

EUCLIDEAN EMBEDDING OF CO-PROVEN QUERIES

PRESENTATION OF MASTER THESIS

Hannes Schulz

Albert-Ludwigs-Universität Freiburg

2009-5-20

THE BLIND AND THE ELEPHANT



THE BLIND AND THE ELEPHANT



OUTLINE

- 1 RELATIONAL DATA
- 2 EMBEDDING OF CO-PROVEN QUERIES AND INTERPRETATIONS
- 3 EXPERIMENTS
- 4 CONCLUSION

1 RELATIONAL DATA

- Representation
- Algorithms (The High Altitude View)
- Semantically Grounded Distances

2 EMBEDDING OF CO-PROVEN QUERIES AND INTERPRETATIONS

3 EXPERIMENTS

4 CONCLUSION

1 RELATIONAL DATA

- Representation

- Algorithms (The High Altitude View)

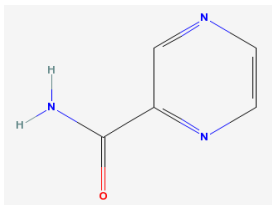
- Semantically Grounded Distances

2 EMBEDDING OF CO-PROVEN QUERIES AND INTERPRETATIONS

3 EXPERIMENTS

4 CONCLUSION

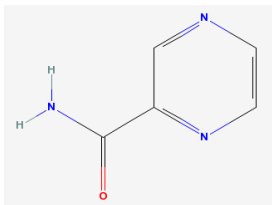
RELATIONAL DATABASES



atom		bond			
id	name	id	a1	a2	type
1	h	1	1	3	1
2	h	2	2	3	1
3	n	3	3	4	1
4	c	4	4	5	2
5	o			...	



RELATIONAL DATABASES



atom		bond			
id	name	id	a1	a2	type
1	h	1	1	3	1
2	h	2	2	3	1
3	n	3	3	4	1
4	c	4	4	5	2
5	o			...	

- Entities are **discrete**
- Elements part of **arbitrary number** of relations
- Implicit information in **views**
- ▶ No mapping to feature vectors

INDUCTIVE LOGIC PROGRAMMING TO THE RESCUE

```
atom(1 , d1_1 , c , 22 , -0.117).  
atom(1 , d1_2 , c , 22 , -0.117).  
atom(1 , d1_3 , c , 22 , -0.117).  
...  
bond(1 , d1_1 , d1_2 , 7).  
bond(1 , d1_2 , d1_3 , 7).  
bond(1 , d1_3 , d1_4 , 7).  
bond(1 , d1_4 , d1_5 , 7).  
...
```

Inductive Logic Programming

- Retains complexity of DB in **interpretations**



INDUCTIVE LOGIC PROGRAMMING TO THE RESCUE

```
ring (K) :-  
  atom (K, A, c, _, _),  
  bond (K, A, B, _),  
  atom (K, B, c, _, _),  
  bond (K, B, C, _),  
  ...  
  bond (K, F, A, _).
```

Inductive Logic Programming

- Retains complexity of DB in **interpretations**
- Uses rules to reflect implicit information



INDUCTIVE LOGIC PROGRAMMING TO THE RESCUE

```
atom(K, A, c, _, _),  
bond(K, A, B, _),  
atom(K, B, c, _, _),  
bond(K, B, C, _)
```

is frequent in organic molecules

Inductive Logic Programming

- Retains complexity of DB in **interpretations**
- Uses rules to reflect implicit information
- Abstracts using variables in **queries**



1 RELATIONAL DATA

- Representation
- Algorithms (The High Altitude View)
- Semantically Grounded Distances

2 EMBEDDING OF CO-PROVEN QUERIES AND INTERPRETATIONS

3 EXPERIMENTS

4 CONCLUSION

ILP ALGORITHMS: LIFTED FROM PROPOSITIONAL CASE

One table to n tables

- Pattern miners (WARMR)
- Rule learners (ALEPH)
- Decision Trees (TILDE)
- ...



