

tu

technische universität dortmund



SFB 876 - Providing Information by Resource-Constrained Data Analysis




Mining Smartphone and Mobility Data Through Graphical Models

Katharina Morik, Artificial Intelligence,

TU Dortmund University, Germany

tu

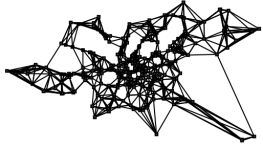
technische universität dortmund



SFB 876 - Providing Information by Resource-Constrained Data Analysis


Overview

- Graphical Models: Markov Random Fields
- Spatio-Temporal Random Fields
  - App Usage Mining
  - Dublin Traffic Prediction
- Integer Markov Random Fields
- Smartphone – The human sensor



tu

technische universität dortmund



SFB 876 - Providing Information by Resource-Constrained Data Analysis

Analyzing Sensor Data

- Select a field of measurements.
- Establish nodes with possible states.
- Indicate dependencies of nodes by edges.
- Look for the joint probability mass function.

Streets in a city (area)

Rain {dry, wet}

Lights {red, yellow, green}

Jam {jam, free}

App usage of a person


Mail {on, off}

Angry Bird {on, off}

Map {on, off}

tu

technische universität dortmund



SFB 876 - Providing Information by Resource-Constrained Data Analysis

Likelihood of joint realizations

- Graph  $G=(V,E)$
- Multi-variate random variable  $\mathbf{X}$  with finite discrete domain  $\mathcal{X}$
- Mapping joint vertex assignment into vector space  $\phi(\mathbf{x}): \mathcal{X} \rightarrow \mathbb{R}^d$
- Parameter to be learned:  $\theta$  in  $\mathbb{R}^d$

MRF :


$$p(\mathbf{x}) = \frac{1}{Z(\theta)} \exp\left(\sum_i \theta \phi_i(\mathbf{x})\right) \frac{1}{a} = \exp(-\ln a)$$

$$= \exp\left[\left(\sum_i \theta \phi_i(\mathbf{x})\right) - \ln Z(\theta)\right]$$

$$= \exp[\langle \theta, \phi(\mathbf{x}) \rangle - A(\theta)]$$

tu

technische universität dortmund



SFB 876 - Providing Information by Resource-Constrained Data Analysis

Probabilistic Graphical Model -- Example

Observation

- $\mathbf{x}_i: \{on, off, off\}$
- $\phi(\mathbf{x}_i) = (1, 0, 0, 1, 0, 1, 1, 0, 0, 0, 1, 0, 0, 0, 0, 0, 1, 0)$
- $d = 18$

Learning

- Likelihood of joint realizations  $p_\theta(\mathbf{x}) = \exp(\langle \theta, \phi(\mathbf{x}) \rangle) - A(\theta)$
- Maximum a posteriori prediction of an unobserved node  $H$ , given observed nodes  $O$

$$\mathbf{x}_H^* = \operatorname{argmax}_{\mathbf{x}_H} p_\theta(\mathbf{x}_H | \mathbf{x}_O)$$

Graph  $G=(V,E)$

Rain {on, off}

Angry Bird {on, off}

Map {on, off}

$\phi(\mathbf{x}):$  (/ \* Nodes \*/)

1. on, /\* dom(angry) \*/

2. off,

3. on, /\* dom(rain) \*/

4. off,

5. map, /\* dom(map) \*/

6. off,

/\* Edges \*/

1. on, off, /\* edge angry-rain \*/

2. on, on,

3. off, off,

4. off, on,

5. on, off, /\* edge angry-map \*/

6. on, on,

7. off, off,

8. off, on,

9. on, off, /\* edge rain-map \*/


10. on, on,

11. off, off,

12. off, on )

tu

technische universität dortmund



SFB 876 - Providing Information by Resource-Constrained Data Analysis

Parameter Estimation

- Maximum Likelihood Estimation (MLE)
- Average loglikelihood
- empirical expectation of the value  $\phi(\mathbf{x})$
- Optimization: partial derivative of  $l(\theta|D)$ ; gradient descent minimizing  $-l(\theta|D)$
- Delivers optimal  $\theta$ .

