

MINING SUBJECTIVELY INTERESTING PATTERNS IN DATA

PART 1/5: MOTIVATION, BACKGROUND & OUTLINE

Jefrey Lijffijt Tijl De Bie
Ghent University

MOTIVATION

SUBJECTIVITY = KEY

Three motivating examples:

1. **Frequent itemset mining**

– Individually frequent items = probably frequent together

2. **Graph embedding**

– High degree nodes = probably embedded centrally

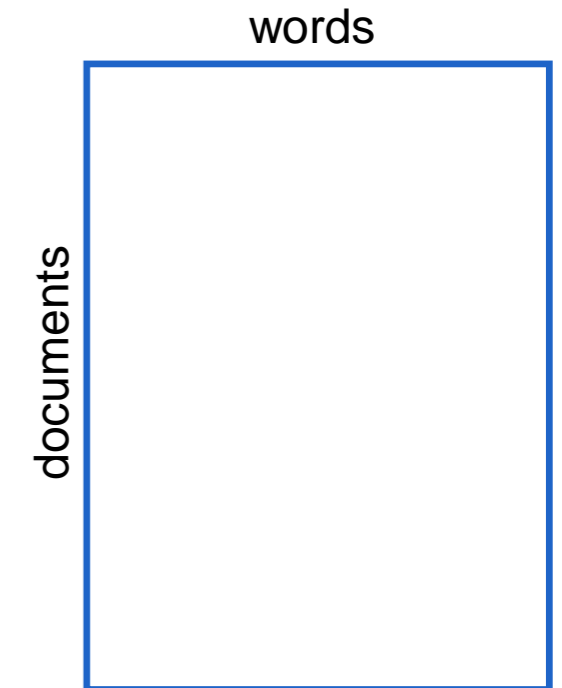
3. **Privacy-preserving data publishing**

– Accounting for background knowledge is key

ASSOCIATION ANALYSIS / ITEMSET MINING

KDD abstracts dataset:

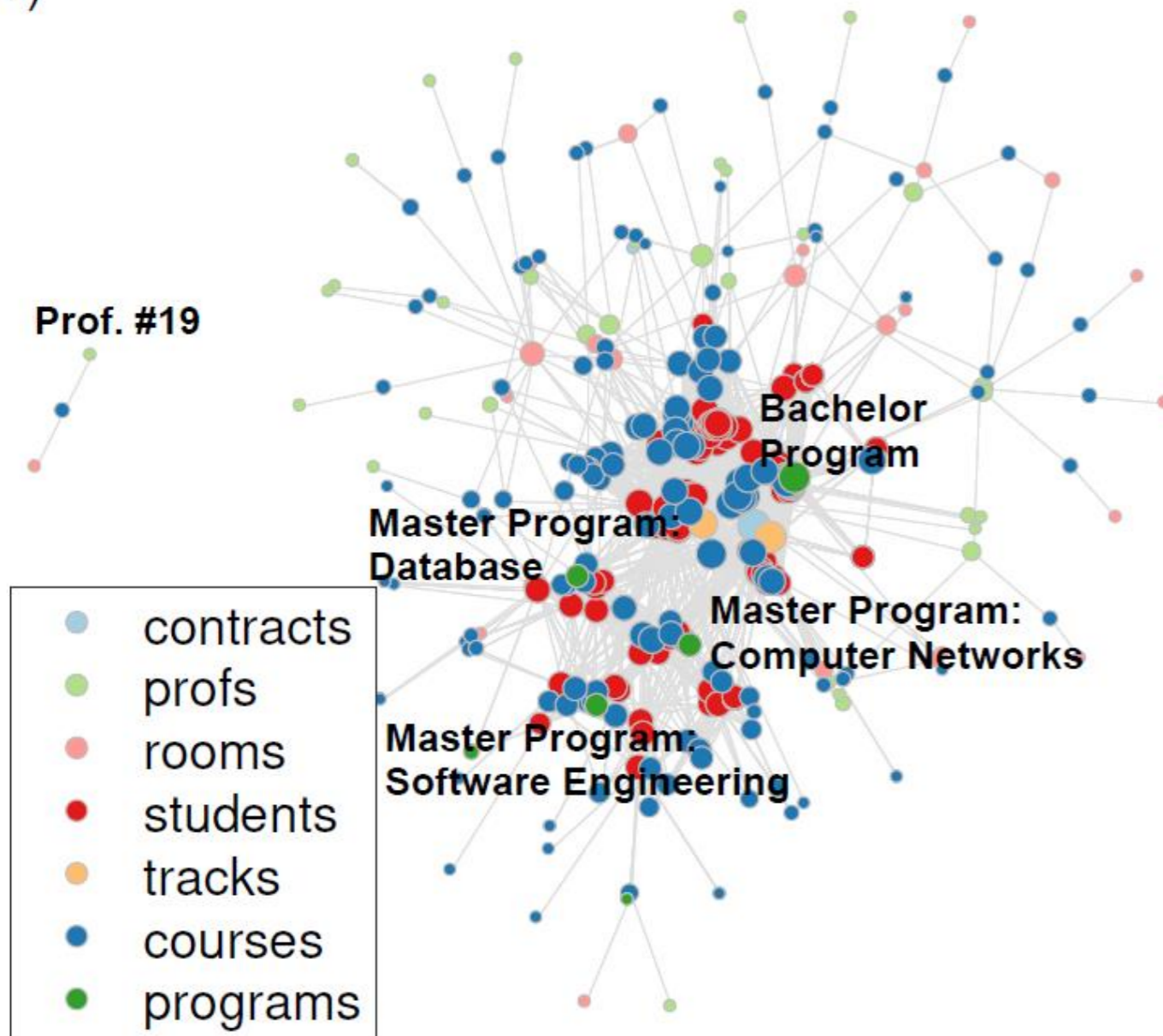
Support x size (area)	#docs
data, paper	389
algorithm, propose	246
data, mine	312
base, method	202
result, show	196
problem	373



