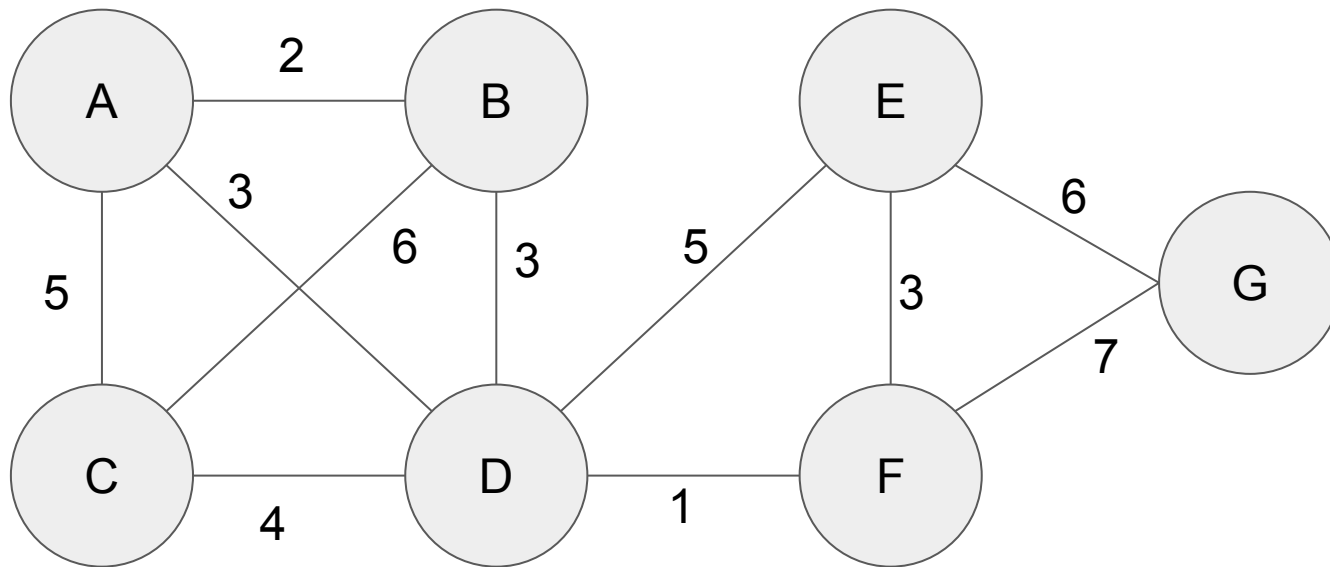


# Greedy algorithms

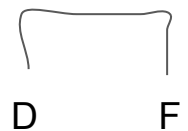
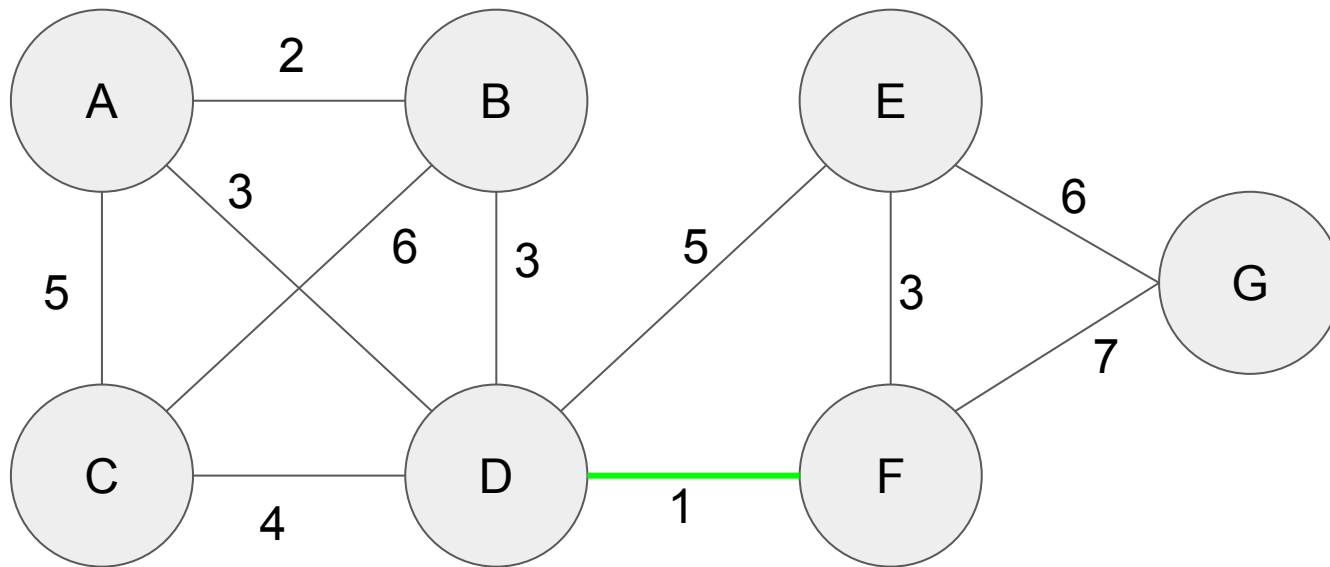
CMPSC 465 - Yana Safonova

# MST and clustering

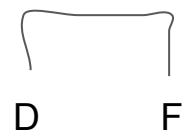
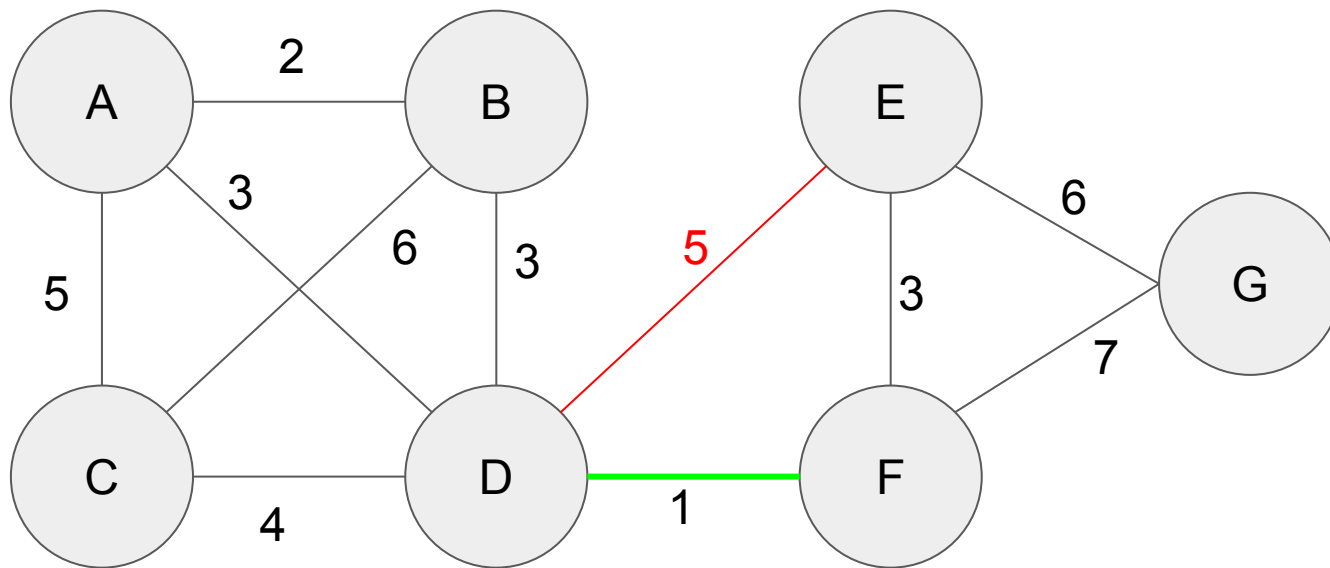
# Kruskal's algorithm aka single-linkage clustering



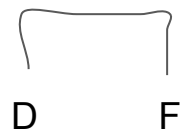
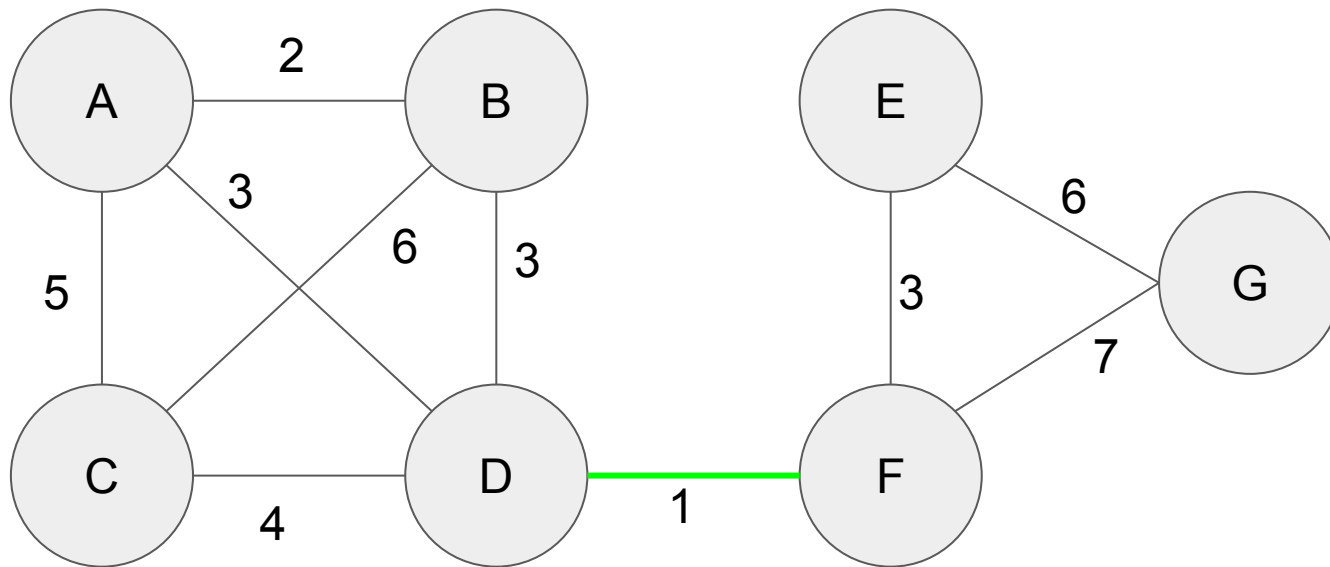
# Kruskal's algorithm - weight 1