$L(\theta|D) = \prod_{x \in D} p(x)$

$$l(\theta|D) = \frac{1}{|D|} \sum_{x \in D} \log p(x)$$

$$= \frac{1}{|D|} \sum_{x \in D} [\langle \theta, \phi(x) \rangle - A(\theta)]$$

$$= \left\langle \theta, \frac{1}{|D|} \sum_{x \in D} \phi(x) \right\rangle - A(\theta)$$

$$= \langle \theta, E(\phi(x)|D) \rangle - A(\theta)$$

1

tu technische universität dortmund SFB 876 - Providing information by Resource-Constrained Data Analysis

### Loopy Belief Propagation

- What is the expectation of  $\phi(x)$ ?
- All nodes need to fit together!
- For all pairs of linked nodes  $u, v$  send messages  $m_{vu}(x_u)$

$$m_{vu}(x_u) = \sum_{x_v \in X} \varphi_{v,u}(x_v, x_u) \varphi(x_v) \prod_{w \in \text{neighbor}_v} m_{vw}(x_w)$$

$$M_v(x) := \prod_{u \in \text{neighbor}_v} m_{uv}(x)$$

$$p_v(x_v) = \frac{\varphi_v(x_v) M_v(x_v)}{\sum_{x \in X} \varphi_v(x_v) M_v(x)}$$

tu technische universität dortmund SFB 876 - Providing information by Resource-Constrained Data Analysis

### Probabilistic Graphical Model – Example

Observation

- $x_1$ : {on, off, off}
- $\phi(x_1) = (1, 0, 0, 1, 0, 1, 1, 0, 0, 1, 0, 0, 0, 1)$
- $d = 18$

Learning

- Likelihood of joint realizations
- $p_\theta(x) = \exp(\langle \theta, \phi(x) \rangle) - A(\theta)$
- Maximum a posteriori prediction of an unobserved node  $H$ , given observed nodes  $O$
- $x_H^* = \text{argmax } p_\theta(x_H | x_O)$

Graph  $G=(V,E)$

Nodes: Rain {on, off}, Angry Bird {on, off}, Map {on, off}

Edges: (Rain, Angry Bird), (Rain, Map), (Angry Bird, Map)

Parameters:

- $\phi(x)$ : (/\* Nodes\*/
- 1. on, /\*dom(angry)\*/
- 2. off,
- 3. on, /\*dom(rain)\*/
- 4. off,
- 5. on, /\*dom(map)\*/
- 6. off,
- /\*Edges\*/
- 1. on, off, /\*edge angry-rain\*/
- 2. on, on,
- 3. off, off,
- 4. off, on,
- 5. on, off, /\*edge angry-map\*/
- 6. on, on,
- 7. off, off,
- 8. off, on,
- 9. on, off, /\*edge rain-map\*/
- 10. on, on,
- 11. off, off,
- 12. off, on,
- )

tu technische universität dortmund SFB 876 - Providing information by Resource-Constrained Data Analysis

### Spatio-temporal Random Fields

- The spatio-temporal graph is trained to predict each node's maximum a posteriori probability with the marginal probabilities.
- Generative model predicting all nodes.
- Maximum a posteriori prediction of an unobserved node  $H$ , given observed nodes  $O$ .
- The learned model answers diverse questions – you only need to choose  $H$ .

tu technische universität dortmund SFB 876 - Providing information by Resource-Constrained Data Analysis

### Spatio-temporal Models for App Usage Mining

- Spatial graph with Apps as nodes and edges based on covariance matrix.
- Windows of 30 minutes arranged for a day.
- Likelihood of an app at a particular time of day.
- Data from 8 users, public: [http://sfb876.tu-dortmund.de/auto?self=\\$e675o3g0ow](http://sfb876.tu-dortmund.de/auto?self=$e675o3g0ow)

tu technische universität dortmund SFB 876 - Providing information by Resource-Constrained Data Analysis

### STRF modeling usage of apps on Android phones

- Probability of an app being switched on
- Ranking according to probability
- Show the first three ranks with their probability.

