# Frequent Generalized Subgraph Mining via Graph Edit Distances

Richard Palme and Pascal Welke

IoT Streams for Predictive Maintenance | SeDaMi

- Farmers track their livestock and have to treat them well

- Farmers track their livestock and have to treat them well
  - before they get sick

UNIVERSITÄT BONN

- Farmers track their livestock and have to treat them well
  - before they get sick
  - certainly before they die

UNIVERSITÄT BONN

- Farmers track their livestock and have to treat them well
  - before they get sick
  - certainly before they die
  - so that they perform well for a long time

- Farmers track their livestock and have to treat them well
  - before they get sick
  - certainly before they die
  - so that they perform well for a long time
- Breeders track whole populations of cows and need to decide

UNIVERSITÄT BONN

- Farmers track their livestock and have to treat them well
  - before they get sick
  - certainly before they die
  - so that they perform well for a long time
- Breeders track whole populations of cows and need to decide
  - if a breed is a good fit for a changing market or changing climate
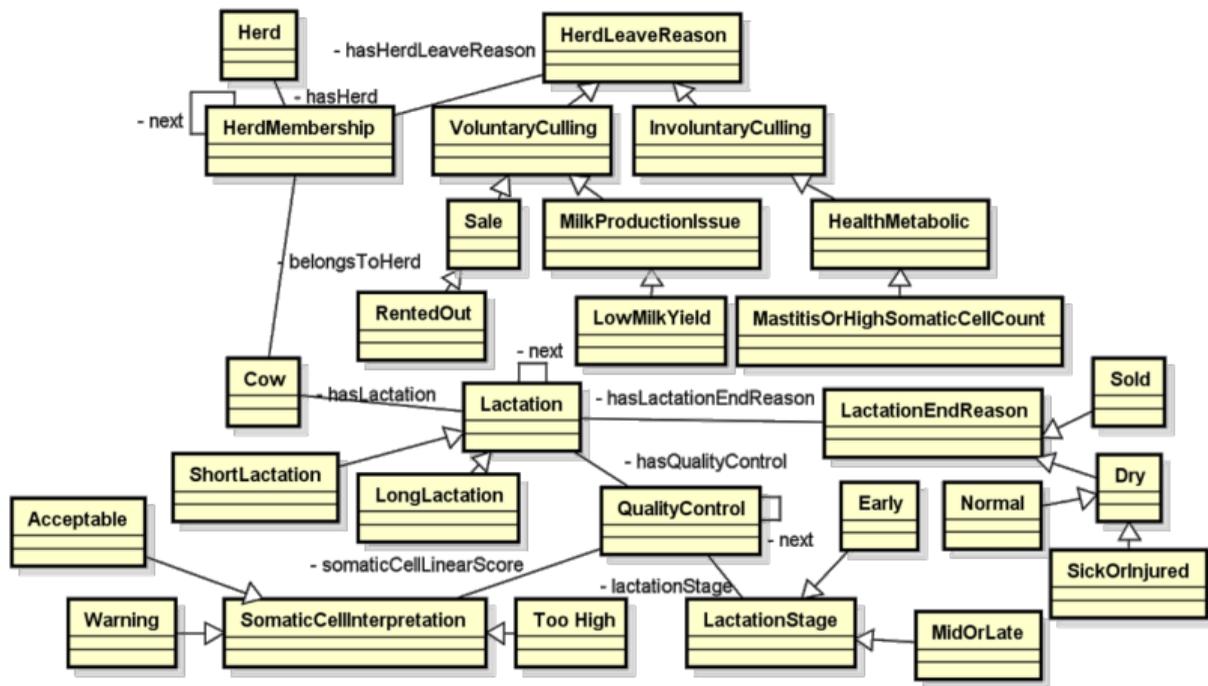
UNIVERSITÄT BONN