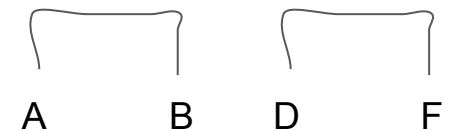
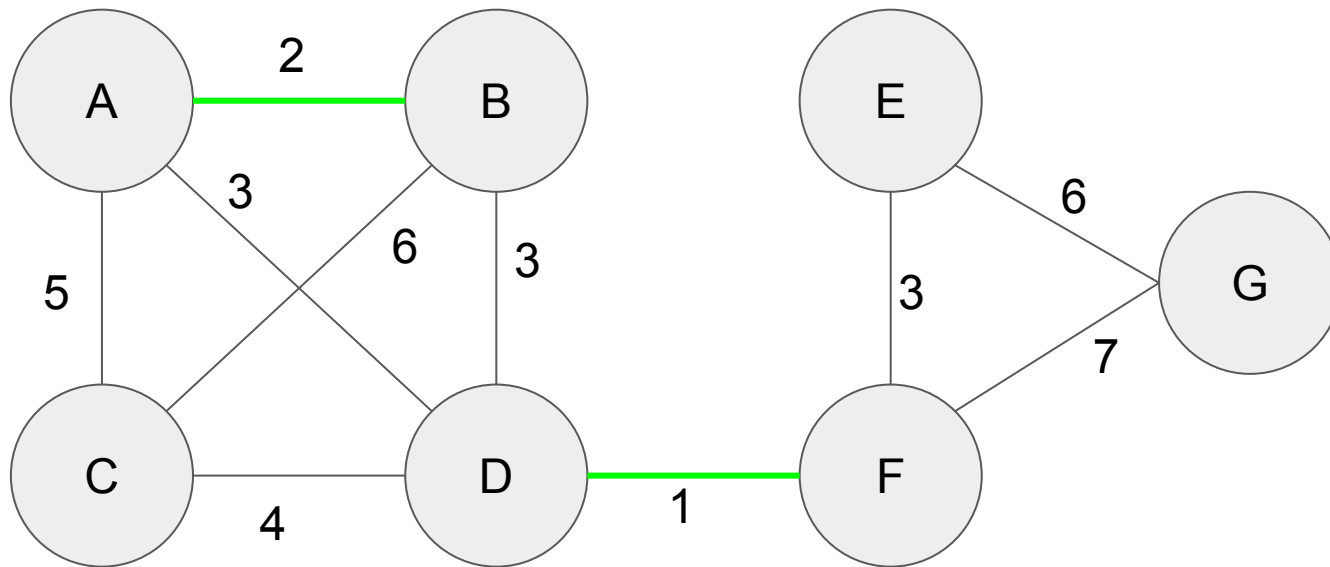
# Kruskal's algorithm - weight 1



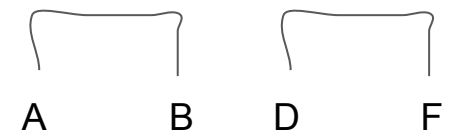
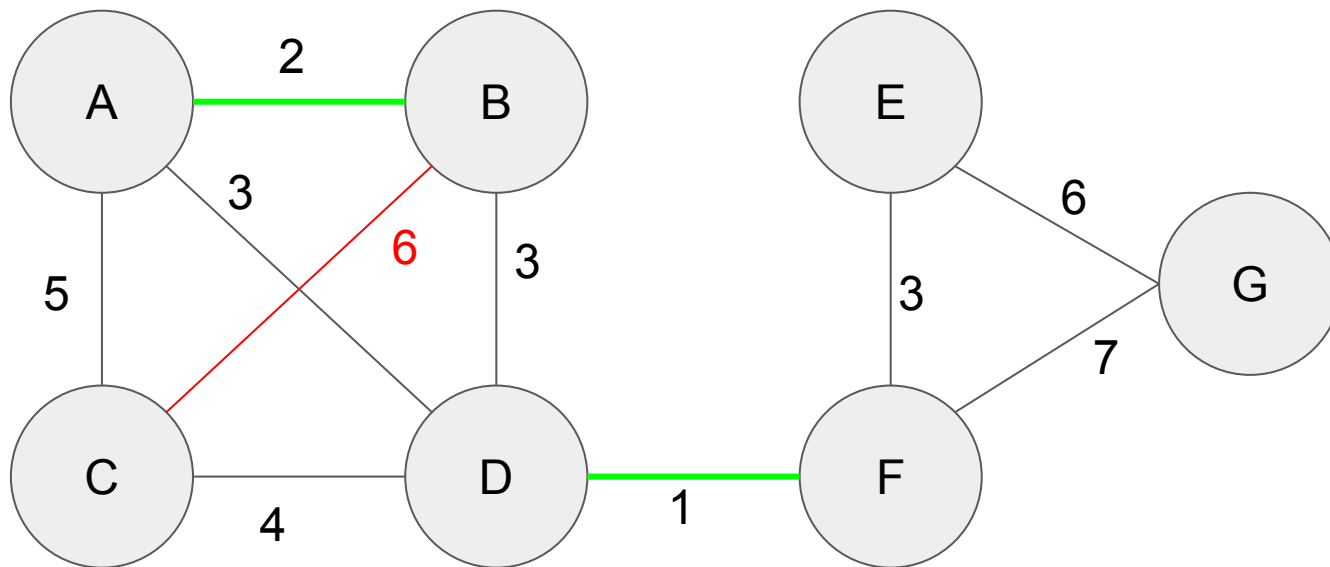
# Kruskal's algorithm - weight 1



# Kruskal's algorithm - weight 2

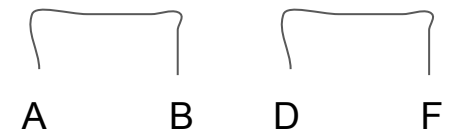
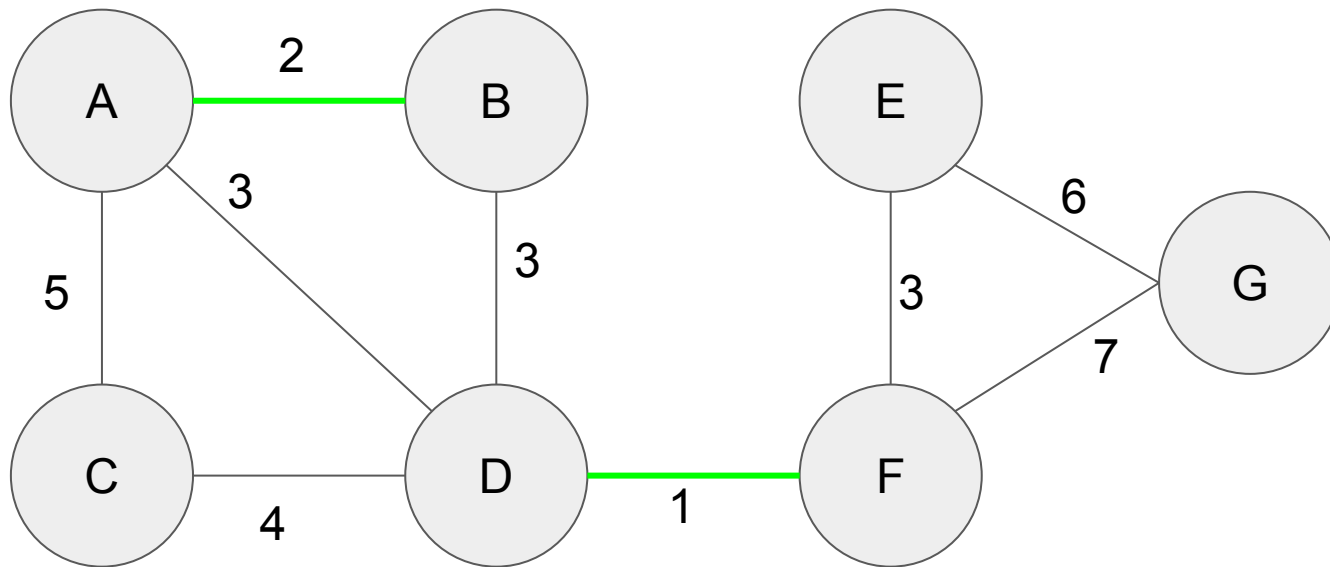


# Kruskal's algorithm - weight 2

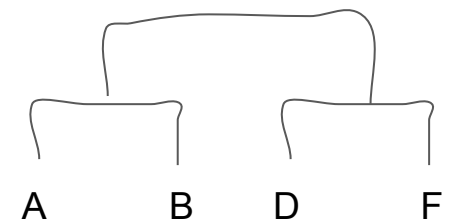
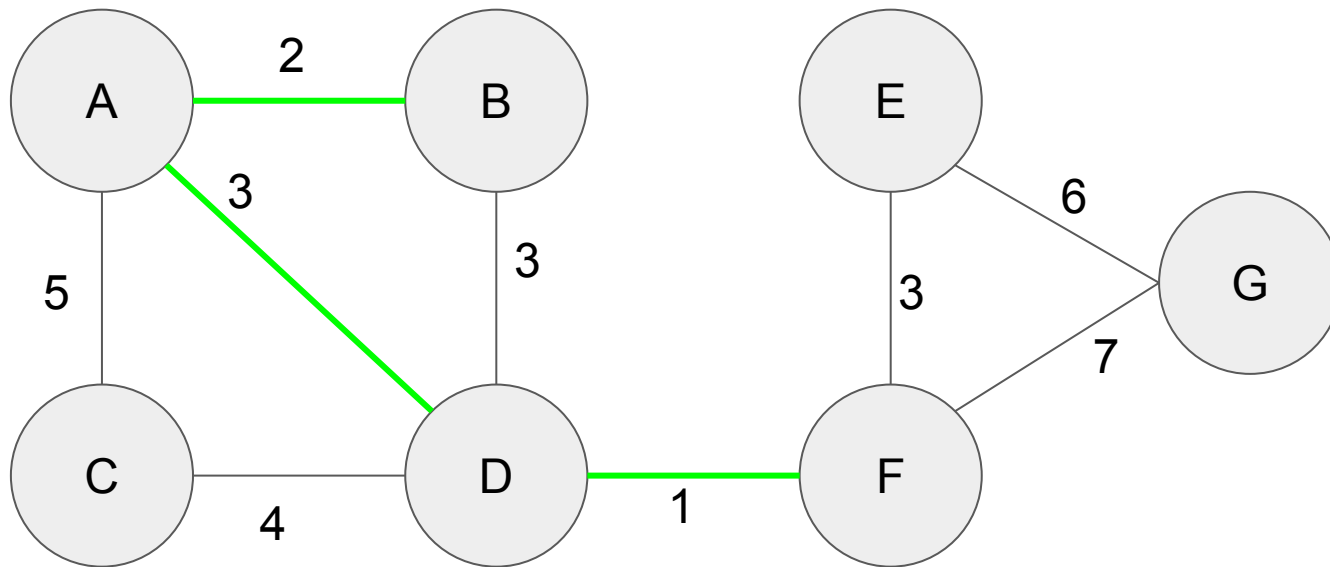