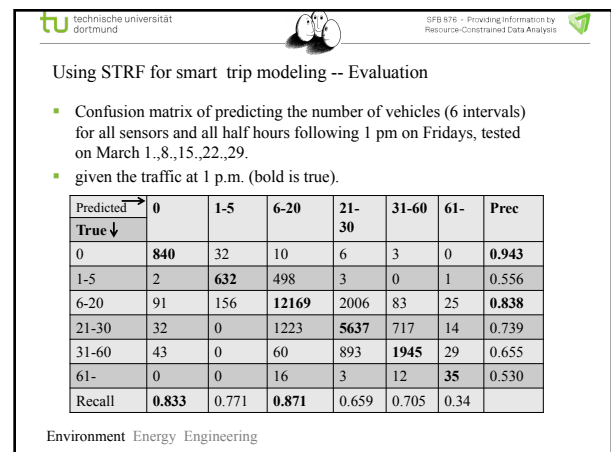
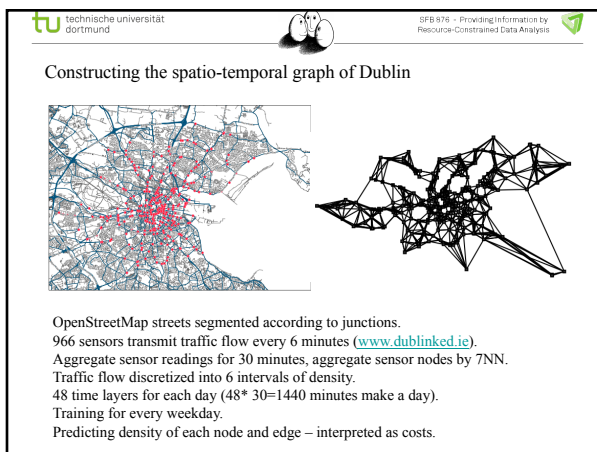
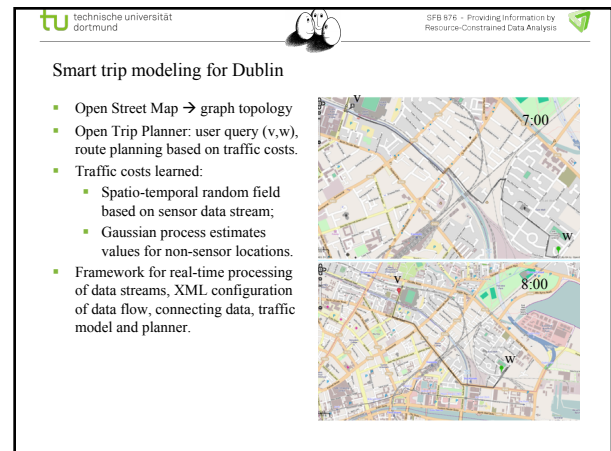
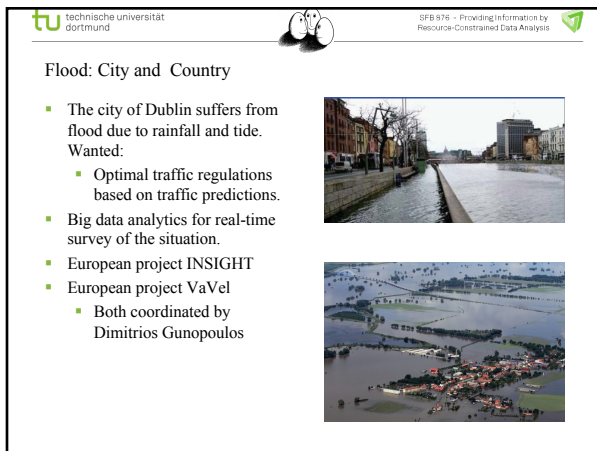
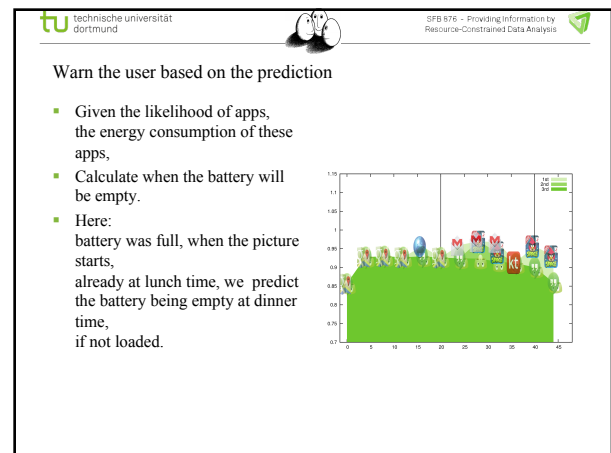
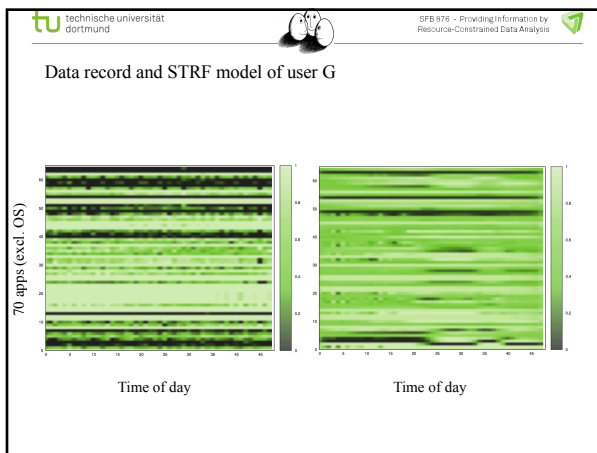
User G keeps battery on, plays Angry Bird, and looks up the map and communicates using hangout.