ILP ALGORITHMS: LIFTED FROM PROPOSITIONAL CASE

One table to n tables

- Pattern miners (*WARMR*)
- Rule learners (*ALEPH*)
- Decision Trees (*TILDE*)
- ...

Usually,

- Input: interpretations, rules (“background knowledge”)
- Search (sub-)space of logical formulæ
- Output:
 - Pattern miners: interesting patterns
 - Classification: concept definitions



BACK TO THE ELEPHANT

Input and output of ILP algorithms feel like elephant parts.

MUTAGENESIS PATTERN MINING

```

key(A) , attp(A,B,28) , attp(A,C,28) , sbond(A,B,C,1)
key(A) , atel(A,B,c) , atel(A,C,c) , attp(A,D,27) , sbond(A,
    B,C,1) , sbond(A,B,D,7) , carbon_6_rings(A)
key(A) , attp(A,B,22) , methyls(A)
key(A) , atel(A,B,c) , atel(A,C,c) , atel(A,D,h) , atel(A,E,n
    ) , sbond(A,B,C,7) , sbond(A,B,D,1) , sbond(A,C,E,1) ,
    benzenes(A) , ring_size_5s(A)
key(A) , atel(A,B,c) , atel(A,C,c) , atel(A,D,h) , atel(A,E,n
    ) , attp(A,F,10) , sbond(A,B,C,7) , sbond(A,B,D,1) ,
    sbond(A,C,E,1)
key(A) , atel(A,B,c) , attp(A,C,10) , attp(A,D,27) , attp(A,
    E,27) , sbond(A,B,C,1) , sbond(A,D,E,7)
key(A) , atel(A,B,c) , attp(A,C,21) , attp(A,D,26) , attp(A,
    E,26) , sbond(A,B,C,7) , sbond(A,B,D,7) , sbond(A,C,
    E,7)
key(A) , atel(A,B,c) , atel(A,C,c) , atel(A,D,n) , attp(A,E
    ,22) , sbond(A,B,D,1) , sbond(A,B,E,7) , sbond(A,C,E
    ,1)
    
```

BACK TO THE ELEPHANT

Input and output of ILP algorithms feel like elephant parts.



MUTAGENESIS PATTERN MINING

```

key(A) , attp(A,B,28) , attp(A,C,28) , sbond(A,B,C,1)
key(A) , atel(A,B,c) , atel(A,C,c) , attp(A,D,27) , sbond(A,
  B,C,1) , sbond(A,B,D,7) , carbon_6_rings(A)
key(A) , attp(A,B,22) , methyls(A)
key(A) , atel(A,B,c) , atel(A,C,c) , atel(A,D,h) , atel(A,E,n)
  , sbond(A,B,C,7) , sbond(A,B,D,1) , sbond(A,C,E,1) ,
  benzenes(A) , ring_size_5s(A)
key(A) , atel(A,B,c) , atel(A,C,c) , atel(A,D,h) , atel(A,E,n)
  , attp(A,F,10) , sbond(A,B,C,7) , sbond(A,B,D,1) ,
  sbond(A,C,E,1)
key(A) , atel(A,B,c) , attp(A,C,10) , attp(A,D,27) , attp(A,
  E,27) , sbond(A,B,C,1) , sbond(A,D,E,7)
key(A) , atel(A,B,c) , attp(A,C,21) , attp(A,D,26) , attp(A,
  E,26) , sbond(A,B,C,7) , sbond(A,B,D,7) , sbond(A,C,
  E,7)
key(A) , atel(A,B,c) , atel(A,C,c) , atel(A,D,n) , attp(A,E,
  22) , sbond(A,B,D,1) , sbond(A,B,E,7) , sbond(A,C,E,
  1)

```


SEE THE WHOLE PICTURE

Idea:

- Embed queries and interpretations into common 2D space

MUTAGENESIS PATTERN MINING

```

key(A) , atyp(A,B,28) , atyp(A,C,28) , sbond(A,B,C,1)
key(A) , atel(A,B,c) , atel(A,C,c) , atyp(A,D,27) , sbond(A,
    B,C,1) , sbond(A,B,D,7) , carbon_6_rings(A)
key(A) , atyp(A,B,22) , methyls(A)
key(A) , atel(A,B,c) , atel(A,C,c) , atel(A,D,h) , atel(A,E,n
    ) , sbond(A,B,C,7) , sbond(A,B,D,1) , sbond(A,C,E,1) ,
    benzenes(A) , ring_size_5s(A)
key(A) , atel(A,B,c) , atel(A,C,c) , atel(A,D,h) , atel(A,E,n
    ) , atyp(A,F,10) , sbond(A,B,C,7) , sbond(A,B,D,1) ,
    sbond(A,C,E,1)
key(A) , atel(A,B,c) , atyp(A,C,10) , atyp(A,D,27) , atyp(A,
    E,27) , sbond(A,B,C,1) , sbond(A,D,E,7)
key(A) , atel(A,B,c) , atyp(A,C,21) , atyp(A,D,26) , atyp(A,
    E,26) , sbond(A,B,C,7) , sbond(A,B,D,7) , sbond(A,C,
    E,7)
key(A) , atel(A,B,c) , atel(A,C,c) , atel(A,D,n) , atyp(A,E
    ,22) , sbond(A,B,D,1) , sbond(A,B,E,7) , sbond(A,C,E
    ,1)
  
```

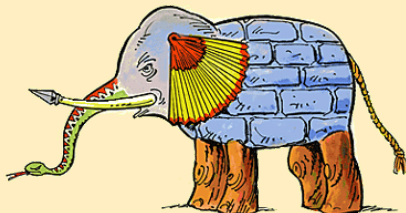


SEE THE WHOLE PICTURE

Idea:

- Embed queries and interpretations into common 2D space
- Show – at a glance! – how they are related

MUTAGENESIS PATTERN MINING



SEE THE WHOLE PICTURE

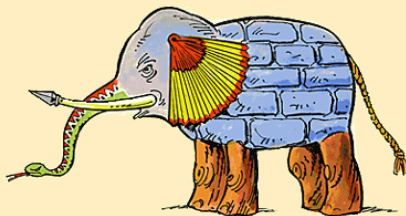
Idea:

- Embed queries and interpretations into common 2D space
- Show – at a glance! – how they are related

However:

- Embedding algorithms need a distance measure

MUTAGENESIS PATTERN MINING



1 RELATIONAL DATA

- Representation
- Algorithms (The High Altitude View)
- **Semantically Grounded Distances**

2 EMBEDDING OF CO-PROVEN QUERIES AND INTERPRETATIONS

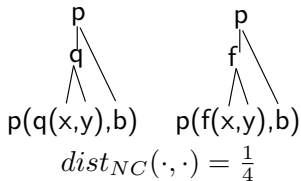
3 EXPERIMENTS

4 CONCLUSION

DISTANCES IN ILP

Typical ILP distances measures

- Based on syntax
- Recursively defined over terms \mathcal{T} and predicates \mathcal{F}



NIENHUYS-CHENG DISTANCE

$$\bigwedge_{t \in \mathcal{T}} dist_{nc}(t, t) = 0$$

$$\bigwedge_{p/n \in \mathcal{F}} \bigwedge_{q/m \in \mathcal{F}} dist_{nc}(p(s_1, \dots, s_n), q(t_1, \dots, t_m)) = 1$$

$$\bigwedge_{p/n \in \mathcal{F}} dist_{nc}(p(s_1, \dots, s_n), p(t_1, \dots, t_n)) = \frac{1}{2n} \sum_{i=1}^n dist_{nc}(s_i, t_i).$$

SYNTAX-BASED DISTANCES NOT “GROUNDED”

- Queries q and p are equivalent with respect to data in E

$$\bigwedge_{e \in E} (q(e) \leftrightarrow p(e)).$$

- ▶ Syntactic difference tells us little about data!



