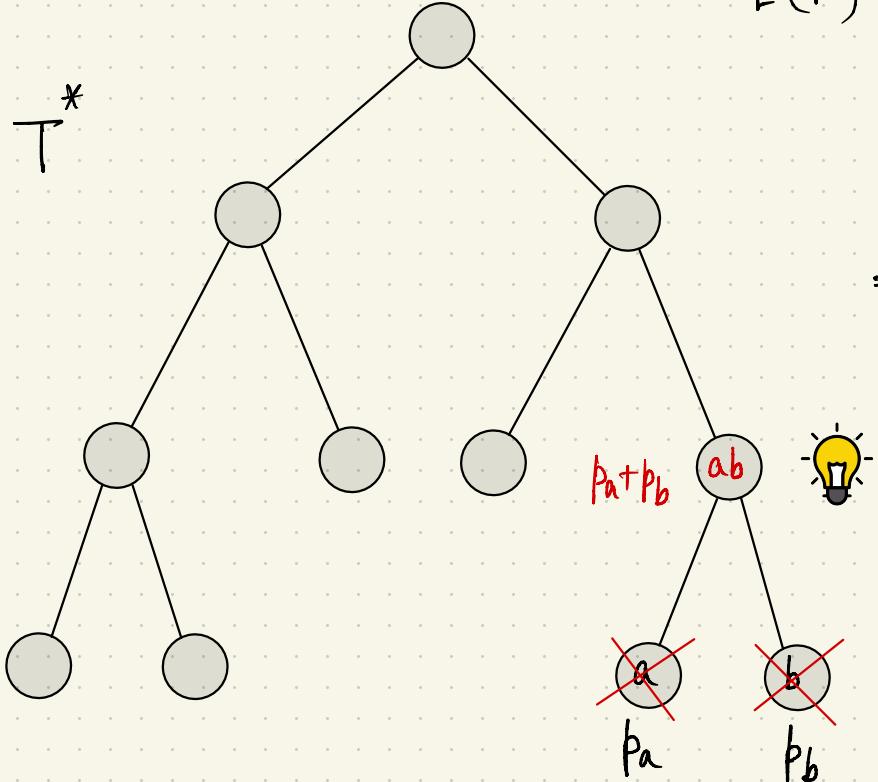
$$\begin{aligned}
 L(T^*) &= \sum_{i \in \Sigma} p_i \cdot \text{depth of leaf labeled } i \\
 &= p_a \cdot \text{depth}(a) + p_b \cdot \text{depth}(b) \\
 &\quad + \text{rest}
 \end{aligned}$$

Two lowest-frequency symbols  
 (safely "locked in")



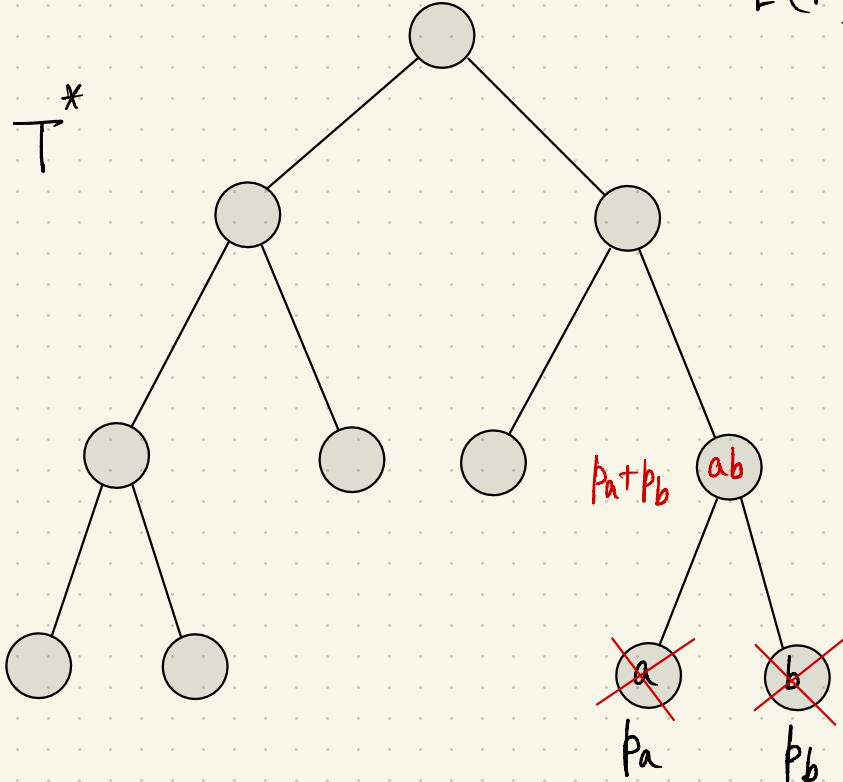
$$\begin{aligned}
 L(T^*) &= \sum_{i \in \Sigma} p_i \cdot \text{depth of leaf labeled } i \\
 &= p_a \cdot \text{depth}(a) + p_b \cdot \text{depth}(b) \\
 &\quad + \text{rest} \\
 &= (p_a + p_b) \cdot \text{depth}(a) + \text{rest}
 \end{aligned}$$

Two lowest-frequency symbols  
(safely "locked in")



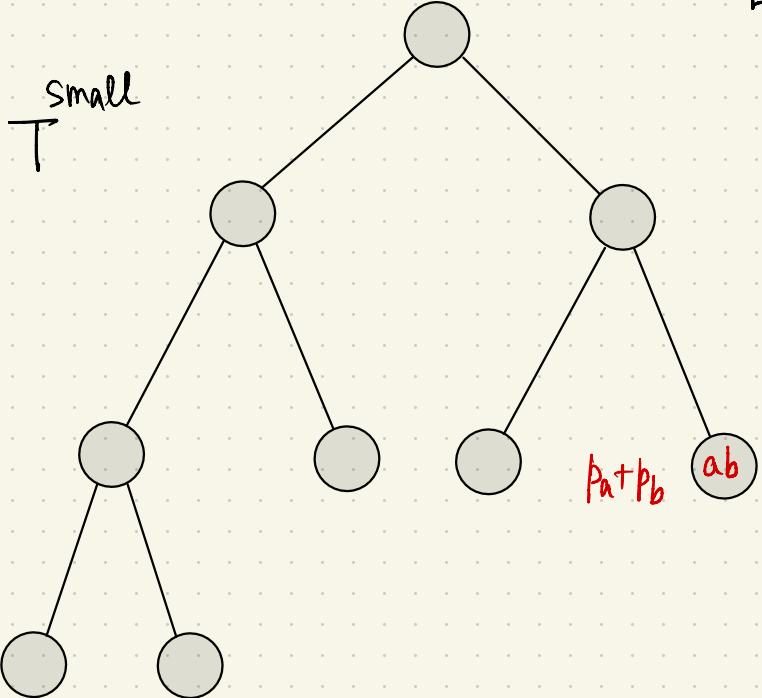
$$\begin{aligned}
 L(T^*) &= \sum_{i \in \Sigma} p_i \cdot \text{depth of leaf labeled } i \\
 &= p_a \cdot \text{depth}(a) + p_b \cdot \text{depth}(b) \\
 &\quad + \text{rest} \\
 &= (p_a + p_b) \cdot \text{depth}(a) + \text{rest}
 \end{aligned}$$

Two lowest-frequency symbols  
(safely "locked in")

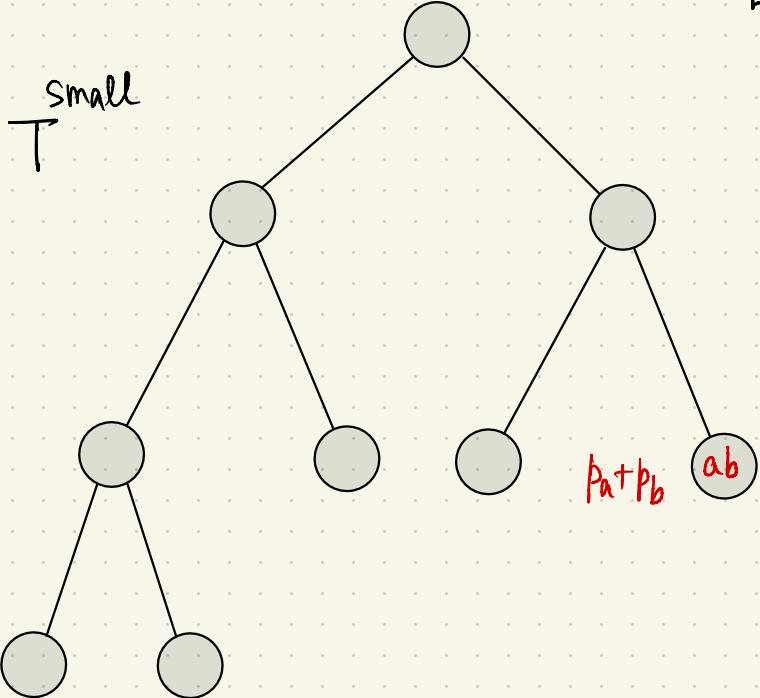


$$\begin{aligned}
 L(T^*) &= \sum_{i \in \Sigma} p_i \cdot \text{depth of leaf labeled } i \\
 &= p_a \cdot \text{depth}(a) + p_b \cdot \text{depth}(b) \\
 &\quad + \text{rest} \\
 &= (p_a + p_b) \cdot \text{depth}(a) + \text{rest} \\
 &= p_{ab} \cdot \text{depth}(ab) + p_a + p_b + \text{rest}
 \end{aligned}$$

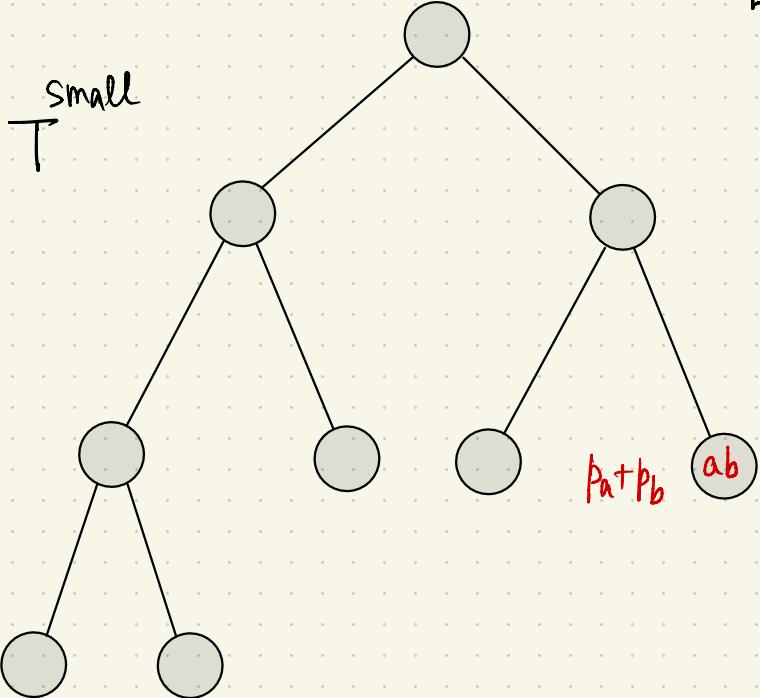
Two lowest-frequency symbols  
(safely "locked in")



$$\begin{aligned}
 L(T^*) &= \sum_{i \in \Sigma} p_i \cdot \text{depth of leaf labeled } i \\
 &= p_a \cdot \text{depth}(a) + p_b \cdot \text{depth}(b) \\
 &\quad + \text{rest} \\
 &= (p_a + p_b) \cdot \text{depth}(a) + \text{rest} \\
 &= p_{ab} \cdot \text{depth}(ab) + p_a + p_b + \text{rest} \\
 &= L(T) + p_a + p_b
 \end{aligned}$$



$$\begin{aligned}
 L(T^*) &= \sum_{i \in \Sigma} p_i \cdot \text{depth of leaf labeled } i \\
 &= p_a \cdot \text{depth}(a) + p_b \cdot \text{depth}(b) \\
 &\quad + \text{rest} \\
 &= (p_a + p_b) \cdot \text{depth}(a) + \text{rest} \\
 &= p_{ab} \cdot \text{depth}(ab) + p_a + p_b + \text{rest} \\
 &= L(T^{\text{small}}) + \underbrace{p_a + p_b}_{\text{fixed}}
 \end{aligned}$$



$$\begin{aligned}
 L(T^*) &= \sum_{i \in \Sigma} p_i \cdot \text{depth of leaf labeled } i \\
 &= p_a \cdot \text{depth}(a) + p_b \cdot \text{depth}(b) \\
 &\quad + \text{rest} \\
 &= (p_a + p_b) \cdot \text{depth}(a) + \text{rest} \\
 &= p_{ab} \cdot \text{depth}(ab) + p_a + p_b + \text{rest} \\
 &= L(T^{\text{small}}) + p_a + p_b
 \end{aligned}$$

*fixed*

reurse!

