

PERSONALISED PATTERN MINING

Jefrey Lijffijt

FWO [Pegasus]² Marie Skłodowska-Curie Fellow @ Ghent University

OUTLINE

Formalizing pattern mining

Case studies

- Graph patterns
- Subgroups
- Visualizations

Open challenges

Summary

A FORMALISATION OF PATTERN MINING

WHAT IS A PATTERN?

WHAT IS A PATTERN?

Itemset: co-occurring attributes in binary data

WHAT IS A PATTERN?

Itemset: co-occurring attributes in binary data

Association rule: locally 'predict' values of other attributes

WHAT IS A PATTERN?

Itemset: co-occurring attributes in binary data

Association rule: locally 'predict' values of other attributes

Subgroup: rule predicting specific target variable

WHAT IS A PATTERN?

Itemset: co-occurring attributes in binary data

Association rule: locally 'predict' values of other attributes

Subgroup: rule predicting specific target variable

Community: densely connected subgraph

WHAT IS A PATTERN?

Itemset: co-occurring attributes in binary data

Association rule: locally 'predict' values of other attributes

Subgroup: rule predicting specific target variable

Community: densely connected subgraph

...

WHAT IS A PATTERN?

Itemset: co-occurring attributes in binary data

Association rule: locally 'predict' values of other attributes

Subgroup: rule predicting specific target variable

Community: densely connected subgraph

...

Clustering? Regression? Probabilistic model?

WHAT IS A PATTERN?

Itemset: co-occurring attributes in binary data

Association rule: locally 'predict' values of other attributes

Subgroup: rule predicting specific target variable

Community: densely connected subgraph

...

Clustering? Regression? Probabilistic model?

Patterns specify some aspect of the data

UNIFIED VIEW OF PATTERNS

Lijffijt et al. 2010, De Bie 2011

Data \hat{X}

Data \hat{X}

Data space Ω

- Every $X \in \Omega$ is a value combination what could be \hat{X}
- \hat{X} is constant, but we do not *know* the values

Data \hat{X}

Data space Ω

- Every $X \in \Omega$ is a value combination what could be \hat{X}
- \hat{X} is constant, but we do not *know* the values

Pattern P is a set $P \subseteq \Omega$ s.t. $\hat{X} \in P$

- That is, it may limit the possibilities for \hat{X}

Data \hat{X} ← $G = (V, E)$

Data space Ω

- Every $X \in \Omega$ is a value combination what could be \hat{X}
- \hat{X} is constant, but we do not *know* the values

Pattern P is a set $P \subseteq \Omega$ s.t. $\hat{X} \in P$

- That is, it may limit the possibilities for \hat{X}

Data \hat{X} ← $G = (V, E)$

Data space Ω ← All possible E for this set V

- Every $X \in \Omega$ is a value combination what could be \hat{X}
- \hat{X} is constant, but we do not *know* the values

Pattern P is a set $P \subseteq \Omega$ s.t. $\hat{X} \in P$

- That is, it may limit the possibilities for \hat{X}

Data \hat{X} ← $G = (V, E)$

In case V is known

Data space Ω ← All possible E for this set V

- Every $X \in \Omega$ is a value combination what could be \hat{X}
- \hat{X} is constant, but we do not *know* the values

Pattern P is a set $P \subseteq \Omega$ s.t. $\hat{X} \in P$

- That is, it may limit the possibilities for \hat{X}

Data \hat{X} ← $G = (V, E)$

In case V is known

Data space Ω ← All possible E for this set V

- Every $X \in \Omega$ is a value combination what could be \hat{X}
- \hat{X} is constant, but we do not *know* the values

Pattern P is a set $P \subseteq \Omega$ s.t. $\hat{X} \in P$

- That is, it may limit the possibilities for \hat{X}

Certain edges that do or do not exist

Data \hat{X}

Data space Ω

- Every $X \in \Omega$ is a value combination what could be \hat{X}
- \hat{X} is constant, but we do not *know* the values

Pattern P is a set $P \subseteq \Omega$ s.t. $\hat{X} \in P$

- That is, it may limit the possibilities for \hat{X}

Data \hat{X} ← Data matrix

Data space Ω

- Every $X \in \Omega$ is a value combination what could be \hat{X}
- \hat{X} is constant, but we do not *know* the values

Pattern P is a set $P \subseteq \Omega$ s.t. $\hat{X} \in P$

- That is, it may limit the possibilities for \hat{X}

Data \hat{X} ← Data matrix

Data space Ω ← All possible value combinations

- Every $X \in \Omega$ is a value combination what could be \hat{X}
- \hat{X} is constant, but we do not *know* the values

Pattern P is a set $P \subseteq \Omega$ s.t. $\hat{X} \in P$

- That is, it may limit the possibilities for \hat{X}

Data \hat{X} ← Data matrix

Data space Ω ← All possible value combinations

- Every $X \in \Omega$ is a value combination what could be \hat{X}
- \hat{X} is constant, but we do not *know* the values

Pattern P is a set $P \subseteq \Omega$ s.t. $\hat{X} \in P$

- That is, it may limit the possibilities for \hat{X}

Range, moment, correlation, clustering, ...

INTERESTINGNESS OF PATTERNS

THE AIM OF PATTERN MINING

Suppose: aim is to **inform** the user about the data

THE AIM OF PATTERN MINING

Suppose: aim is to **inform** the user about the data

Form of communication

– Computer \xrightarrow{P} user

THE AIM OF PATTERN MINING

Suppose: aim is to **inform** the user about the data

Form of communication

– Computer \xrightarrow{P} user

Central premise of pattern mining:
Patterns can provide *information*
in a compressed form

THE AIM OF PATTERN MINING

Suppose: aim is to **inform** the user about the data

Form of communication

– Computer \xrightarrow{P} user

Central premise of pattern mining:
Patterns can provide *information*
in a compressed form

How to optimize this communication?

≈ Maximal structure in minimal time

Pattern contains *information*

- What we really learn about data

Pattern has a *descriptive complexity*

- How time consuming is it to internalize

Learning is being surprised

- What we learn depends on what we know/expect

In particular, need to define $\Pr(P)$

- Typically from mass or density function over Ω

Used a lot in data mining, but often implicitly

- Models for random graphs, sequences, matrices, ...

BACKGROUND MODEL

$\Pr(P)$ is called the background model (or distribution)

Many options

- Maximum Entropy model Tatti 2008, De Bie 2011
- Randomization Gionis et al. 2006, Lijffijt et al. 2014
- Other probabilistic models

BACKGROUND MODEL

$\Pr(P)$ is called the background model (or distribution)

Explicitly neutral model
given certain expectations

Many options

– Maximum Entropy model Tatti 2008, De Bie 2011

– Randomization Gionis et al. 2006, Lijffijt et al. 2014

– Other probabilistic models

BACKGROUND MODEL: EXAMPLE

$\Pr(P)$ is called the background model (or distribution)

Given a graph $G = (V, E)$, and Ω all possible E

- Without knowledge: MaxEnt dist. is Uniform over Ω
- Given degree for every vertex: find MaxEnt dist.
through optimization

BACKGROUND MODEL: EXAMPLE

$\Pr(P)$ is called the background model (or distribution)

Given a graph $G = (V, E)$, and Ω all possible E

– Without knowledge: MaxEnt dist. is Uniform over Ω

– Given degree for every vertex: find MaxEnt dist.

through optimization

Cannot do arbitrary constraints, but many types of expectations lead to tractable parameter inference

