

Foundations of Data Mining

Homework Assignment 1C

Frequent Subsequence and Subgraph Mining

Homework group 4:

Árni Ólafur Einarsson

Stefán G. Jónsson

Gudmundur Ómi Pálsson

Nick Geerjens

Troy Maasland

Frequent Subsequence Mining

We impose a minimal support constraint of

$$s_{\min} = 0.75$$

Homework 2C

Data:

Sequence ID

1

Sequence
 $\langle a(abc)(ac)d(cf) \rangle$

2

$\langle (ad)c(bc)(ae) \rangle$

3

$\langle (ef)(ab)(df)cb \rangle$

4

$\langle eg(ab)cbc \rangle$

Minimal support: $s_{min} = 0.75$

GSP

→ Scan dataset once
 find all frequent items

Repeat until done:

→ feed seed set of
 sequential patterns to join
 operator, which determines
 new set of candidates

→ Scan dataset once to determine
 support of new candidate set
 remove those that do not satisfy
 minimal support constraint.

Resulting set: seed for next
 iteration

Candidate	Support
$\langle a \rangle$	1
$\langle b \rangle$	1
$\langle c \rangle$	1
$\langle d \rangle$	0.75
$\langle e \rangle$	0.75
$\langle f \rangle$	0.5
$\langle g \rangle$	0.25

frequent
 length-1
 sequences

→ candidate
 sequences

B. All combinations of the frequent length-1 candidate subsequences are length-2 candidate subsequences

Also all combinations of the frequent length-1 candidate subsequences in a single event are length-2 candidate subsequences (for which the order does not matter)

If N is the amount of frequent length-1 candidate subsequences, then this leaves us with ~~N^2~~ N^2 combinations

$N(N-1)/2$ evens

$\rightarrow N^2 + N(N-1)/2$ length-2 candidate subsequences

length-2

Candidate subsequences:

$\langle a \rangle$	$\langle b \rangle$	$\langle c \rangle$	$\langle d \rangle$	$\langle e \rangle$
$\langle a \rangle$	$\langle (aa) \rangle$	$\langle ab \rangle$	$\langle ac \rangle$	$\langle ad \rangle$
$\langle b \rangle$	$\langle ba \rangle$	$\langle bb \rangle$	$\langle bc \rangle$	$\langle bd \rangle$
$\langle c \rangle$	$\langle ca \rangle$	$\langle cb \rangle$	$\langle cc \rangle$	$\langle cd \rangle$
$\langle d \rangle$	$\langle da \rangle$	$\langle db \rangle$	$\langle dc \rangle$	$\langle dd \rangle$
$\langle e \rangle$	$\langle ea \rangle$	$\langle eb \rangle$	$\langle ec \rangle$	$\langle ed \rangle$
$\langle fa \rangle$				$\langle ee \rangle$

$\langle a \rangle$	$\langle b \rangle$	$\langle c \rangle$	$\langle d \rangle$	$\langle e \rangle$
$\langle a \rangle$		$\langle (ab) \rangle$	$\langle (ac) \rangle$	$\langle (ad) \rangle$
$\langle b \rangle$			$\langle (bc) \rangle$	$\langle (bd) \rangle$
$\langle c \rangle$				$\langle (ce) \rangle$
$\langle d \rangle$				$\langle (de) \rangle$
$\langle e \rangle$				

Now for all candidates we have to check the support.

Here, only the candidates with enough support (0.75) are written down: (so that the following are the most frequent).

$\langle ab \rangle: s = 0.75$ $\langle ac \rangle: s = 1$ $\langle cc \rangle: s = 0.75$

$\langle cb \rangle: s = 0.75$ $\langle bc \rangle: s = 0.75$ ~~$\langle ad \rangle: s = 0.75$~~

$\langle (ab) \rangle: s = 0.75$

$\langle dc \rangle: s = 0.75$

$$7. N_f = 5, \quad N_{2c} = N(N-1)/2 + N^2 \\ = (5 \cdot 4)/2 + 5^2 \\ = 35 \text{ length-2 candidate subsequences}$$