# HUFFMAN'S ALGORITHM

input: alphabet  $\Sigma$ , frequencies  $\{p_a\}_{a \in \Sigma}$

output:  $\Sigma$ -tree with minimum average leaf depth

## HUFFMAN'S ALGORITHM

input: alphabet  $\Sigma$ , frequencies  $\{p_a\}_{a \in \Sigma}$

output:  $\Sigma$ -tree with minimum average leaf depth

if  $|\Sigma| = 2$

    └ encode one letter using 0 and other using 1

## HUFFMAN'S ALGORITHM

input: alphabet  $\Sigma$ , frequencies  $\{P_a\}_{a \in \Sigma}$

output:  $\Sigma$ -tree with minimum average leaf depth

if  $|\Sigma| = 2$

    └ encode one letter using 0 and other using 1

else

$x, y \leftarrow$  two lowest-frequency letters

## HUFFMAN'S ALGORITHM

input: alphabet  $\Sigma$ , frequencies  $\{p_a\}_{a \in \Sigma}$

output:  $\Sigma$ -tree with minimum average leaf depth

if  $|\Sigma| = 2$

    └ encode one letter using 0 and other using 1

else

$x, y \leftarrow$  two lowest-frequency letters

    Create new alphabet  $\Sigma'$  by merging  $x$  and  $y$  into  $z$  with  $p_z := p_x + p_y$ .

## HUFFMAN'S ALGORITHM

input: alphabet  $\Sigma$ , frequencies  $\{p_a\}_{a \in \Sigma}$

output:  $\Sigma$ -tree with minimum average leaf depth

if  $|\Sigma| = 2$

    [ encode one letter using 0 and other using 1

else

$x, y \leftarrow$  two lowest-frequency letters

    Create new alphabet  $\Sigma'$  by merging  $x$  and  $y$  into  $z$  with  $p_z := p_x + p_y$ .

    Recursively construct tree  $T'$  for  $(\Sigma', p)$

## HUFFMAN'S ALGORITHM

input: alphabet  $\Sigma$ , frequencies  $\{p_a\}_{a \in \Sigma}$

output:  $\Sigma$ -tree with minimum average leaf depth

if  $|\Sigma| = 2$

    [ encode one letter using 0 and other using 1

else

$x, y \leftarrow$  two lowest-frequency letters

    Create new alphabet  $\Sigma'$  by merging  $x$  and  $y$  into  $z$  with  $p_z := p_x + p_y$ .

    Recursively construct tree  $T'$  for  $(\Sigma', p)$

return  $T'$  after replacing leaf labeled  $z$  with unlabeled node  
and two children labeled  $x$  and  $y$ .

## HUFFMAN'S ALGORITHM

input: alphabet  $\Sigma$ , frequencies  $\{p_a\}_{a \in \Sigma}$

output:  $\Sigma$ -tree with minimum average leaf depth

if  $|\Sigma| = 2$

    [ encode one letter using 0 and other using 1

else

$x, y \leftarrow$  two lowest-frequency letters

    Create new alphabet  $\Sigma'$  by merging  $x$  and  $y$  into  $z$  with  $p_z := p_x + p_y$ .

    Recursively construct tree  $T'$  for  $(\Sigma', p)$

return  $T'$  after replacing leaf labeled  $z$  with unlabeled node  
and two children labeled  $x$  and  $y$ .

# HUFFMAN'S ALGORITHM

(A)

60 %.

(B)

25 %.

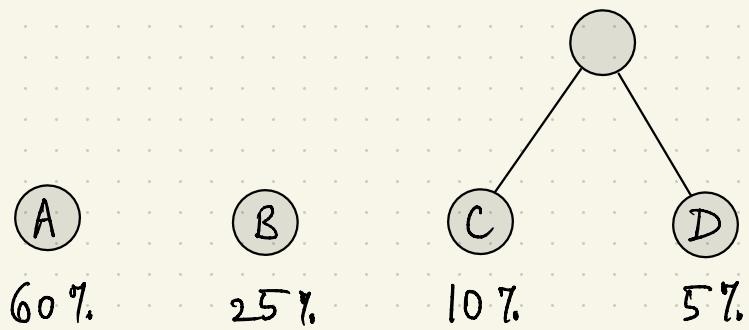
(C)

10 %.

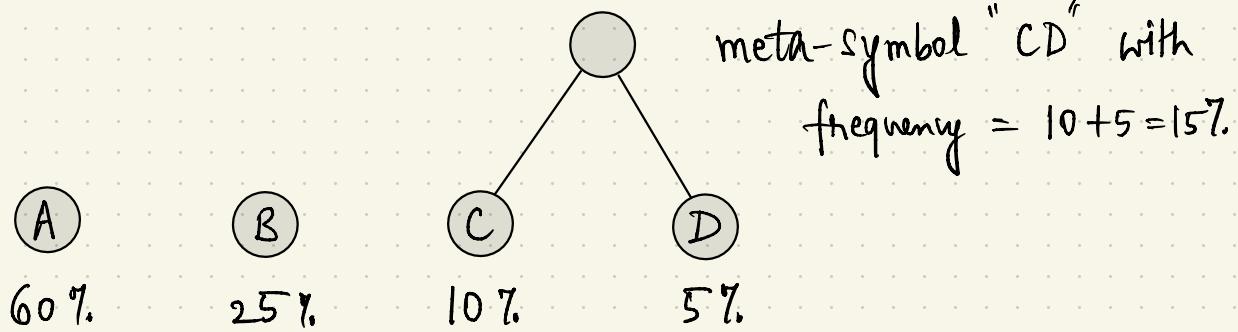
(D)

5 %.

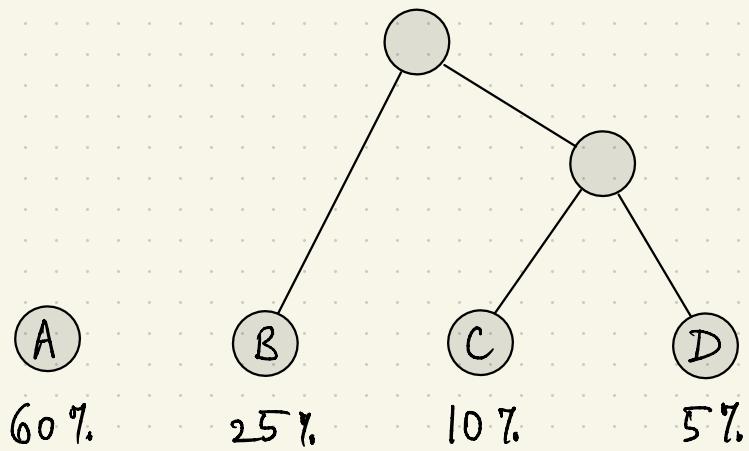
# HUFFMAN'S ALGORITHM



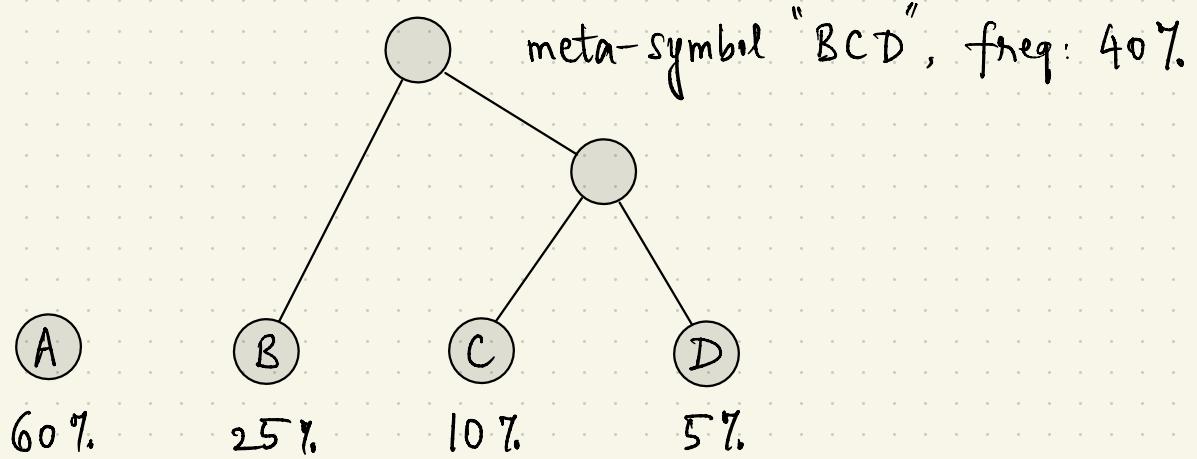
# HUFFMAN'S ALGORITHM



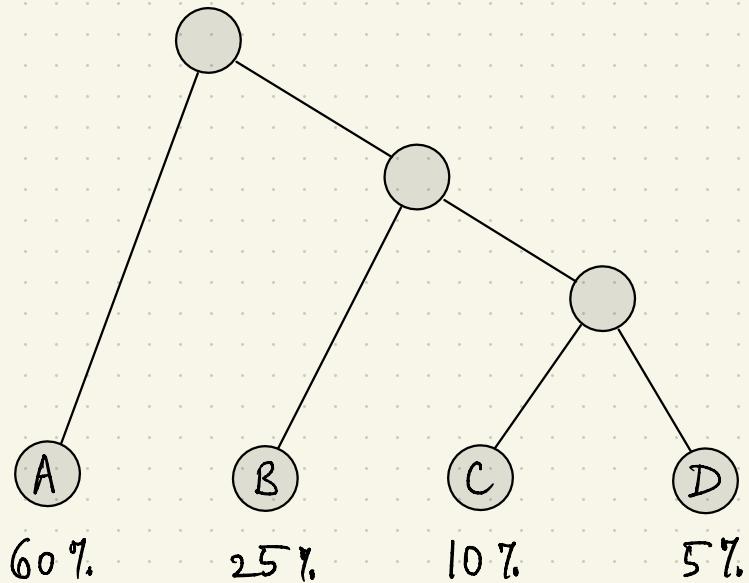
# HUFFMAN'S ALGORITHM



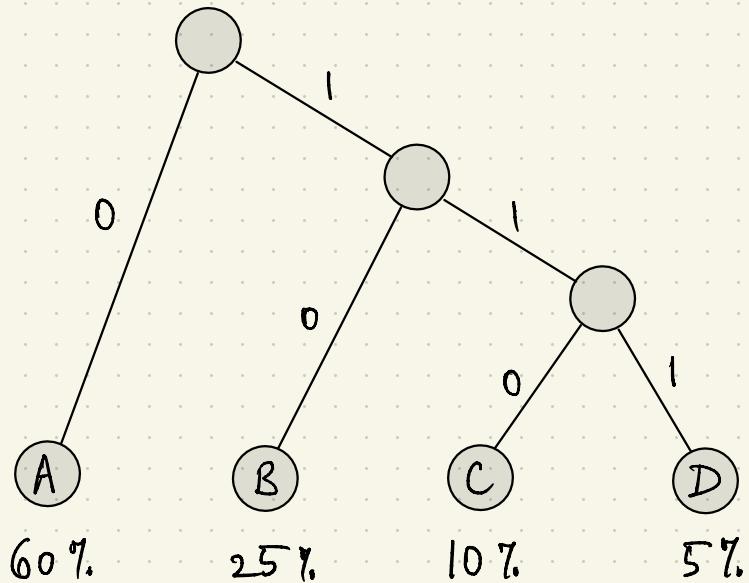
# HUFFMAN'S ALGORITHM



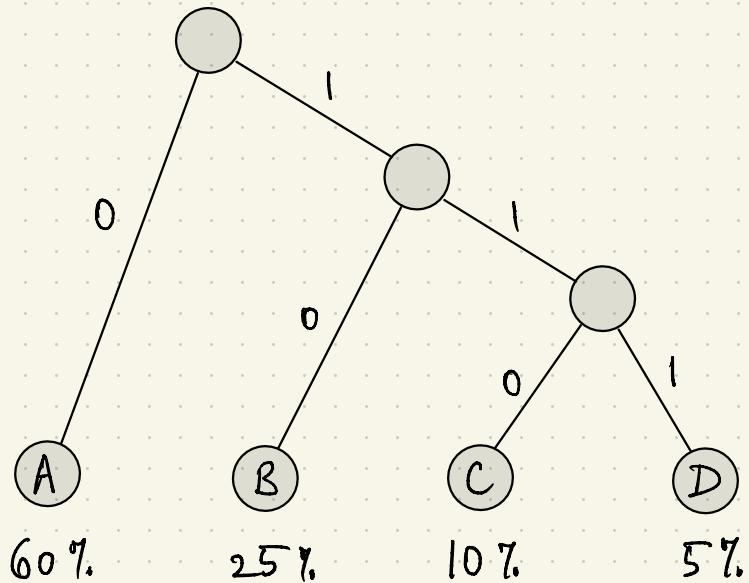
# HUFFMAN'S ALGORITHM



# HUFFMAN'S ALGORITHM



# HUFFMAN'S ALGORITHM



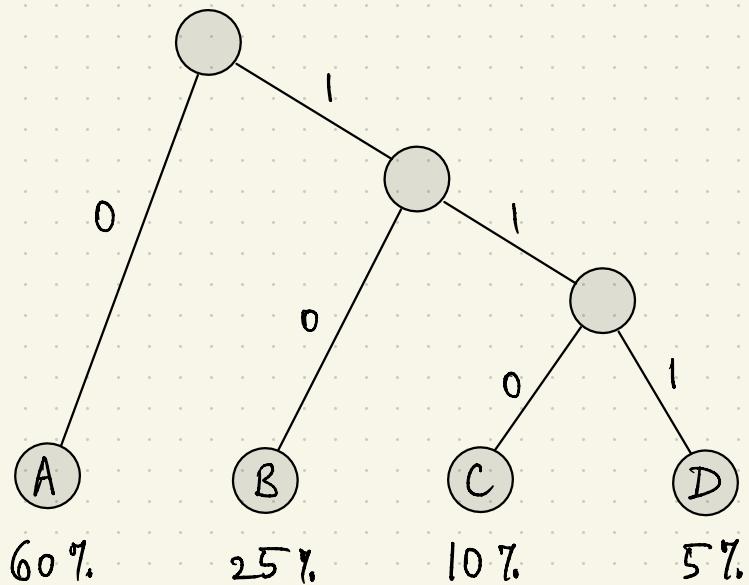
A : 0

B : 01

C : 110

D : 111

# HUFFMAN'S ALGORITHM



Huffman code

A : 0

B : 01

C : 110

D : 111

# HUFFMAN'S ALGORITHM

(A)

70 %.