tu technische universität dortmund SFB 876 - Providing information by Resource-Constrained Data Analysis

### The Google addict

YouTube, hangout, music, maps

Facebook in the evening




tu technische universität dortmund

SFB 876 - Providing Information by Resource-Constrained Data Analysis

### Smart traffic for smart cities

- Spatio-temporal modeling is natural for traffic prediction.
- Several questions can be answered using the same learned model.
- The answers come along with their probabilities. This might be helpful for decision makers.
- Integration of Spatio-temporal random fields into the Open Trip Planner and Gaussian information completion resulted in an excellent navigation system.



Platkowski, Lee, Morik (2013) Spatio-temporal random fields: compressible representation and distributed estimation, *Machine Learning Journal* 93:1, 115 – 140.  
 Ludwig, Platkowski, Bockermann, Morik (2014) Predictive Trip Planning – Smart Routing in Smart Cities, *Mining Urban Data Workshop at 17th Intern. Conf. on Extending Database Technology*.

tu technische universität dortmund

SFB 876 - Providing Information by Resource-Constrained Data Analysis

### Graphical models on resource-restricted processors

- Spatio-Temporal Random Fields on the phone. but
- Power consumption is restricted!
- “The most obvious technique to conserve power is to reduce the number of cycles it takes to complete a workload.”  
 (Intel 64, IA-32 architectures optimization reference manual, guidelines for extending battery life)
- Restrict the parameter space of the Markov Random Field  
 $\theta \in \{0, 1, \dots, K\} \subset \mathbb{N}$

	Sandy Bridge		ARM 11	
	Real	Int	Real	Int
+	3	1	8	1
*	5	3	8	4-5
/	14	13-15	19	-
Bit shift	-	3	-	2

Clock cycles for arithmetics on different processors:  
Real vs. integer.

tu technische universität dortmund

SFB 876 - Providing Information by Resource-Constrained Data Analysis

### Parameter space transformation

- Graph model tree-structured
- Transform the parameter space:  
 $\eta_i(\theta) = \theta_i \ln 2$