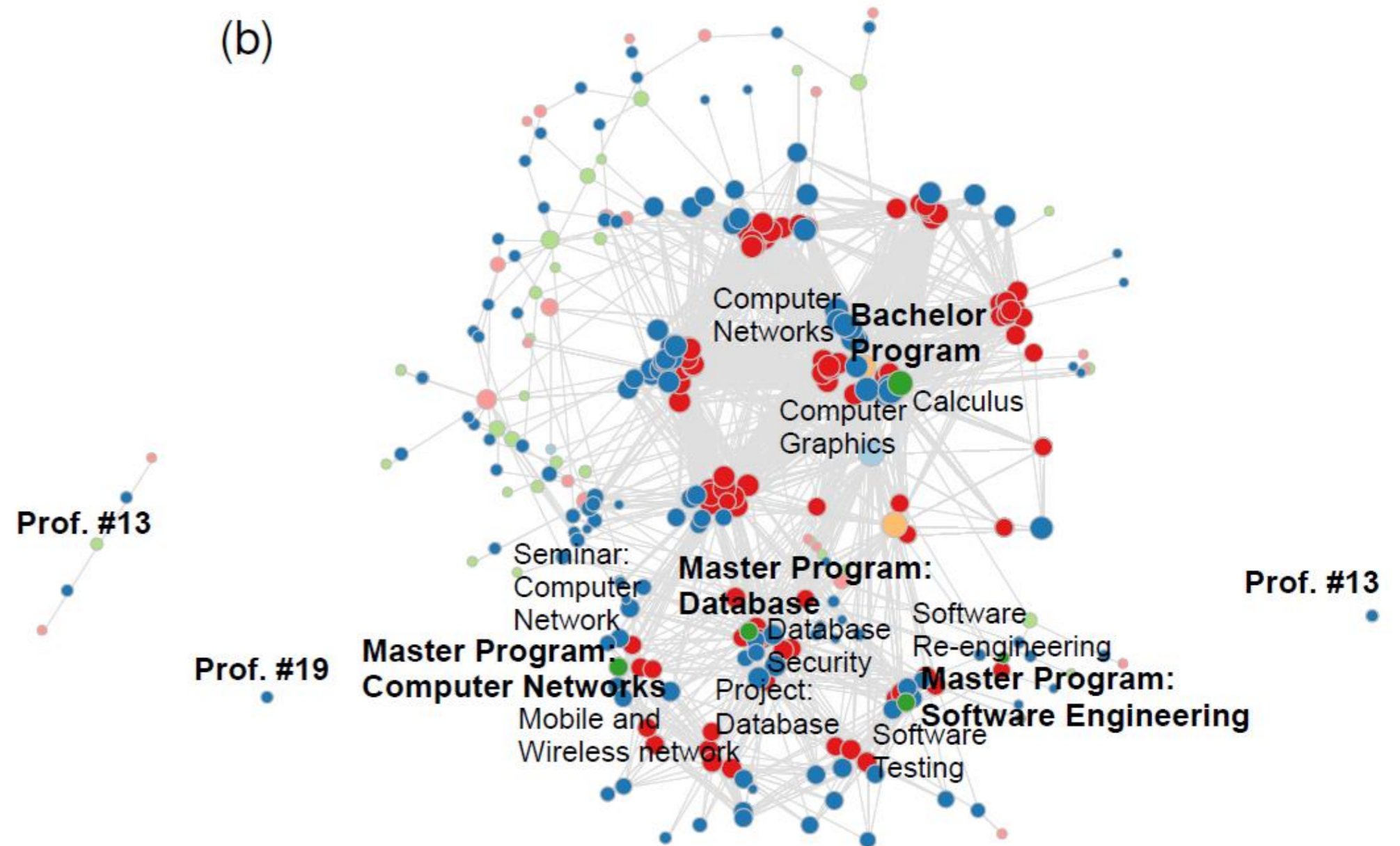
Subjective interestingness ranking	#docs
Prior info on: Row & column sums	
svm, support , machin, vector	25
state, art	39
unlabelled, labelled, supervised , learn	10
associ , rule, mine	36
gene , express	25

