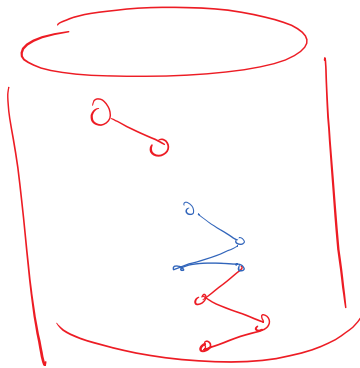
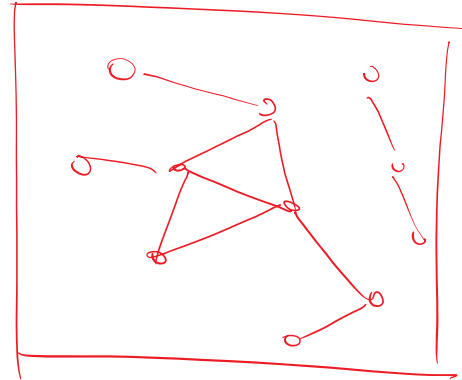


graph pattern mining



Database of graphs

vs.



One large graph
(social graph)

labeled graphs

$$G = (V, E)$$

$L(u)$ = label of vertex u

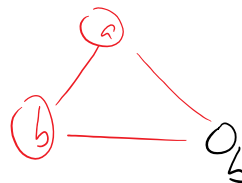
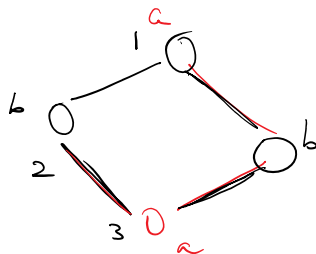
$L(u, v)$ = label of an edge (u, v)

Task: find commonly occurring subgraphs

frequent subgraphs pattern

$$\begin{matrix} \text{freq}(P) \\ \uparrow \\ \text{frequency} \end{matrix} = \begin{matrix} \text{sup}(P) \\ \uparrow \\ \text{support} \end{matrix} = \left\{ \begin{array}{l} \# \text{ of distinct graphs } G_i \text{ that contain } P \\ \# \text{ of occurrences of } P \text{ over all } G_i \in D \end{array} \right.$$

graph database



P:

$$\text{sup}(P) = 2 \quad (\text{distinct graphs})$$

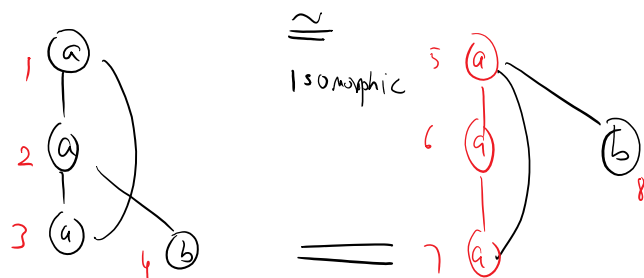
$$\text{sup}(P) = 6 \quad (\text{total \# of occurrences})$$

Mining Task: given a minimum support threshold θ find all

frequent subgraphs, i.e. $\text{sup}(P) \geq \theta$



potentially exponential space



	1	2	3	4	$\leftarrow x \in G_1$
$\phi(x) \in G_2$	6	5	7	8	
	7	5	6	8	

Graph G_1 is isomorphic to G_2 iff

$\Rightarrow \phi$: 1-1 correspondence between G_1 & G_2
such that both into and onto

$$x \in G_1$$

$$\phi(x) \in G_2$$

① Structure has to be preserved

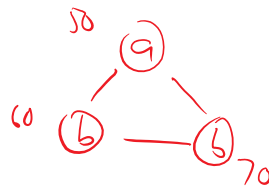
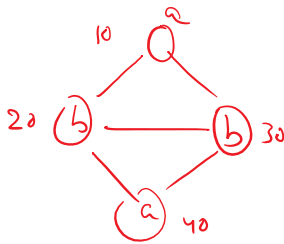
$$(x, y) \in G_1 \iff (\phi(x), \phi(y)) \in G_2$$