**MRF :**

$$p(\vec{x}) = \frac{1}{Z(\theta)} \exp\left(\sum_i \theta_i \phi_i(\vec{x})\right)$$

$$= \exp[\langle \theta, \phi(\vec{x}) \rangle - A(\theta)]$$

**IntegerMRF :**

$$p(\vec{x}) = \exp[\langle \eta(\theta), \phi(\vec{x}) \rangle]$$

$$= 2^{[\langle \theta, \phi(\vec{x}) \rangle - A(\eta(\theta))]}$$

$$= \frac{2^{\langle \theta, \phi(\vec{x}) \rangle}}{\sum_{y \in \mathbb{N}} 2^{\langle \theta, \phi(y) \rangle}}$$

tu technische universität dortmund

SFB 876 - Providing Information by Resource-Constrained Data Analysis

### Integer belief propagation

- Simply replacing the  $\exp(\cdot)$  by  $2^{(\cdot)}$  is not sufficient
  - Overflows are normally avoided by normalization.
  - Normalization is impossible in integer division.
- Magnitude of messages corresponds to probability
  - Use the length of each message
  - Bit-length is similar to log

$$m_{vu}(y) = \sum_{x \in \mathbb{N}_v} \exp(\theta_{vu=xy} + \theta_{v=x}) \prod_{w \in N_v - \{u\}} m_{vw}(x)$$

$$\tilde{m}_{vu}(y) = \sum_{x \in \mathbb{N}_v} 2^{(\theta_{vu=xy} + \theta_{v=x})} \prod_{w \in N_v - \{u\}} \tilde{m}_{vw}(x)$$

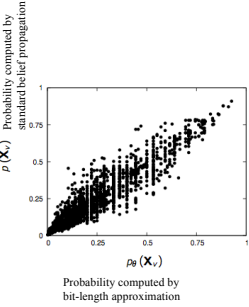
$$\beta_{vu}(y) = \max_{x \in \mathbb{N}_v} \theta_{vu=xy} + \theta_{v=x} + \sum_{w \in N_v - \{u\}} \beta_{vw}(x)$$

tu technische universität dortmund

SFB 876 - Providing Information by Resource-Constrained Data Analysis

### Discretized probability space

- Belief propagation is now bit-length propagation, i.e. the MAP and marginals are computed using the bit-length.
- The approximation error depends on the number of neighboring nodes and the space of states.
- Some true probabilities (y axis) cannot be expressed by the integer approximation (x axis).



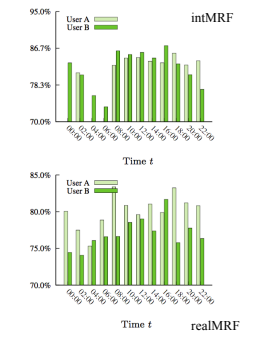
Probability computed by bit-length approximation

tu technische universität dortmund

SFB 876 - Providing Information by Resource-Constrained Data Analysis

### Experiments

- User A
  - 2064 vertices (app at time = {on, off})
  - 23368 spatio-temporal edges
  - 38 training instances, 32 test instances
- User B
  - 2064 vertices
  - 14998 spatio-temporal edges
  - 125 training instances, 119 test instances
- Comparing integer MRF and real MRF
  - Average of 50 trainings
  - Runtime: intMRF 156 (136) s, realMRF 466 (419) s
  - Predicting apps at t+1, given apps at t





technische universität dortmund

SFB 876 - Providing Information by Resource-Constrained Data Analysis

## Smartphones – the Human Sensor

- Rayid Ghani: chief scientist at the 2012 Obama campaign: fundraising, voter mobilization, volunteer (invited talk ECML PKDD 2013 Prague)
- Personality testing on smart phones: B. Rammstedt and O. P. John. Measuring personality in one minute or less: A 10-item short version of the big five inventory in English and German. Journal of Research in Personality, 2007
- MyPersonality (started 2008 at Cambridge, UK, by Michael Kosinski) publicly available data for unregistered users: Facebook status updates of 250 users + personality + Facebook properties; 3120 hand-labeled status updates for registered users; 3,100,000 Big 5 scores; Facebook friendship graphs, activity records, topics, dictionaries, demographics <https://applymagicsauce.com>
- Cambridge Analytica (Alexander J. A. Nix) claims to have won the election for Trump using 5000 data points on 220,000,000 Americans, predictive models with 100 variables <https://ocean.cambridgeanalytica.org>

technische universität dortmund

SFB 876 - Providing Information by Resource-Constrained Data Analysis

## Big Five (OCEAN) from Questionnaire

OCEAN score

Your unique OCEAN score

High

Average

Low

O C E A N

High in Openness

Average in Conscientiousness

High in Extraversion

Average in Agreeableness

Average in Neuroticism

Psychomedia L.Satow (German)