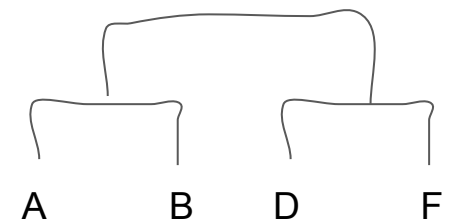
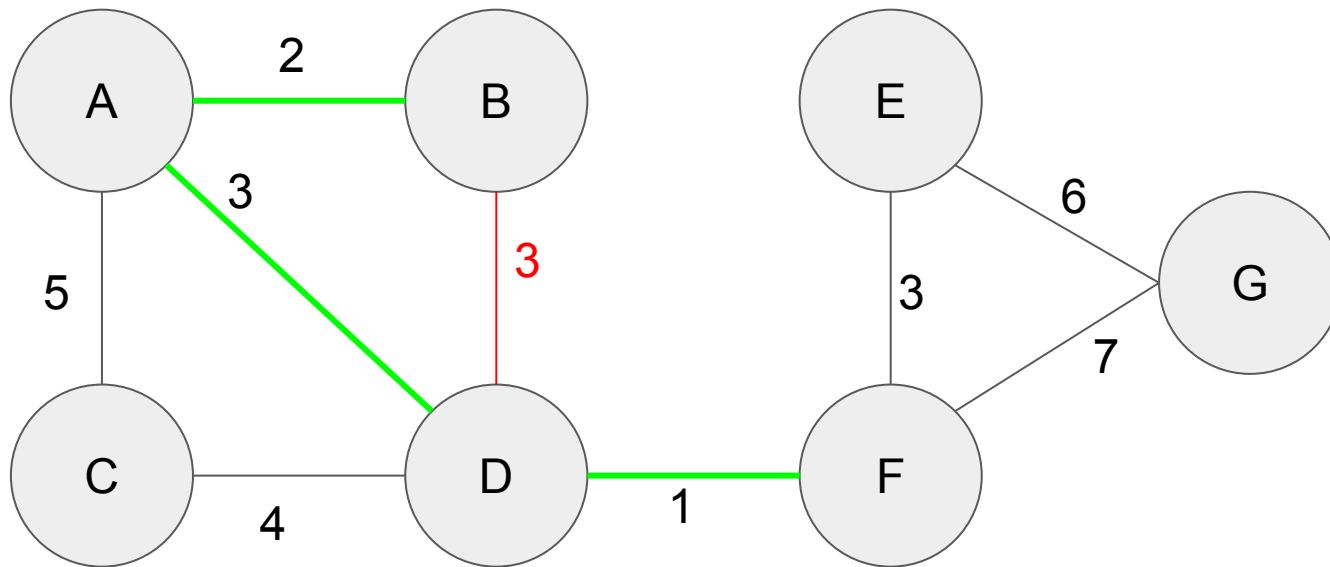
# Kruskal's algorithm - weight 2



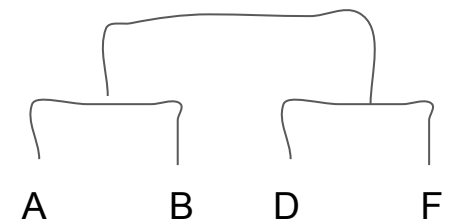
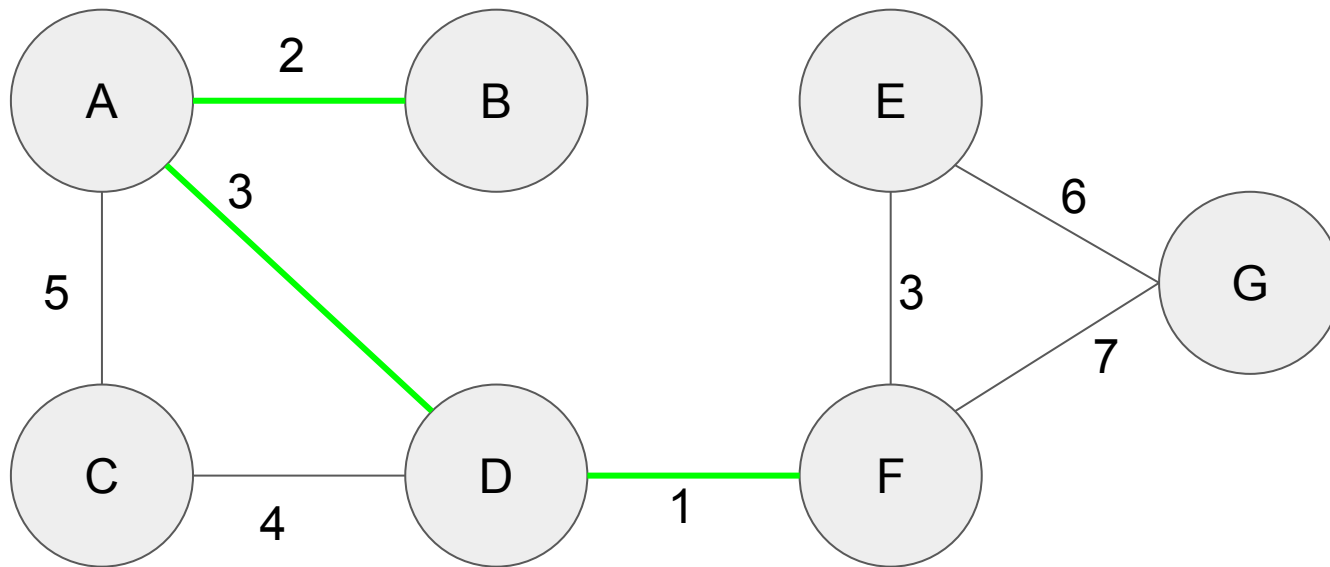
# Kruskal's algorithm - weight 3



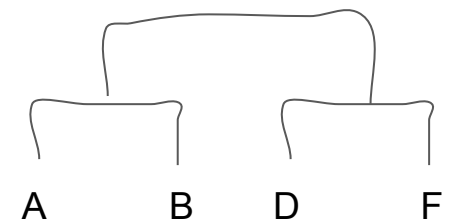
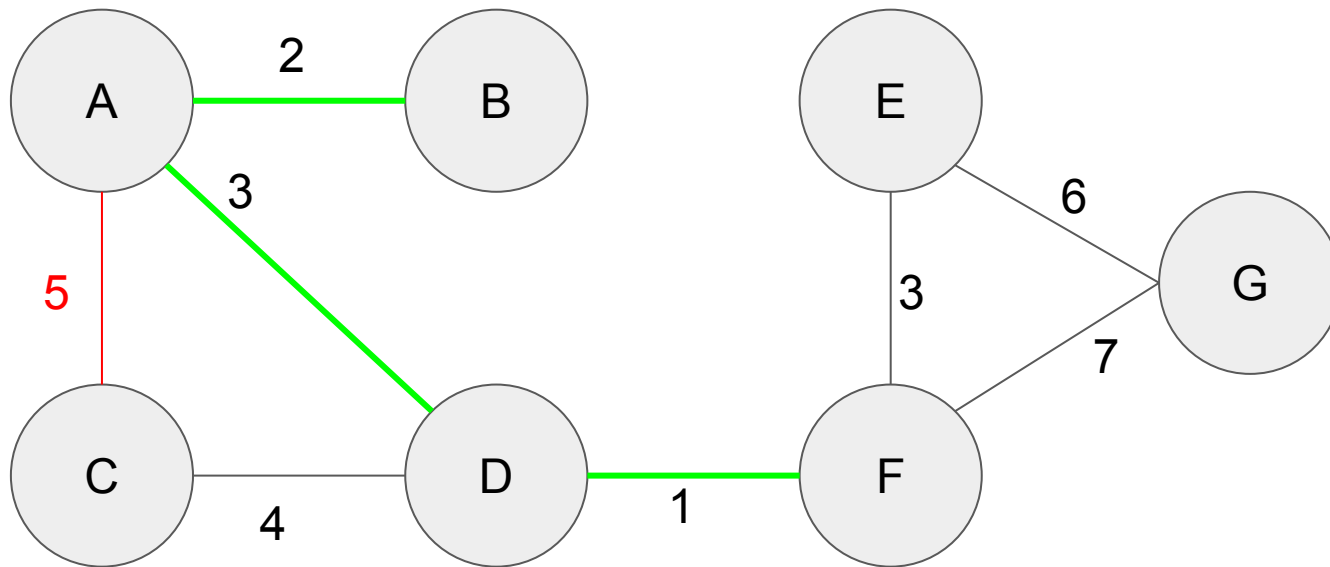
# Kruskal's algorithm - weight 3



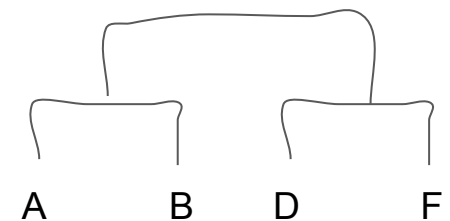
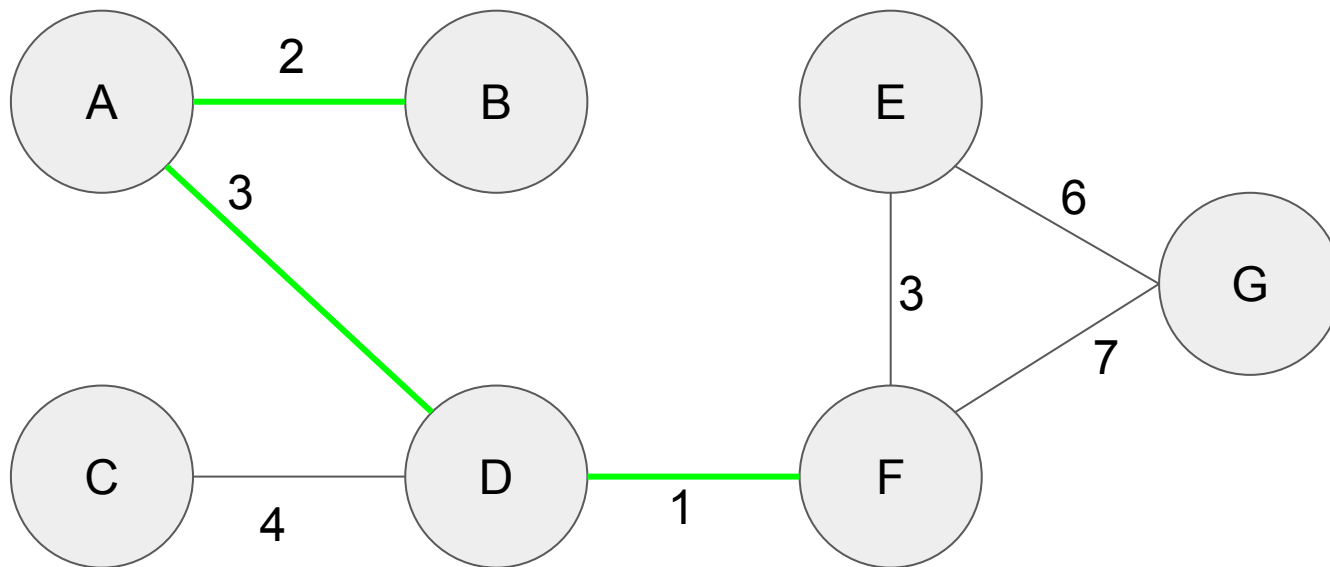
# Kruskal's algorithm - weight 3