SYNTAX-BASED DISTANCES NOT “GROUNDED”

- Queries q and p are equivalent with respect to data in E

$$\bigwedge_{e \in E} (q(e) \leftrightarrow p(e)).$$

- ▶ Syntactic difference tells us little about data!

- Two queries q and p are not related with respect to the data

$$\bigwedge_{e \in E} ((q(e) \rightarrow \neg p(e)) \wedge (p(e) \rightarrow \neg q(e))).$$

- ▶ $p \equiv \neg q$? Probably not. Property of the data? Probably.



DEFINE SIMILARITY USING DATABASE

QUERY-QUERY SIMILARITY

$$\text{sim}(q_1, q_2) =$$

$$\left| \{e \mid e \in E \wedge q_1(e) \wedge q_2(e)\} \right|$$

- ▶ Co-Proven queries are similar



DEFINE SIMILARITY USING DATABASE

QUERY-QUERY SIMILARITY

$$\text{sim}(q_1, q_2) = \left| \left\{ e \mid e \in E \wedge q_1(e) \wedge q_2(e) \right\} \right|$$

QUERY-INTERPRETATION SIMILARITY

$$\text{sim}(q, e) = \begin{cases} 1 & \text{if } q(e) \\ 0 & \text{sonst.} \end{cases}$$

- ▶ Co-Proven queries are similar
- ▶ Queries are similar to interpretations in which they are true



DEFINE SIMILARITY USING DATABASE

QUERY-QUERY SIMILARITY

$$\text{sim}(q_1, q_2) = \frac{|\{e \mid e \in E \wedge q_1(e) \wedge q_2(e)\}|}{|E|}$$

QUERY-INTERPRETATION SIMILARITY

$$\text{sim}(q, e) = \begin{cases} 1 & \text{if } q(e) \\ 0 & \text{sonst.} \end{cases}$$

- ▶ Co-Proven queries are similar
- ▶ Queries are similar to interpretations in which they are true
- ▶ Can be seen as joint probability when normalized:

$$p_{QQ}(q_1, q_2) = \eta \cdot \text{sim}(q_1, q_2) \text{ and } p_{QE}(q, e) = \nu \cdot \text{sim}(q, e)$$



DEFINE SIMILARITY USING DATABASE

QUERY-QUERY SIMILARITY

$$\text{sim}(q_1, q_2) = \frac{|\{e \mid e \in E \wedge q_1(e) \wedge q_2(e)\}|}{|E|}$$

QUERY-INTERPRETATION SIMILARITY

$$\text{sim}(q, e) = \begin{cases} 1 & \text{if } q(e) \\ 0 & \text{sonst.} \end{cases}$$

- ▶ Co-Proven queries are similar
- ▶ Queries are similar to interpretations in which they are true
- ▶ Can be seen as joint probability when normalized:
 $p_{QQ}(q_1, q_2) = \eta \cdot \text{sim}(q_1, q_2)$ and $p_{QE}(q, e) = \nu \cdot \text{sim}(q, e)$
- ▶ Queries with same co-occurrence can be **removed**



OUTLINE

- 1 RELATIONAL DATA
- 2 EMBEDDING OF CO-PROVEN QUERIES AND INTERPRETATIONS
 - CODE
- 3 EXPERIMENTS
- 4 CONCLUSION

OUTLINE

- 1 RELATIONAL DATA
- 2 EMBEDDING OF CO-PROVEN QUERIES AND INTERPRETATIONS
 - CODE
- 3 EXPERIMENTS
- 4 CONCLUSION

CODE

- CODE (Globerson et al. 2007): template for Co-Occurrence Data Embedding algorithms
- We use instance of CODE