(B)

8 %.

(C)

10 %.

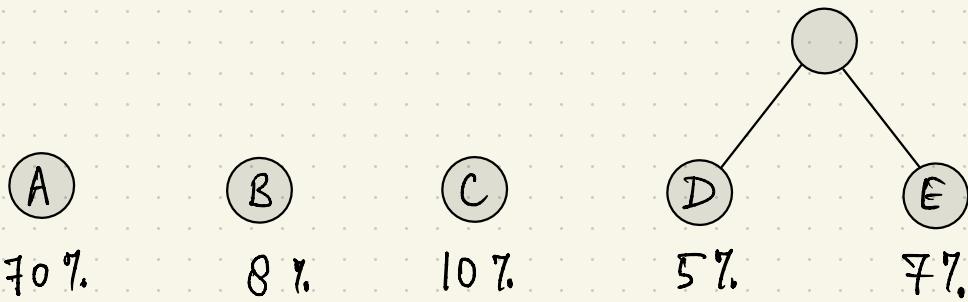
(D)

5 %.

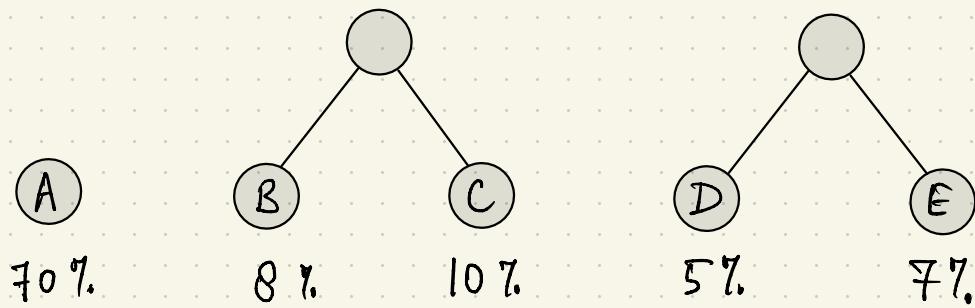
(E)

7 %.

# HUFFMAN'S ALGORITHM

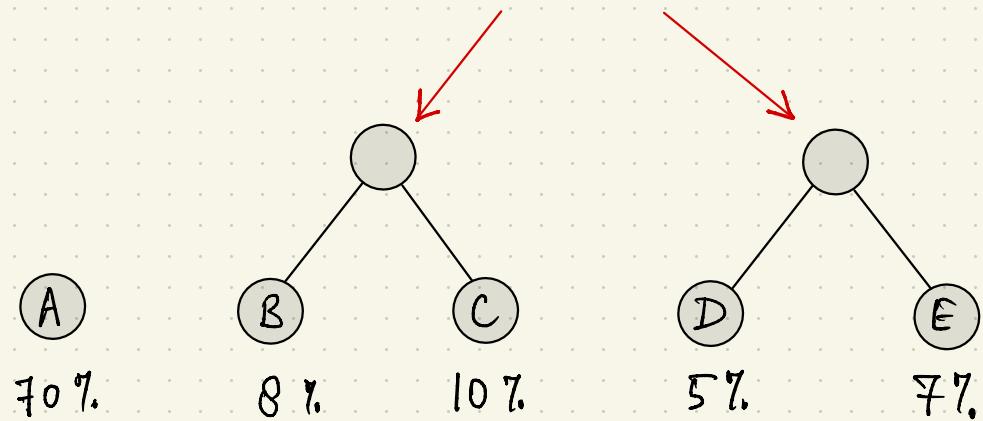


# HUFFMAN'S ALGORITHM

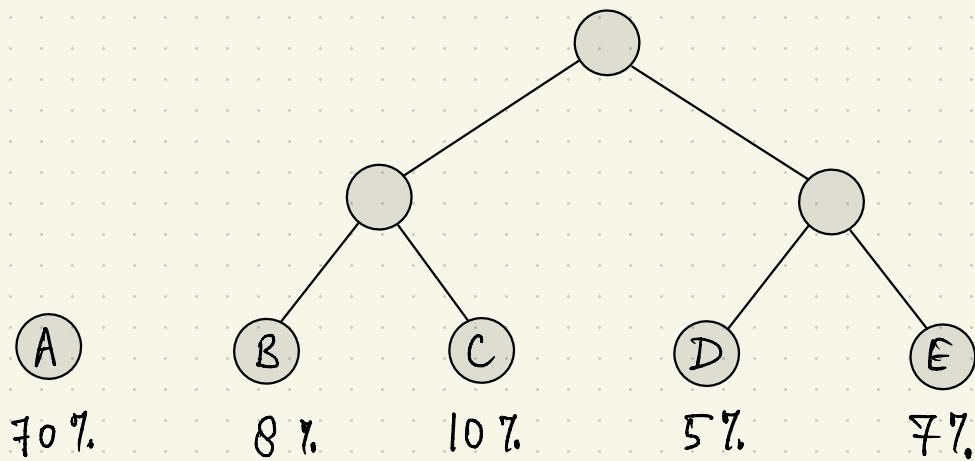


# HUFFMAN'S ALGORITHM

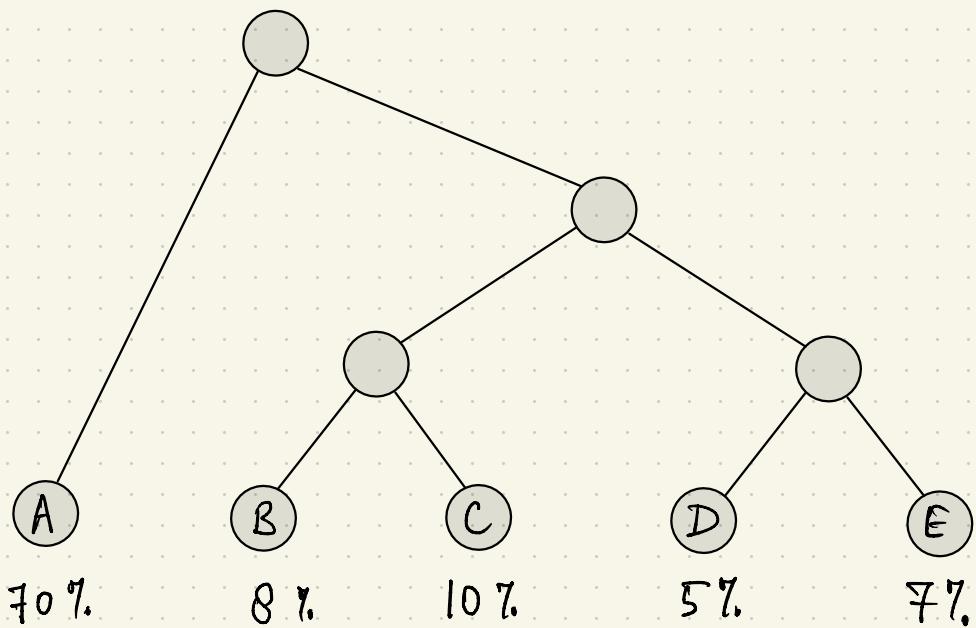
Combine trees



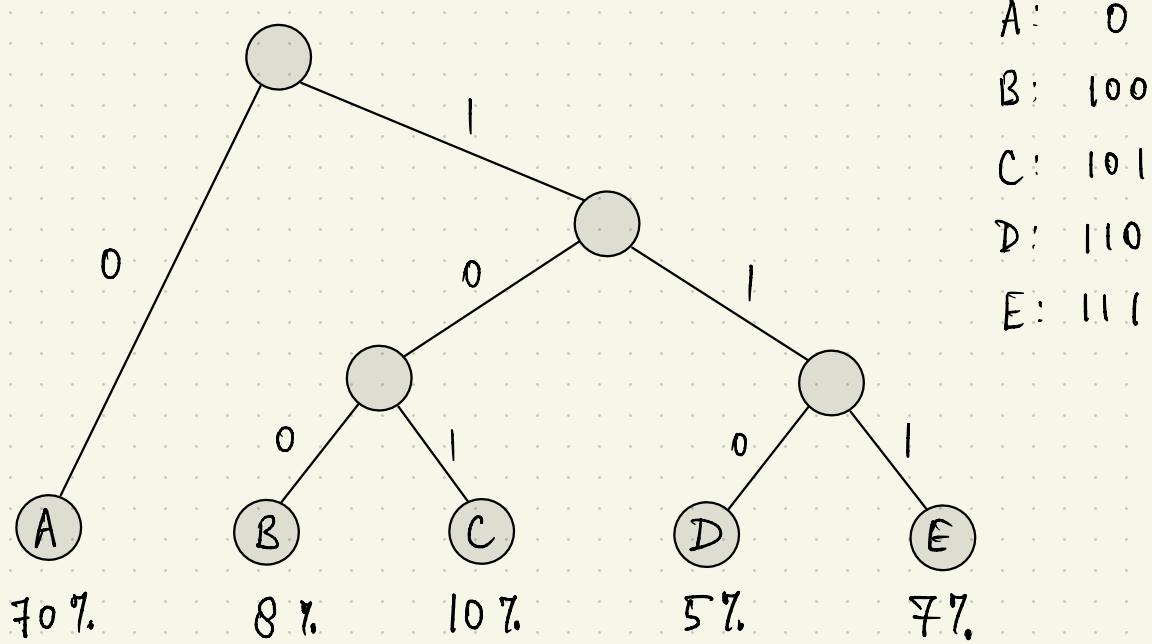
# HUFFMAN'S ALGORITHM



# HUFFMAN'S ALGORITHM

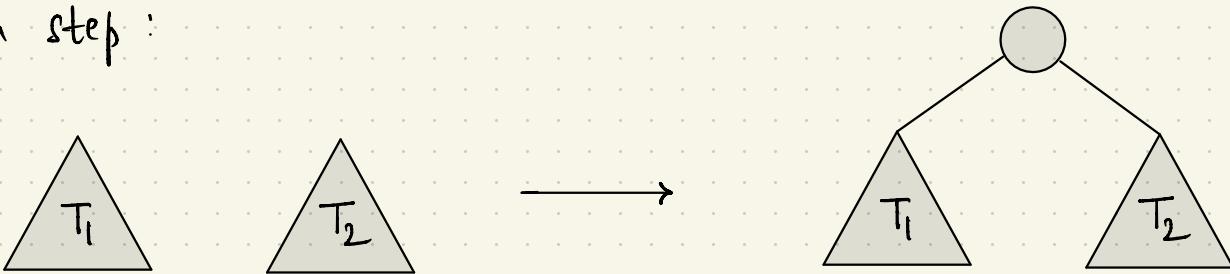


# HUFFMAN'S ALGORITHM



# AN ITERATIVE VIEW OF HUFFMAN'S ALGORITHM

at each step :



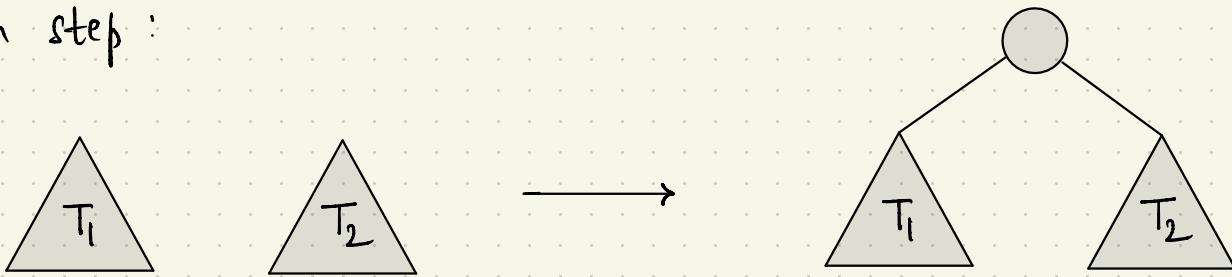
after each merge

👍 # trees in the forest decrease by 1

👎 average depth of leaves increases

# AN ITERATIVE VIEW OF HUFFMAN'S ALGORITHM

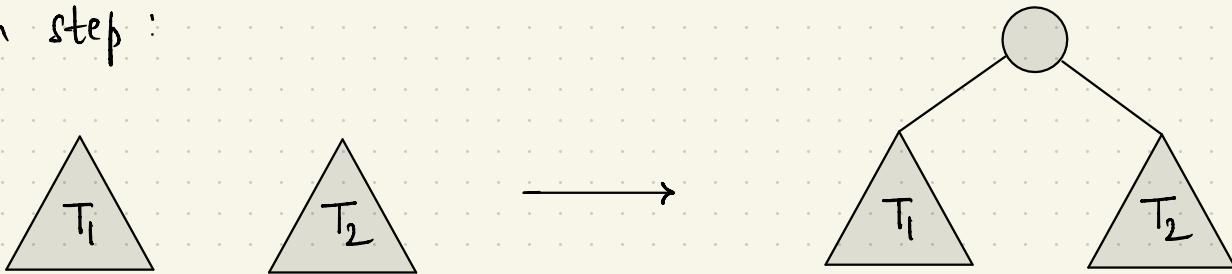
at each step :



$$\text{increase in average leaf-depth} = \sum_{a \in T_1} p_a + \sum_{b \in T_2} p_b$$

# AN ITERATIVE VIEW OF HUFFMAN'S ALGORITHM

at each step :



$$\text{increase in average leaf-depth} = \sum_{a \in T_1} p_a + \sum_{b \in T_2} p_b$$

Huffman's greedy criterion

Merge the pair of trees that causes  
minimum possible increase in average leaf depth.