CONDITIONAL NETWORK EMBEDDINGS

(a)

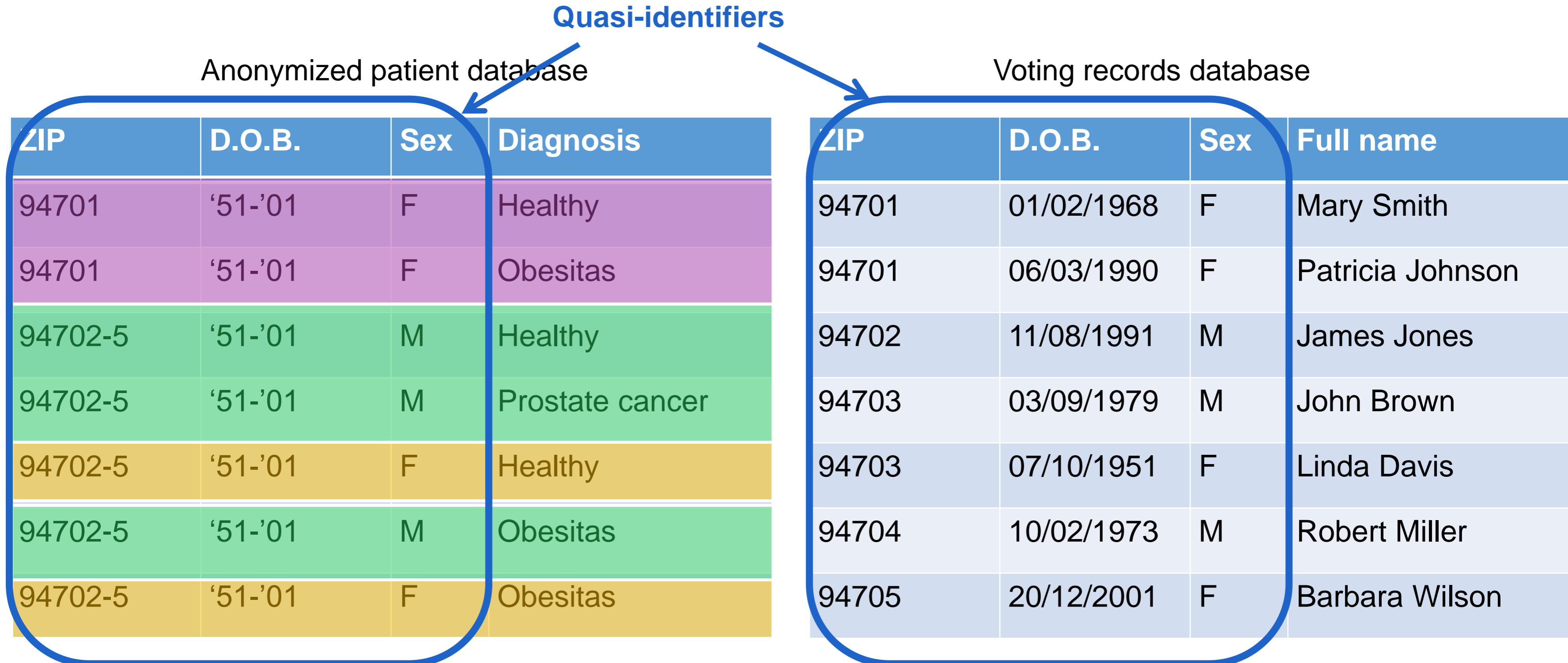


(b)



PRIVACY-PRESERVING DATA PUBLISHING

- **Anonymization** insufficient to protect sensitive attributes (**linkage attack**)
- **Generalization!**



EXPLORING DATA

– The search for interesting *patterns* in *data*

- Association analysis

Frequency, lift, confidence, leverage, coverage,...

- Graph embedding

Node2Vec, Path2Vec, MetaPath2Vec,...

- Privacy-preserving data publishing

Discernibility, generalization height, average group size,...

- Dimensionality reduction

PCA, ICA, projection pursuit, Laplacian Eigenmaps, tSNE, LLE,...

- Clustering

K-means clustering, hierarchical clustering, Mixture of Gaussians, spectral clustering,...

- Community detection

Stochastic block modelling, modularity, k-cores, quasi-cliques, dense subgraphs,...

- ...