Recall: pattern P is a set $\hat{X} \in P \subseteq \Omega$

Information Content (IC) of P is $-\log(\text{Pr}(P))$

Description Length (DL) of P depends on syntax

Subjective Interestingness (SI) of P is $\frac{\text{IC}(P)}{\text{DL}(P)}$

Recall: pattern P is a set $\hat{X} \in P \subseteq \Omega$

Derived from
information gain

Information Content (IC) of P is $-\log(\text{Pr}(P))$

Description Length (DL) of P depends on syntax

Subjective Interestingness (SI) of P is $\frac{\text{IC}(P)}{\text{DL}(P)}$

THE FORMALIZATION

De Bie 2011, 2013

Recall: pattern P is a set $\hat{X} \in P \subseteq \Omega$

Derived from
information gain

Information Content (IC) of P is $-\log(\text{Pr}(P))$

Description Length (DL) of P depends on syntax

Subjective Interestingness (SI) of P is $\frac{\text{IC}(P)}{\text{DL}(P)}$

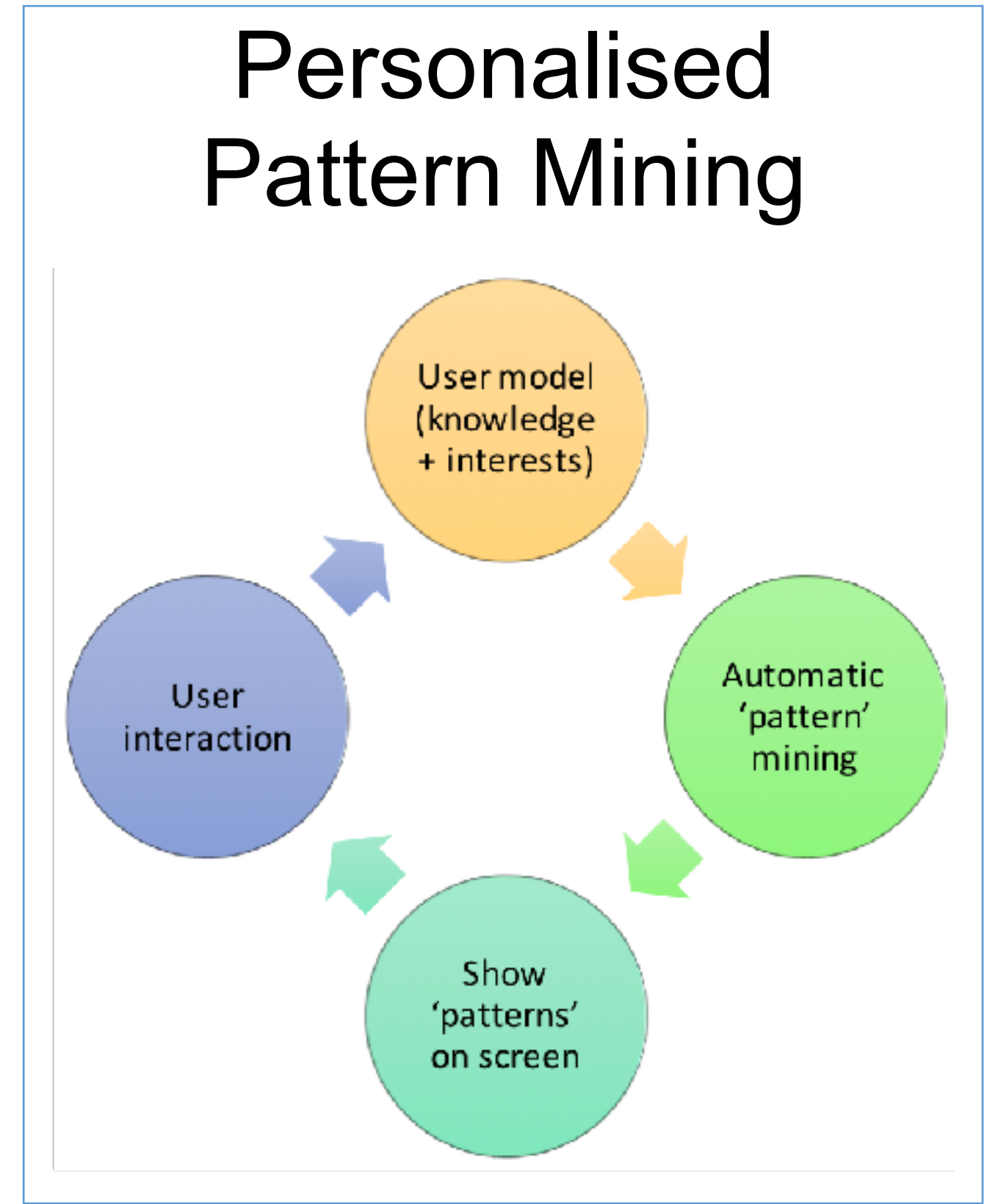
Rate at which we gain information

THE FORSIED PROCESS

1. Background model
2. Pattern syntax
 - IC & DL of patterns
3. Mining
4. Update background model
5. Iterative mining

THE FORSIED PROCESS

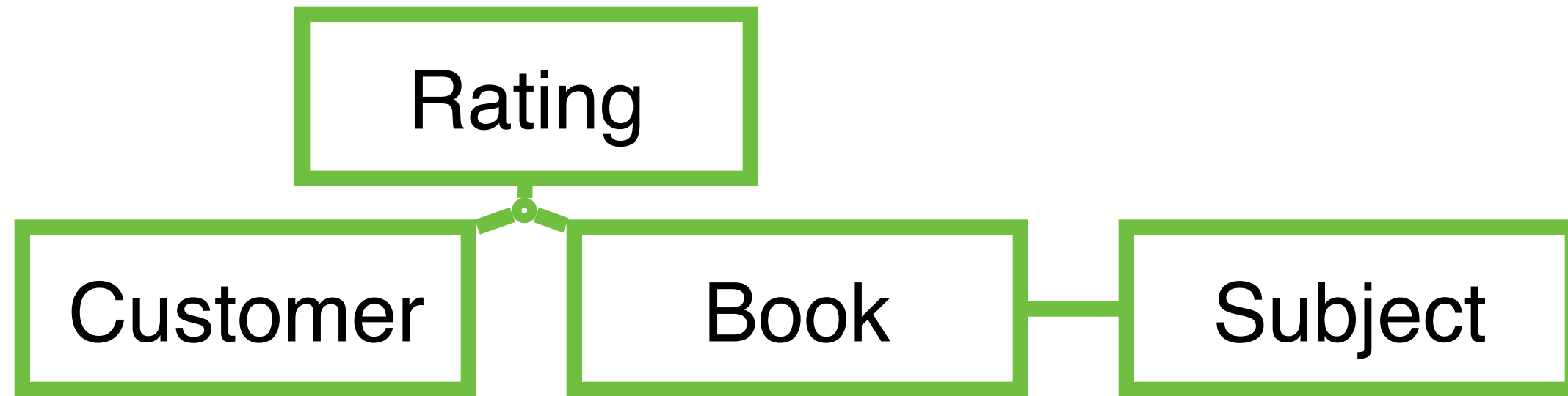
1. Background model
2. Pattern syntax
 - IC & DL of patterns
3. Mining
4. Update background model
5. Iterative mining



CONNECTIONS IN RELATIONAL DATA

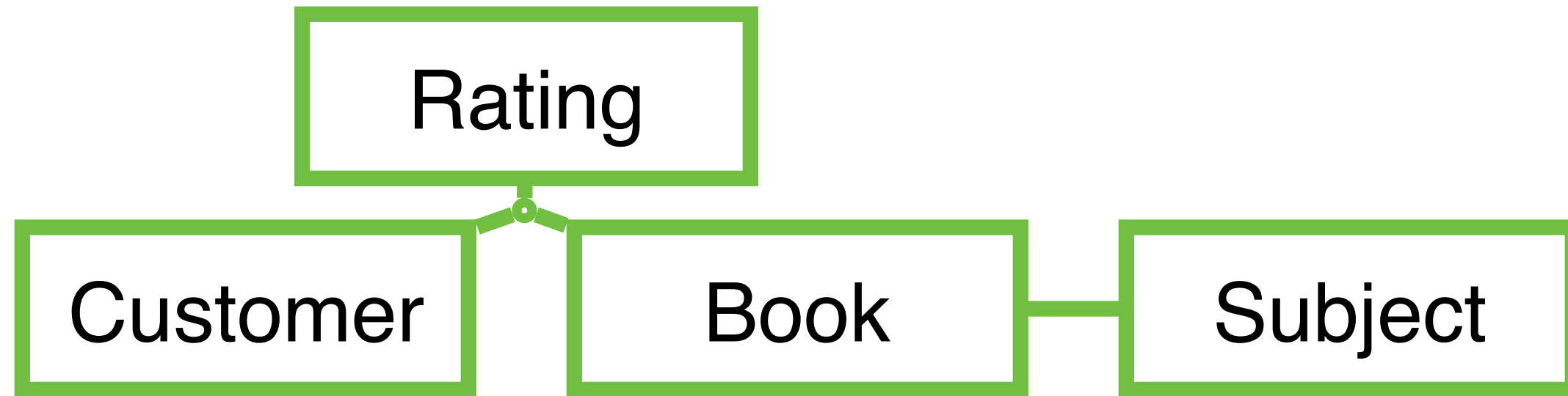
RELATIONAL DATA

Example: Amazon ratings



RELATIONAL DATA

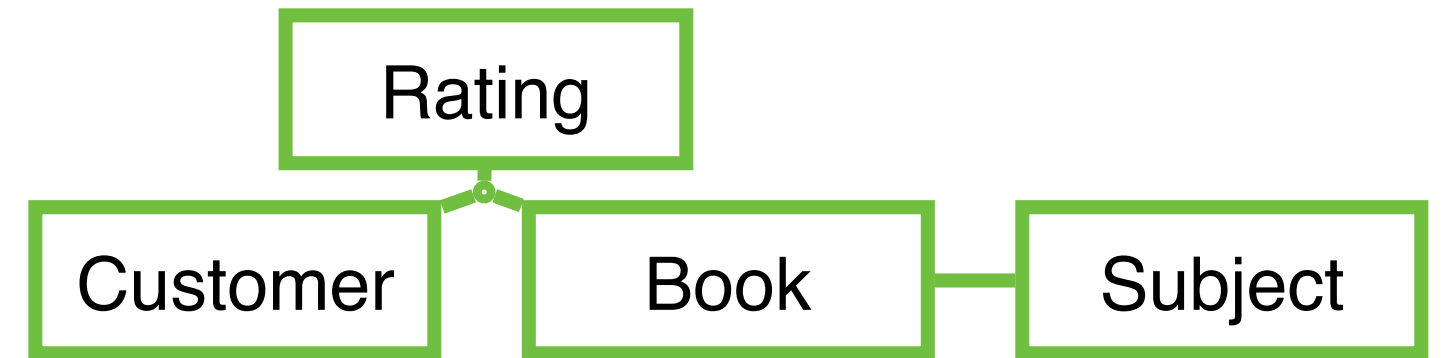
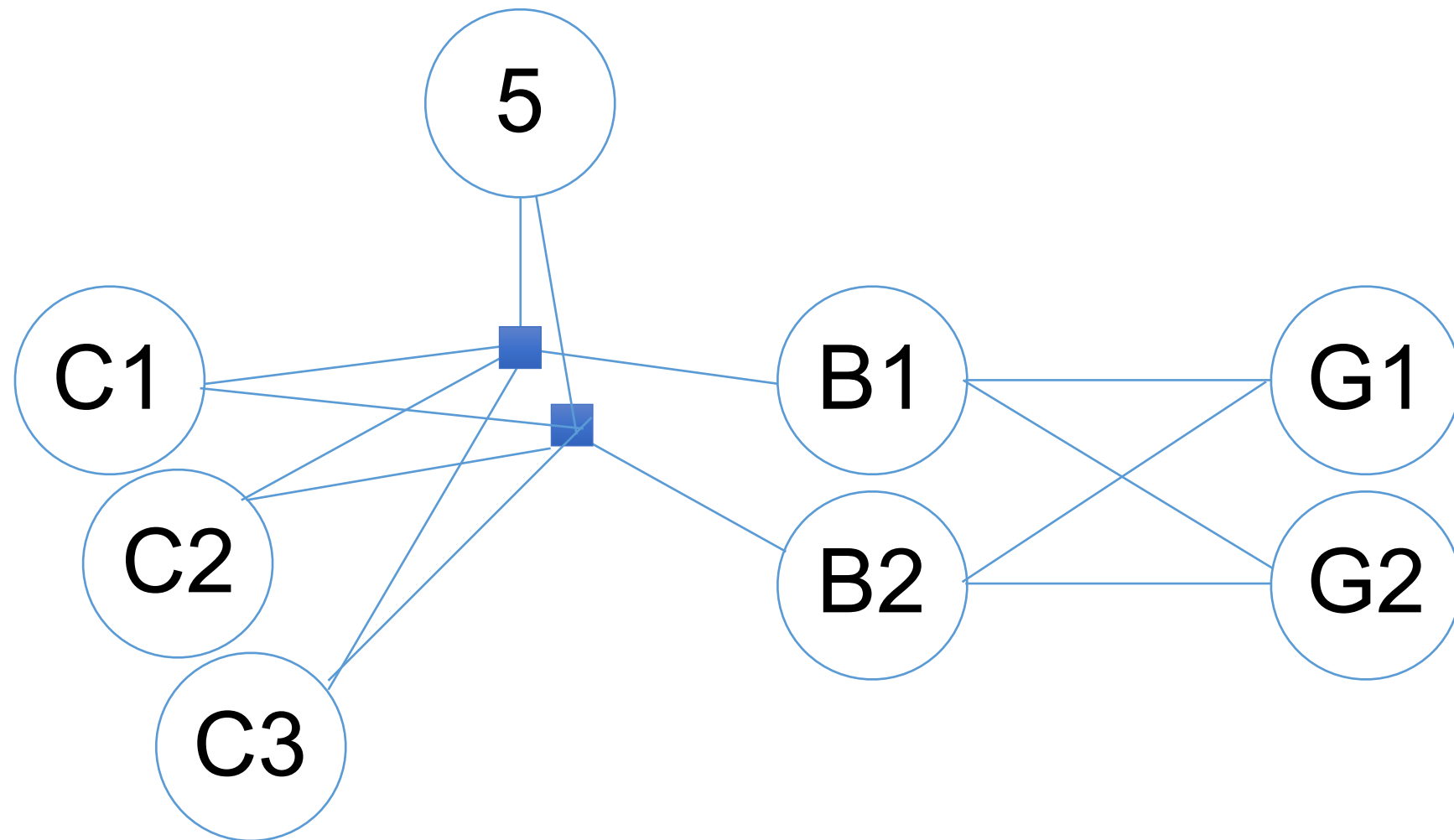
Example: Amazon ratings