# AN ITERATIVE VIEW OF HUFFMAN'S ALGORITHM

input: alphabet  $\Sigma$ , frequencies  $\{p_a\}_{a \in \Sigma}$

output:  $\Sigma$ -tree with minimum average leaf depth

# AN ITERATIVE VIEW OF HUFFMAN'S ALGORITHM

input: alphabet  $\Sigma$ , frequencies  $\{p_a\}_{a \in \Sigma}$

output:  $\Sigma$ -tree with minimum average leaf depth

// initialize

for each  $a \in \Sigma$

$T_a :=$  tree with one node labeled "a"

$P(T_a) := p_a$

# AN ITERATIVE VIEW OF HUFFMAN'S ALGORITHM

input: alphabet  $\Sigma$ , frequencies  $\{p_a\}_{a \in \Sigma}$

output:  $\Sigma$ -tree with minimum average leaf depth

// initialize

for each  $a \in \Sigma$

$T_a :=$  tree with one node labeled "a"

$P(T_a) := p_a$

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

# AN ITERATIVE VIEW OF HUFFMAN'S ALGORITHM

input: alphabet  $\Sigma$ , frequencies  $\{p_a\}_{a \in \Sigma}$

output:  $\Sigma$ -tree with minimum average leaf depth

// initialize

for each  $a \in \Sigma$

$T_a :=$  tree with one node labeled "a"

$P(T_a) := p_a$

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

// invariant:

$$\#T \in F \quad P(T) := \sum_{a \in T} p_a$$

# AN ITERATIVE VIEW OF HUFFMAN'S ALGORITHM

input: alphabet  $\Sigma$ , frequencies  $\{p_a\}_{a \in \Sigma}$

output:  $\Sigma$ -tree with minimum average leaf depth

// initialize

for each  $a \in \Sigma$

$T_a :=$  tree with one node labeled "a"

$P(T_a) := p_a$

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

// main loop

// invariant:

$$\#T \in F \quad P(T) := \sum_{a \in T} p_a$$

# AN ITERATIVE VIEW OF HUFFMAN'S ALGORITHM

input: alphabet  $\Sigma$ , frequencies  $\{p_a\}_{a \in \Sigma}$

output:  $\Sigma$ -tree with minimum average leaf depth

// initialize

for each  $a \in \Sigma$

$T_a :=$  tree with one node labeled "a"

$P(T_a) := p_a$

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

// invariant:

$$\#T \in F \quad P(T) := \sum_{a \in T} p_a$$

// main loop

while  $F$  has at least two trees

$$T_1 := \underset{T \in F}{\operatorname{argmin}} P(T), \quad T_2 := \underset{T \in F \setminus T_1}{\operatorname{argmin}} P(T)$$

# AN ITERATIVE VIEW OF HUFFMAN'S ALGORITHM

**input:** alphabet  $\Sigma$ , frequencies  $\{p_a\}_{a \in \Sigma}$

**output:**  $\Sigma$ -tree with minimum average leaf depth

// initialize

for each  $a \in \Sigma$

$T_a :=$  tree with one node labeled "a"

$P(T_a) := p_a$

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

// invariant:

$$\#T \in F \quad P(T) := \sum_{a \in T} p_a$$

// main loop

while  $F$  has at least two trees

$$T_1 := \underset{T \in F}{\operatorname{argmin}} P(T), \quad T_2 := \underset{T \in F \setminus T_1}{\operatorname{argmin}} P(T)$$

$T_3 :=$  merger of  $T_1$  and  $T_2$

$$P(T_3) := P(T_1) + P(T_2)$$

# AN ITERATIVE VIEW OF HUFFMAN'S ALGORITHM

**input:** alphabet  $\Sigma$ , frequencies  $\{p_a\}_{a \in \Sigma}$

**output:**  $\Sigma$ -tree with minimum average leaf depth

// initialize

for each  $a \in \Sigma$

$T_a :=$  tree with one node labeled "a"

$P(T_a) := p_a$

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

// invariant:

$$\#T \in F \quad P(T) := \sum_{a \in T} p_a$$

// main loop

while  $F$  has at least two trees

$$T_1 := \underset{T \in F}{\operatorname{argmin}} P(T), \quad T_2 := \underset{\substack{T \in F \\ T \neq T_1}}{\operatorname{argmin}} P(T)$$

$T_3 :=$  merge of  $T_1$  and  $T_2$

$$P(T_3) := P(T_1) + P(T_2)$$

$$F := F \cup T_3 \setminus \{T_1, T_2\}$$

# AN ITERATIVE VIEW OF HUFFMAN'S ALGORITHM

**input:** alphabet  $\Sigma$ , frequencies  $\{p_a\}_{a \in \Sigma}$

**output:**  $\Sigma$ -tree with minimum average leaf depth

// initialize

for each  $a \in \Sigma$

$T_a :=$  tree with one node labeled "a"

$P(T_a) := p_a$

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

// invariant:

$\forall T \in F \quad P(T) := \sum_{a \in T} p_a$

// main loop

while  $F$  has at least two trees

$T_1 := \underset{T \in F}{\operatorname{argmin}} P(T), \quad T_2 := \underset{\substack{T \in F \setminus T_1}}{\operatorname{argmin}} P(T)$

$T_3 :=$  merge of  $T_1$  and  $T_2$

$P(T_3) := P(T_1) + P(T_2)$

$F := F \cup T_3 \setminus \{T_1, T_2\}$

return unique tree in  $F$ .

# RUNNING TIME OF HUFFMAN'S ALGORITHM

# RUNNING TIME OF HUFFMAN'S ALGORITHM

if  $|\Sigma| = 2$

    encode one letter using 0 and other using 1

else

$x, y \leftarrow$  two lowest-frequency letters

    Create new alphabet  $\Sigma'$  by merging  $x$  and  $y$  into  $z$  with  $p_z := p_x + p_y$ .

    Recursively construct tree  $T'$  for  $(\Sigma', p)$

# RUNNING TIME OF HUFFMAN'S ALGORITHM

if  $|\Sigma| = 2$

    encode one letter using 0 and other using 1

else

$x, y \leftarrow$  two lowest-frequency letters

    Create new alphabet  $\Sigma'$  by merging  $x$  and  $y$  into  $z$  with  $p_z := p_x + p_y$ .

    Recursively construct tree  $T'$  for  $(\Sigma', p)$

$O(n)$

# RUNNING TIME OF HUFFMAN'S ALGORITHM

if  $|\Sigma| = 2$

  | encode one letter using 0 and other using 1

else

$x, y \leftarrow$  two lowest-frequency letters

  Create new alphabet  $\Sigma'$  by merging  $x$  and  $y$  into  $z$  with  $p_z := p_x + p_y$ .

  Recursively construct tree  $T'$  for  $(\Sigma', p)$

$O(n)$

$n-1$  calls

# RUNNING TIME OF HUFFMAN'S ALGORITHM

if  $|\Sigma| = 2$

    encode one letter using 0 and other using 1

else

$x, y \leftarrow$  two lowest-frequency letters

    Create new alphabet  $\Sigma'$  by merging  $x$  and  $y$  into  $z$  with  $p_z := p_x + p_y$ .

    Recursively construct tree  $T'$  for  $(\Sigma', p)$

$O(n)$

$O(n^2)$  overall

$n-1$  calls

# RUNNING TIME OF HUFFMAN'S ALGORITHM

if  $|\Sigma| = 2$

  | encode one letter using 0 and other using 1

else

$x, y \leftarrow$  two lowest-frequency letters

  Create new alphabet  $\Sigma'$  by merging  $x$  and  $y$  into  $z$  with  $p_z := p_x + p_y$ .

  Recursively construct tree  $T'$  for  $(\Sigma', p)$

$O(n^2)$  overall

$O(\log n)$  with  
heap

~~$O(n)$~~

$n-1$  calls

# RUNNING TIME OF HUFFMAN'S ALGORITHM

if  $|\Sigma| = 2$

  | encode one letter using 0 and other using 1

else

$x, y \leftarrow$  two lowest-frequency letters

  Create new alphabet  $\Sigma'$  by merging  $x$  and  $y$  into  $z$  with  $p_z := p_x + p_y$ .

  Recursively construct tree  $T'$  for  $(\Sigma', p)$

$O(n \log n)$

~~$O(n^2)$~~  overall

$O(\log n)$  with  
heap

~~$O(n)$~~

$n-1$  calls

# CORRECTNESS OF HUFFMAN'S ALGORITHM

## CORRECTNESS OF HUFFMAN'S ALGORITHM

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

## CORRECTNESS OF HUFFMAN'S ALGORITHM

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

**Proof:** (by strong induction on  $|\Sigma| = n$ )

## CORRECTNESS OF HUFFMAN'S ALGORITHM

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

**Proof:** (by strong induction on  $|\Sigma| = n$ )

Base case :

## CORRECTNESS OF HUFFMAN'S ALGORITHM

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

**Proof:** (by strong induction on  $|\Sigma| = n$ )

**Base case :**  $n = 2$

Huffman uses one bit per symbol — **optimal!**

## CORRECTNESS OF HUFFMAN'S ALGORITHM

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

**Proof :** Induction step : Huffman is optimal for all  $\Sigma$  with  $|\Sigma| \leq n$ .

## CORRECTNESS OF HUFFMAN'S ALGORITHM

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

**Proof :** Induction step : Huffman is optimal for all  $\Sigma$  with  $|\Sigma| \leq n$ .  
Consider  $\Sigma$  such that  $|\Sigma| = n+1$ .

## CORRECTNESS OF HUFFMAN'S ALGORITHM

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

**Proof :** Induction step : Huffman is optimal for all  $\Sigma$  with  $|\Sigma| \leq n$ .  
Consider  $\Sigma$  such that  $|\Sigma| = n+1$ .

Huffman's algo merges  $x, y \in \Sigma$  into meta-symbol  $z$ .

## CORRECTNESS OF HUFFMAN'S ALGORITHM

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

**Proof :** Induction step : Huffman is optimal for all  $\Sigma$  with  $|\Sigma| \leq n$ .

Consider  $\Sigma$  such that  $|\Sigma| = n+1$ .

Huffman's algo merges  $x, y \in \Sigma$  into meta-symbol  $z$ .

Let  $\Sigma'$  be the new alphabet. Thus,  $|\Sigma'| = n$ .

## CORRECTNESS OF HUFFMAN'S ALGORITHM

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

**Proof :** Induction step : Huffman is optimal for all  $\Sigma$  with  $|\Sigma| \leq n$ .

Consider  $\Sigma$  such that  $|\Sigma| = n+1$ .

Huffman's algo merges  $x, y \in \Sigma$  into meta-symbol  $z$ .

Let  $\Sigma'$  be the new alphabet. Thus,  $|\Sigma'| = n$ .

Let  $T'$  be algorithm's output for  $\Sigma' \Rightarrow T'$  is Optimal !

# CORRECTNESS OF HUFFMAN'S ALGORITHM

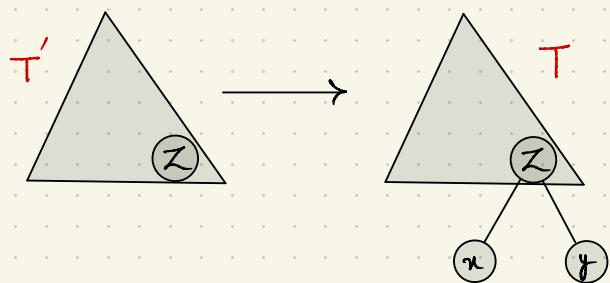
**Theorem:** For every alphabet  $\Sigma$  and non-negative frequencies  $\{p_a\}_{a \in \Sigma}$ , Huffman's algorithm returns a prefix-free code with minimum possible average encoding length.