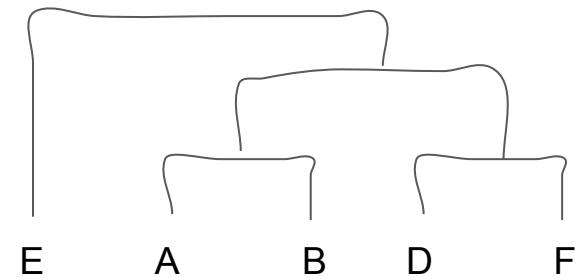
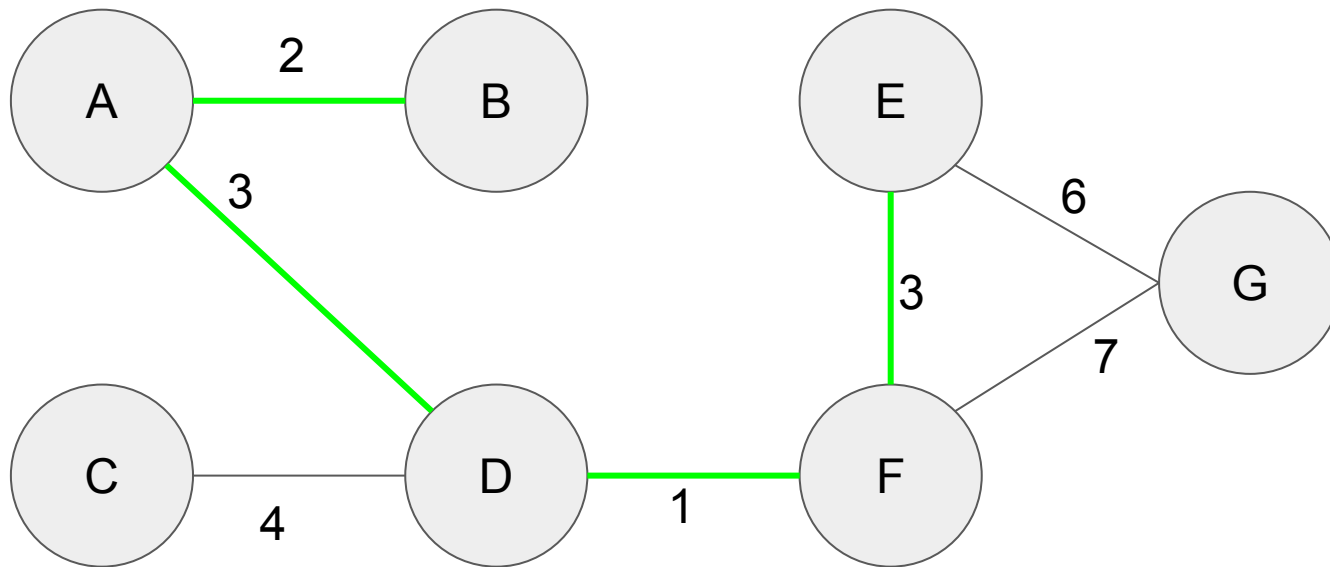
# Kruskal's algorithm - weight 3



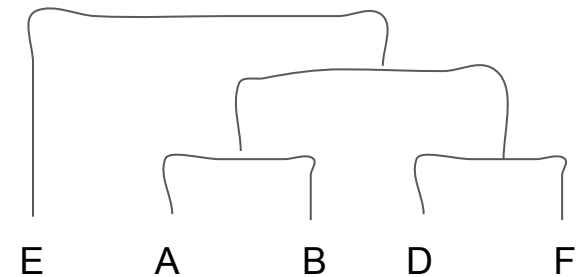
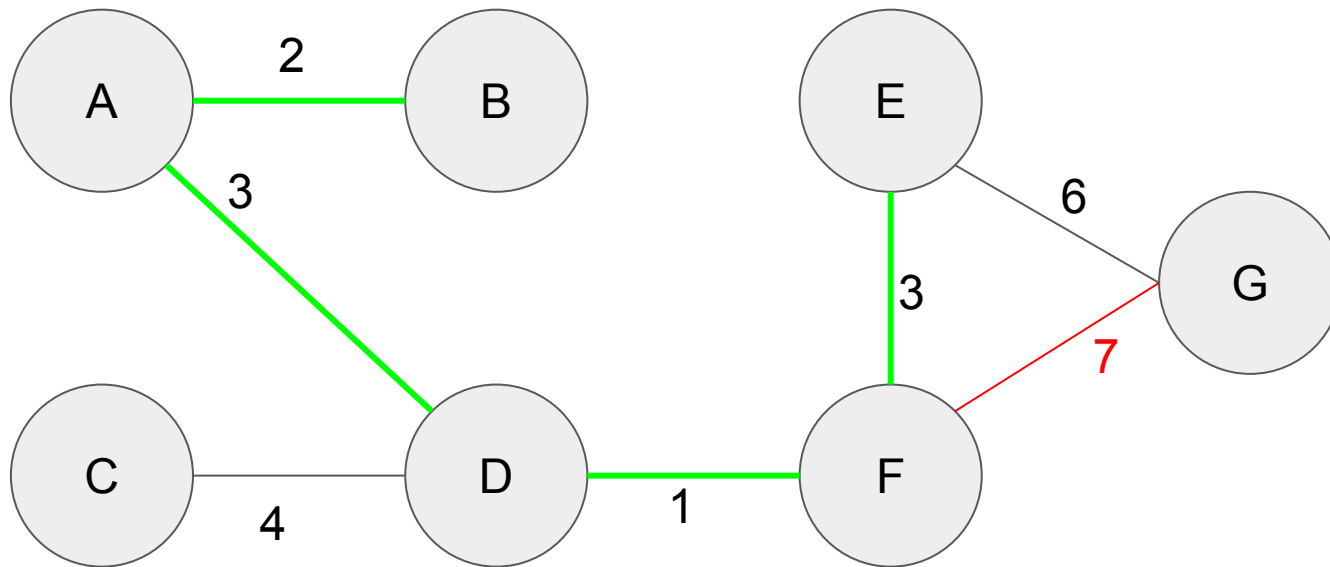
# Kruskal's algorithm - weight 3



# Kruskal's algorithm - weight 3

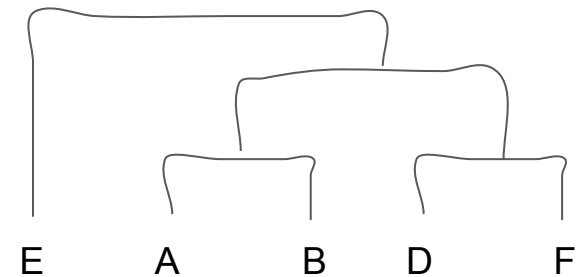
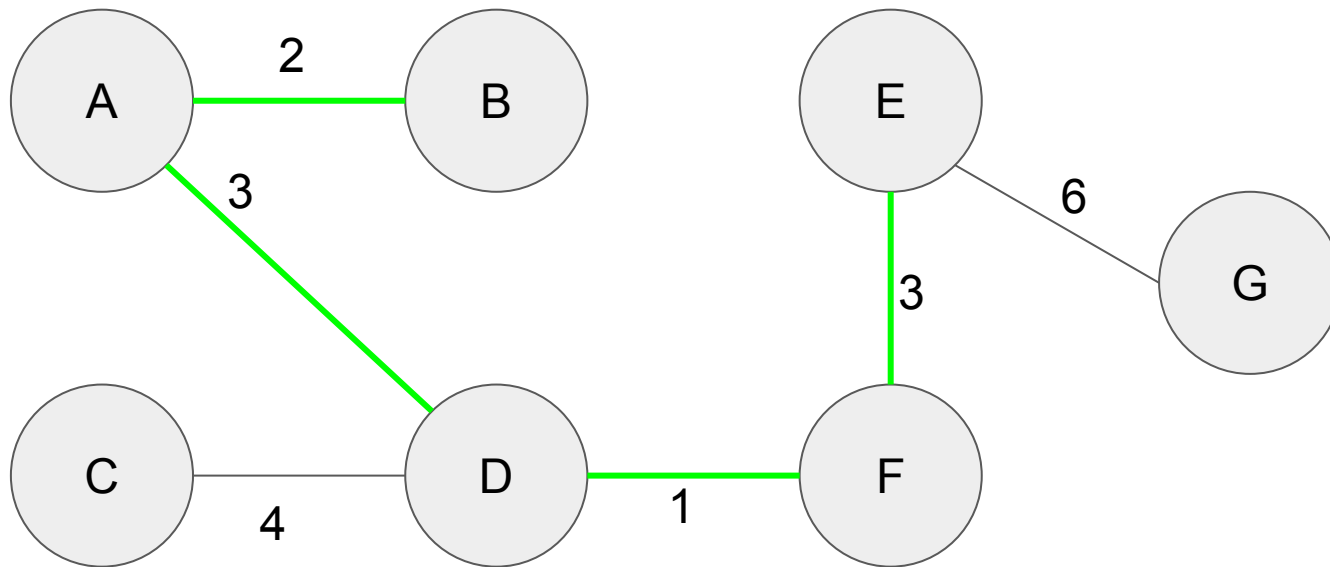