Patterns like “Many customers rate all seven volumes of Harry Potter highly, which all belong to Fantasy and Fiction”

PATTERN SYNTAX

Fully connected set of entities



INFORMATION CONTENT

‘Compression’ due to fully connectedness

– Only have to communicate vertices: $P \subseteq V$

If background model factorizes over edges:

$$\text{IC}(P) = \sum_{(v,w) \in P \times P} \Pr((v,w))$$

DESCRIPTION LENGTH

Proportional to number of vertices:

$$DL(P) = \gamma|P| + \eta$$

NB. Absolute values for SI are irrelevant, only ranking matters. Hence either parameter can be constant

MINING

RMiner: exhaustive candidate enumeration Spyropoulou et al. 2014

Constraint programming version: branch and bound Guns et al. 2016

CP works well, but scalability remains an issue

No heuristic algorithms developed yet

CONNECTING TREES IN GRAPHS

EXAMPLE

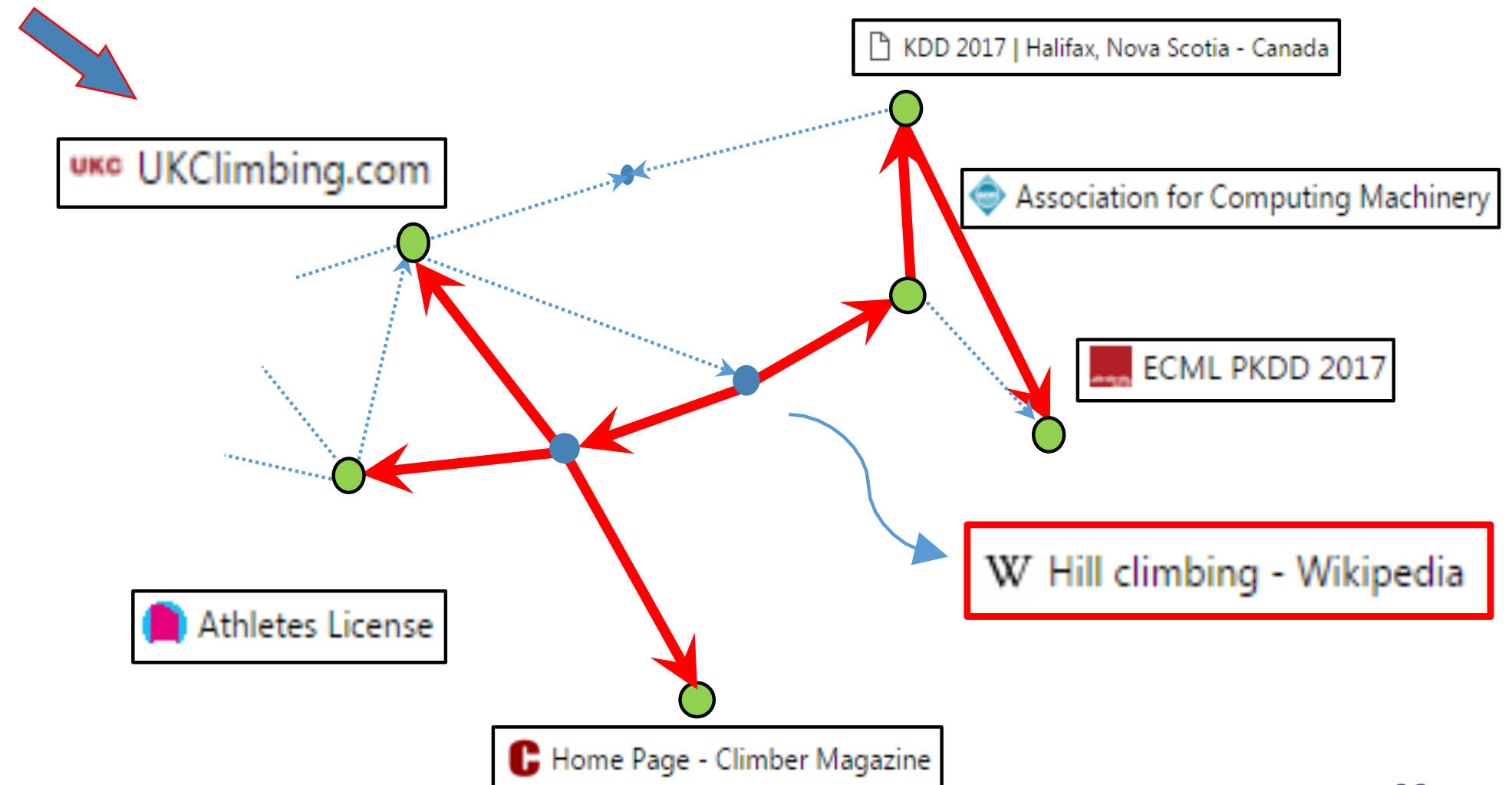
Bookmark Manager

Folders ▾

- Bookmarks bar
- Other bookmarks

Organize ▾

- ECML PKDD 2017
- Florian Adriaens - Outlook Web App
- Athletes License
- KDD 2017 | Halifax, Nova Scotia - Canada
- Association for Computing Machinery
- UKC UKClimbing.com
- Home Page - Climber Magazine



Given a query Q , integer k

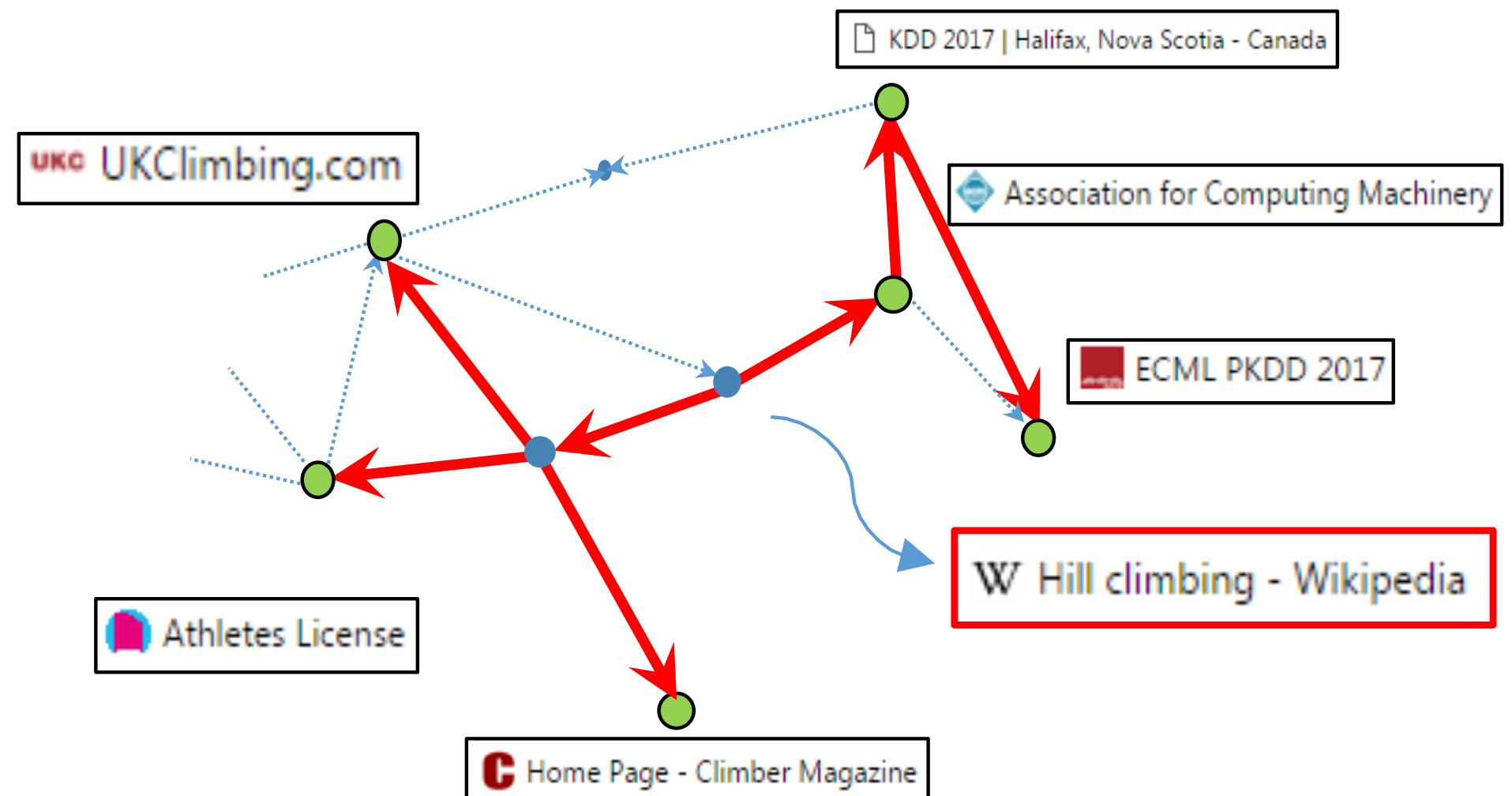
Pattern: arborescence of height $\leq k$ that connects Q