– Zillions of

Objective functions

Quality functions

Utility functions

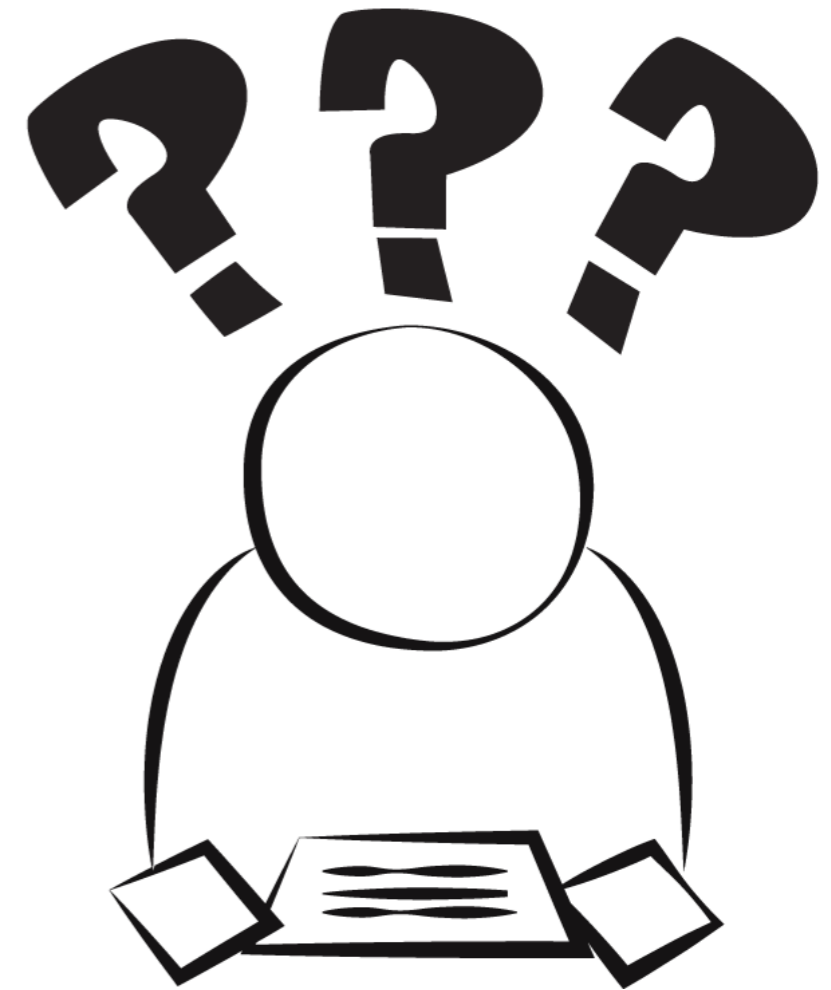
Cost functions

...

← 'Interestingness measures'

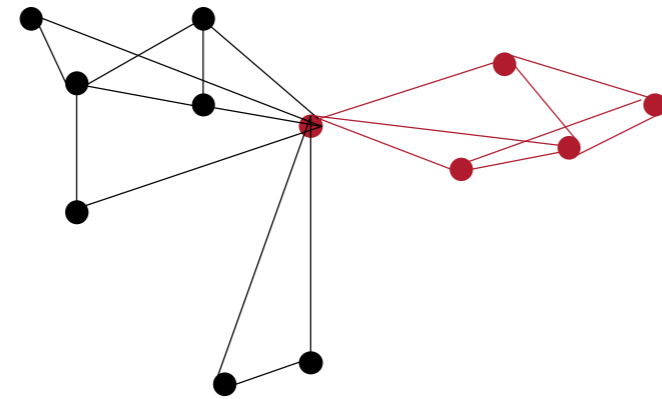
THE CHALLENGE

- Zillions of interestingness measures = good & bad
 - Good: more options!
 - Bad: the trees & the forest...
- Challenge:
 - **Formalise *true* interestingness!**
 - With minimal user interaction
 - Without requiring user expertise



MOTIVATING EXAMPLE

- Community detection:



- What makes for an interesting community?
 - Densely connected?
 - Large?
 - Few neighbours outside community?
 - Unrelated to certain known ‘affiliations’?
 - ...

THE FORSIED APPROACH

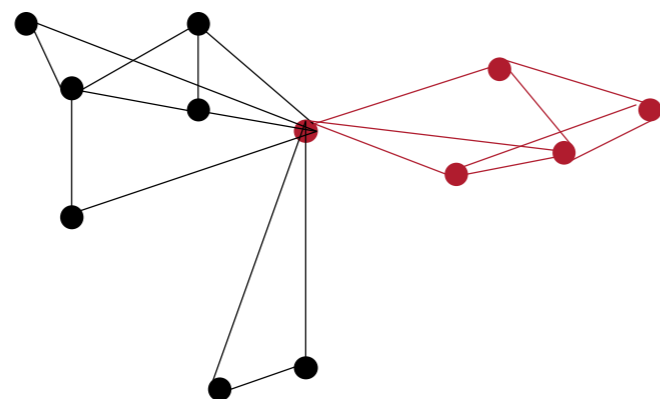


Interestingness(pattern) → Interestingness(pattern, analyst)

Interestingness = **subjective**

MOTIVATING EXAMPLE

- Community detection:



- User states expectations / beliefs
 - Formalized as a ‘background distribution’
- Any ‘pattern’ that contrasts with this and is easy to describe
 - = *subjectively interesting*

OUTLINE

OUTLINE

Part 1: Introduction and motivation

15mins

Part 2: The FORSIED framework

40mins

Part 3: Binary matrices, graphs, and relational data

45mins

COFFEE BREAK

Part 4: Numeric and mixed data

55mins

Part 5: Advanced topics, outlook & conclusions

30mins

Q&A

15mins

Feel free to interrupt
for questions anytime