We would have generated $(7 \cdot 6)/2 + 7^2 = 70$ length-2 candidate subsequences if we did not prune away infrequent length-1 candidate subsequences

S $\langle ab \rangle, \langle cb \rangle, \langle ac \rangle, \langle bc \rangle, \langle cc \rangle, \langle dc \rangle, \langle (ab) \rangle$

Join operator

Given a seed set of sequential patterns
the join operator

	$\langle ab \rangle$	$\langle cb \rangle$	$\langle ac \rangle$	$\langle bc \rangle$	$\langle cc \rangle$	$\langle dc \rangle$	$\langle (ab) \rangle$
$\langle ab \rangle$	X	X	X	$\langle abc \rangle$	X	X	X
$\langle cb \rangle$	X	X	X	$\langle cbc \rangle$	X	X	X
$\langle ac \rangle$	X	$\langle acb \rangle$	X	X	$\langle acc \rangle$	X	X
$\langle bc \rangle$	X	$\langle bcb \rangle$	X	X	$\langle bcc \rangle$	X	X
$\langle cc \rangle$	X	$\langle ccb \rangle$	X	X	$\langle ccc \rangle$	X	X
$\langle dc \rangle$	X	$\langle dc b \rangle$	X	X	$\langle dcc \rangle$	X	X
$\langle (ab) \rangle$	X	X	X	$\langle (abc) \rangle$	X	X	X

So, the length-3 candidate sequences are

$\langle acb \rangle, \langle bcb \rangle, \langle ccb \rangle, \langle dc b \rangle, \langle abc \rangle, \langle cbc \rangle,$
 $\langle (ab)c \rangle, \langle acc \rangle, \langle bcc \rangle, \langle ccc \rangle, \langle dcc \rangle$

E.

E.

~~s is closed within S~~

A sequential pattern s is closed if no strict superpattern s' exists such that s' and s have the same support

$\langle acb \rangle$: support: ~~s~~: 0.75

~~Superpatterns~~

Superpatterns: $\langle acba \rangle, \langle acbb \rangle, \langle acbc \rangle, \langle acbd \rangle,$
 $\langle acbe \rangle, \langle acbf \rangle, \langle acbg \rangle$

b, g support < 0.75, so those are left out

$\langle acba \rangle$: support 0.25

$\langle acbb \rangle$: support 0

$\langle acbc \rangle$: support 0.25

$\langle acbd \rangle$: support 0

$\langle acbe \rangle$: support 0.25

So, yes, $\langle acb \rangle$ is closed within S

$\langle adc \rangle$: support 0.5

$\langle adca \rangle$: $s = 0$

$\langle adcb \rangle$: $s = 0.25$ ~~(w)~~

$\langle addc \rangle$: $s =$

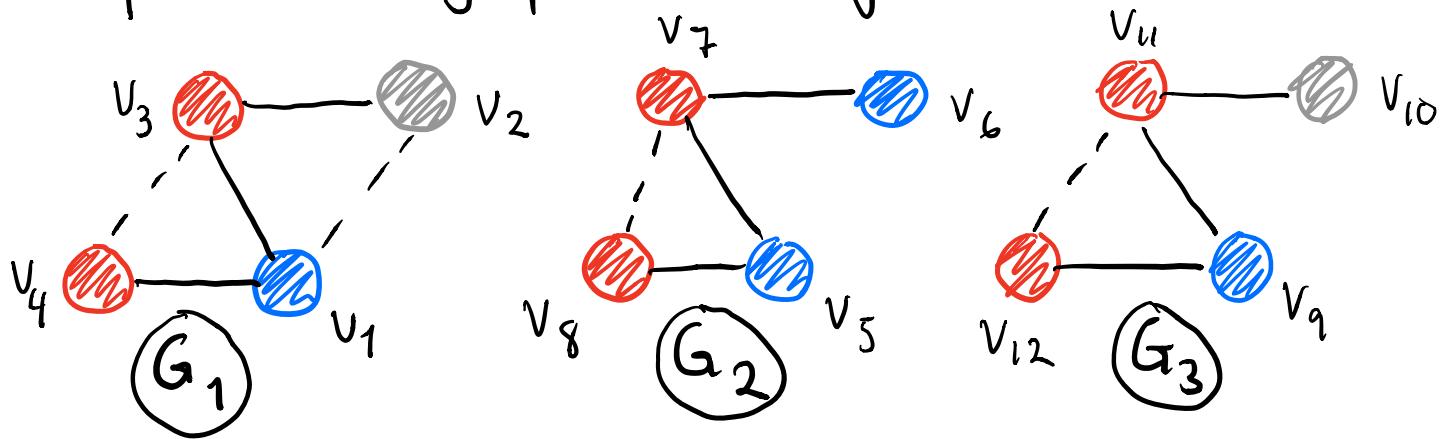
The rest does not have to be checked. Only