- All leafs must be query nodes
- May not exist

Given a query Q , integer k

Pattern: arborescence of height $\leq k$ that connects Q

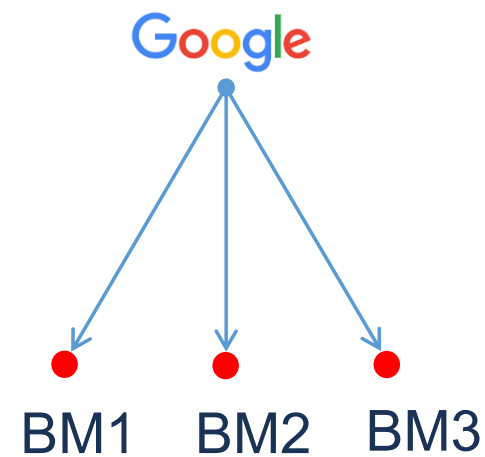
- All leafs must be query nodes
- May not exist



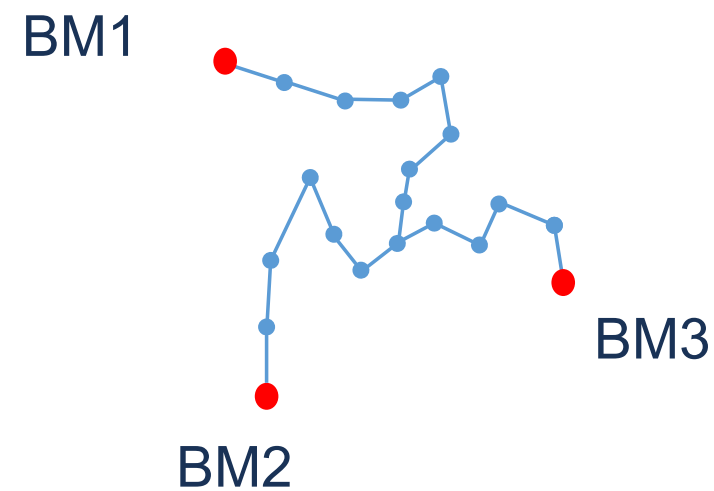
QUANTIFICATION

Idea: tree is informative if

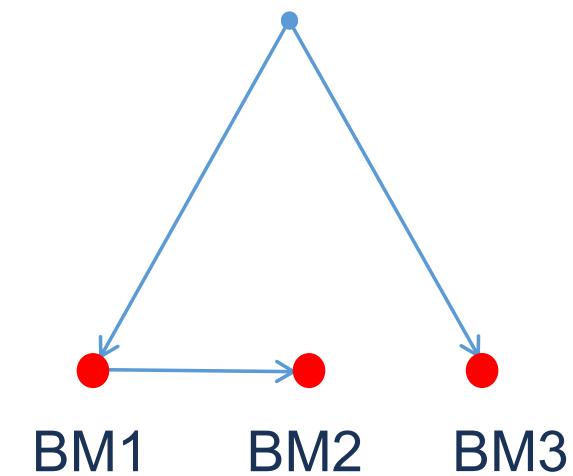
- (links to) vertices are unexpected
- vertices share parents



Non-informative



Non-minimal



Minimal and informative

QUANTIFICATION

Information Content like before

$$\text{IC}(P) = \sum_{(v,w) \in P \times P} \Pr((v,w))$$

Communicate (new) vertices + parent for each vertex

$$\text{DL}(T) = (|V_T| - |Q| + 1) \log(|V| - |Q| + 1) + |V_T| \log(|V_T| + 1)$$

Heuristic construction of trees

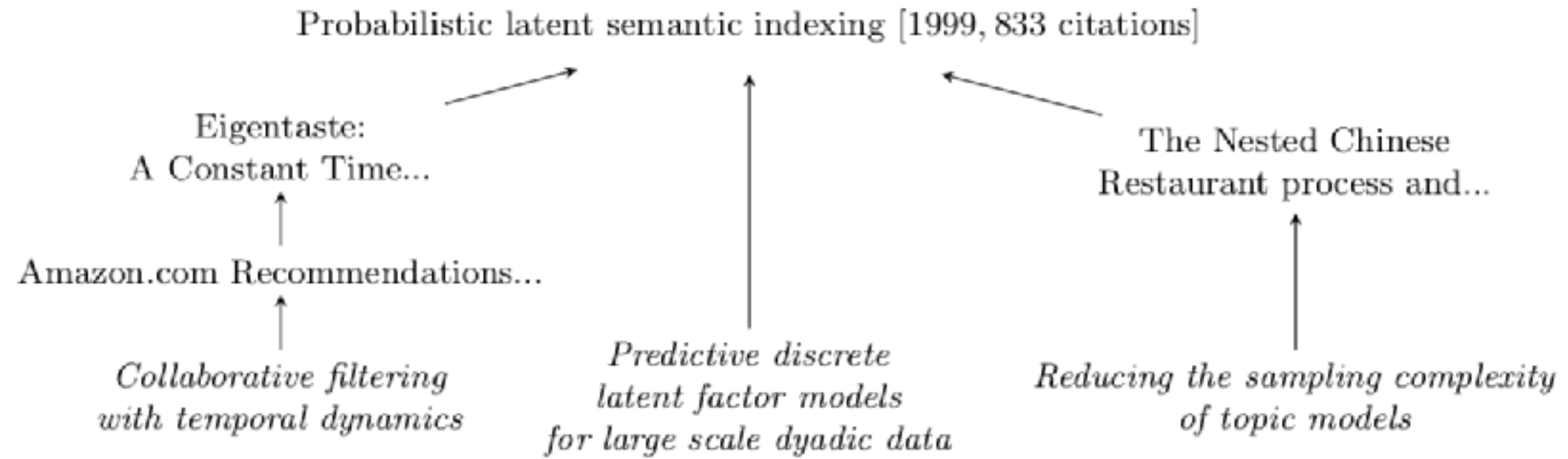
- Not so straightforward to greedily construct a tree

Background distribution for growing graph:

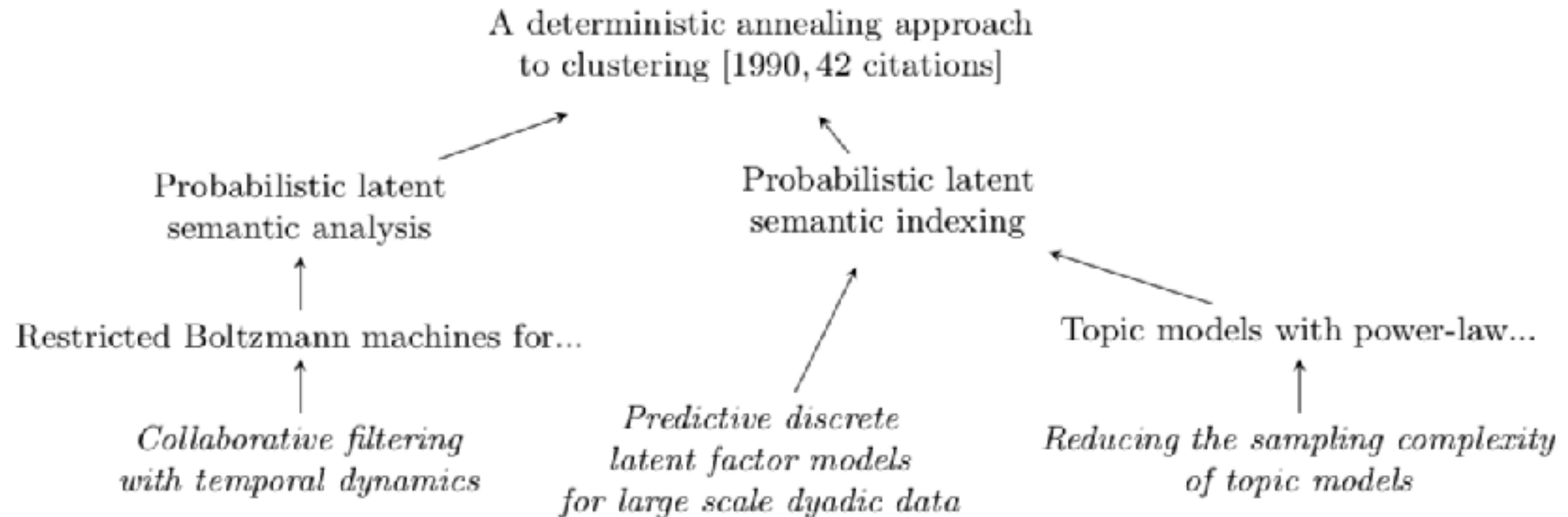
- Vertices can only have links to older vertices
- Parameter inference through intelligent grouping

EXAMPLE RESULT: KDD BEST PAPERS

Without
prior:

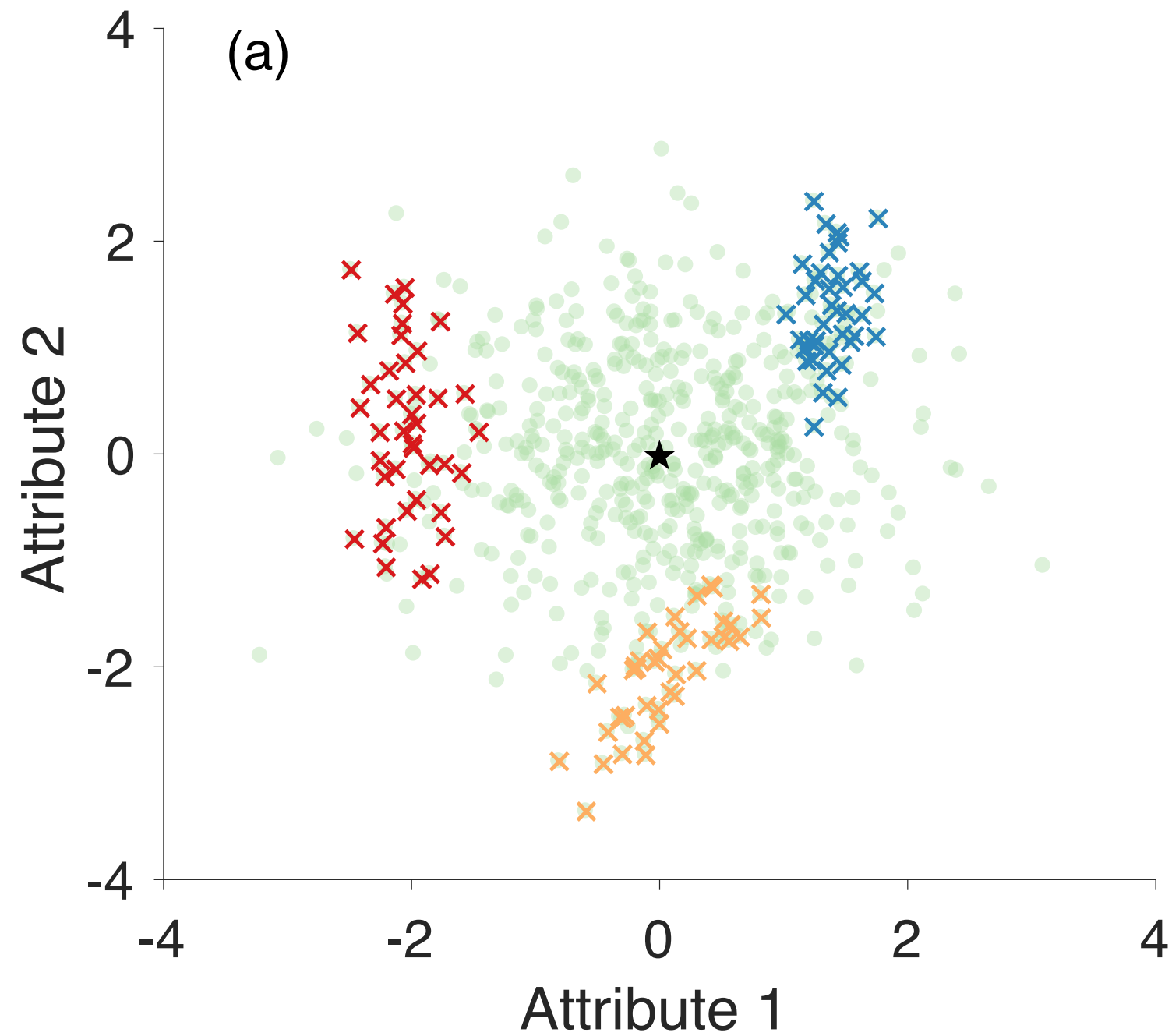


With degree
and time:

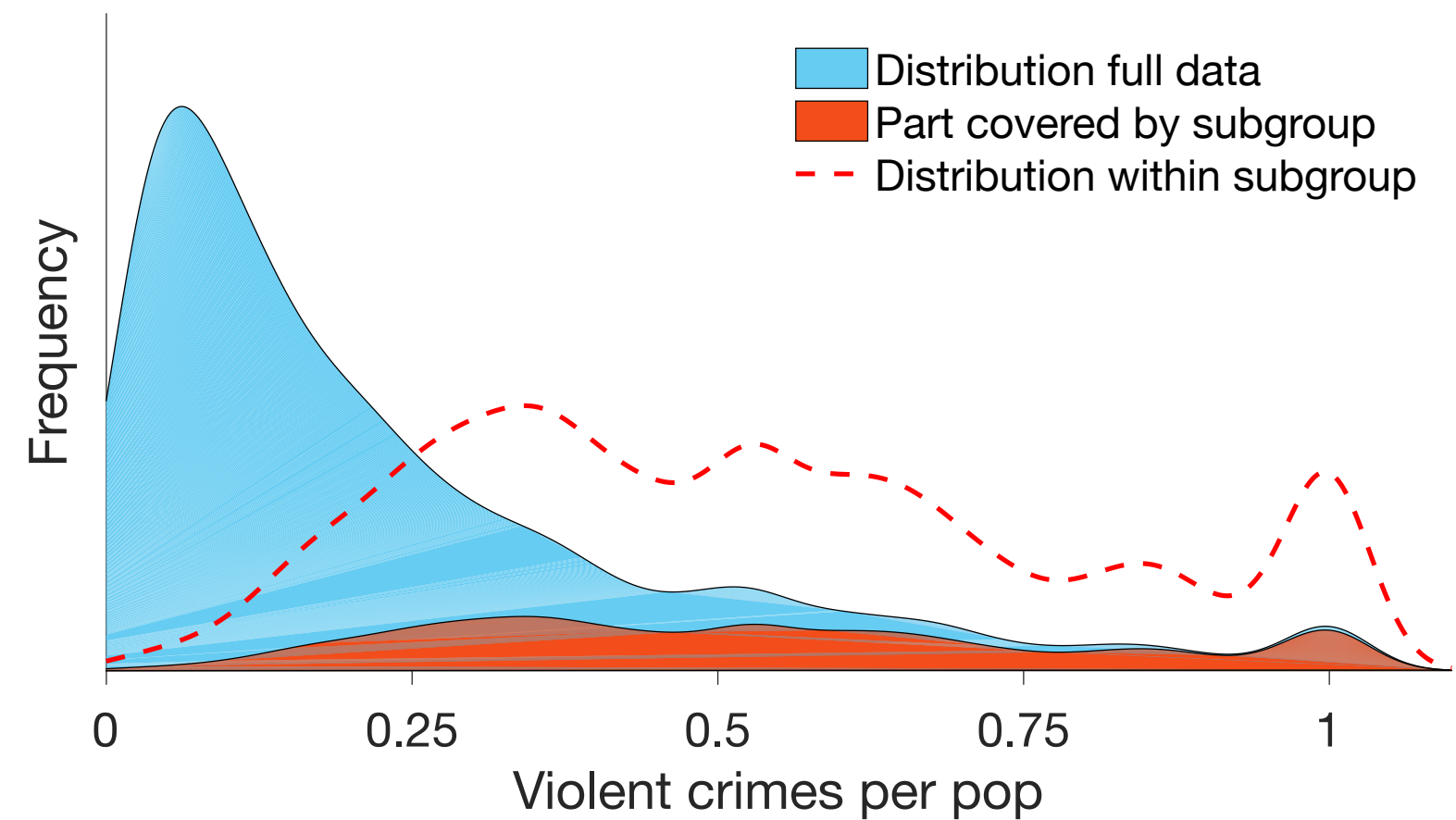


SUBGROUP DISCOVERY IN NUMERIC DATA

EXAMPLE



Condition $Z \rightarrow$ High crime rate

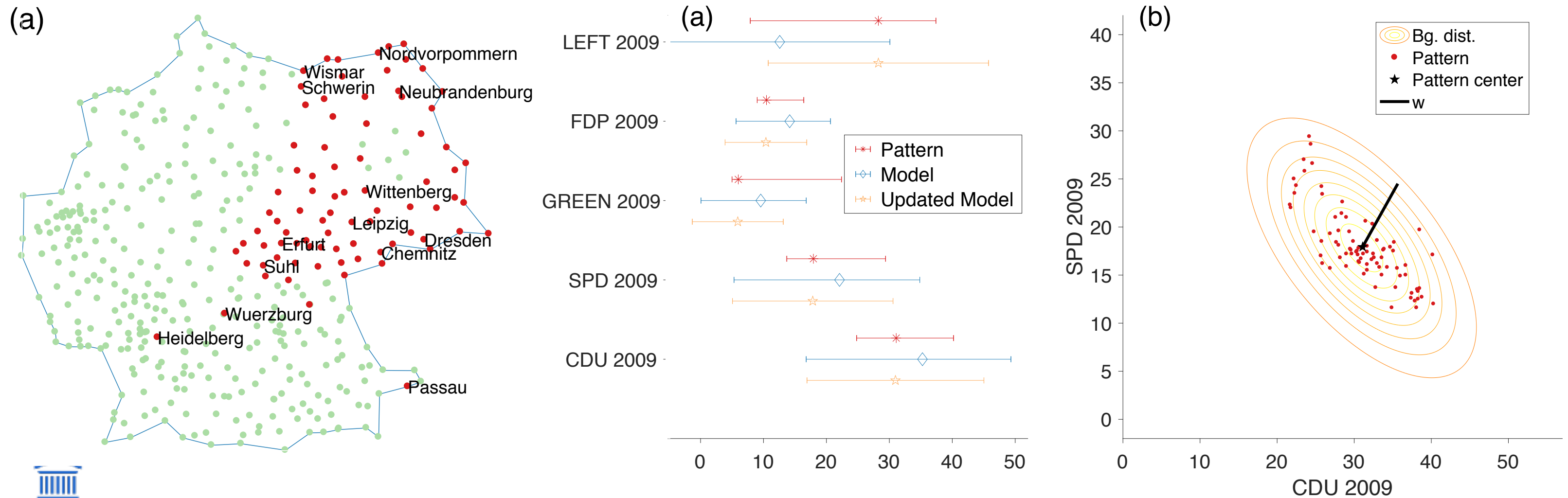


Condition $Z \rightarrow$ means are $f(\hat{X}_Z)$

Condition $Z \rightarrow$ variance in direction w is $g(\hat{X}_Z, w)$

Condition $Z \rightarrow$ means are $f(\hat{X}_Z)$

Condition $Z \rightarrow$ variance in direction w is $g(\hat{X}_Z, w)$



Condition: areas with few children

QUANTIFICATION

Surprisal of the location

$$\text{IC}(f(\hat{X}_Z)) = \log((2\pi)^d |\Sigma_Z|) / 2 + \left(f(\hat{X}_Z - \mu_Z) \right)^T \Sigma_Z^{-1} \left(f(\hat{X}_Z - \mu_Z) \right) / 2$$

Surprisal of the spread

$$\text{IC}(g(\hat{X}_Z, w)) \approx \log(2^{\frac{m}{2}} \Gamma(m/2)) + \alpha - (m/2 - 1) \log\left(\left(g(\hat{X}_Z, w - \beta) / \alpha\right) + \left(g(\hat{X}_Z, w - \beta) / (2\alpha)\right)\right)$$

Description length

$$\text{DL}(\cdot) = \gamma |\text{Cond}|(+\gamma) + \eta$$

QUANTIFICATION

Surprisal of the location

$$\text{IC}(f(\hat{X}_Z)) =$$

$$\log \left((2\pi)^d |\Sigma_Z| \right) / 2 + \left(f(\hat{X}_Z - \mu_Z) \right)^T \Sigma_Z^{-1} \left(f(\hat{X}_Z - \mu_Z) \right) / 2$$

Too much algebra to explain

Surprisal of the spread

$$\text{IC}(g(\hat{X}_Z, w)) \approx \log \left(2^{\frac{m}{2}} \Gamma(m/2) \right)$$

$$+ \alpha - (m/2 - 1) \log \left(\left(g(\hat{X}_Z, w - \beta) \right) / \alpha \right) + \left(g(\hat{X}_Z, w - \beta) \right) / (2\alpha)$$