② Labels have to be preserved

$$L(x) = L(\phi(x))$$

$$L(x, y) = L(\phi(x), \phi(y))$$

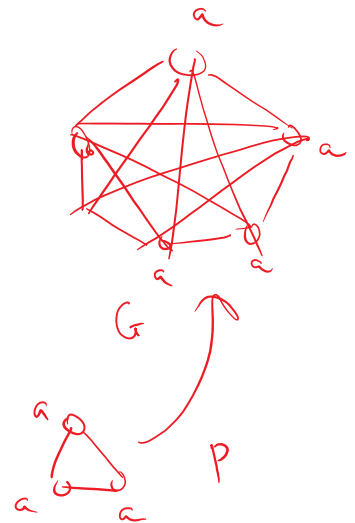
graph isomorphism problem
is it in P?



	a	b	b
	0	1	2
G_1	10	20	30
	10	30	20
	40	20	30
	40	30	20
G_2	50	60	70
	50	70	60

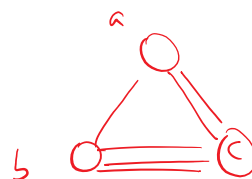
ϕ : subgraph isomorphism
 $x_1 \quad x_2 \quad \dots \quad x_n \leftarrow P$
 embedding: $(\phi(x_1), \phi(x_2), \dots, \phi(x_n)) \leftarrow G_i$

subgraph isomorphism
NP-hard problem



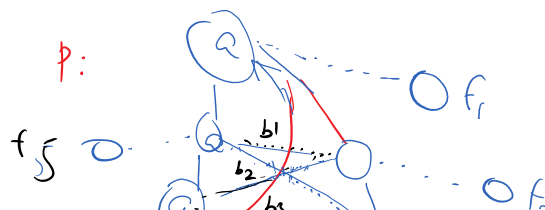
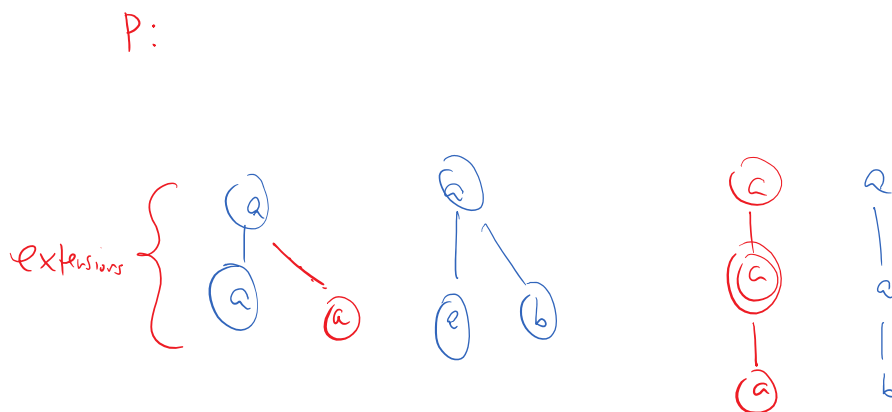
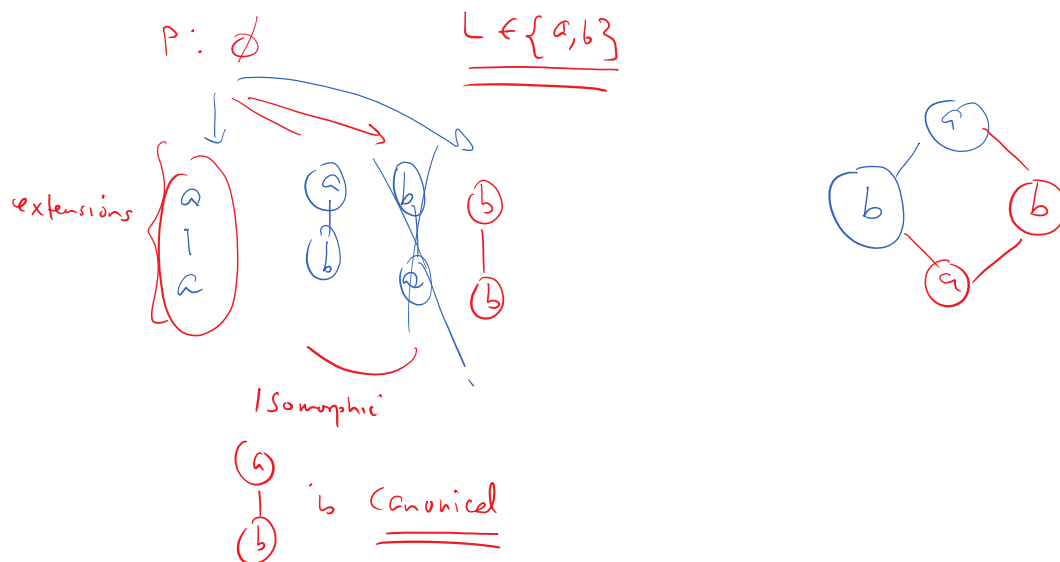
① Systematic search