two sequences contain $\langle adc \rangle$ within S

The first has only an f after the c, within the same event (which does not count)

→ support is at maximum 0.25 of

Frequency Subgraph Mining



Q. Characterize each of the graphs G_i by writing down their vertex set V_i , their edge set E_i , their label set L_i , and their label function $l_i : V_i \cup E_i \rightarrow L_i$. For their edge sets, you may choose yourself if you want to write them down by the adjacency matrix or the adjacency list

Solution: We get

$$G_1 = (V(G_1), E(G_1))$$

$$V(G_1) = \{v_1, v_2, v_3, v_4\}$$

$$E(G_1) = \{E_1 = (v_1, v_2), E_2 = (v_1, v_3), E_3 = (v_2, v_3),$$

$$E_4 = (v_1, v_4), E_5 = \{v_3, v_4\}\}$$

$$\ell(v_1) = \text{BLUE} \quad \ell(E_1) = \text{dotted}$$

$$\ell(v_2) = \text{GRAY} \quad \ell(E_2) = \text{straight}$$

$$\ell(v_3) = \text{RED} \quad \ell(E_3) = \text{straight}$$

$$\ell(v_4) = \text{RED} \quad \ell(E_4) = \text{straight}$$

$$\ell(E_5) = \text{dotted}$$

$$\mathcal{L}_1 = \{\text{BLUE}, \text{RED}, \text{GRAY}, \text{dotted}, \text{straight}\}$$