Description length

$$\text{DL}(\cdot) = \gamma |\text{Cond}| (+\gamma) + \eta$$

QUANTIFICATION

Surprisal of the location

$$\text{IC}(f(\hat{X}_Z)) =$$

$$\log \left((2\pi)^d |\Sigma_Z| \right) / 2 + \left(f(\hat{X}_Z - \mu_Z) \right)^T \Sigma_Z^{-1} \left(f(\hat{X}_Z - \mu_Z) \right) / 2$$

Too much algebra to explain

Surprisal of the spread

$$\text{IC}(g(\hat{X}_Z, w)) \approx \log \left(2^{\frac{m}{2}} \Gamma(m/2) \right)$$

$$+ \alpha - (m/2 - 1) \log \left(\left(g(\hat{X}_Z, w - \beta) \right) / \alpha \right) + \left(g(\hat{X}_Z, w - \beta) \right) / (2\alpha)$$

Description length

$$\text{DL}(\cdot) = \gamma |\text{Cond}| (+\gamma) + \eta$$

Weight times number of conditions + constant
(+spread has one term more than location)

Finding subgroups:

Beam-search through off-the-shelf toolbox (Cortana)

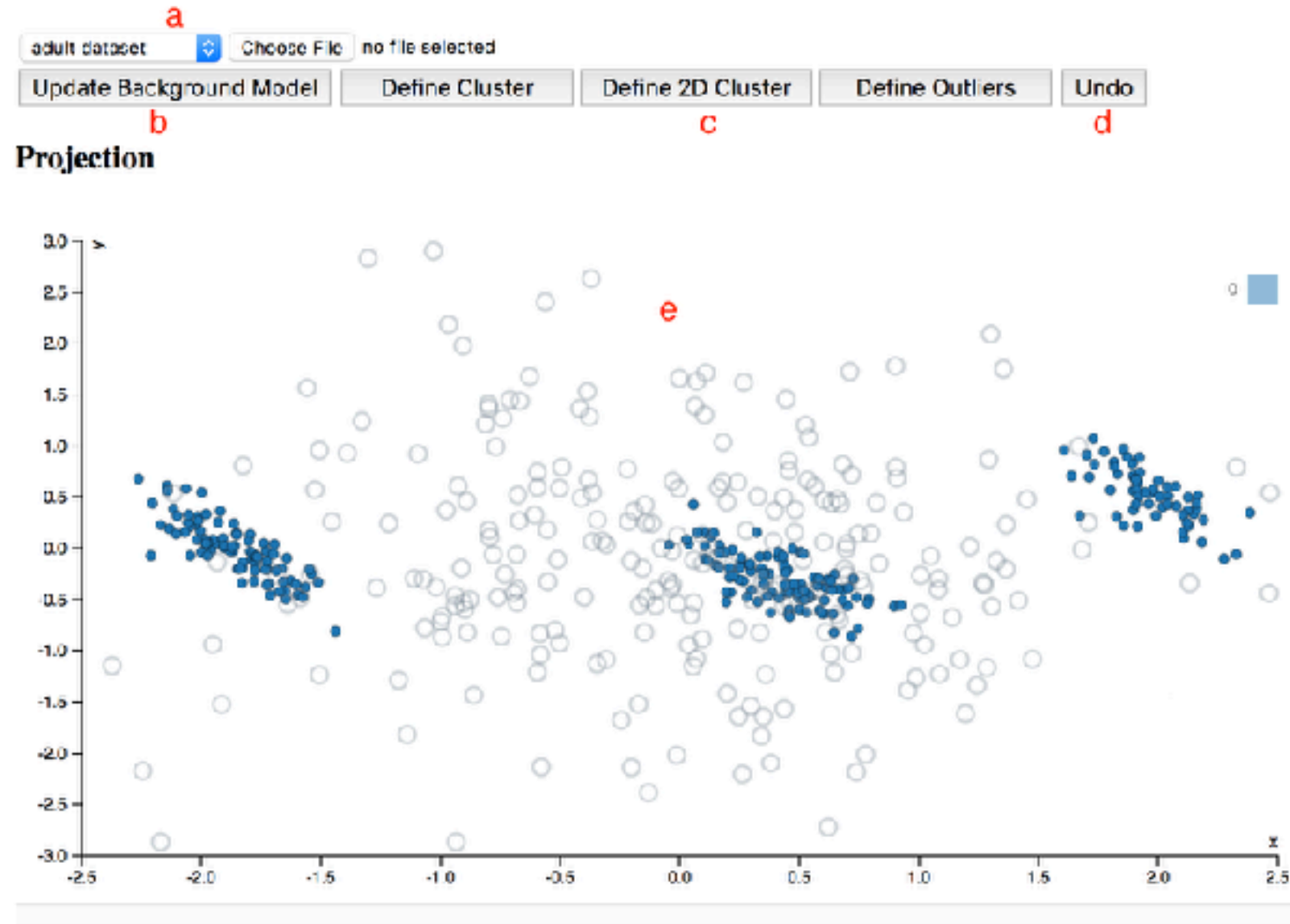
Finding weight vector for spread pattern:

Manifold learning toolbox (ManOpt)

