# Finding representatives
# in a heterogeneous network

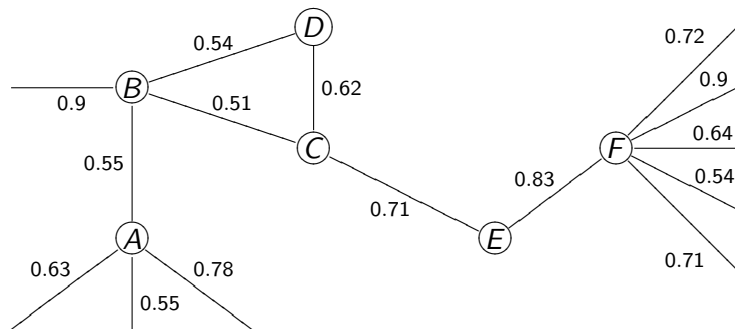Laura Langohr

Department of Computer Science
University of Helsinki

May 19, 2009

## Motivation

- **Finding representative vertices**
- Given a list of 100 vertices
- But only resources to study 10 vertices
- Cluster 100 vertices in 10 clusters
- For each cluster suggest a vertex as representative
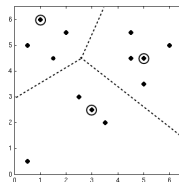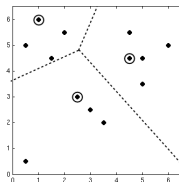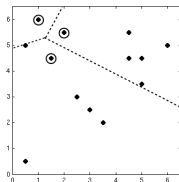
# Example graph

# *K*-medoids

- Clustering method
- Objects are partitioned into *k* clusters
- First, an initial partitioning is created
- The partition is then iteratively improved
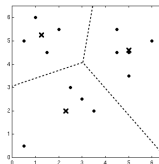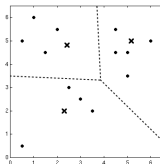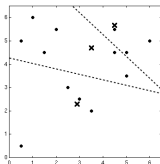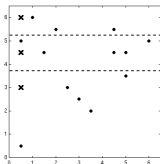- Cluster centers are objects $\rightarrow$ medoids

# Algorithm

1. *K* objects are randomly chosen as medoids
2. Assign remaining objects to the medoid that is the nearest
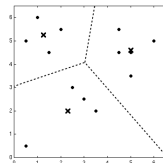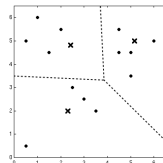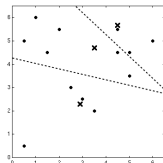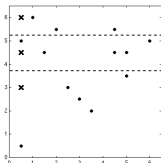3. Calculate new medoid for each cluster

# *K*-means

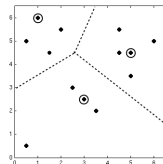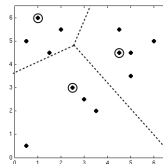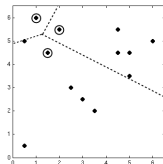- *K*-medoids is similar to *k*-means
- *K*-means uses mean value as cluster center

# $K$-medoids vs $k$-means

# *K*-medoids in a heterogeneous network

- Select few representatives from a large set of vertices
- Representatives should be independent of each other
- Relations between two vertices in a graph $\rightarrow$ link
- Including undiscovered relations
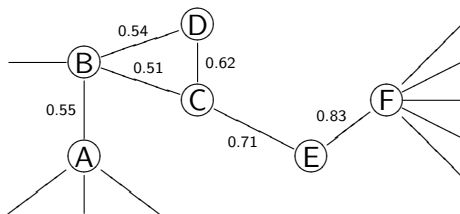- Undiscovered relations are manifested as path(s)

## Measure for link strength

- Probability of a path is the product of the probabilities of the edges along the path

$$g(\mathbf{p}) = \prod_{i=1}^{k} w(e_i)$$

- Probability of the best path between two vertices

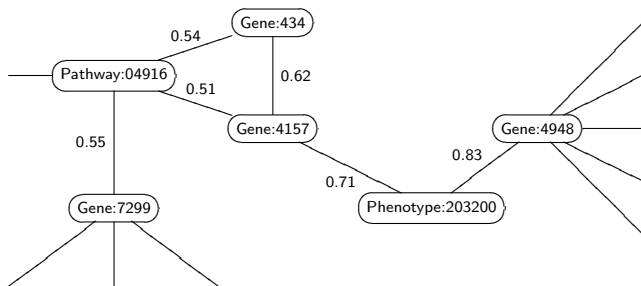$$P_{bp} = \max_{p \in Pa(G,o,o')} g(\mathbf{p})$$

## Algorithm

1. Calculate similarity matrix
2. Choose $k$ objects randomly as initial medoids
3. Assign each remaining object to the most similar medoid
4. Calculate new medoid for each cluster

$$medoid(C_j) = \underset{o \in C_j}{\operatorname{argmax}} \prod_{\substack{o' \in C_j \\ o' \neq o}} P_{bp}(G, o, o')$$

Repeat steps 3. and 4. until clustering converges

# Biomine

- 12 biological databases are integrated
- Over 1 million vertices
- Over 9 million edges



http://biomine.cs.helsinki.fi

# Artificial example

- Three phenotypes, for each three genes
- $k$-medoids with nine genes, and $k = 3$

# Result

# Future Work

- Hierarchical clustering
- Statistical evaluation
- Comparison to an existing method

# Conclusion

- Finding representative vertices, e.g. genes
- *K*-medoids on Biomine
- Example with nine genes is promising