- Farmers track their livestock and have to treat them well
  - before they get sick
  - certainly before they die
  - so that they perform well for a long time
- Breeders track whole populations of cows and need to decide
  - if a breed is a good fit for a changing market or changing climate
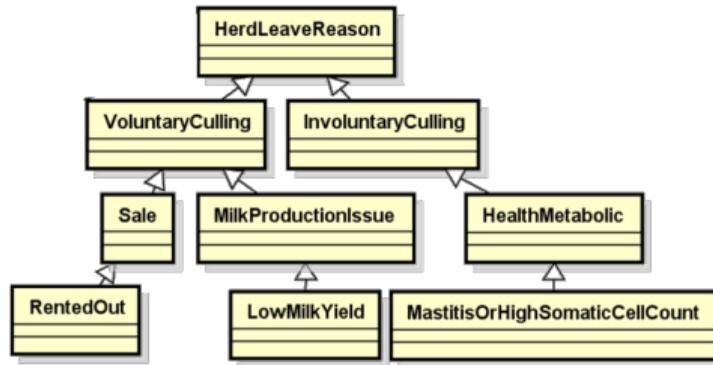  - which traits to improve by selective breeding

# A Dairy Cattle Performance Ontology

UNIVERSITÄT BONN

- For each cow, there is a (knowledge) graph recording events

- For each cow, there is a (knowledge) graph recording events
- Are there some reoccurring patterns?

UNIVERSITÄT BONN

A pattern graph $H$ is a *generalized subgraph* (with respect to an ontology $O$) of a graph $G$ if

UNIVERSITÄT BONN

A pattern graph $H$ is a *generalized subgraph* (with respect to an ontology $O$) of a graph $G$ if

- there is a subgraph isomorphism from $H'$ to $G$

A pattern graph $H$ is a *generalized subgraph* (with respect to an ontology $O$) of a graph $G$ if

- there is a subgraph isomorphism from $H'$ to $G$

- $H'$ can be constructed by replacing vertex labels in $H$ with *more specific labels* (according to our ontology)

A pattern graph $H$ is a *generalized subgraph* (with respect to an ontology $O$) of a graph $G$ if

- there is a subgraph isomorphism from $H'$ to $G$
- $H'$ can be constructed by replacing vertex labels in $H$ with *more specific labels* (according to our ontology)

A pattern graph $H$ is a *generalized subgraph* (with respect to an ontology $O$) of a graph $G$ if

- there is a subgraph isomorphism from $H'$ to $G$
- $H'$ can be constructed by replacing vertex labels in $H$ with *more specific labels* (according to our ontology)

UNIVERSITÄT BONN

A pattern graph $H$ is a *generalized subgraph* (with respect to an ontology $O$) of a graph $G$ if
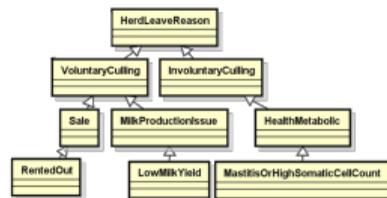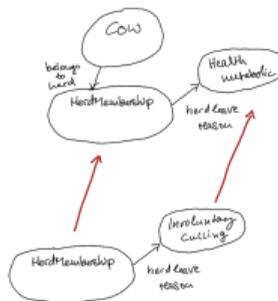
- there is a subgraph isomorphism from $H'$ to $G$
- $H'$ can be constructed by replacing vertex labels in $H$ with *more specific labels* (according to our ontology)

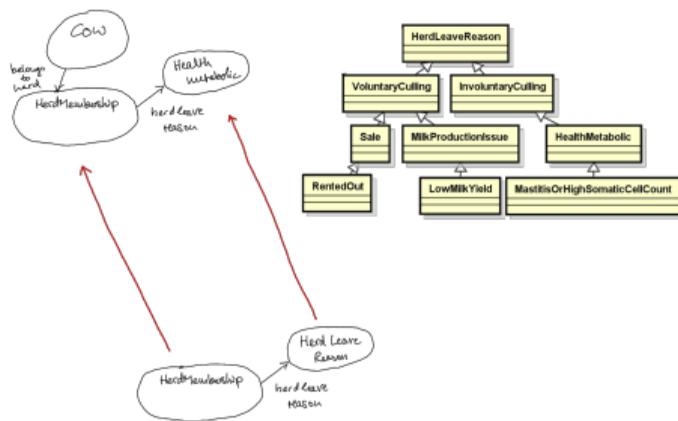UNIVERSITÄT BONN

The *Frequent Generalized Subgraph Mining Problem* is then:

The *Frequent Generalized Subgraph Mining Problem* is then:

Given: A database $D$ of graphs, an ontology $O$ and a frequency threshold $t$



$t = 4$

UNIVERSITÄT BONN

# Generalized Subgraph Mining

The *Frequent Generalized Subgraph Mining Problem* is then:

Given: A database $D$ of graphs, an ontology $O$ and a frequency threshold $t$

Compute: The set of all graphs that are generalized subgraphs of at least $t$ graphs in $D$



$t = 4$

UNIVERSITÄT BONN

- Not automatically...