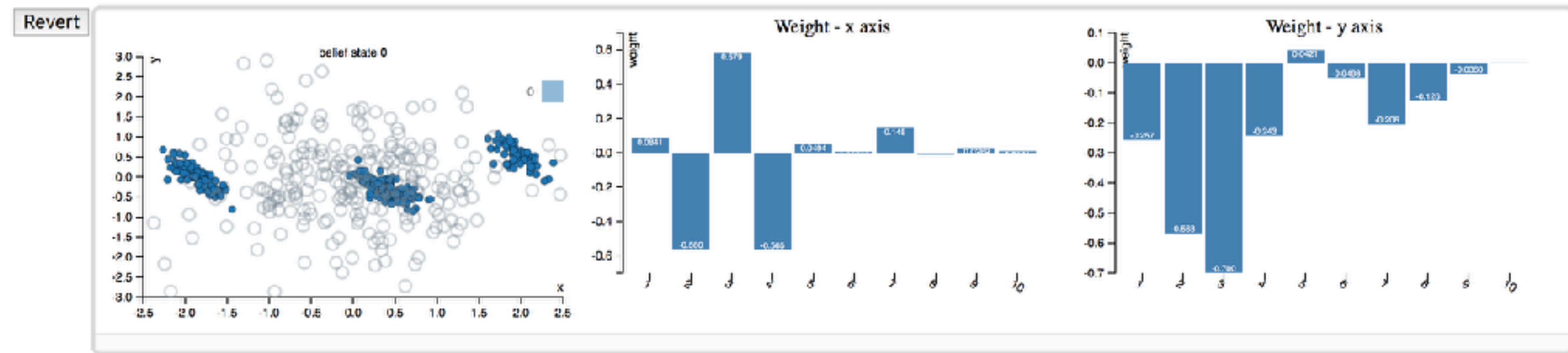
VISUALIZATIONS AS PATTERNS

Synthetic Data Case Study

De Bie et al. 2016
Kang et al. 2016
Puolamäki et al. 2016



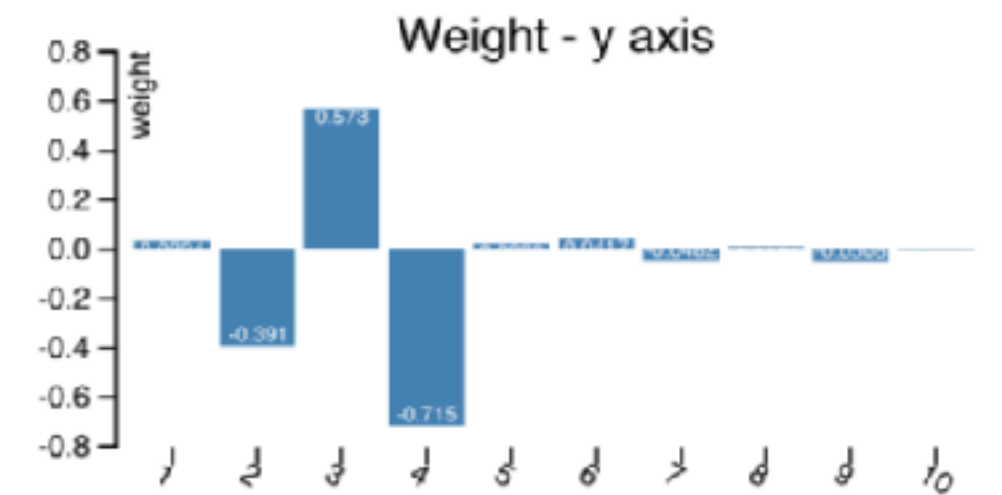
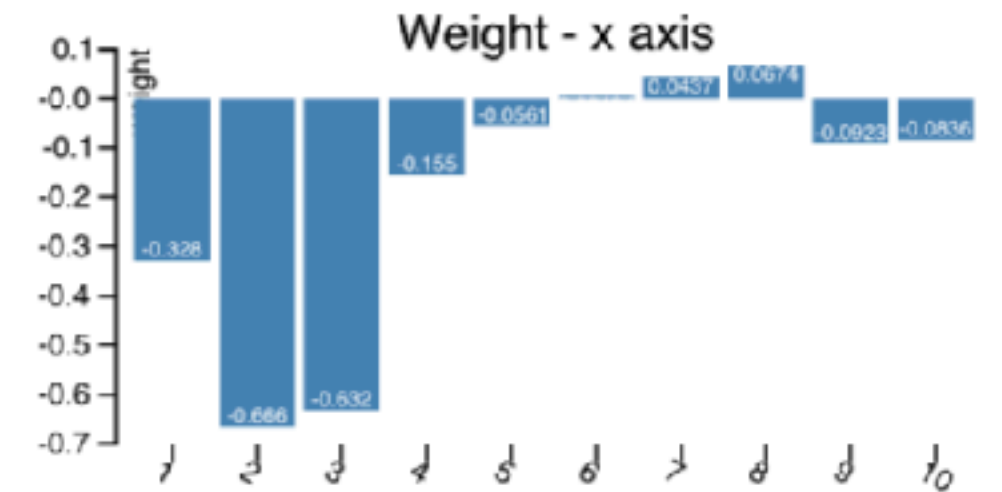
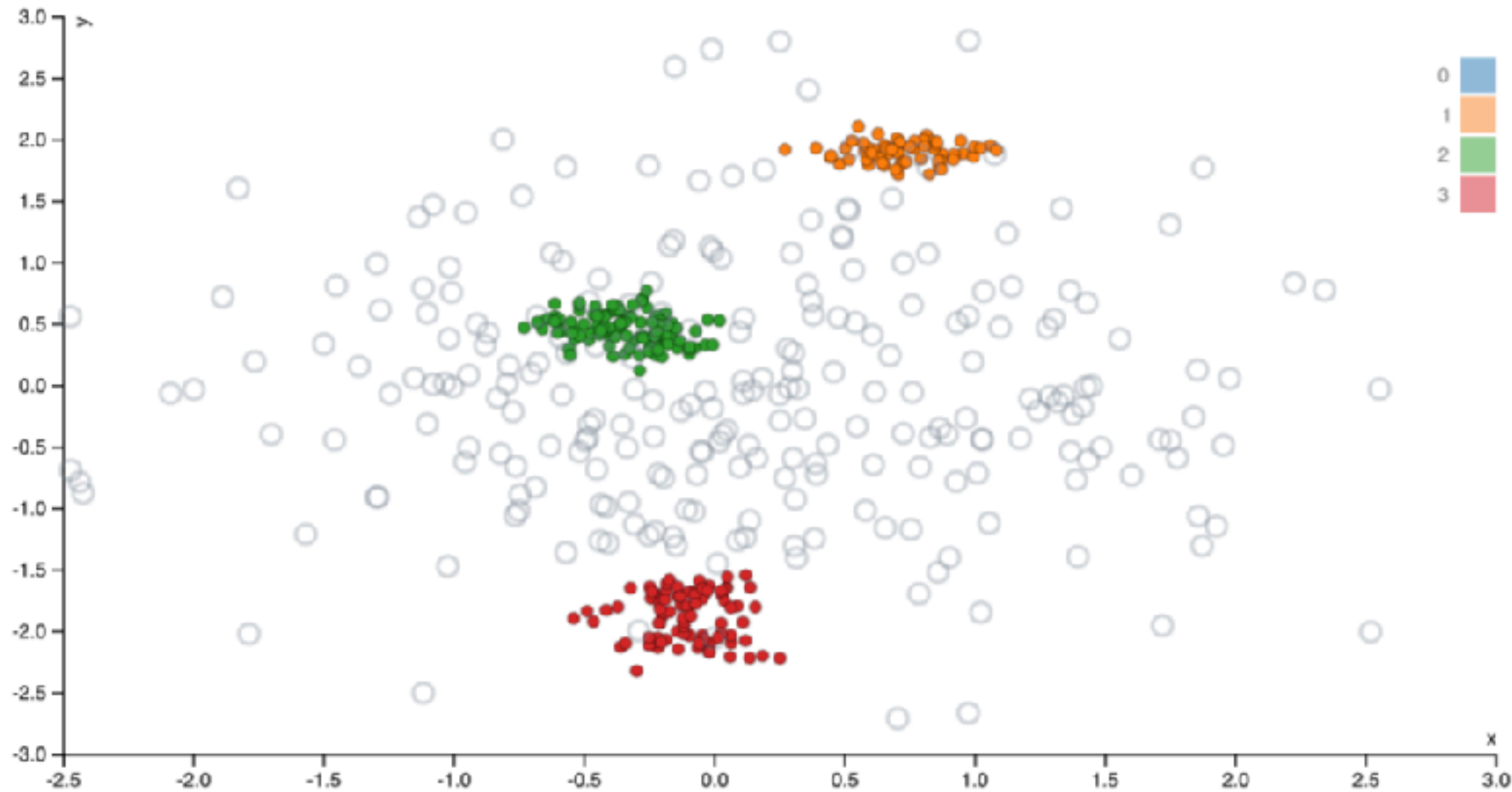
Snapshots



PATTERN SYNTAX

De Bie et al. 2016
Kang et al. 2016
Puolamäki et al. 2016

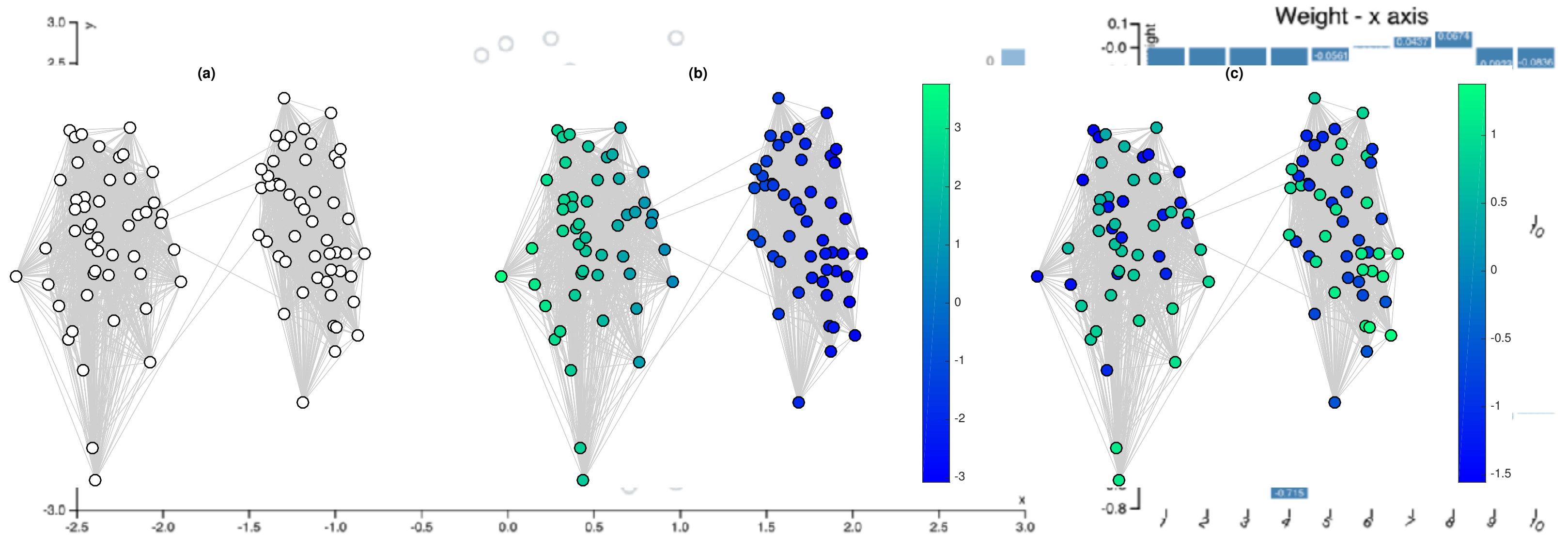
Weight vector and (approximate) projection



PATTERN SYNTAX

De Bie et al. 2016
Kang et al. 2016
Puolamäki et al. 2016

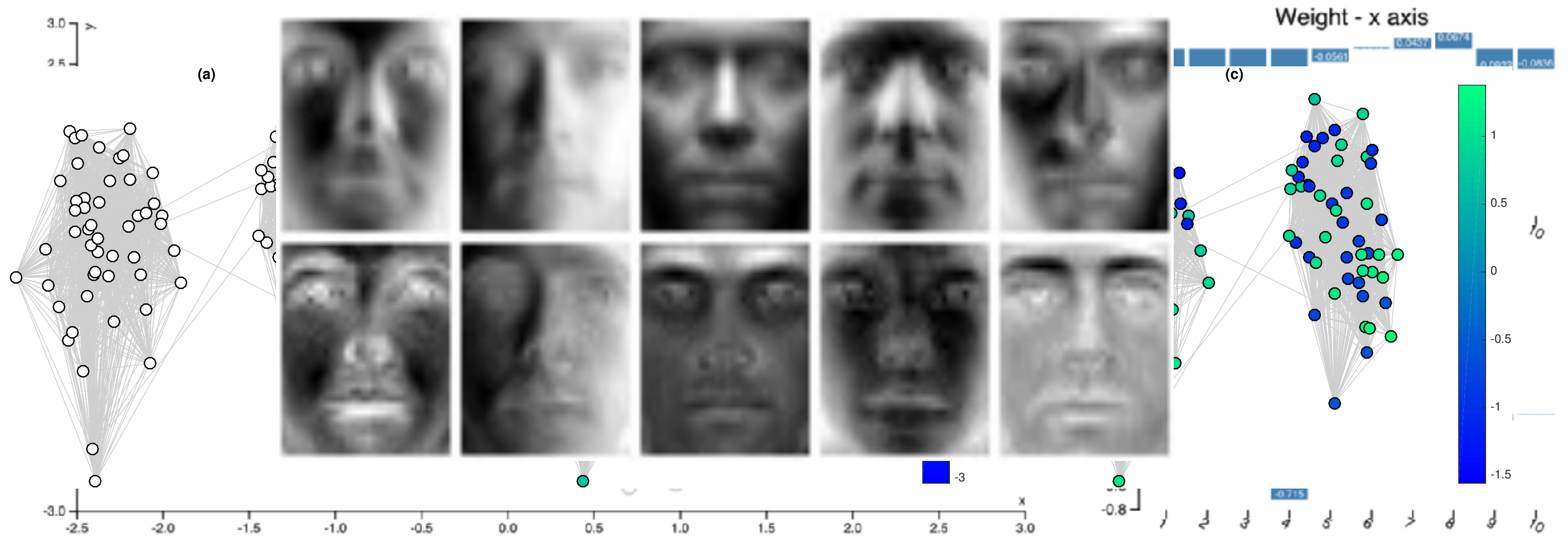
Weight vector and (approximate) projection