**Proof :** Algorithm's output for  $\Sigma' \rightarrow T'$

# CORRECTNESS OF HUFFMAN'S ALGORITHM

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

**Proof :** Algorithm's output for  $\Sigma' \rightarrow T'$



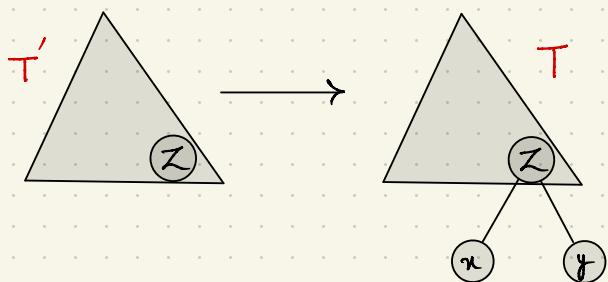
# CORRECTNESS OF HUFFMAN'S ALGORITHM

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

**Proof:** Algorithm's output for  $\Sigma' \rightarrow T'$

" " " " "  $\Sigma \rightarrow T$

$$L(T') = L(T) - (p_x + p_y) \text{ (exercise)}$$



# CORRECTNESS OF HUFFMAN'S ALGORITHM

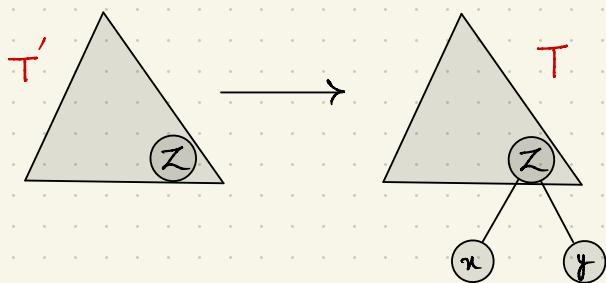
**Theorem:** For every alphabet  $\Sigma$  and non-negative frequencies  $\{p_a\}_{a \in \Sigma}$ , Huffman's algorithm returns a prefix-free code with minimum possible average encoding length.

**Proof:** Algorithm's output for  $\Sigma' \rightarrow T'$

" " " " " "  $\Sigma \rightarrow T$

$$L(T') = L(T) - (p_x + p_y) \text{ (exercise)}$$

If  $T$  not optimal, then  $L(T) > L(S)$  for some tree  $S$



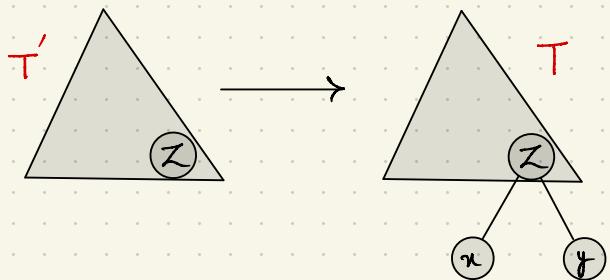
# CORRECTNESS OF HUFFMAN'S ALGORITHM

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

**Proof :** Algorithm's output for  $\Sigma' \rightarrow T'$

" " " " "  $\Sigma \rightarrow T$

$$L(T') = L(T) - (p_x + p_y) \text{ (exercise)}$$



If  $T$  not optimal, then  $L(T) > L(S)$  for some tree  $S$

**Corollary**  $\Rightarrow x$  and  $y$  are siblings in  $S$ .

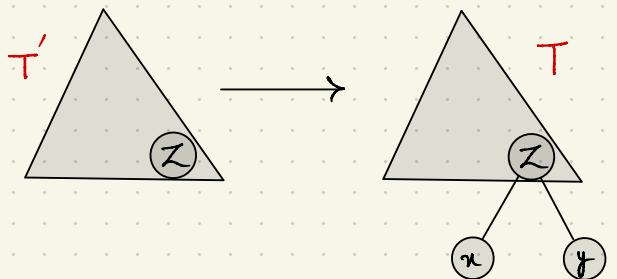
# CORRECTNESS OF HUFFMAN'S ALGORITHM

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

**Proof:** Algorithm's output for  $\Sigma' \rightarrow T'$

" " " " "  $\Sigma \rightarrow T$

$$L(T') = L(T) - (p_x + p_y) \text{ (exercise)}$$



If  $T$  not optimal, then  $L(T) > L(S)$  for some tree  $S$

**Corollary**  $\Rightarrow x$  and  $y$  are siblings in  $S$ .

Create  $S'$  by "merging"  $x$  and  $y$ .

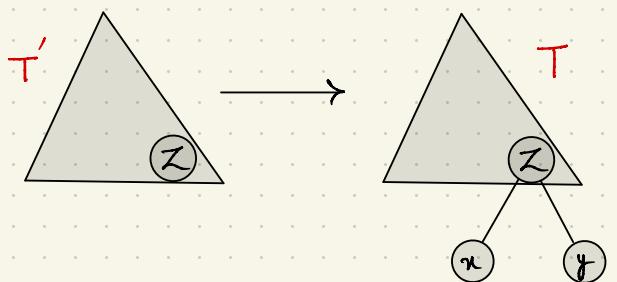
# CORRECTNESS OF HUFFMAN'S ALGORITHM

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

**Proof:** Algorithm's output for  $\Sigma' \rightarrow T'$

" " " " "  $\Sigma \rightarrow T$

$$L(T') = L(T) - (p_x + p_y) \text{ (exercise)}$$



If  $T$  not optimal, then  $L(T) > L(S)$  for some tree  $S$ .

**Corollary**  $\Rightarrow x$  and  $y$  are siblings in  $S$ .

Create  $S'$  by "merging"  $x$  and  $y$ .

Then,  $L(S') < L(T')$ . Contradiction!

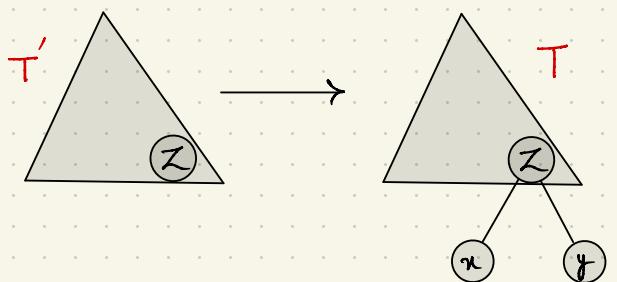
# CORRECTNESS OF HUFFMAN'S ALGORITHM

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

**Proof:** Algorithm's output for  $\Sigma' \rightarrow T'$

" " " " "  $\Sigma \rightarrow T$

$$L(T') = L(T) - (p_x + p_y) \text{ (exercise)}$$



If  $T$  not optimal, then  $L(T) > L(S)$  for some tree  $S$ .

**Corollary**  $\Rightarrow x$  and  $y$  are siblings in  $S$ .

Create  $S'$  by "merging"  $x$  and  $y$ .

Then,  $L(S') < L(T')$ . Contradiction!



# SUMMARY

Goal : Prefix-free code with minimum average encoding length

- Key ideas :
- ① Reimagining codes as trees
  - ② Two lowest-frequency symbols are "locked in" as deepest leaves.
  - ③ "Merging" into meta-symbol facilitates recursion

MINIMUM SPANNING TREES

# SPANNING TREES

# SPANNING TREES

A spanning tree of a connected graph

# SPANNING TREES

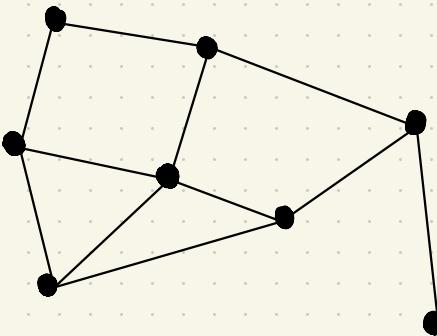
A spanning tree of a connected graph is a subgraph that is itself a tree and

# SPANNING TREES

A spanning tree of a connected graph is a subgraph that is itself a tree and whose vertex set is the same as that of the original graph.

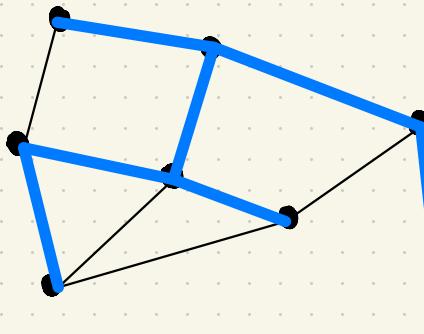
# SPANNING TREES

A spanning tree of a connected graph is a subgraph that is itself a tree and whose vertex set is the same as that of the original graph.



# SPANNING TREES

A spanning tree of a connected graph is a subgraph that is itself a tree and whose vertex set is the same as that of the original graph.



**Theorem:** Every connected graph has a spanning tree.

**Theorem:** Every connected graph has a spanning tree.

**Proof:** By contradiction.

Suppose, for contradiction, that some connected graph  $G$  has no spanning tree.

**Theorem:** Every connected graph has a spanning tree.

**Proof:** By contradiction.

Suppose, for contradiction, that some connected graph  $G$  has no spanning tree.

Let  $T$  be a connected subgraph of  $G$  with the same set of vertices and the **smallest** number of edges possible.

**Theorem:** Every connected graph has a spanning tree.

**Proof:** By contradiction.

Suppose, for contradiction, that some connected graph  $G$  has no spanning tree.

Let  $T$  be a connected subgraph of  $G$  with the same set of vertices and the **smallest** number of edges possible.



$T$  is well-defined because of well ordering principle.

**Theorem:** Every connected graph has a spanning tree.

**Proof:** By contradiction.

Suppose, for contradiction, that some connected graph  $G$  has no spanning tree.

Let  $T$  be a connected subgraph of  $G$  with the same set of vertices and the **smallest** number of edges possible.

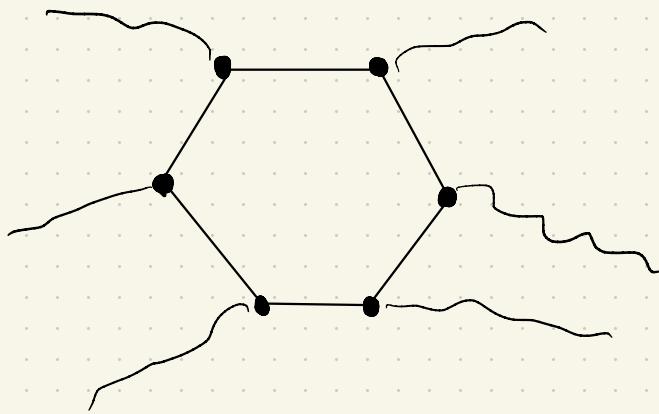
Then, by contradiction assumption,  $T$  is **NOT** a tree.

**Theorem:** Every connected graph has a spanning tree.

**Proof:** Therefore,  $T$  must have a cycle.

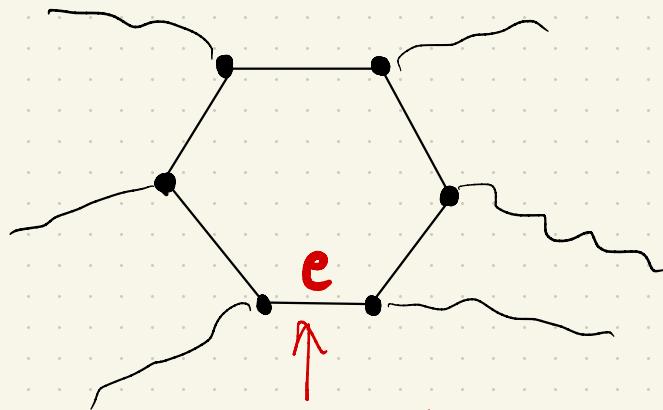
**Theorem:** Every connected graph has a spanning tree.

**Proof:** Therefore,  $T$  must have a cycle.



**Theorem:** Every connected graph has a spanning tree.

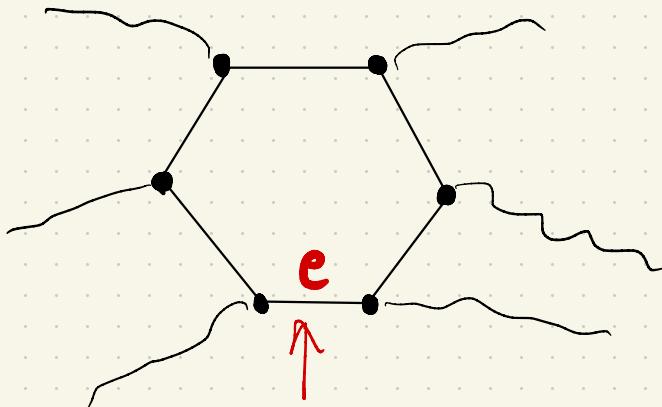
**Proof:** Therefore,  $T$  must have a cycle.



remove this

**Theorem:** Every connected graph has a spanning tree.

**Proof:** Therefore,  $T$  must have a cycle.



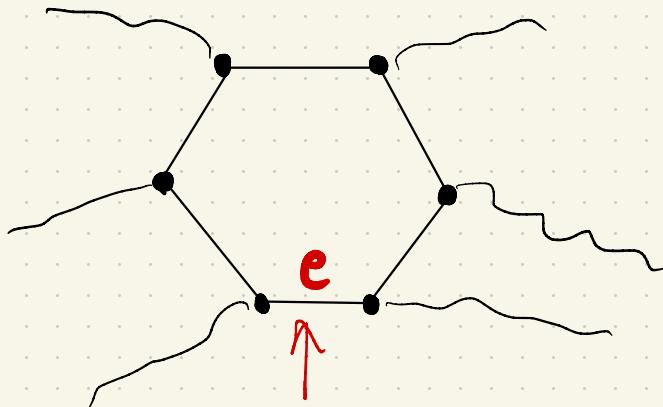
Case I :



does not contain  $e$ .