$$G_2 = (V(G_2), E(G_2))$$

$$V(G_2) = \{v_5, v_6, v_7, v_8\}$$

$$E(G_2) = \{E_1 = (v_5, v_8), E_2 = (v_5, v_7), E_3 = (v_6, v_7),$$

$$E_4 = (v_7, v_8)\}$$

$$\ell(v_5) = \text{BLUE} \quad \ell(E_1) = \text{straight}$$

$$\ell(v_6) = \text{BLUE} \quad \ell(E_2) = \text{straight}$$

$$\ell(v_7) = \text{RED} \quad \ell(E_3) = \text{straight}$$

$$\ell(v_8) = \text{RED} \quad \ell(E_4) = \text{dotted}$$

$$\mathcal{L}_2 = \{\text{BLUE}, \text{RED}, \text{dotted}, \text{straight}\}$$

$$G_3 = (V(G_3), E(G_3))$$

$$V(G_3) = \{v_9, v_{10}, v_{11}, v_{12}\}$$

$$E(G_3) = \{E_1 = (v_9, v_{12}), E_2 = (v_9, v_{11}), E_3 = (v_{10}, v_{11}),$$

$$E_4 = (v_{11}, v_{12})\}$$

$$\ell(v_9) = \text{BLUE} \quad \ell(E_1) = \text{straight}$$

$$\ell(v_{10}) = \text{GRAY} \quad \ell(E_2) = \text{straight}$$

$$\ell(v_{11}) = \text{RED} \quad \ell(E_3) = \text{straight}$$

$$\ell(v_{12}) = \text{RED} \quad \ell(E_4) = \text{dotted}$$

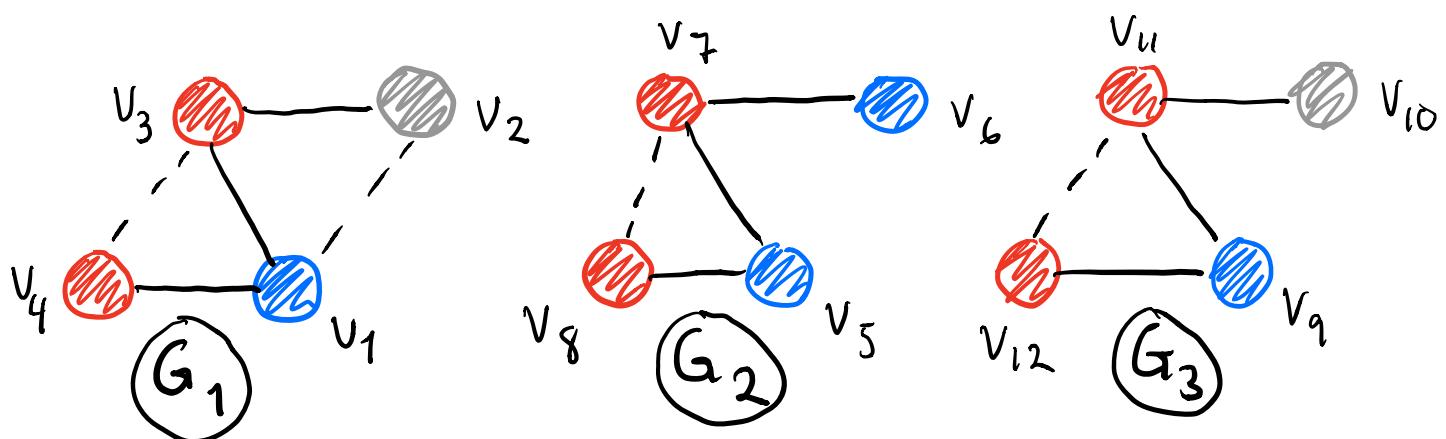
$$\mathcal{L}_3 = \{\text{BLUE}, \text{RED}, \text{GRAY}, \text{dotted}, \text{straight}\}$$

η. Give all frequent 1-subgraphs

Solution: We have $s_{\min} = 2/3$

Three small graphs illustrating frequent 2-subgraphs:

- Graph 1: A red node connected to a blue node. Below it is the formula $\beta = \frac{3}{3}$.
- Graph 2: A red node connected to a blue node, which is then connected to a grey node. Below it is the formula $\beta = \frac{4}{3}$.
- Graph 3: A red node connected to a grey node. Below it is the formula $\beta = \frac{2}{3}$.