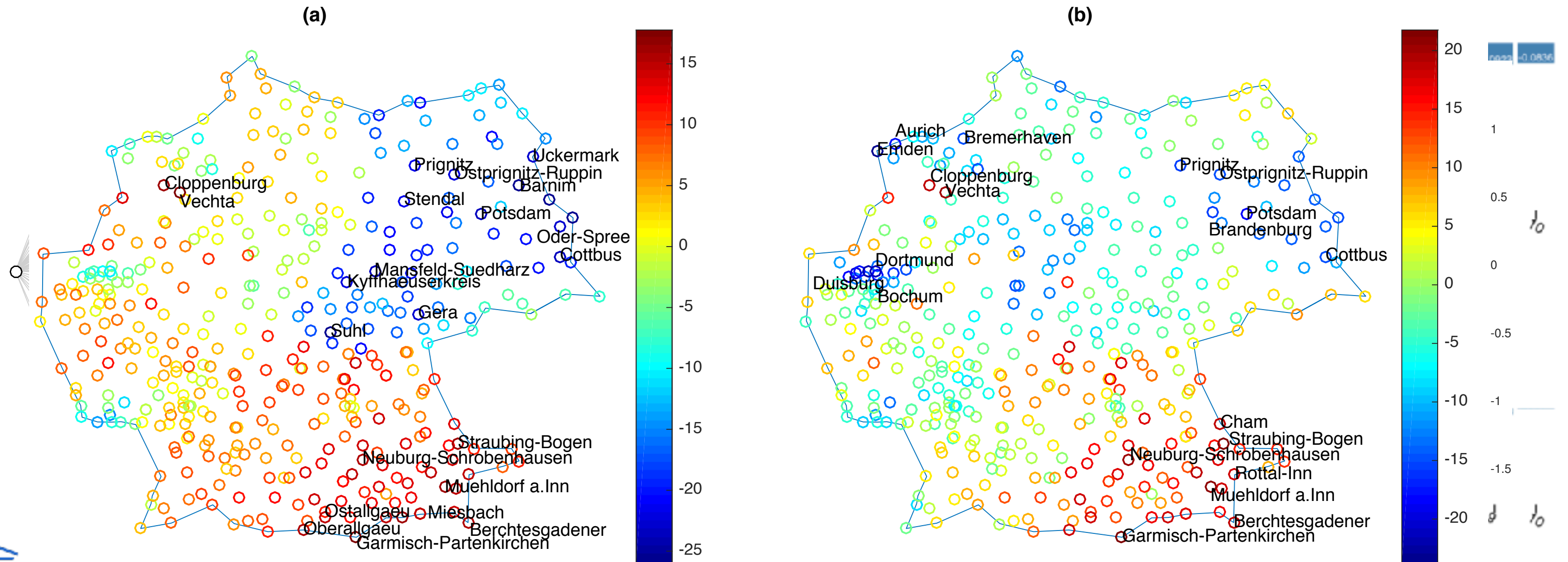
PATTERN SYNTAX

De Bie et al. 2016
Kang et al. 2016
Puolamäki et al. 2016

Weight vector and (approximate) projection



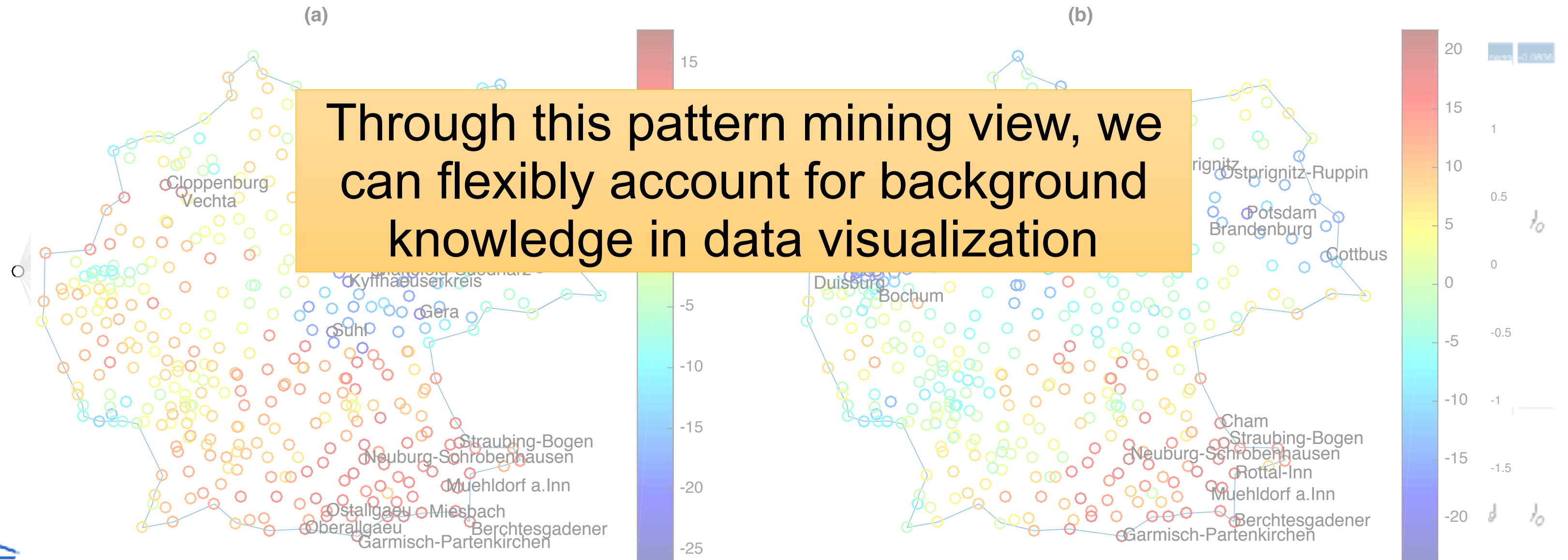
Weight vector and (approximate) projection



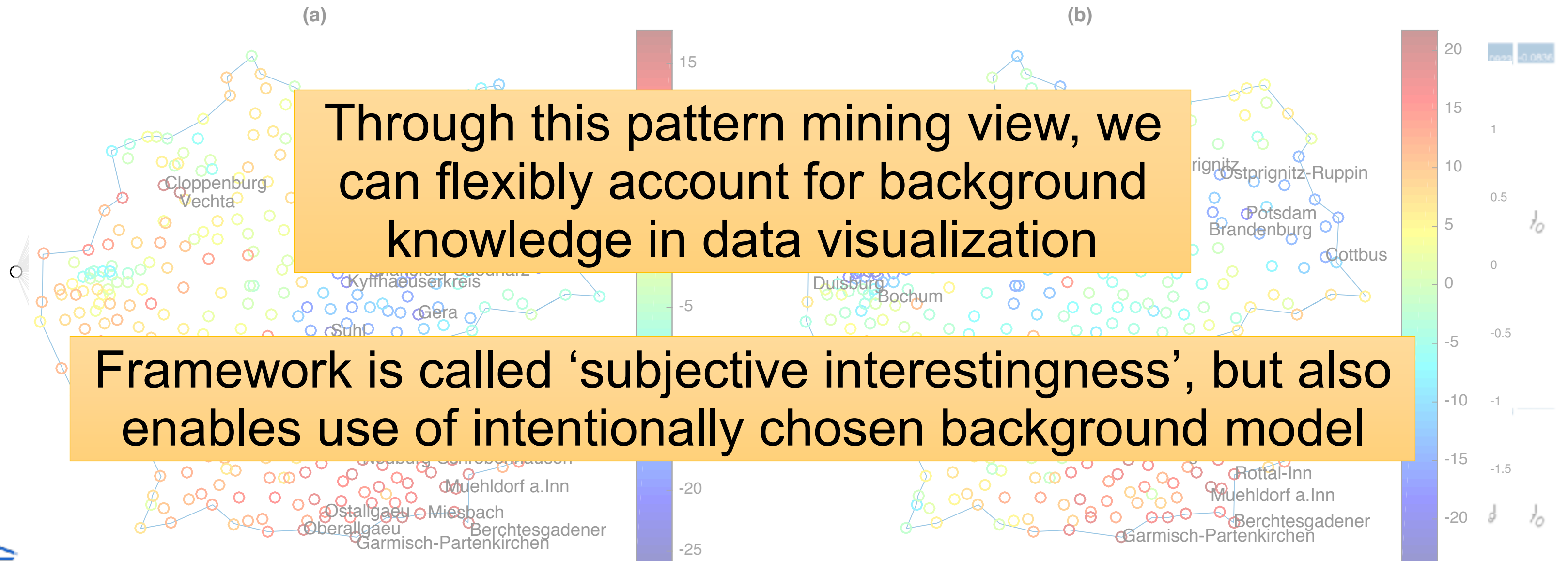
PATTERN SYNTAX

De Bie et al. 2016
Kang et al. 2016
Puolamäki et al. 2016

Weight vector and (approximate) projection



Weight vector and (approximate) projection



Background model inference with gradient descent

Mining patterns is a manifold learning problem

– Optimize through toolbox (ManOpt)

OPEN CHALLENGES

THE FORSIED PROCESS

1. Background model
2. Pattern syntax
 - IC & DL of patterns
3. Mining
4. Update background model
5. Iterative mining

BACKGROUND MODEL

How to provide specification as end-user

- Objective
- Subjective

Parameter inference (of MaxEnt model)

PATTERN SYNTAX

Should correspond to the task

- Probably constructed by researchers
- Users need to be able to select appropriate syntax

Not yet investigated: syntaxes with mixed relevance

- E.g., for graphs with various vertex types

IC & DL OF PATTERNS

Requires extensive knowledge of maths

Not related to cognition yet in any way

MINE GOOD/BEST PATTERNS

All the usual algorithmic challenges

- Typically NP-hard
 - Depends on b.g. model and pattern syntax
- Search strategies: greedy, beam search, branch & bound, SGD, black-box optimization
 - Consider bounds / approximability

UPDATE BACKGROUND MODEL

Insert what the user has learned

How about forgetting?

What kind of inference capabilities do humans have?

- Track remaining degree of vertices in graph, degrees of freedom for projections?

ITERATIVE MINING

May be possible to re-use state of previous mining step

- E.g., in branch-and-bound
- Non-overlapping patterns are generally unaffected

Iterative scheme has an $1-1/e$ bound for greedy

- For the budget used so far

SUMMARY

SUMMARY

A generic approach to define patterns

Information Theory perspective on interestingness

Examples of mining interesting/surprising

- Dense and connecting vertex sets in graphs
- Subgroups in numeric data
- Projections for visualization

Many open challenges remaining

Jefrey Lijffijt, Dr. Sc. (Tech.)

FWO [Pegasus]² MSC Fellow

Internet Technology and Data Science Lab
Dept. of Electronics and Information Systems

E jefrey.lijffijt@ugent.be

M +32 474 53 94 64

users.ugent.be/~jlijffij/

 [jlijffijt](#)

 [@jlijffijt](#)

 [jefrey](#)

Jefrey Lijffijt, Dr. Sc. (Tech.)

FWO [Pegasus]² MSC Fellow

Internet Technology and Data Science Lab
Dept. of Electronics and Information Systems

E jefrey.lijffijt@ugent.be

M +32 474 53 94 64

users.ugent.be/~jlijffij/

 [jlijffijt](#)

 [@jlijffijt](#)

 [jefrey](#)

**Interested? Visit us!
And we are hiring :-)**