UNIVERSITÄT BONN

- Not automatically...
- But domain experts can interpret (smaller) frequent patterns

- Not automatically...
- But domain experts can interpret (smaller) frequent patterns
- And they can interpret changes in frequent patterns

UNIVERSITÄT BONN

Frequent Generalized Subgraph Mining has a long history

- ⏱
- ⏱

UNIVERSITÄT BONN

Frequent Generalized Subgraph Mining has a long history

- ?
- ?

And recently gained more traction

- ?
- ?
- ?

Frequent Generalized Subgraph Mining has a long history

- ⌛
- ⌛

And recently gained more traction

- ⌛
- ⌛
- ⌛

These papers all modify classical frequent subgraph mining algorithms.

UNIVERSITÄT BONN

Frequent Generalized Subgraph Mining has a long history

- ⌀
- ⌀

And recently gained more traction

- ⌀
- ⌀
- ⌀

These papers all modify classical frequent subgraph mining algorithms.

As a subroutine, they use subgraph isomorphism algorithms.

UNIVERSITÄT BONN

# Our Approach

UNIVERSITÄT BONN

We propose to replace subgraph isomorphism computations by *graph edit distance* computations

We propose to replace subgraph isomorphism computations by *graph edit distance* computations

- This makes the problem *harder*

UNIVERSITÄT BONN

We propose to replace subgraph isomorphism computations by *graph edit distance* computations

- This makes the problem *harder*
- But it allows some nice *freedom and simplicity* in modeling

We propose to replace subgraph isomorphism computations by *graph edit distance* computations

- This makes the problem *harder*
- But it allows some nice *freedom and simplicity* in modeling
- And it *simplifies* rather intricate mining algorithms

UNIVERSITÄT BONN

The *graph edit distance* between $G$ and $H$ is the smallest cost of a sequence of edits transforming $G$ into $H$.

UNIVERSITÄT BONN

# Graph Edit Distance

| Edit operation | Edit cost |
|---|---|
| Insert an isolated vertex with label $\alpha \in \Sigma$ | $c(\varepsilon, \alpha)$ |
| Delete an isolated vertex $u$ | $c(\lambda(u), \varepsilon)$ |
| Substitute the label of a vertex $u$ by $\alpha \in \Sigma$ | $c(\lambda(u), \alpha)$ |
| Insert an edge with label $\alpha \in \Sigma$ | $c(\varepsilon, \alpha)$ |
| Delete an edge $e$ | $c(\lambda(e), \varepsilon)$ |
| Substitute the label of an edge $e$ by $\alpha \in \Sigma$ | $c(\lambda(e), \alpha)$ |

The *graph edit distance* between $G$ and $H$ is the smallest cost of a sequence of edits transforming $G$ into $H$.

UNIVERSITÄT BONN

The GED can be used to solve the *subgraph isomorphism problem (SGI)* by imposing the following three constraints on the edit cost function:

The GED can be used to solve the *subgraph isomorphism problem (SGI)* by imposing the following three constraints on the edit cost function:

$$\forall \beta \in \Sigma_\varepsilon \colon c(\varepsilon, \beta) = 0 \qquad \text{(free insertions)}$$

$$\forall \alpha \in \Sigma \colon c(\alpha, \varepsilon) > 0 \qquad \text{(paid deletions)}$$

$$\forall \alpha, \beta \in \Sigma \colon c(\alpha, \beta) > 0 \iff \alpha \neq \beta \qquad \text{(paid substitutions)}$$

The GED can be used to solve the *subgraph isomorphism problem (SGI)* by imposing the following three constraints on the edit cost function:

$$\forall \beta \in \Sigma_{\varepsilon} \colon c(\varepsilon, \beta) = 0 \qquad \text{(free insertions)}$$

$$\forall \alpha \in \Sigma \colon c(\alpha, \varepsilon) > 0 \qquad \text{(paid deletions)}$$

$$\forall \alpha, \beta \in \Sigma \colon c(\alpha, \beta) > 0 \iff \alpha \neq \beta \qquad \text{(paid substitutions)}$$