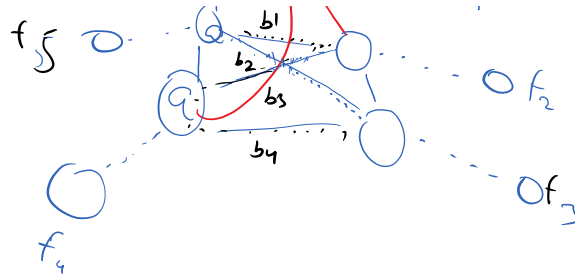
→ extend a pattern by one extra edge
 → pattern extension step



- check for "duplicates"; list/count only the
distinct subgraphs/patterns
- 2) (graph isomorphism)
- (2) Collect the frequency for each pattern P
- 3) subgraph isomorphism step

Step 1: Candidate generation / Pattern extension





$S \times |L|$
forward extensions
 $f_1 \dots f_s$
 add new nodes.

backward extensions
 $b_1 \dots b_n$
 add cycles

Every pattern will be represented by its
 Canonical code
 DFS code
 → DFS tree

extend

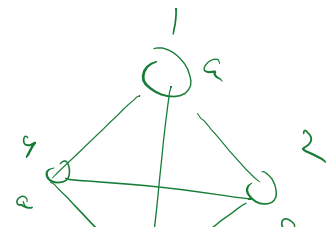
① go depth-first
 do a forward extension from the current node

② add backwards edges closer to root first from the last vertex.

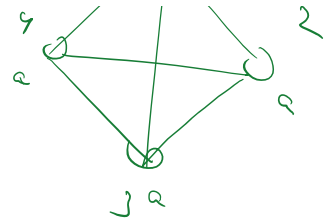
③ start branching bottom to top

Only to right most path

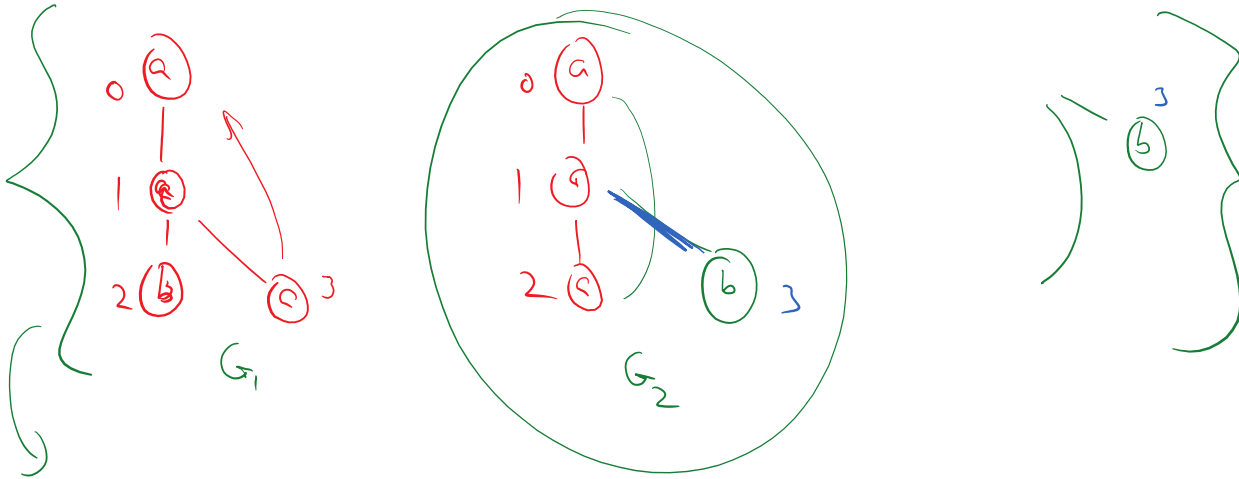
$S = \{ G_1, G_2, \dots, G_n \}$
isomorphic graphs