**Theorem:** Every connected graph has a spanning tree.

**Proof:** Therefore,  $T$  must have a cycle.

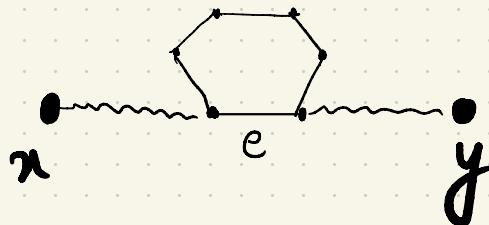


Case I :



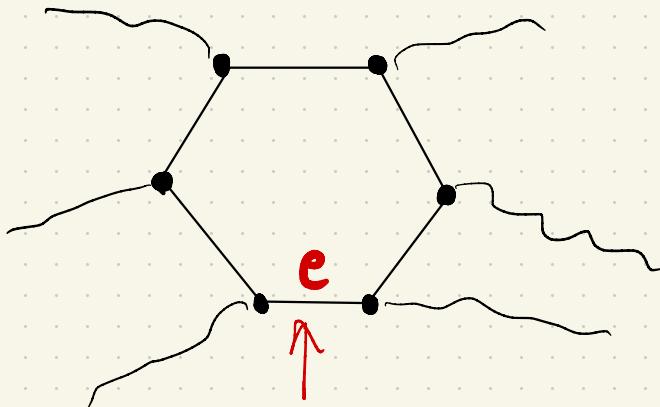
does not contain  $e$ .

Case II :



**Theorem:** Every connected graph has a spanning tree.

**Proof:** Therefore, T must have a cycle.

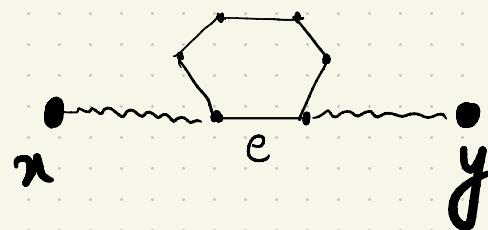


Case I :



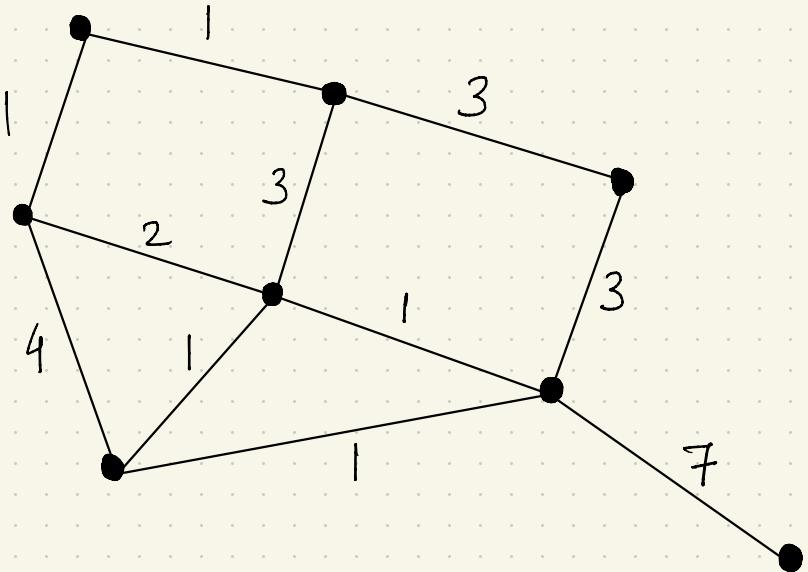
does not contain e.

Case II :

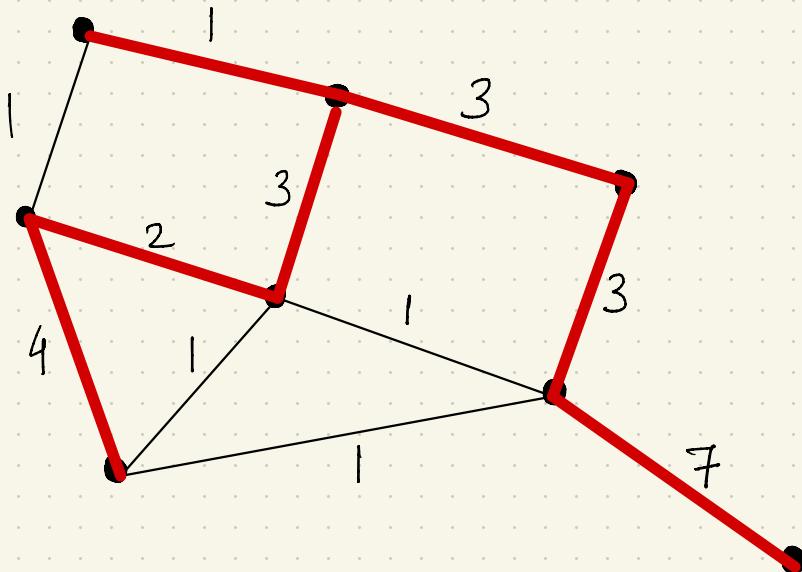


All vertices in T are still connected after removing e.

# MINIMUM WEIGHT SPANNING TREE

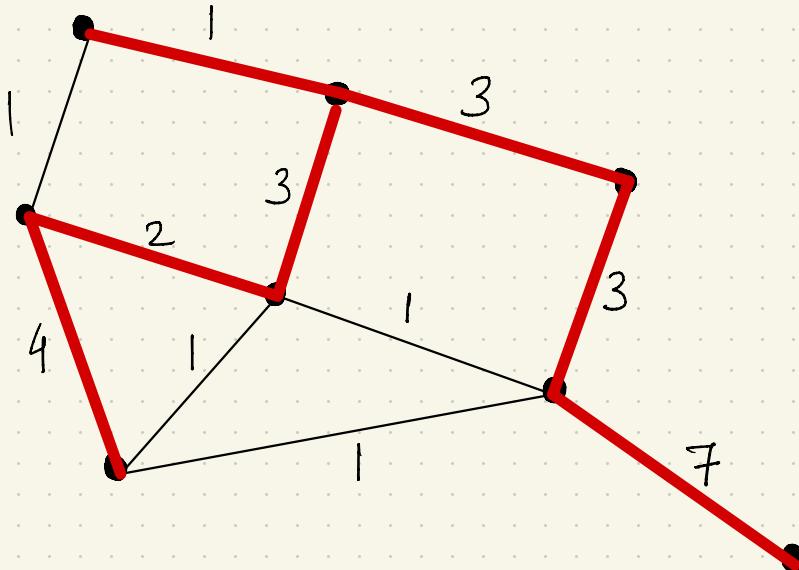


# MINIMUM WEIGHT SPANNING TREE



$$\text{Weight} = 1 + 3 + 3 + 3 + 2 + 4 + 7 = 23$$

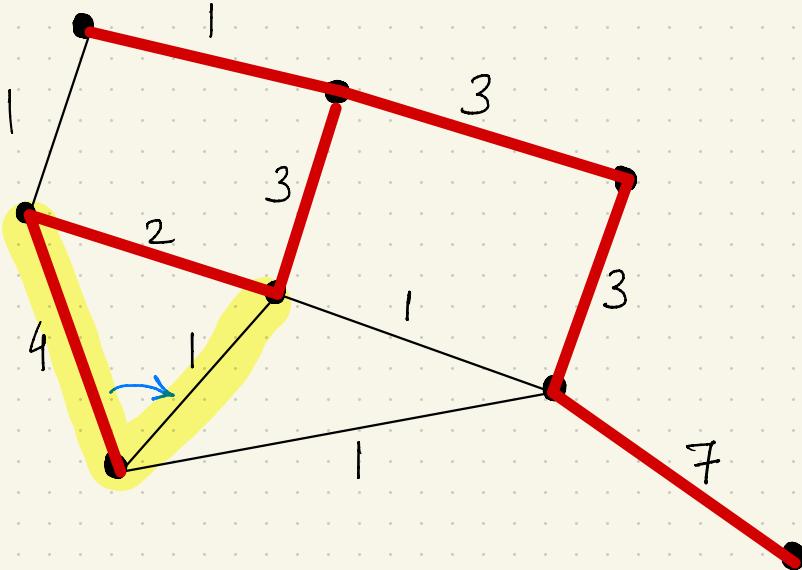
# MINIMUM WEIGHT SPANNING TREE



$$\text{Weight} = 1 + 3 + 3 + 3 + 2 + 4 + 7 = 23$$

Can we do better?

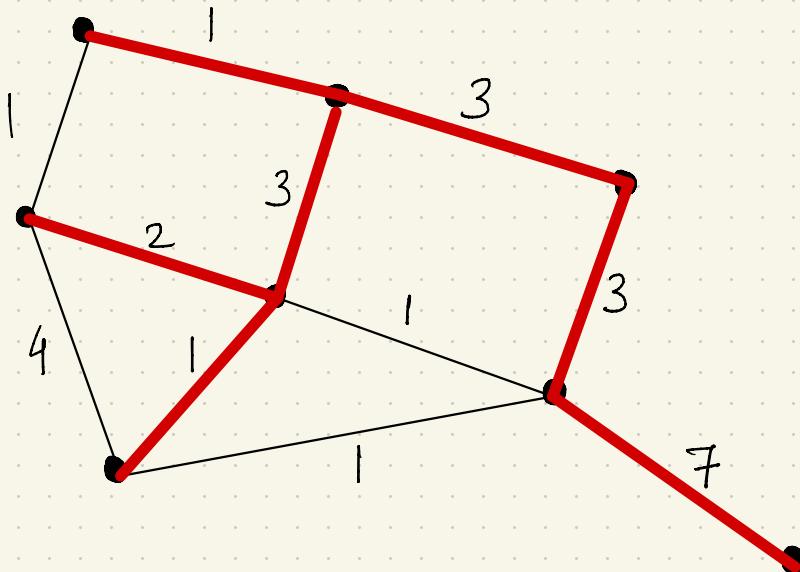
# MINIMUM WEIGHT SPANNING TREE



$$\text{Weight} = 1 + 3 + 3 + 3 + 2 + 4 + 7 = 23$$

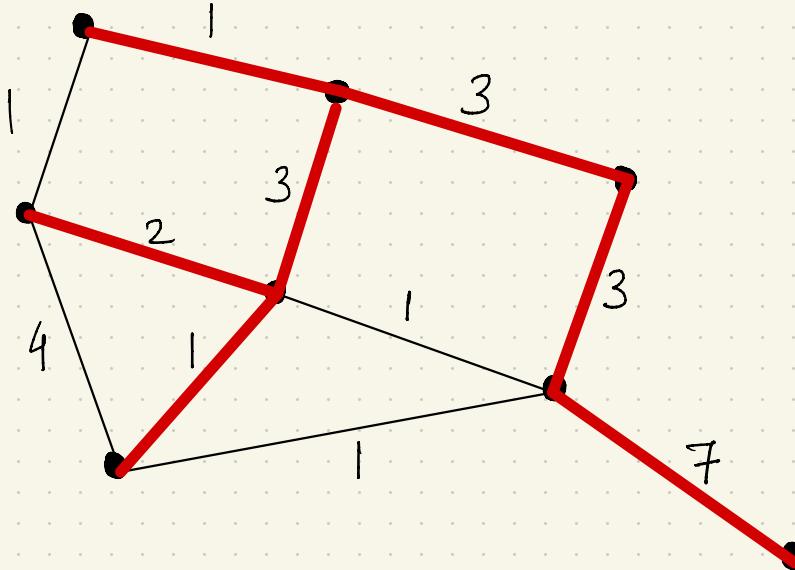
Can we do better?