UNIVERSITÄT BONN

The GED can be used to solve the *subgraph isomorphism problem (SGI)* by imposing the following three constraints on the edit cost function:

$$\forall \beta \in \Sigma_\varepsilon \colon c(\varepsilon, \beta) = 0 \qquad \text{(free insertions)}$$
$$\forall \alpha \in \Sigma \colon c(\alpha, \varepsilon) > 0 \qquad \text{(paid deletions)}$$
$$\forall \alpha, \beta \in \Sigma \colon c(\alpha, \beta) > 0 \iff \alpha \neq \beta \qquad \text{(paid substitutions)}$$

$$\mathrm{SGI}(H, G) = \text{true} \iff \mathrm{GED}(H, G) = 0.$$

UNIVERSITÄT BONN

The GED can be used to solve the *subgraph isomorphism problem (SGI)* by imposing the following three constraints on the edit cost function:

$$\forall \beta \in \Sigma_\varepsilon \colon c(\varepsilon, \beta) = 0 \qquad \text{(free insertions)}$$

$$\forall \alpha \in \Sigma \colon c(\alpha, \varepsilon) > 0 \qquad \text{(paid deletions)}$$

$$\forall \alpha, \beta \in \Sigma \colon c(\alpha, \beta) > 0 \iff \alpha \neq \beta \qquad \text{(paid substitutions)}$$
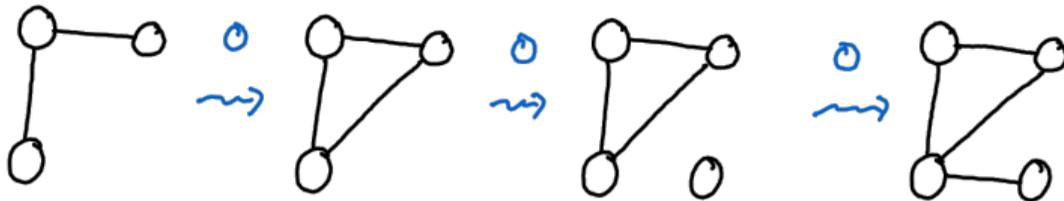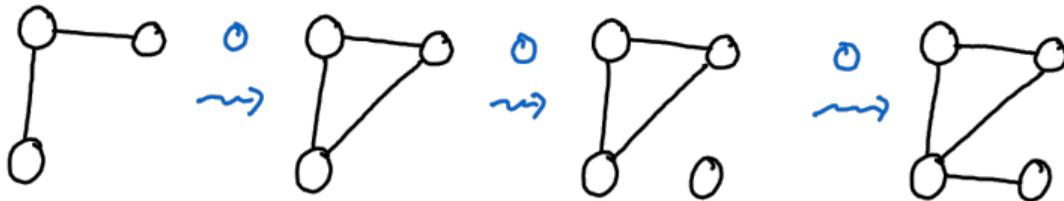
$$\mathrm{SGI}(H, G) = \text{true} \iff \mathrm{GED}(H, G) = 0.$$

To solve the *generalized subgraph isomorphism problem (GSGI)*, we impose the following four constraints on the edit cost function:

To solve the *generalized subgraph isomorphism problem (GSGI)*, we impose the following four constraints on the edit cost function:

$$\forall \beta \in \Sigma_\varepsilon \colon c(\varepsilon, \beta) = 0 \qquad \text{(free insertions)}$$

$$\forall \alpha \in \Sigma \colon c(\alpha, \varepsilon) > 0 \qquad \text{(paid deletions)}$$

$$\forall \alpha, \beta \in \Sigma \colon c(\alpha, \beta) > 0 \iff \alpha \neq \beta \text{ and } \alpha \text{ is not more general than } \beta \qquad \text{(paid substitutions)}$$

$$\forall \alpha, \beta \in \Sigma \colon c(\alpha, \beta) = 0 \iff \alpha = \beta \text{ or } \alpha \text{ is more general than } \beta \qquad \text{(free specializations)}$$

To solve the *generalized subgraph isomorphism problem (GSGI)*, we impose the following four constraints on the edit cost function:

$$\forall \beta \in \Sigma_\varepsilon \colon c(\varepsilon, \beta) = 0 \qquad \text{(free insertions)}$$