# Kruskal's algorithm - weight 3

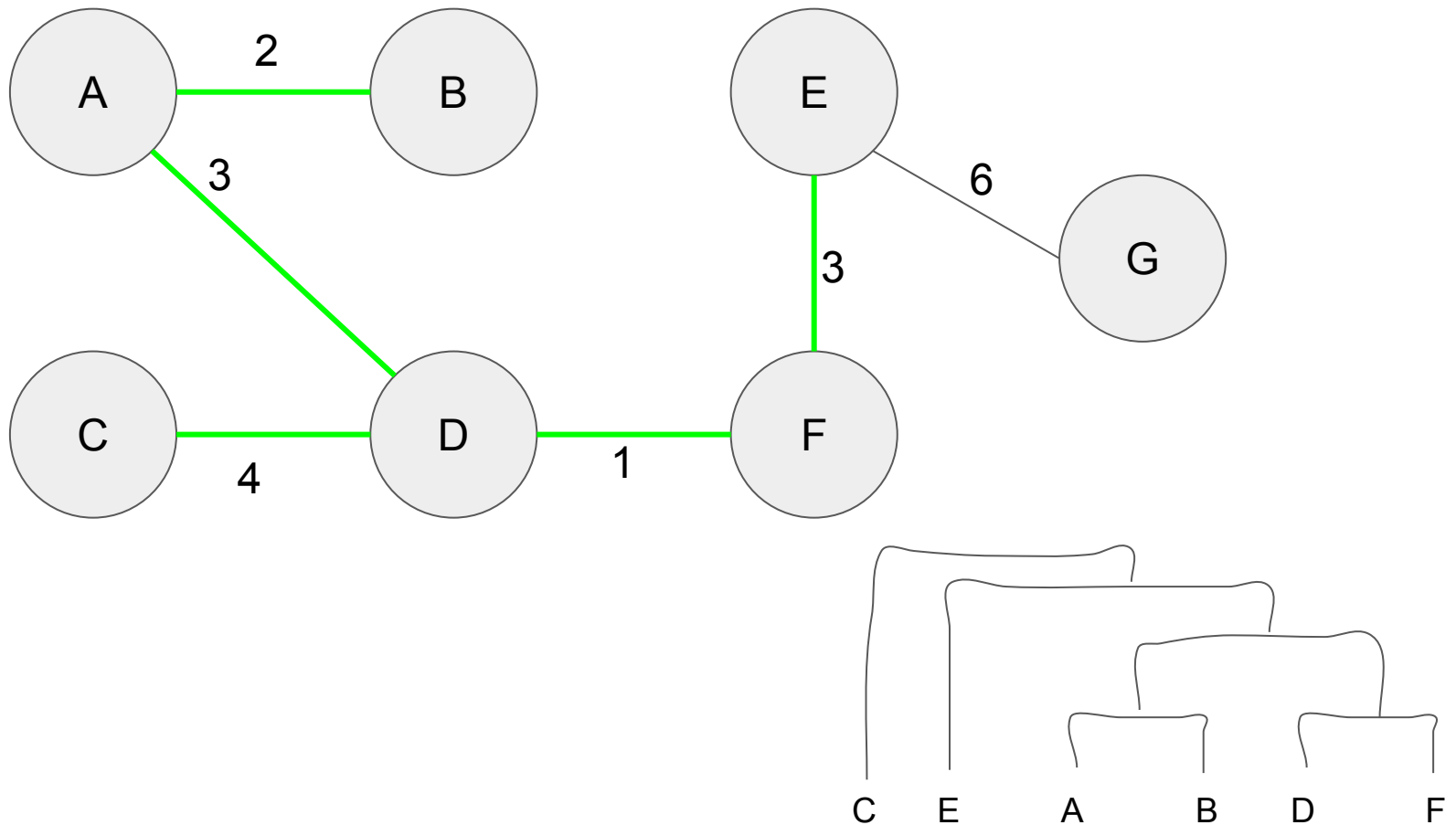




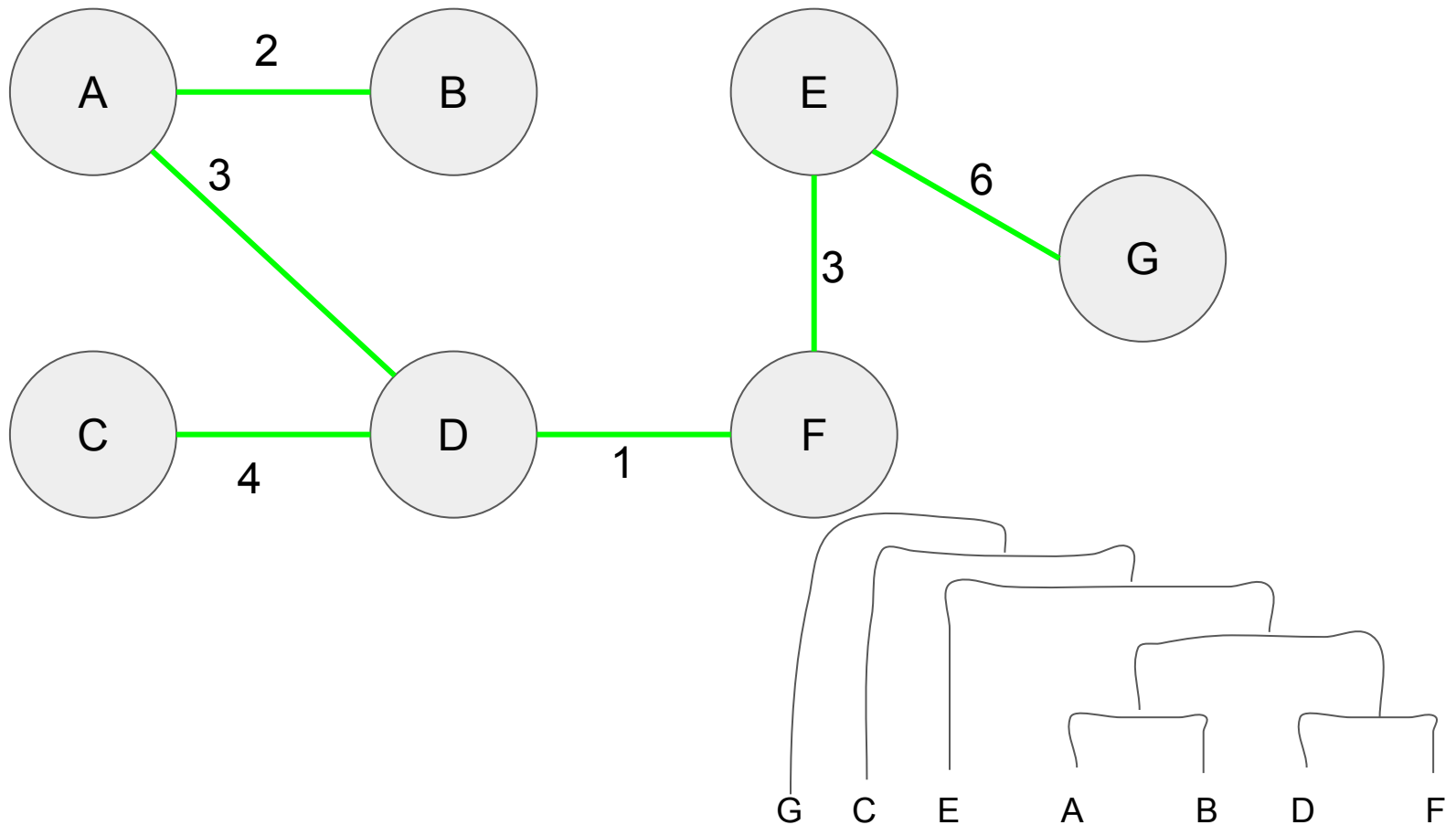
# Kruskal's algorithm - weight 3



# Kruskal's algorithm - weight 4



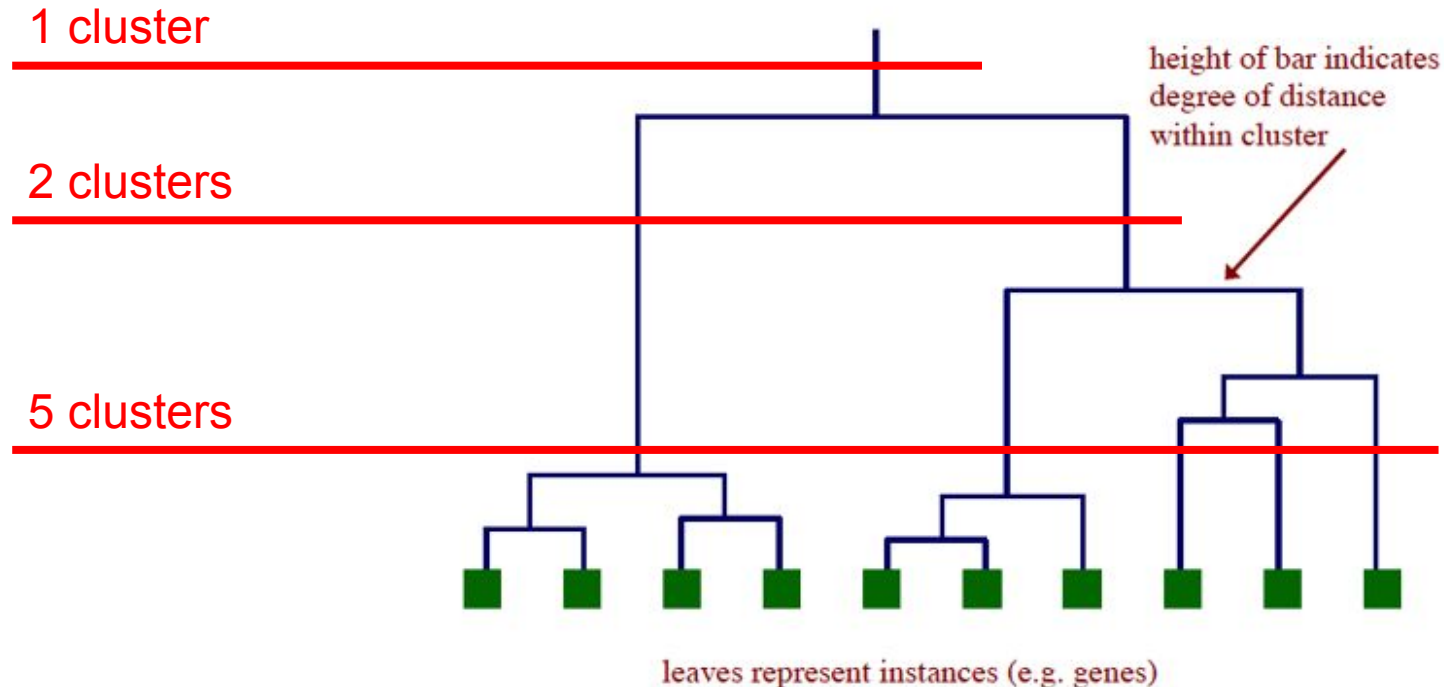
# Kruskal's algorithm - weight 6



# Dendrogram

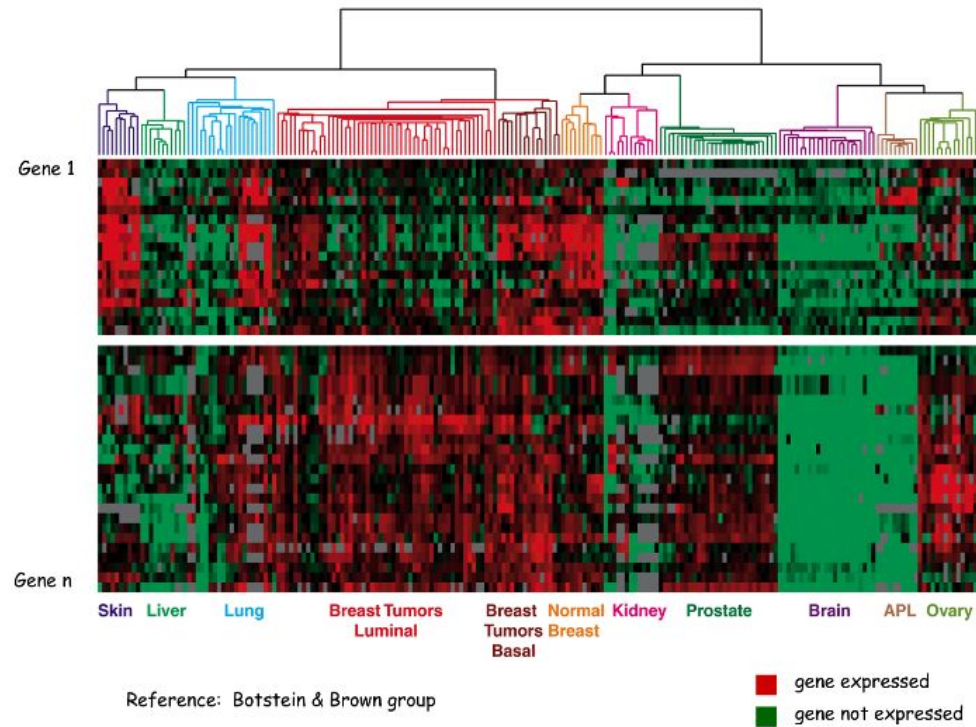
**Dendrogram.** Scientific visualization of hypothetical sequence of evolutionary events.

- Leaves = genes.
- Internal nodes = hypothetical ancestors.

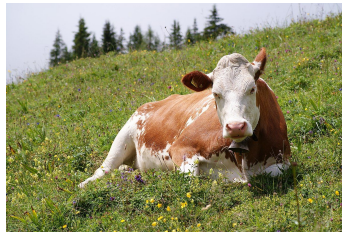


# Applications in data analysis

Tumors in similar tissues cluster together.



# Applications: taxonomy



## Morphological features

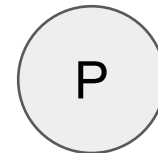
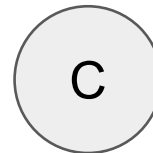
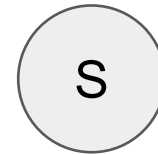
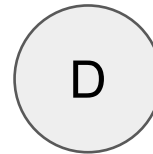
Dolphin (D): aquatic, eats fish

Seal (S): semi-aquatic, eats fish

Cow (C): terrestrial, eats grass

Panda (P): terrestrial, eats grass

Morphological MST:



# Applications: taxonomy



## Morphological features

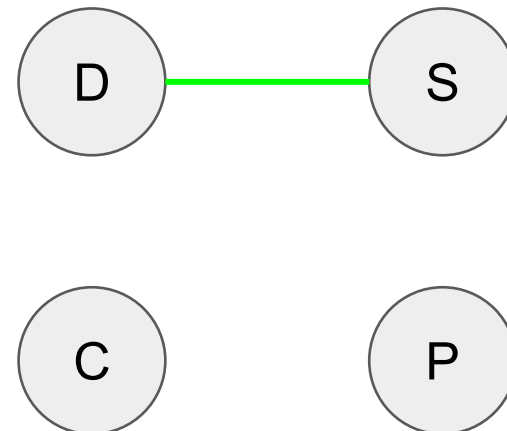
Dolphin (D): aquatic, eats fish

Seal (S): semi-aquatic, eats fish

Cow (C): terrestrial, eats grass

Panda (P): terrestrial, eats grass

Morphological MST:



# Applications: taxonomy



## Morphological features

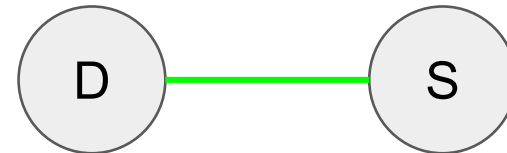
Dolphin (D): aquatic, eats fish

Seal (S): semi-aquatic, eats fish

Cow (C): terrestrial, eats grass

Panda (P): terrestrial, eats grass

Morphological MST:





# Applications: taxonomy

## Morphological features

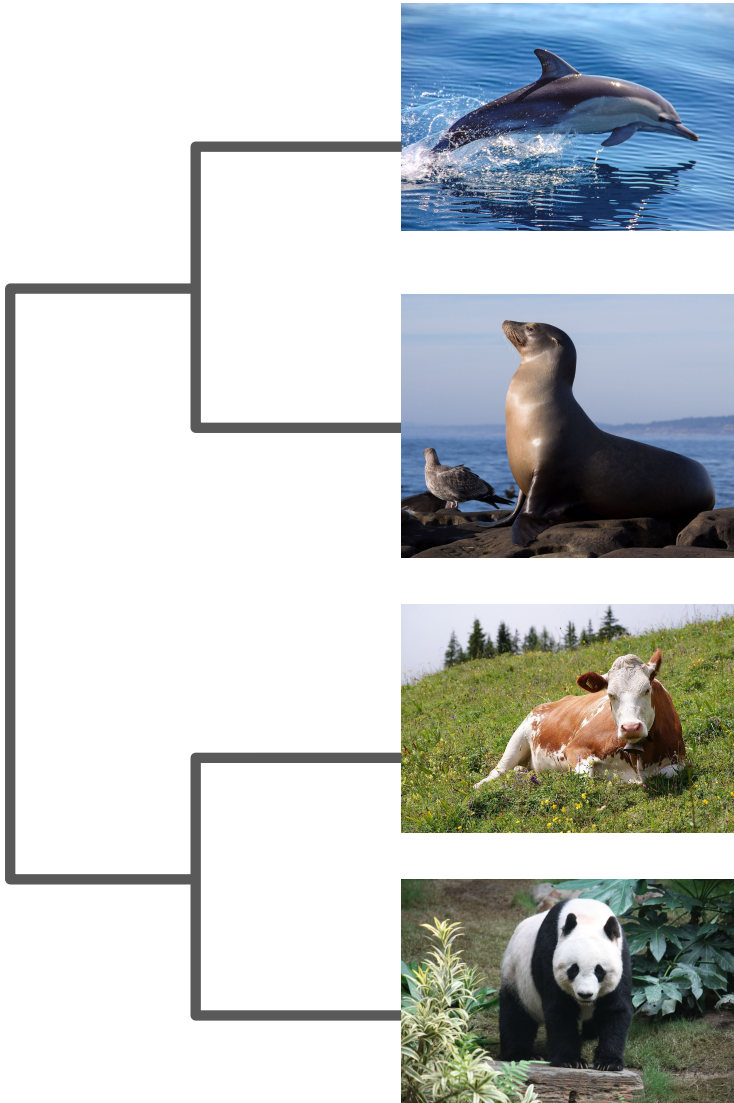
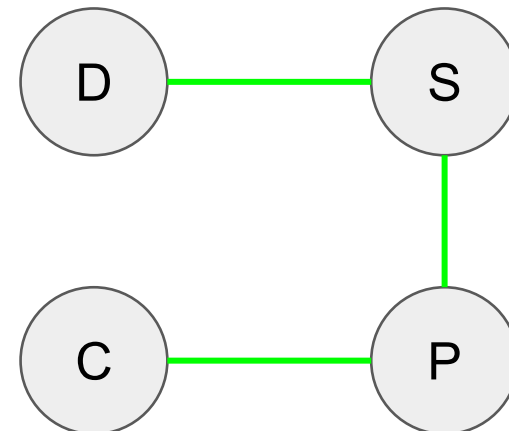
Dolphin (D): aquatic, eats fish

Seal (S): semi-aquatic, eats fish

Cow (C): terrestrial, eats grass

Panda (P): terrestrial, eats grass

Morphological MST:



# Applications: taxonomy

## Morphological features

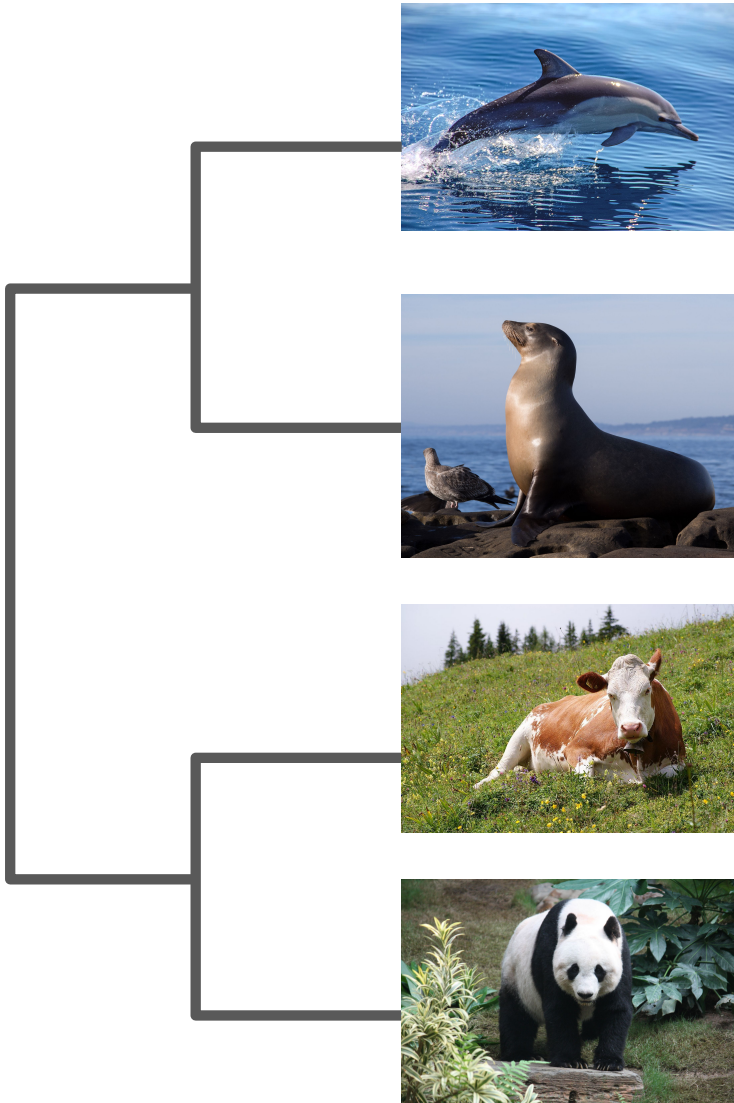
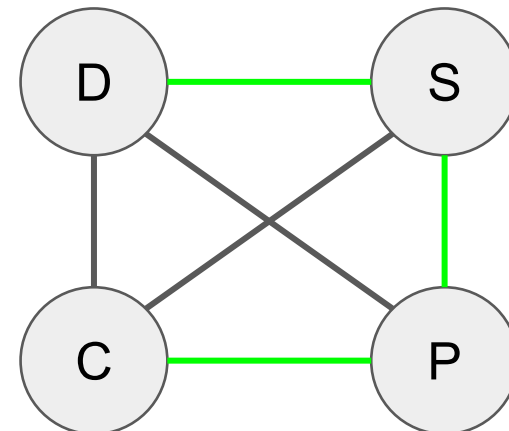
Dolphin (D): aquatic, eats fish