Cambridge Analytica online test

technische universität dortmund

SFB 876 - Providing Information by Resource-Constrained Data Analysis

## Big 5 (OCEAN) from Facebook Likes

- University of Cambridge, Psychometrics Centre Method: ApplyMagicSauce
- Michael Kosinski (now Stanford)

Big 5 Personality: individuals are represented in personality traits

Openness

Conscientiousness

Extraversion

Agreeableness

Neuroticism

Using 12 Fb LIKES

technische universität dortmund

SFB 876 - Providing Information by Resource-Constrained Data Analysis

## App Usage Mining

- Smartphones produce big data.
- Each user generates about 60 gigabytes of data per year. This is big data!
- Data curation is time-consuming.
- Basic Data representation:
  - Did user<sub>i</sub> start app<sub>j</sub>?

Binary	app <sub>1</sub>	...	app <sub>m</sub>
usr <sub>1</sub>	1		0
...			
usr <sub>n</sub>	0		1

technische universität dortmund

SFB 876 - Providing Information by Resource-Constrained Data Analysis

## App Usage Mining

- Smartphones produce big data.
- Each user generates about 60 gigabytes of data per year. This is big data!
- Data curation is time-consuming.
- Basic Data representation:
  - Did user<sub>i</sub> start app<sub>j</sub>?
  - How often did user<sub>i</sub> start app<sub>j</sub>?
- Probability distributions deliver an overview of the data.

Frequent	app <sub>1</sub>	...	app <sub>m</sub>
usr <sub>1</sub>	#starts		
...			
usr <sub>n</sub>			#starts(n,m)

technische universität dortmund

SFB 876 - Providing Information by Resource-Constrained Data Analysis

## Probability Distributions of Apps

App usage follows the Zipf' or Mandelbrot distribution:

- The probability of an app is inverse proportional to its rank:  $p(n) = 1/n$

Zipf-Mandelbrot

tu technische universität dortmund

SFB 876 - Providing Information by Resource-Constrained Data Analysis

### App Usage Mining

- Smartphones produce big data.
- Each user generates about 60 gigabytes of data per year. This is big data!
- Data curation is time-consuming.
- Basic Data representation:
  - Did user<sub>i</sub> start app<sub>j</sub>?
  - How often did user<sub>i</sub> start app<sub>j</sub>?
- Probability distributions deliver an overview of the data.
- Additional attributes of users can be predicted.

Binary	app <sub>1</sub>	...	app <sub>m</sub>	target
usr <sub>1</sub>	0		1	
usr <sub>n</sub>	1		0	

tu technische universität dortmund

SFB 876 - Providing Information by Resource-Constrained Data Analysis

### Predicting Attributes -- Gender

- Menthal data from Alexander Markowetz <https://menthal.org>
- 42,482 men and 31,700 women were logged.
- Top 99 apps of women and top 99 apps of men overlap by 82 apps.
- We can predict the gender by app usage data using logistic regression and the L1 norm: accuracy 91.2%.