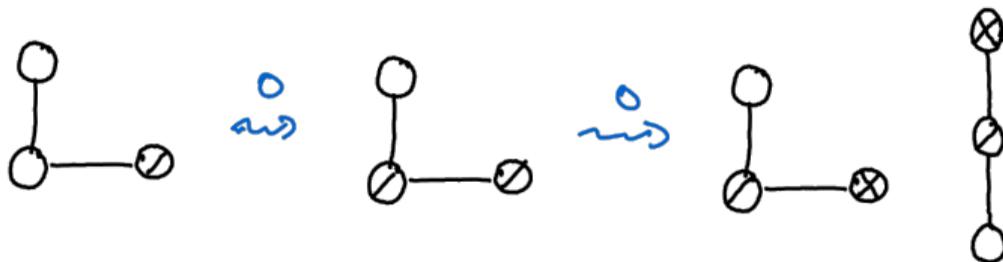$$\forall \alpha \in \Sigma \colon c(\alpha, \varepsilon) > 0 \qquad \text{(paid deletions)}$$

$$\forall \alpha, \beta \in \Sigma \colon c(\alpha, \beta) > 0 \iff \alpha \neq \beta \text{ and } \alpha \text{ is not more general than } \beta \qquad \text{(paid substitutions)}$$

$$\forall \alpha, \beta \in \Sigma \colon c(\alpha, \beta) = 0 \iff \alpha = \beta \text{ or } \alpha \text{ is more general than } \beta \qquad \text{(free specializations)}$$

# Another Nice Polynomial Reduction

To solve the *generalized subgraph isomorphism problem (GSGI)*, we impose the following four constraints on the edit cost function:

$$\forall \beta \in \Sigma_\varepsilon \colon c(\varepsilon, \beta) = 0 \qquad \text{(free insertions}$$

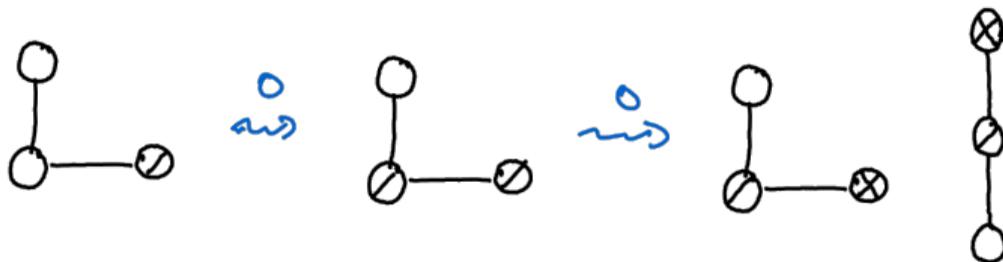$$\forall \alpha \in \Sigma \colon c(\alpha, \varepsilon) > 0 \qquad \text{(paid deletions}$$

$$\forall \alpha, \beta \in \Sigma \colon c(\alpha, \beta) > 0 \iff \alpha \neq \beta \text{ and } \alpha \text{ is not more general than } \beta \qquad \text{(paid substitutions}$$

$$\forall \alpha, \beta \in \Sigma \colon c(\alpha, \beta) = 0 \iff \alpha = \beta \text{ or } \alpha \text{ is more general than } \beta \qquad \text{(free specializations}$$

UNIVERSITÄT BONN

# Another Nice Polynomial Reduction

To solve the *generalized subgraph isomorphism problem (GSGI)*, we impose the following four constraints on the edit cost function:

$$\forall \beta \in \Sigma_\varepsilon \colon c(\varepsilon, \beta) = 0 \qquad \text{(free insertions)}$$

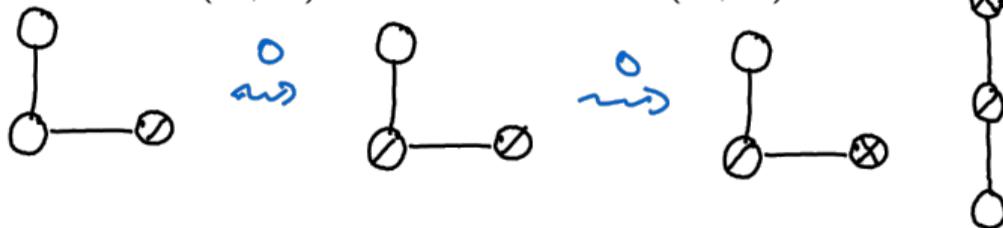$$\forall \alpha \in \Sigma \colon c(\alpha, \varepsilon) > 0 \qquad \text{(paid deletions)}$$