Q. Give all frequent 2-subgraphs

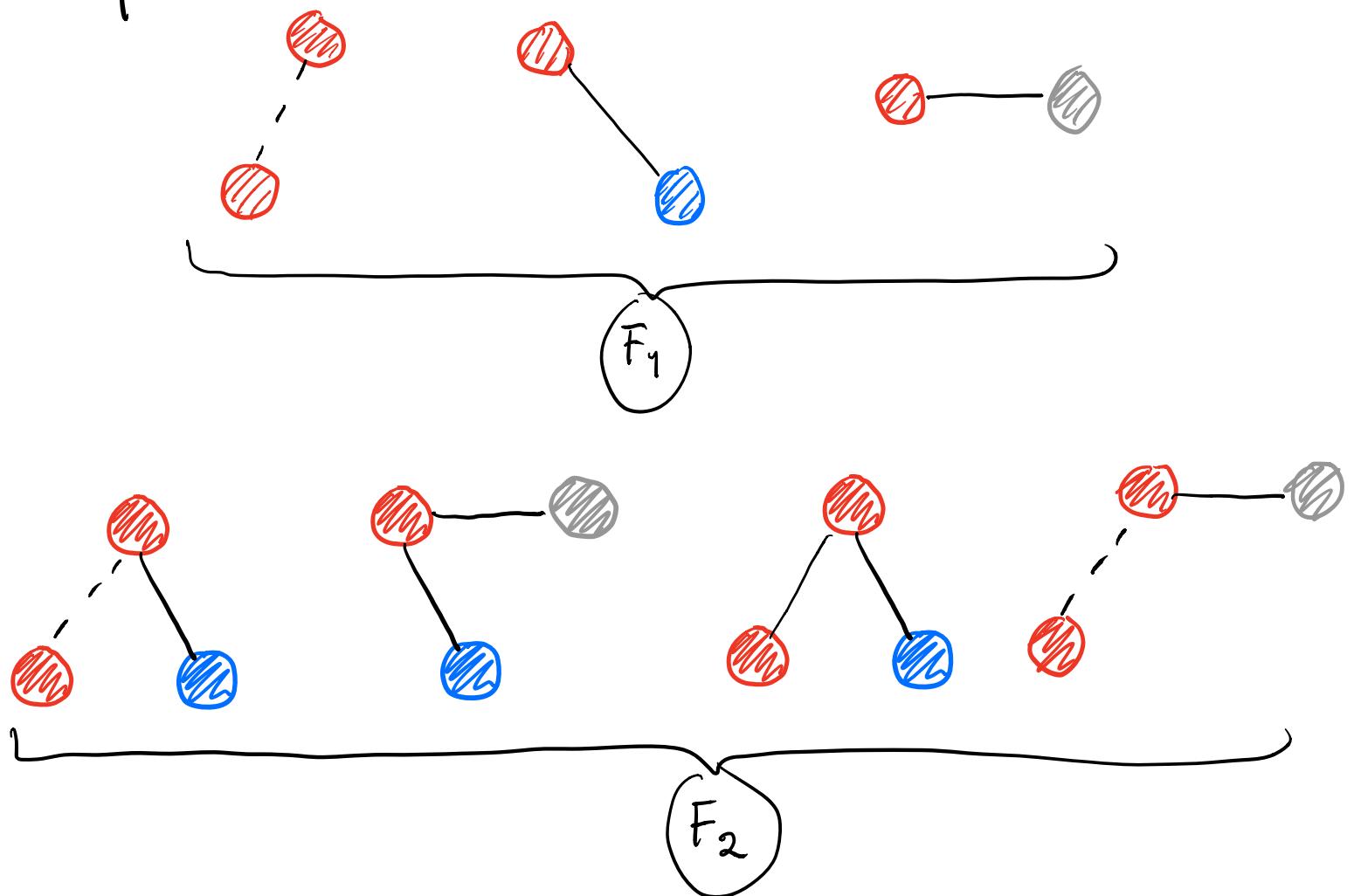
Solution: We get

Four small graphs corresponding to the ones above:

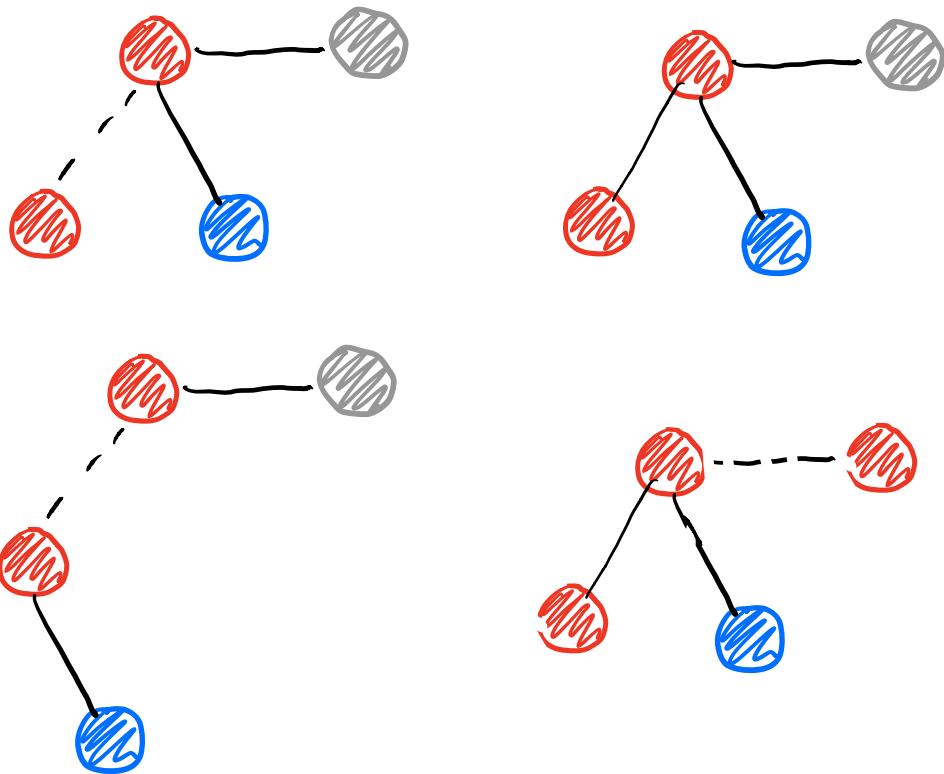
- Graph 1: A red node connected to a blue node. Below it is the formula $\beta = \frac{4}{3}$.
- Graph 2: A red node connected to a blue node, which is then connected to a grey node. Below it is the formula $\beta = \frac{2}{3}$.
- Graph 3: A red node connected to a blue node, which is then connected to another red node. Below it is the formula $\beta = \frac{3}{3}$.
- Graph 4: A red node connected to a grey node. Below it is the formula $\beta = \frac{2}{3}$.