CODE FOR RELATIONAL DATA AND QUERIES

- 1 Define $\Phi(\cdot)$, $\Psi(\cdot)$ placing queries, interpretations in 2D space
- 2 Distance in 2D space reflects co-occurrence probability

$$p_{QQ}(q_1, q_2) \propto \exp(-\|\Phi(q_1) - \Phi(q_2)\|^2)$$
$$\frac{p_{QE}(q, e)}{p(e)} \propto \exp(-\|\Phi(q) - \Psi(e)\|^2)$$

- 3 Maximize log-likelihood of embedding

$$l(\Phi, \Psi) = \sum_{q, e} p_{QE}(q, e) \log p_{QE}(q, e) +$$
$$\eta \sum_{q_1, q_2} p_{QQ}(q_1, q_2) \log p_{QQ}(q_1, q_2)$$



- 1 RELATIONAL DATA
- 2 EMBEDDING OF CO-PROVEN QUERIES AND INTERPRETATIONS
- 3 EXPERIMENTS**
 - Datasets and Pattern Miners
 - Distance vs. Co-Occurrence in the Embedding
 - From Embedding to Visualization
- 4 CONCLUSION

OUTLINE

- 1 RELATIONAL DATA
- 2 EMBEDDING OF CO-PROVEN QUERIES AND INTERPRETATIONS
- 3 EXPERIMENTS**
 - **Datasets and Pattern Miners**
 - Distance vs. Co-Occurrence in the Embedding
 - From Embedding to Visualization
- 4 CONCLUSION

DATASETS AND PATTERN MINERS

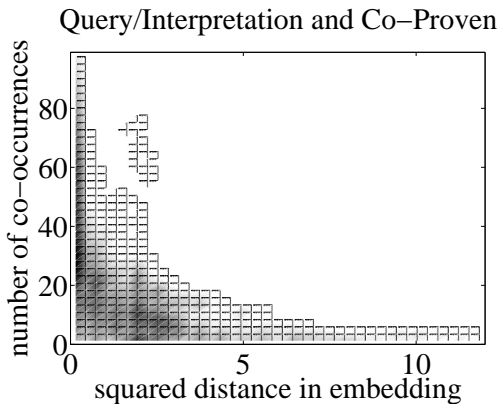
In molecular databases: Interpretations=molecules;
Queries=properties

- Mutagenesis: 188 molecules, ca. 30 atoms per molecule. C-ARMR finds 16 Mio pattern, 505 semantically different.
- AIDS: Sampled 800 molecules, MOLFEA finds 3310 linear fragments (e. g. C-C=C-N-C).
- Estrogen: 232 chemicals, MOLFEA finds 843 different linear fragments.



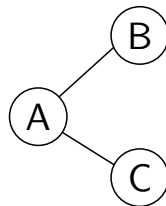
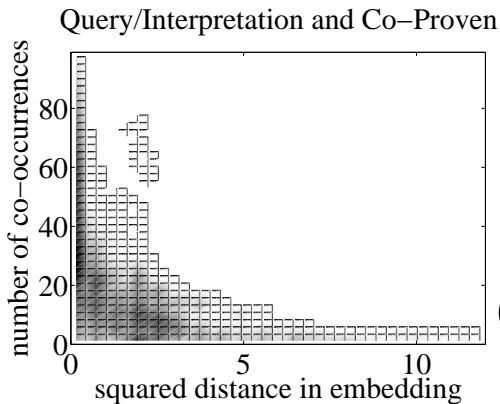
- 1 RELATIONAL DATA
- 2 EMBEDDING OF CO-PROVEN QUERIES AND INTERPRETATIONS
- 3 EXPERIMENTS**
 - Datasets and Pattern Miners
 - Distance vs. Co-Occurrence in the Embedding
 - From Embedding to Visualization
- 4 CONCLUSION

CODE BEST RETAINS CO-OCCURRENCE STATISTICS



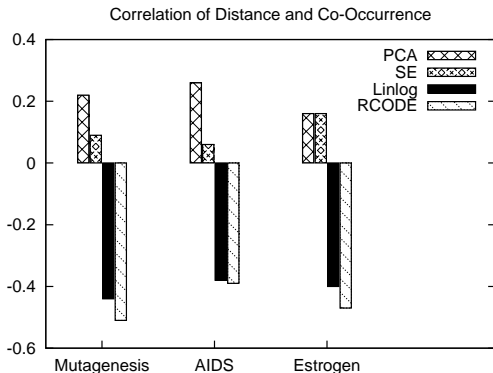
- Cannot be perfect because of intransitivity of co-occurrence