$$\forall \alpha, \beta \in \Sigma \colon c(\alpha, \beta) > 0 \iff \alpha \neq \beta \text{ and } \alpha \text{ is not more general than } \beta \qquad \text{(paid substitutions)}$$

$$\forall \alpha, \beta \in \Sigma \colon c(\alpha, \beta) = 0 \iff \alpha = \beta \text{ or } \alpha \text{ is more general than } \beta \qquad \text{(free specializations)}$$

$$\text{GSGI}(H, G) = \text{true} \iff \text{GED}(H, G) = 0.$$

- There are fast heuristics ⌀

– There are fast heuristics  ⟳
– There was some nice work on lower bounds at ECMLPKDD this year which might be adapted  ⟳

UNIVERSITÄT BONN

- There are fast heuristics ⟳
- There was some nice work on lower bounds at ECMLPKDD this year which might be adapted ⟳

- We say that a graph is a generalized subgraph if it has a small generalized subgraph edit distance

UNIVERSITÄT BONN

- We implemented a proof of concept graph mining algorithm

UNIVERSITÄT BONN

- We implemented a proof of concept graph mining algorithm
- It is available on `https://github.com/RichardPalme/fasm`

UNIVERSITÄT BONN
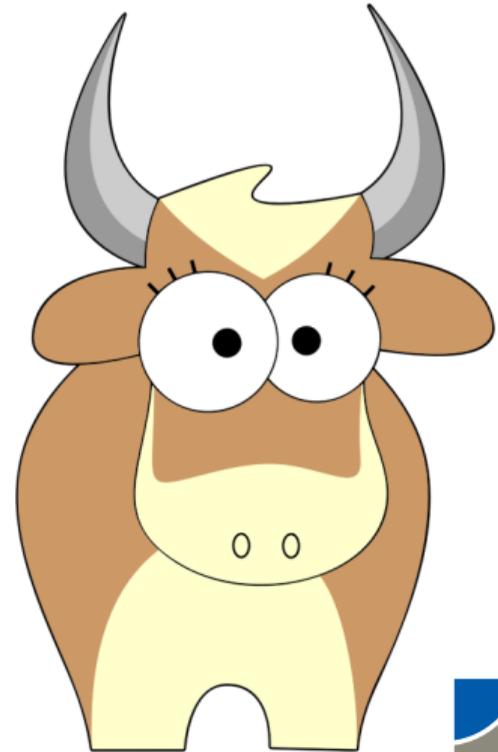
- Generalized Subgraphs can be identified with a specialized variant of the graph edit distance
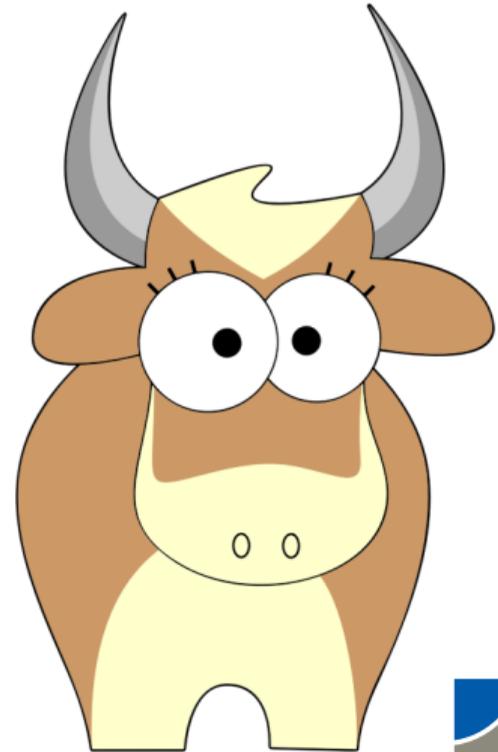
UNIVERSITÄT BONN

- Generalized Subgraphs can be identified with a specialized variant of the graph edit distance
- This allows to mine generalized patterns in an elegant way

- Generalized Subgraphs can be identified with a specialized variant of the graph edit distance
- This allows to mine generalized patterns in an elegant way
- We can include interesting costs (checkout the paper) to make the mining practically better