L. Use the FSG algorithm, without any optimizations, to find the set of all frequent subgraphs.

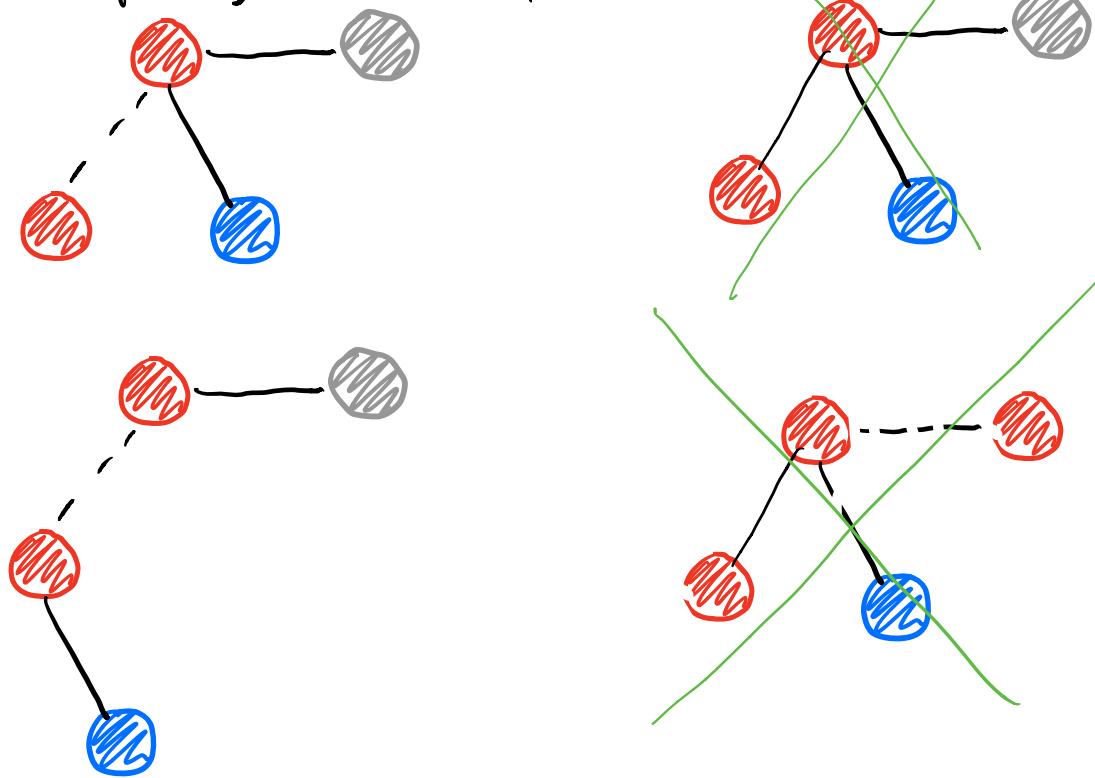
Solution. We have F_1 and F_2 from previous steps