Binary	app <sub>1</sub>	...	app <sub>m</sub>	gender
usr <sub>1</sub>	0		1	f
usr <sub>n</sub>	1		0	m

tu technische universität dortmund

SFB 876 - Providing Information by Resource-Constrained Data Analysis

### Inspecting the Results -- Gender

- Menthal data from Alexander Markowetz <https://menthal.org>
- 42,482 men and 31,700 women were logged.
- Ranking apps according to their frequency, compare rankings of men and women.
- Largest differences:
  - ClashOfClans men 80, women 219
  - FarmHeroes men 249, women 78
  - GoogleDocs men 52, women 92

Frequent	app <sub>1</sub>	...	app <sub>m</sub>
usr <sub>1</sub>	#starts		
usr <sub>n</sub>			#starts(n,m)

tu technische universität dortmund

SFB 876 - Providing Information by Resource-Constrained Data Analysis

### App Usage Mining

- Smartphones produce big data.
- Each user generates about 60 gigabytes of data per year. This is big data!
- Data curation is time-consuming.
- Individual data:
  - Where did user<sub>i</sub> use app<sub>j</sub>?
  - Include only places that occur more than t times.

usr <sub>n</sub>	app <sub>1</sub>	...	app <sub>m</sub>
usr <sub>1</sub>	app <sub>1</sub>	...	app <sub>m</sub>
place <sub>1</sub>	#starts		
usr <sub>1</sub>	app <sub>1</sub>	...	app <sub>m</sub>
place <sub>1</sub>	#starts		
place <sub>n</sub>			#starts(n, m)

tu technische universität dortmund

SFB 876 - Providing Information by Resource-Constrained Data Analysis

### App Usage Mining

- Smartphones produce big data.
- Each user generates about 60 gigabytes of data per year. This is big data!
- Data curation is time-consuming.
- Individual data:
  - Where did user<sub>i</sub> use app<sub>j</sub>?
  - Include only places that occur more than t times.
  - When did user<sub>i</sub> use app<sub>j</sub>?
  - Discretized time of day: 30 minutes

usr <sub>n</sub>	app <sub>1</sub>	...	app <sub>m</sub>
usr <sub>1</sub>	app <sub>1</sub>	...	app <sub>m</sub>
0	#starts		
usr <sub>1</sub>	app <sub>1</sub>	...	app <sub>m</sub>
0	#starts		
1			
...			
48			#starts(n, m)

tu technische universität dortmund

SFB 876 - Providing Information by Resource-Constrained Data Analysis

### References

- Nico Piatkowski, Sangyun Lee, Katharina Morik (2016) Integer Undirected Graphical Models for Resource-Constrained Systems. In: Neurocomputing, 173:1, 9 – 23
- Nico Piatkowski, Katharina Morik (2016) Stochastic Discrete Clenshaw-Curtis Quadrature. In: 33<sup>rd</sup> Int. Conference on Machine Learning, Proceedings JMLR
- Marco Stolpe, Thoms Liebig, Katharina Morik (2015) Communication-efficient learning of traffic flow in a network of wireless presence sensors. In: Workshop Parallel and Distributed Computing for KDD, CEUR-WS
- Thomas Liebig, Nico Piatkowski, Christian Bockermann, Katharina Morik (2014) Predictive Trip Planning – Smart Routing in Smart Cities. In: EDBT/ICDT Workshops CEUR-WS
- François Schnitzler, Thomas Liebig, Shie Mannor, Katharina Morik (2014) Combining a Gaussian Markov model and Gaussian process for traffic prediction in Dublin city center. In: EDBT/ICDT, CEUR-WS
- Jochen Streicher, Nico Piatkowski, Katharina Morik, Olaf Spinczyk (2013) Open Smartphone Data for Mobility and Utilization Analysis in Ubiquitous Environments. In: Atzmüller, Scholz (eds) 4<sup>th</sup> Int. Workshop on Mining Ubiquitous and Social Environments (MUSE)
- Nico Piatkowski, Sangyun Lee, Katharina Morik (2013) Spatio-temporal Random Fields: Compressible Representation and Distributed Estimation. In: Machine Learning Journal, 93:1, 115 – 140.
- Thomas Liebig, Nico Piatkowski, Christian Bockermann, Katharina Morik (2014) Predictive Trip Planning – Smart Routing in Smart Cities. In: Mining Urban Data Workshop at 17th Intern. Conf. on Extending Database Technology

Code and local data available at <http://sfb876.tu-dortmund.de/index.htm>