# MINIMUM WEIGHT SPANNING TREE



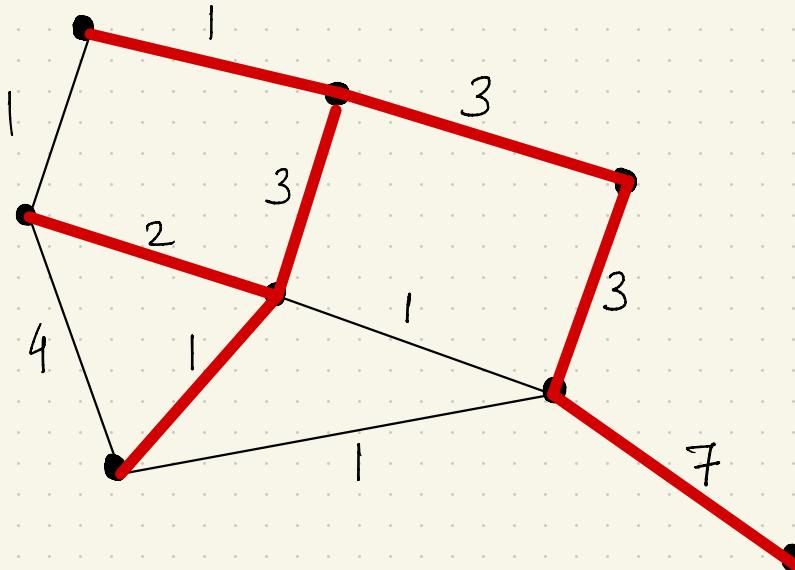
$$\text{Weight} = 1 + 3 + 3 + 3 + 2 + \cancel{4} + 7 = 2\cancel{3} \ 20$$

# MINIMUM WEIGHT SPANNING TREE



$$\text{Weight} = 1 + 3 + 3 + 3 + 2 + 1 + 7 = 20$$

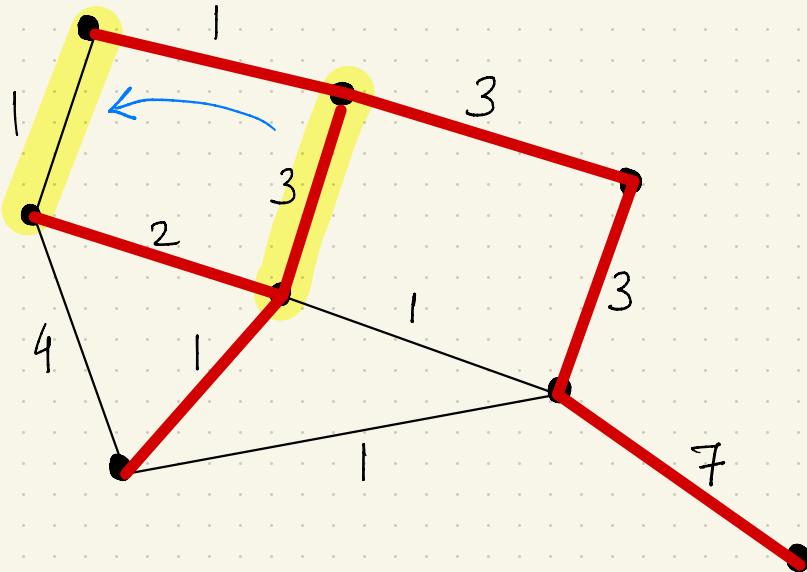
# MINIMUM WEIGHT SPANNING TREE



$$\text{Weight} = 1 + 3 + 3 + 3 + 2 + 1 + 7 = 20$$

Can we do better?

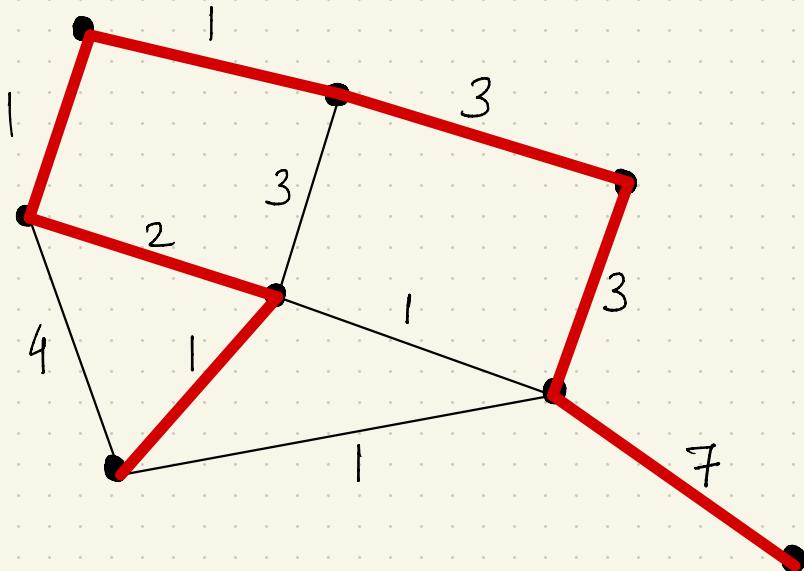
# MINIMUM WEIGHT SPANNING TREE



$$\text{Weight} = 1 + 3 + 3 + 3 + 2 + 1 + 7 = 20$$

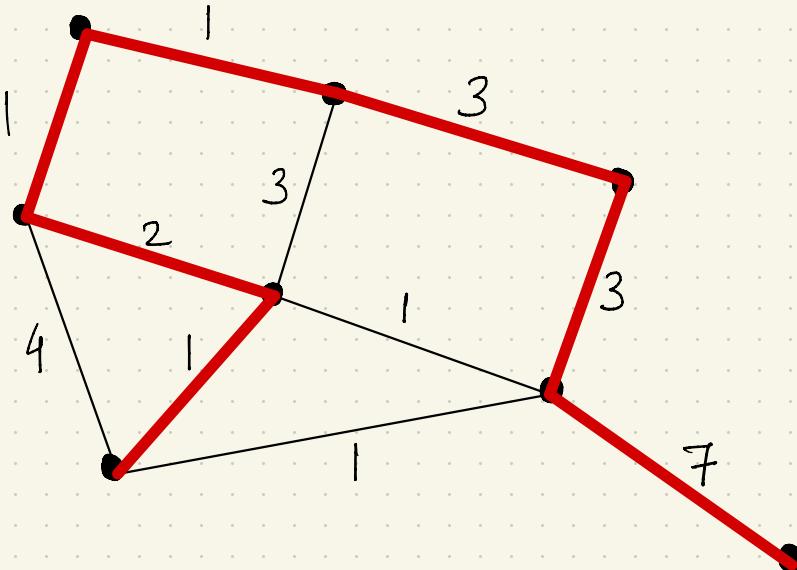
Can we do better?

# MINIMUM WEIGHT SPANNING TREE



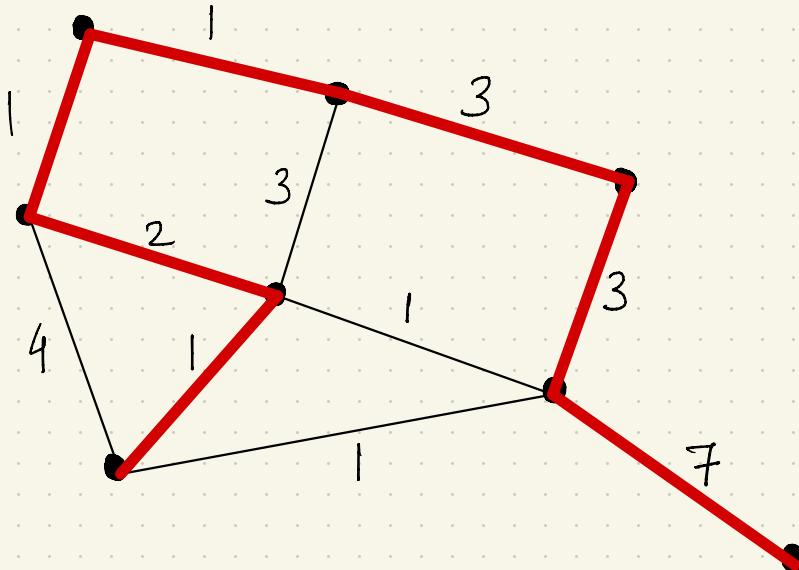
$$\text{Weight} = 1 + \cancel{8} + 3 + 3 + 2 + 1 + 7 = \cancel{20} 18$$

# MINIMUM WEIGHT SPANNING TREE



$$\text{Weight} = 1 + 1 + 3 + 3 + 2 + 1 + 7 = 18$$

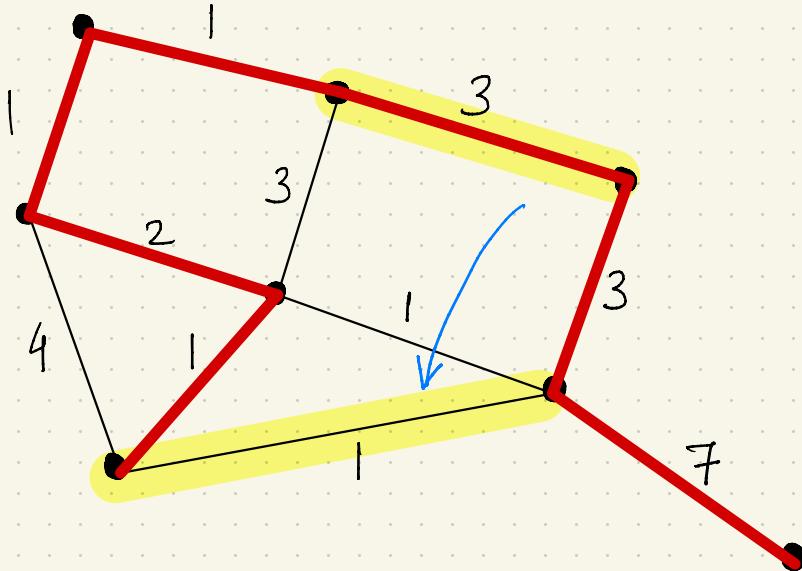
# MINIMUM WEIGHT SPANNING TREE



$$\text{Weight} = 1 + 1 + 3 + 3 + 2 + 1 + 7 = 18$$

Can we do better?

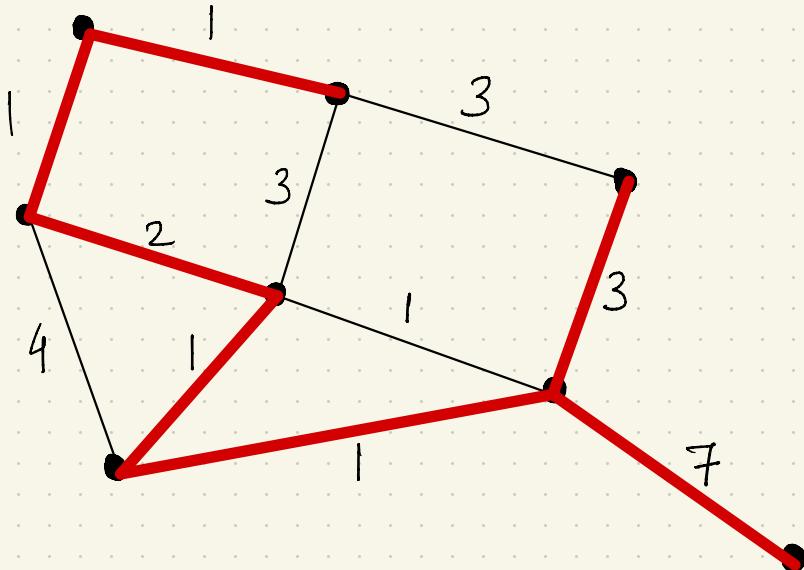
# MINIMUM WEIGHT SPANNING TREE



$$\text{Weight} = 1 + 1 + 3 + 3 + 2 + 1 + 7 = 18$$

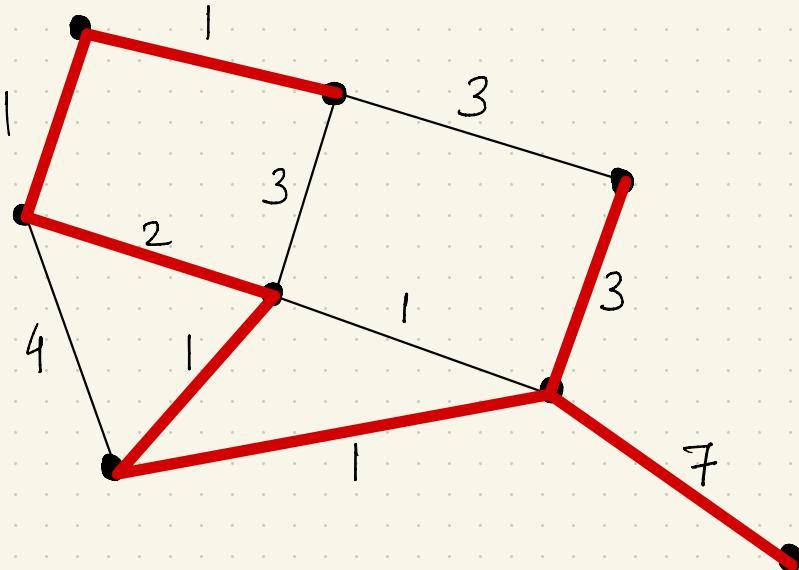
Can we do better?