Seal (S): semi-aquatic, eats fish

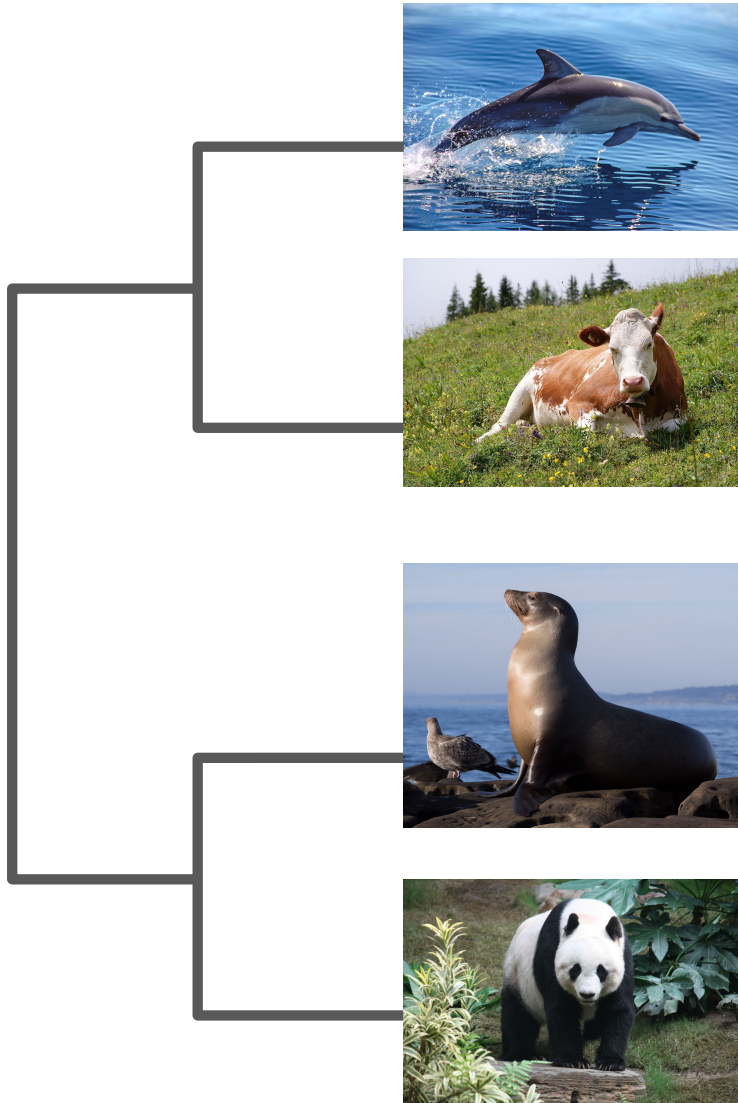
Cow (C): terrestrial, eats grass

Panda (P): terrestrial, eats grass

Morphological MST:



# Applications: taxonomy



## Morphological features

Dolphin (D): aquatic, eats fish

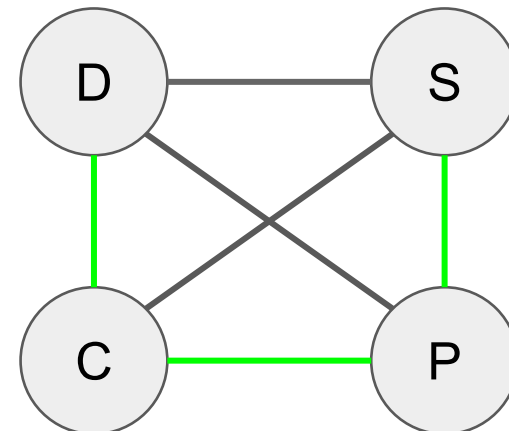
Seal (S): semi-aquatic, eats fish

Cow (C): terrestrial, eats grass

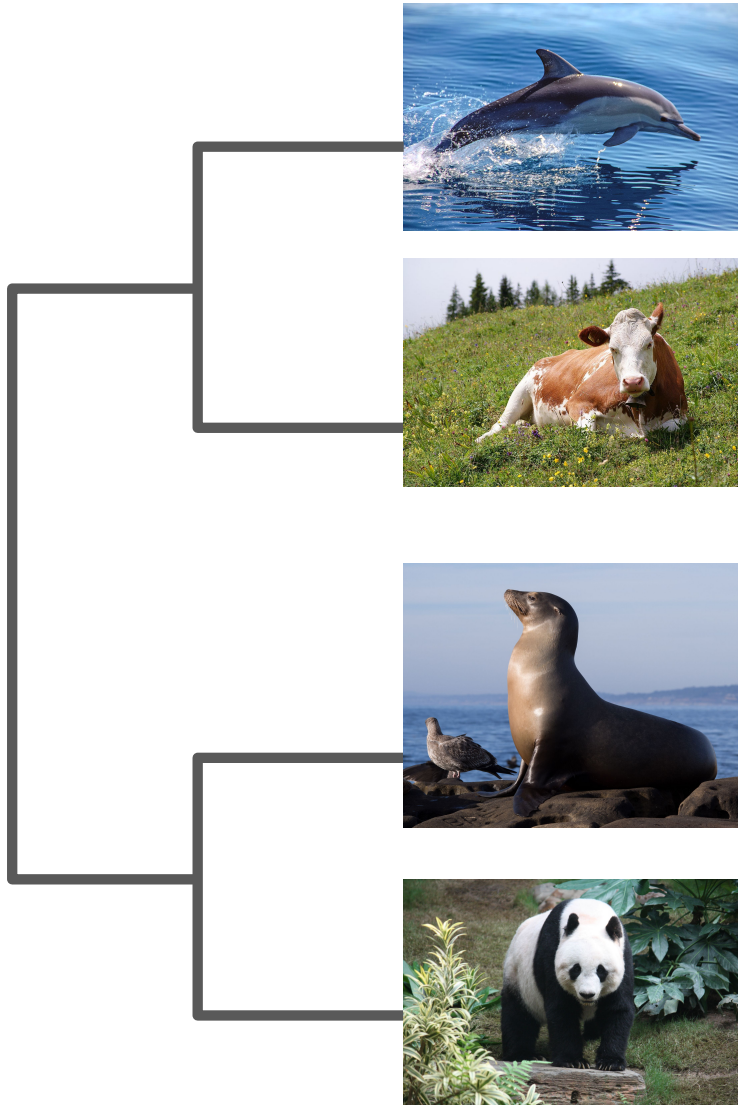
Panda (P): terrestrial, eats grass

## Genetic MST

(the actual evolutionary path):



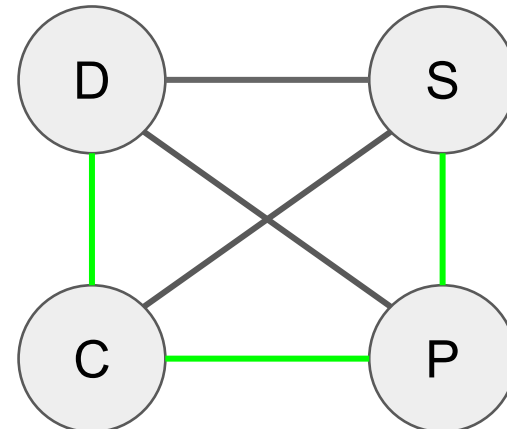
# Applications: taxonomy



Computing genetic distances  
will be discussed in the  
dynamic programming section

Genetic MST

(the actual evolutionary path):



# Huffman coding

# Huffman encoding

An encoding scheme used in, e.g., MP3 encoding

**Data:** a string  $S$  of symbols over an alphabet  $\Gamma$

**Goal:** find a binary encoding  $e$  of  $\Gamma$  resulting in minimum encoded length of  $S$

Denote the encoded string by  $S_e$

	$a$	01100001
Example: ASCII encoding	$b$	01100010
	$\vdots$	$\vdots$

# Different encodings

Consider  $\Gamma = \{a, b, c\}$

Stats on  $S$ :  $a$  appears 45 times,  $b$  16 times, and  $c$  twice

# Different encodings

Consider  $\Gamma = \{a, b, c\}$

Stats on  $S$ :  $a$  appears 45 times,  $b$  16 times, and  $c$  twice

- Fixed-length encoding

$a \rightarrow 00$

$e_1 : b \rightarrow 01 \quad |S_{e_1}| = 45 \times 2 + 16 \times 2 + 2 \times 2 = 126$

$c \rightarrow 10$



# Different encodings

Consider  $\Gamma = \{a, b, c\}$

Stats on  $S$ :  $a$  appears 45 times,  $b$  16 times, and  $c$  twice

- Fixed-length encoding

$a \rightarrow 00$

$e_1 : b \rightarrow 01 \quad |S_{e_1}| = 45 \times 2 + 16 \times 2 + 2 \times 2 = 126$

$c \rightarrow 10$

- Variable-length encoding

$a \rightarrow 0$

$e_2 : b \rightarrow 10 \quad |S_{e_2}| = 45 \times 1 + 16 \times 2 + 2 \times 2 = 81$

$c \rightarrow 11$

# Different encodings

Consider  $\Gamma = \{a, b, c\}$

Stats on  $S$ :  $a$  appears 45 times,  $b$  16 times, and  $c$  twice

- Fixed-length encoding

$$a \rightarrow 00$$

$$e_1 : b \rightarrow 01 \quad |S_{e_1}| = 45 \times 2 + 16 \times 2 + 2 \times 2 = 126$$

$$c \rightarrow 10$$

- Variable-length encoding

$$a \rightarrow 0$$

$$e_2 : b \rightarrow 10 \quad |S_{e_2}| = 45 \times 1 + 16 \times 2 + 2 \times 2 = 81$$

$$c \rightarrow 11$$

$$a \rightarrow 0$$

- Be careful!  $e_2 : b \rightarrow 1$       Decoding will lead to ambiguity

$$c \rightarrow 01$$

# Prefix-free encoding

$a \rightarrow 0$

Consider the bad encoding  $e_2$  :  $b \rightarrow 1$

$c \rightarrow 01$

*ababba?*, *ccba?*, *abcba?*, or ...?

How to decode 010110?