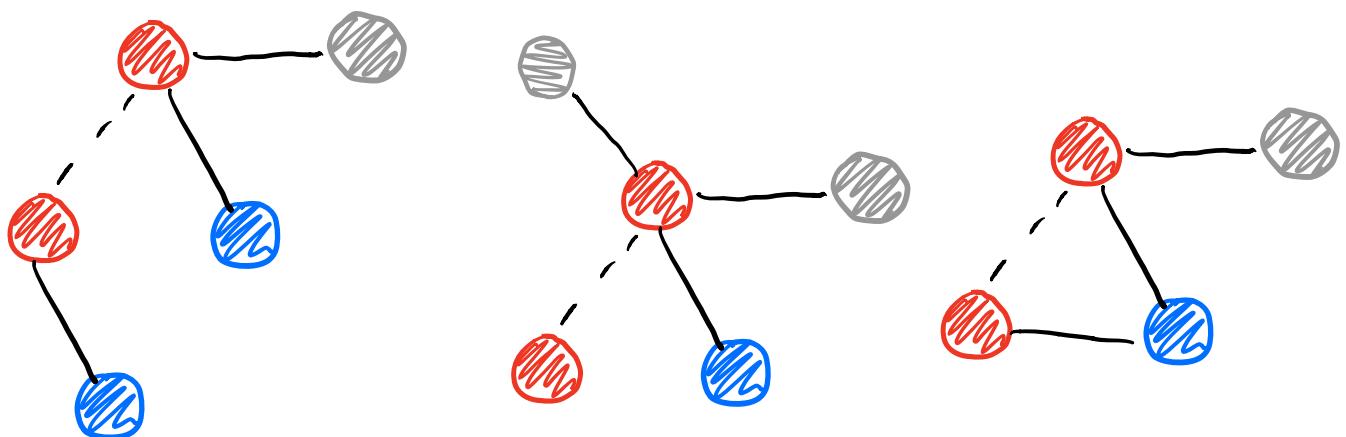
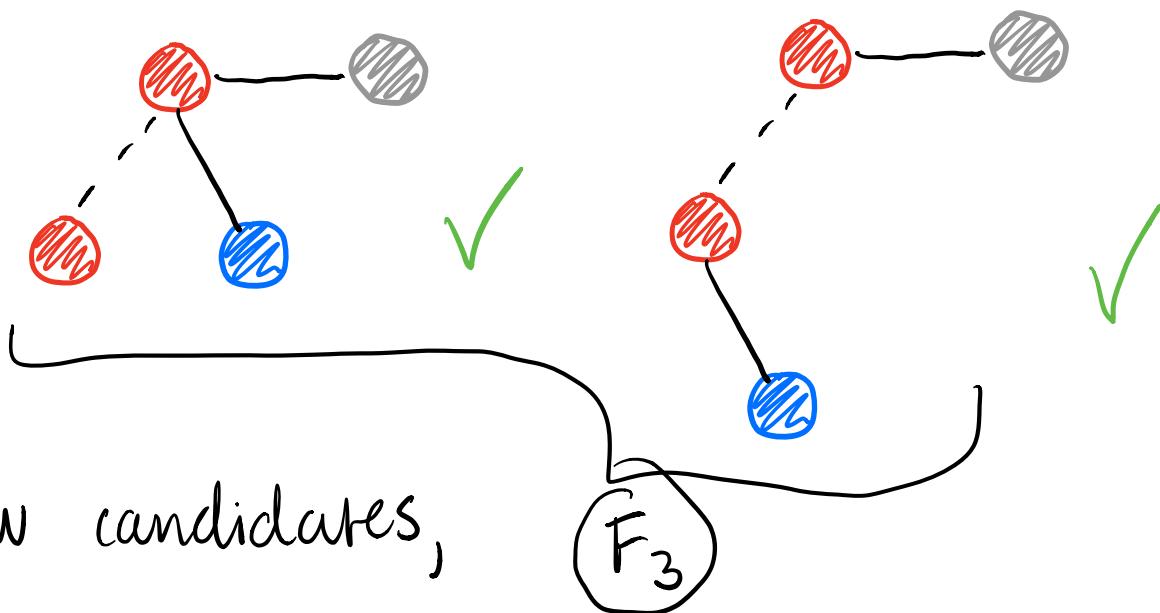
We generate candidates