CODE BEST RETAINS CO-OCCURRENCE STATISTICS



- Cannot be perfect because of intransitivity of co-occurrence

CODE BEST RETAINS CO-OCCURRENCE STATISTICS

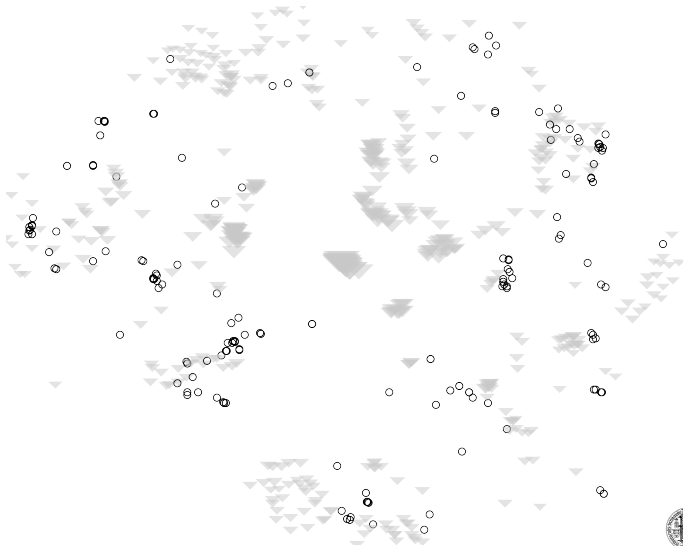


- Cannot be perfect because of intransitivity of co-occurrence



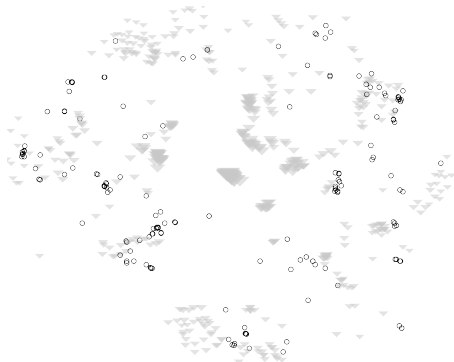
- 1 RELATIONAL DATA
- 2 EMBEDDING OF CO-PROVEN QUERIES AND INTERPRETATIONS
- 3 EXPERIMENTS**
 - Datasets and Pattern Miners
 - Distance vs. Co-Occurrence in the Embedding
 - From Embedding to Visualization**
- 4 CONCLUSION