# Prefix-free encoding

$$a \rightarrow 0$$

Consider the bad encoding  $e_2$  :  $b \rightarrow 1$       How to decode 010110?

$$c \rightarrow 01$$

*ababba?*, *ccba?*, *abcba?*, or ...?

To avoid ambiguity, we need the encoding to be **prefix-free**

## Definition

An encoding is **prefix-free** if no codeword is a prefix of any other codewords

# Tree representation of a prefix-free encoding

## Definition

A **full binary tree** is a binary tree where each node is either a leaf or it has two children

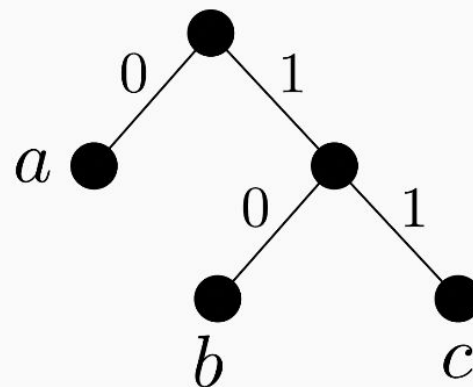
# Tree representation of a prefix-free encoding

## Definition

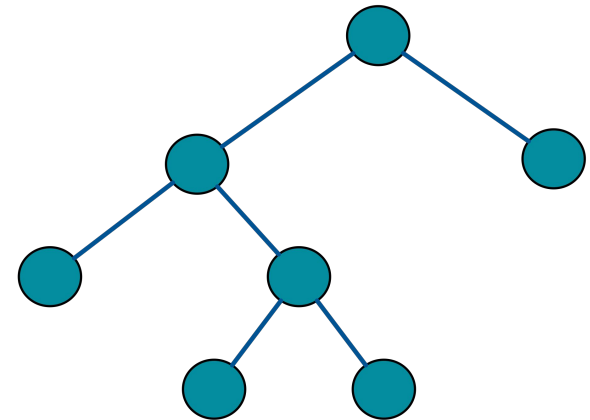
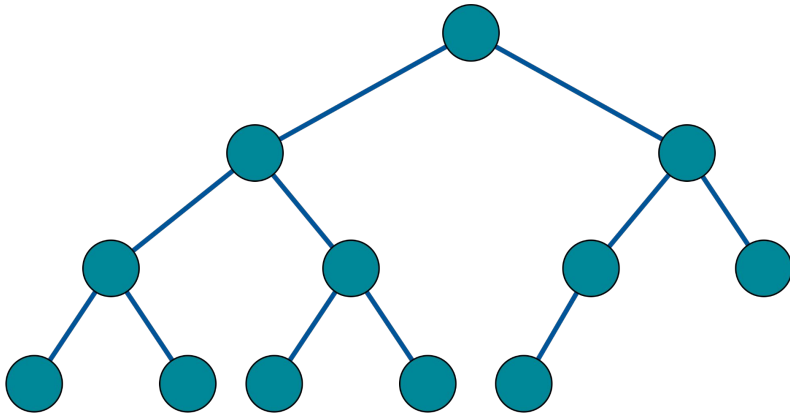
A **full binary tree** is a binary tree where each node is either a leaf or it has two children

We use a full binary tree to represent a prefix-free encoding

- leaves are corresponding to symbols in  $\Gamma$
- label edge to the left child with 0
- label edge to the right child with 1



**Is it a full binary tree?**



# Tree representation of a prefix-free encoding

## Definition

A **full binary tree** is a binary tree where each node is either a leaf or it has two children

We use a full binary tree to represent a prefix-free encoding

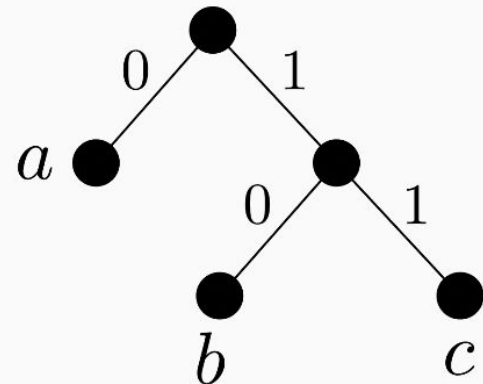
- leaves are corresponding to symbols in  $\Gamma$
- label edge to the left child with 0
- label edge to the right child with 1

To obtain the encoding, read edge labels from root to a symbol

$a \rightarrow 0$ ,  $b \rightarrow 10$ ,  $c \rightarrow 11$ ,

Depth of a leaf  $\equiv$  length of its codeword

It guarantees to be prefix-free





# Tree representation of a prefix-free encoding

Let  $e$  be an encoding represented by a tree

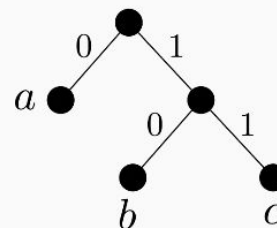
For string  $S$ , let  $f_v$  be the symbol count in  $S$  for each  $v \in \Gamma$

# Tree representation of a prefix-free encoding

Let  $e$  be an encoding represented by a tree

For string  $S$ , let  $f_v$  be the symbol count in  $S$  for each  $v \in \Gamma$

$$|S_e| = \sum_{v \in \Gamma} f_v \cdot \text{depth}(v)$$



$$|S_e| = 45 \times 1 + 16 \times 2 + 2 \times 2 = 81$$

Consider  $\Gamma = \{a, b, c\}$

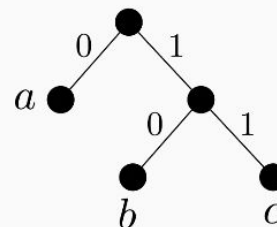
Stats on  $S$ :  $a$  appears 45 times,  $b$  16 times, and  $c$  twice

# Tree representation of a prefix-free encoding

Let  $e$  be an encoding represented by a tree

For string  $S$ , let  $f_v$  be the symbol count in  $S$  for each  $v \in \Gamma$

$$|S_e| = \sum_{v \in \Gamma} f_v \cdot \text{depth}(v)$$



$$|S_e| = 45 \times 1 + 16 \times 2 + 2 \times 2 = 81$$

A useful re-write: label internal nodes with counts of descendants

Consider  $\Gamma = \{a, b, c\}$

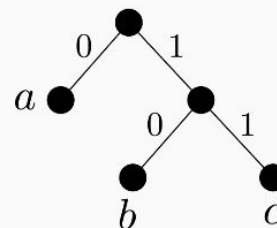
Stats on  $S$ :  $a$  appears 45 times,  $b$  16 times, and  $c$  twice

# Tree representation of a prefix-free encoding

Let  $e$  be an encoding represented by a tree

For string  $S$ , let  $f_v$  be the symbol count in  $S$  for each  $v \in \Gamma$

$$|S_e| = \sum_{v \in \Gamma} f_v \cdot \text{depth}(v)$$



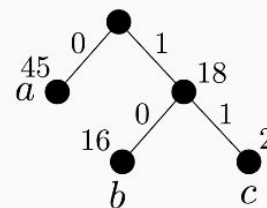
$$|S_e| = 45 \times 1 + 16 \times 2 + 2 \times 2 = 81$$

A useful re-write: label internal nodes with counts of descendants

For all non-root node  $v$ , define

$\text{cost}(v) :=$  sum of leaf node counts descending from  $v$

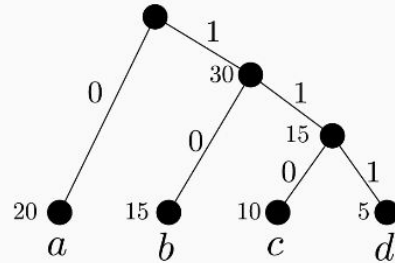
$$|S_e| = \sum_{v \in T - \{\text{root}\}} \text{cost}(v)$$



$$|S_e| = 45 + 16 + 2 + 18 = 81$$

# Constructing the prefix-free encoding tree: examples

- $a : 20, b : 15, c : 10, d : 5$



$a \rightarrow 0$

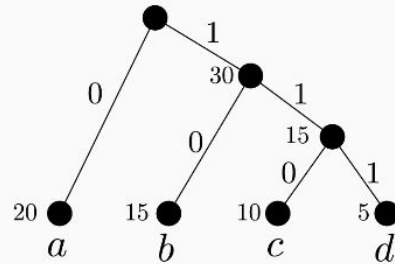
$b \rightarrow 10$

$c \rightarrow 110$

$d \rightarrow 111$

# Constructing the prefix-free encoding tree: examples

- $a : 20, b : 15, c : 10, d : 5$



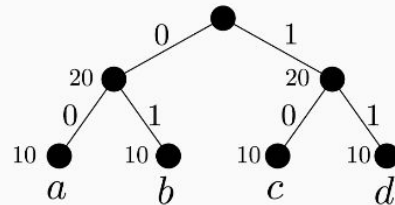
$a \rightarrow 0$

$b \rightarrow 10$

$c \rightarrow 110$

$d \rightarrow 111$

- $a : 10, b : 10, c : 10, d : 10$



$a \rightarrow 00$

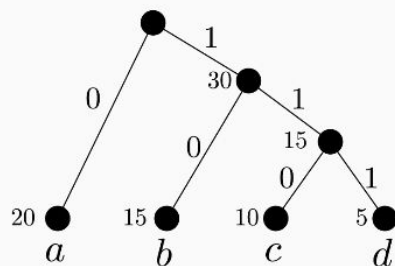
$b \rightarrow 01$

$c \rightarrow 10$

$d \rightarrow 11$

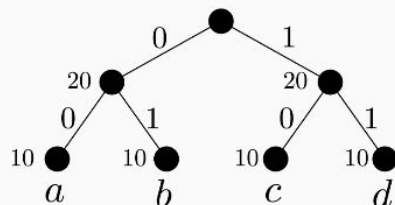
# Constructing the prefix-free encoding tree: examples

- $a : 20, b : 15, c : 10, d : 5$



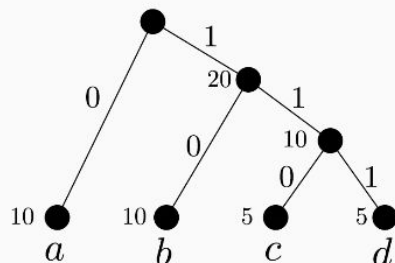
$a \rightarrow 0$   
 $b \rightarrow 10$   
 $c \rightarrow 110$   
 $d \rightarrow 111$

- $a : 10, b : 10, c : 10, d : 10$



$a \rightarrow 00$   
 $b \rightarrow 01$   
 $c \rightarrow 10$   
 $d \rightarrow 11$

- $a : 10, b : 10, c : 5, d : 5$



$a \rightarrow 0$   
 $b \rightarrow 10$   
 $c \rightarrow 110$   
 $d \rightarrow 111$

What are the total costs of these encodings?