# MINIMUM WEIGHT SPANNING TREE



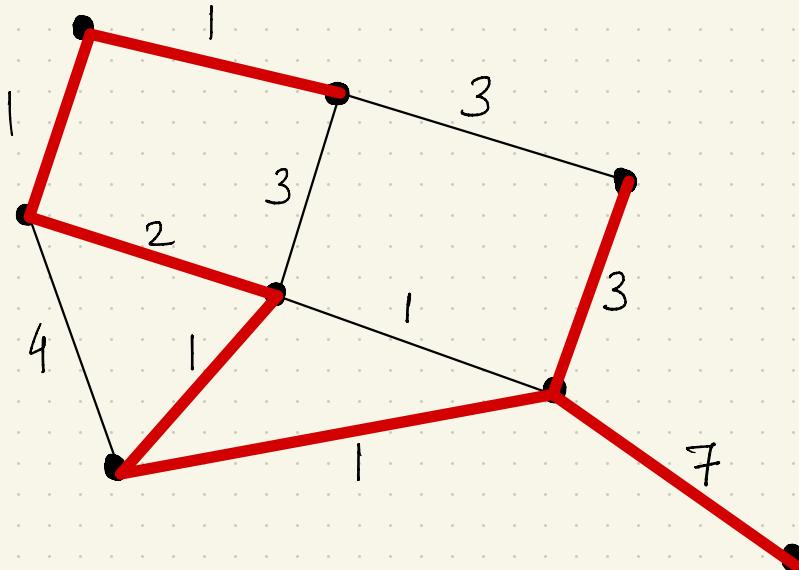
$$\text{Weight} = 1 + 1 + \cancel{X} + 3 + 2 + 1 + 7 = \cancel{18} 16$$

# MINIMUM WEIGHT SPANNING TREE



$$\text{Weight} = 1 + 1 + 1 + 3 + 2 + 1 + 7 = 16$$

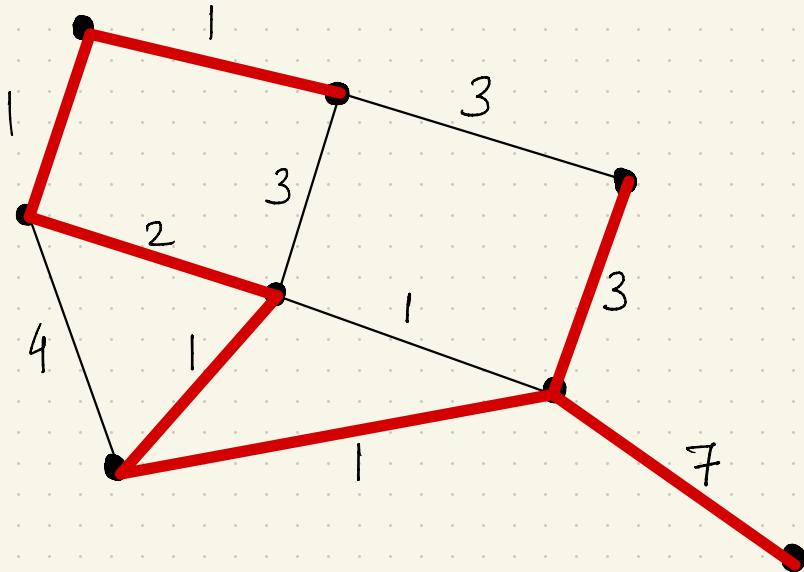
# MINIMUM WEIGHT SPANNING TREE



$$\text{Weight} = 1 + 1 + 1 + 3 + 2 + 1 + 7 = 16$$

Can we do better?

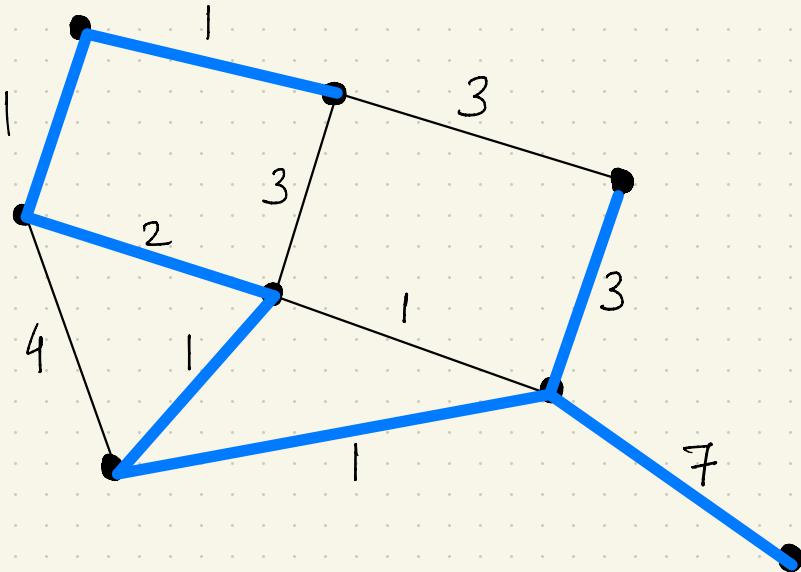
# MINIMUM WEIGHT SPANNING TREE



$$\text{Weight} = 1 + 1 + 1 + 3 + 2 + 1 + 7 = 16$$

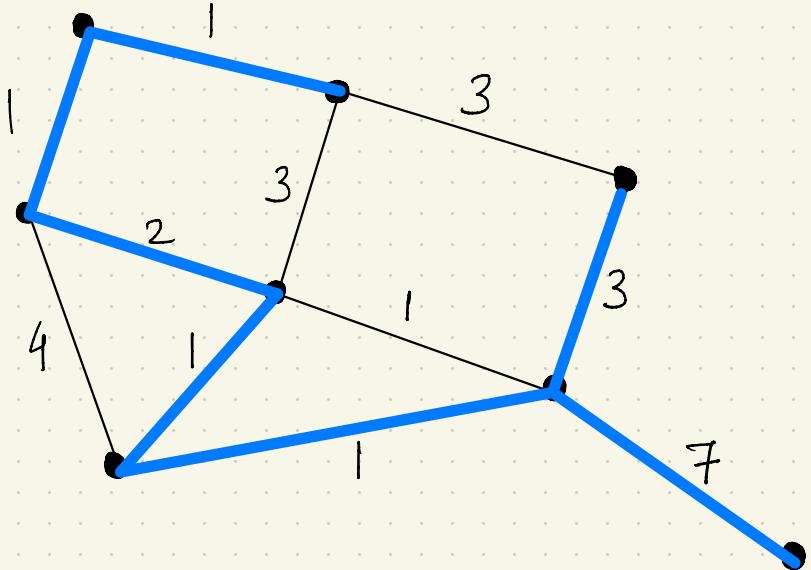
Can we do better? No!

# MINIMUM WEIGHT SPANNING TREE



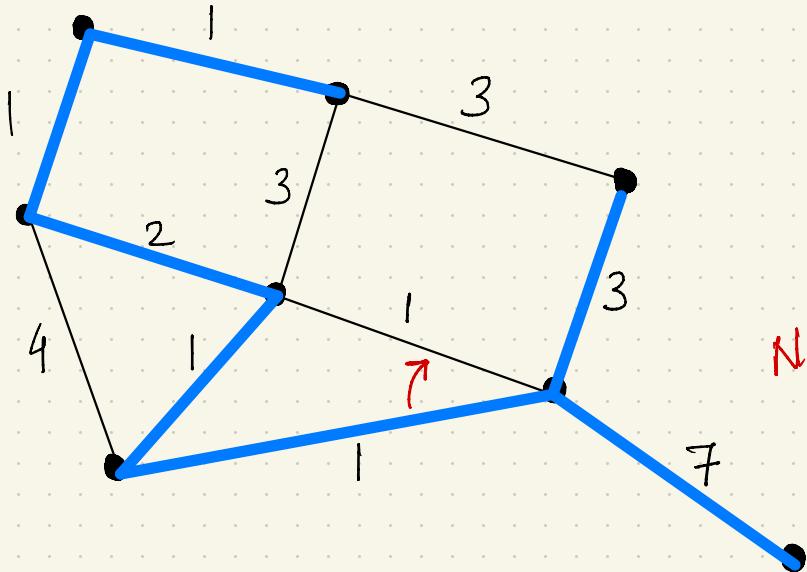
Spanning tree of minimum weight

# MINIMUM WEIGHT SPANNING TREE



$$\text{Weight} = 1 + 1 + 1 + 1 + 2 + 3 + 7 = 16$$

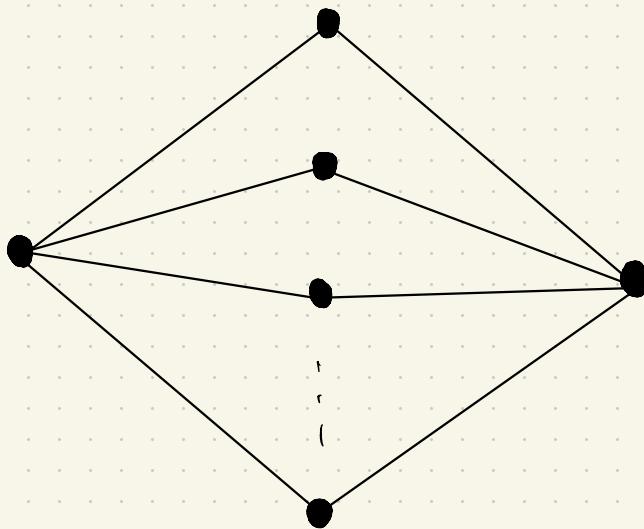
# MINIMUM WEIGHT SPANNING TREE



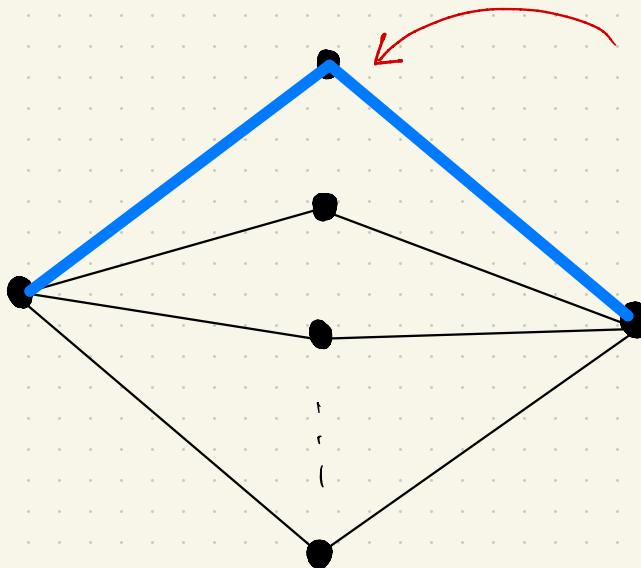
NOT UNIQUE!

$$\text{Weight} = 1 + 1 + 1 + 1 + 2 + 3 + 7 = 16$$

# EXPONENTIALLY MANY SPANNING TREES

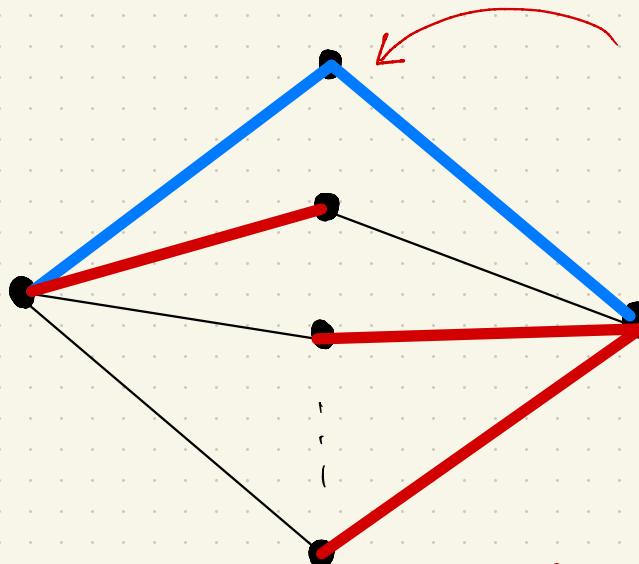


# EXPONENTIALLY MANY SPANNING TREES



Some middle node  
must have both edges  
in spanning tree

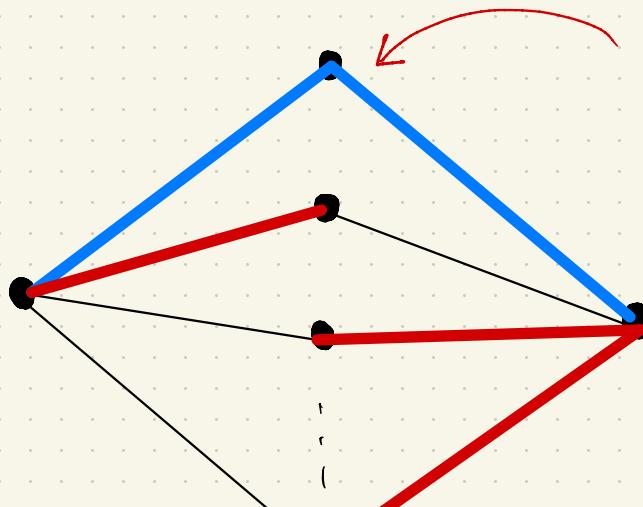
# EXPONENTIALLY MANY SPANNING TREES



Some middle node  
must have both edges  
in spanning tree

for every other node,  
pick either left or right edge

# EXPONENTIALLY MANY SPANNING TREES

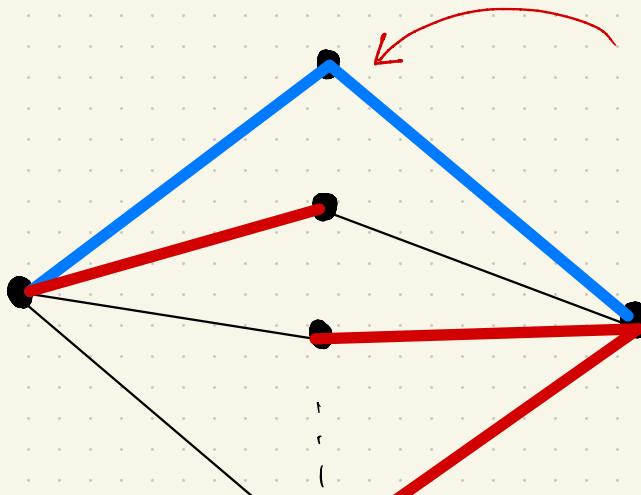


$$\# \text{ Spanning trees} \geq 2^{n-3}$$

Some middle node  
must have both edges  
in spanning tree

for every other node,  
pick either left or right edge

# EXPONENTIALLY MANY SPANNING TREES



$$\# \text{ spanning trees} \geq 2^{n-3}$$

Cannot hope to find an MST  
via brute force enumeration

Some middle node  
must have both edges  
in spanning tree

for every other node,  
pick either left or right edge