STARTING POINT



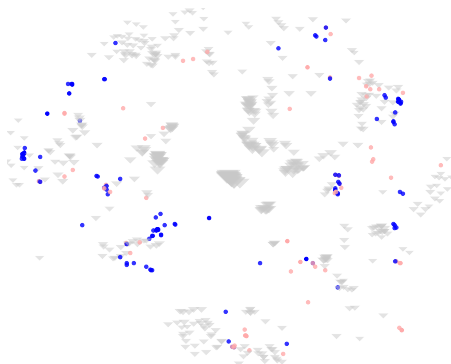
SUPPLYING ADDITIONAL INFORMATION

- Size: Frequency



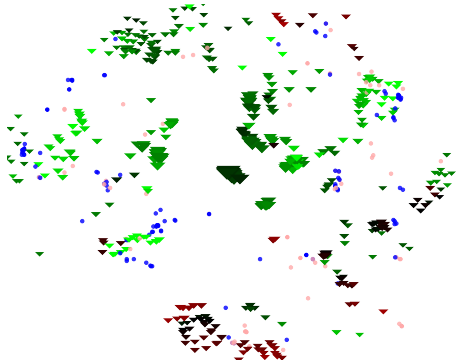
SUPPLYING ADDITIONAL INFORMATION

- Size: Frequency
- Color: interpretation class



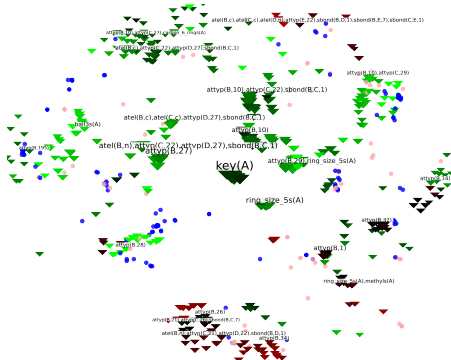
SUPPLYING ADDITIONAL INFORMATION

- Size: Frequency
- Color: interpretation class
- Color: class affinity of queries



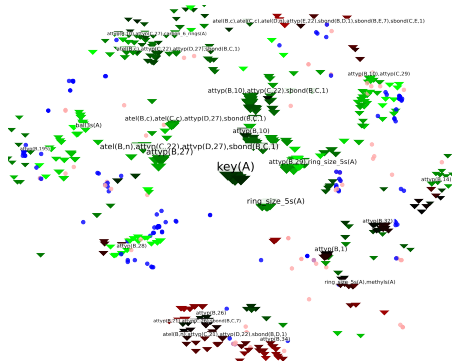
SUPPLYING ADDITIONAL INFORMATION

- Size: Frequency
- Color: interpretation class
- Color: class affinity of queries
- Mark representative queries

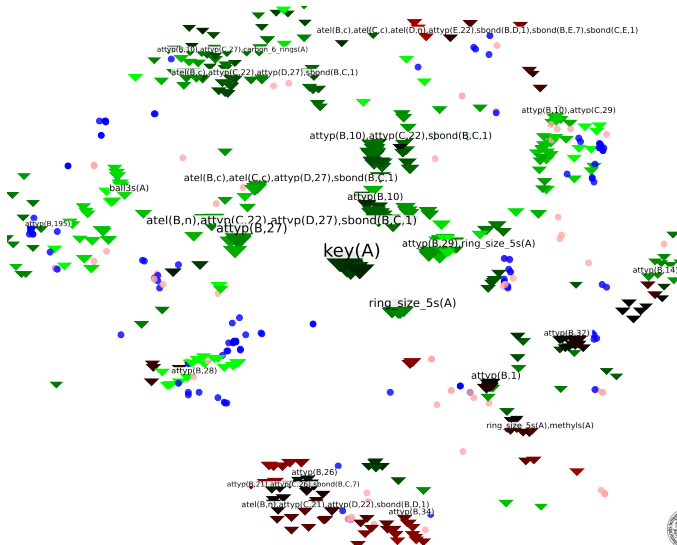


SUPPLYING ADDITIONAL INFORMATION

- Size: Frequency
- Color: interpretation class
- Color: class affinity of queries
- Mark representative queries
- Interactivity to disambiguate embedding



MUTAGENESIS DATASET



OUTLINE

- 1 RELATIONAL DATA
- 2 EMBEDDING OF CO-PROVEN QUERIES AND INTERPRETATIONS
- 3 EXPERIMENTS
- 4 CONCLUSION**

CONCLUSION

- First visualization method for relational data and queries
- Interpretations and queries placed in common Euclidean space
- Use semantically grounded distance measure
- In embedding, co-proven queries are close to each other, interpretations close to their queries
- Developed tools to visualize embedding with side-information



FUTURE WORK

- Improve interactive program to full-fledged Visual Analytics application
- Use graph miner to crawl databases, use for visualization
- Explore use of technique for other applications with large binary feature vectors (Genetic Algorithms, Bag-of-Words, . . .)



Thanks!

(Elephant pictures from http://www.wordinfo.info/words/index/info/view_unit/1/?letter=B&page=3).