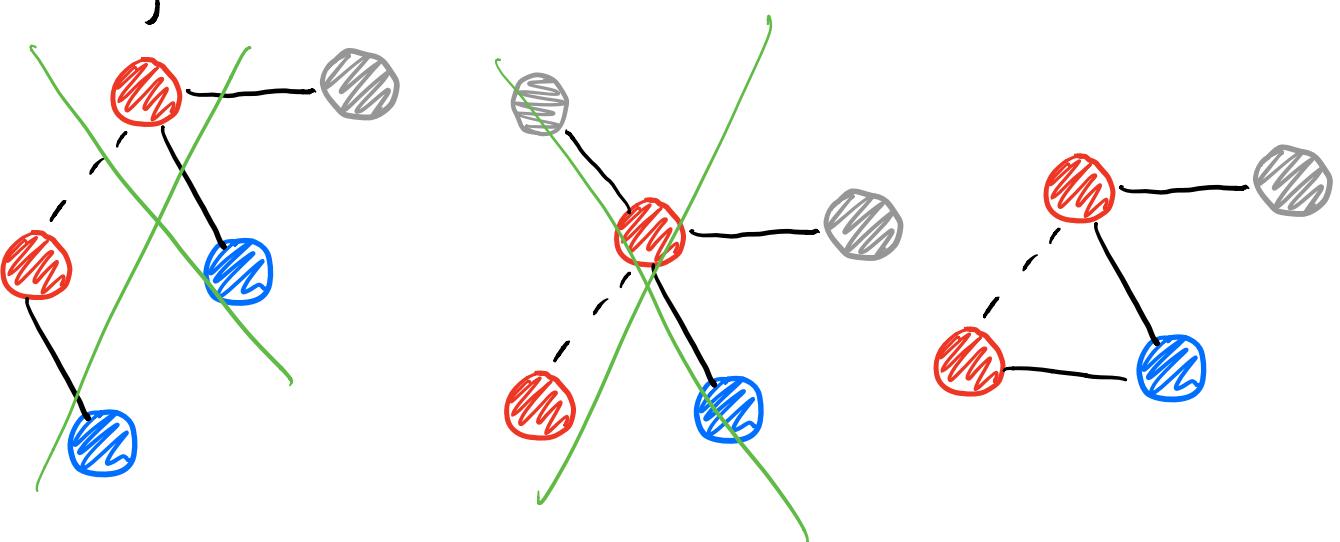
Candidates that do not satisfy antimonotonicity property are pruned



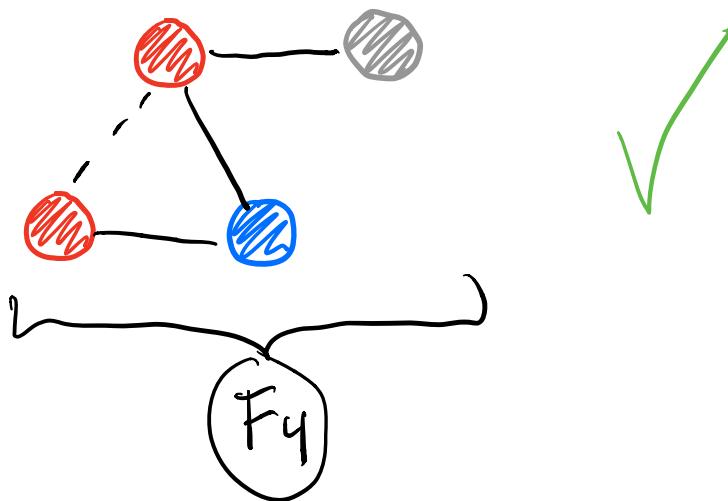
Check if candidates satisfy minimal support constraint, $s_{min} = 2/3$



Prune,



Check minimum support:



We return $F_1 \cup F_2 \cup F_3 \cup F_4$