isomorphic graphs



$\min \text{DFS}_{\text{code}}(G) \rightarrow$ a unique canonical representative in terms of the DFS-tree



$\text{DFS}_{\text{code}}(G_1)$

0	1	a	a	-
1	2	a	b	-
1	3	a	c	-
3	0	a	c	-

X

0	1	a	a	-
1	2	a	a	-
2	0	a	a	-
1	3	a	b	-

$\text{DFS}_{\text{code}}(G_2)$

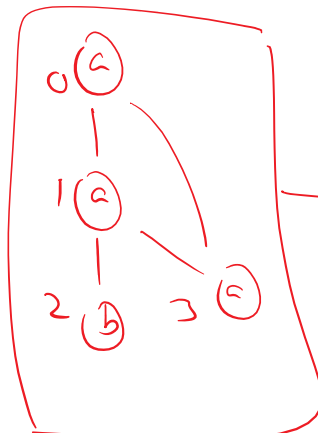
✓✓

0	1	a	a	-
1	2	a	a	-
2	0	a	a	-
0	3	a	b	-

$\text{DFS}_{\text{code}}(G_3)$

X

P:

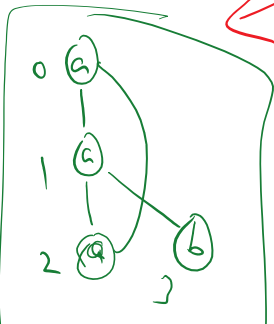


is this canonical?

0	1	a	a	-
1	2	a	b	-
1	3	a	c	-
3	0	a	a	-

DFS_{code}

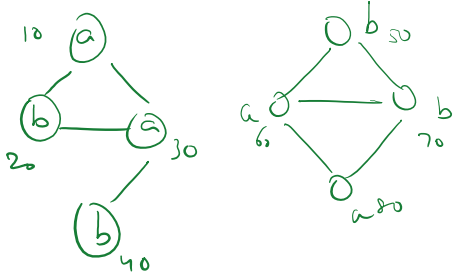
X





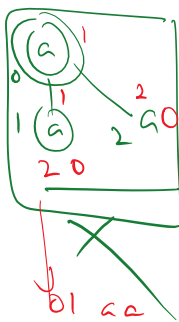
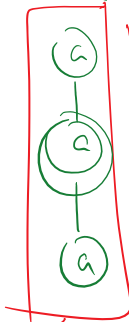
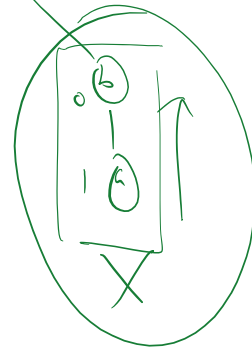
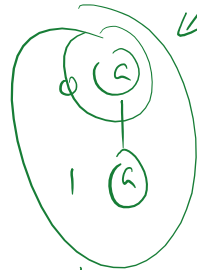
Is this Canonical?

→ is this a duplicate?



	a	b
	0	1
G1	10	30
	30	10
G2	60	80
	80	60

min sup = 2



not canonical

01 aa
02 aa

not Canonical

	0	1	2
	a	a	b
G1	10	30	20
	10	30	40
	30	10	20
G2	60	80	70
	80	60	50
	80	60	70

not yes



