

Unsupervised Neural Net based automatic trend analysis for Weblogs

Johannes Fordemann
Manfred Leisenberg
Timo Timm
Julia Wolff

Motto

*“Irrationality of human
communication
Might convert into a
swarming attitude
after a while”*

[Lehmann2005]

.... Let us discover and understand *swarming attitudes*

Agenda

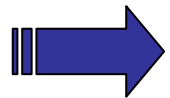
- Introduction
 - Initial situation
 - Trend analysis based on Neural Networks
 - Experimental results
 - Summary
-

Trends are reflected by Weblogs

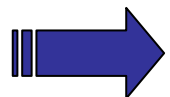
- Weblog based communication produces and reflects socio- cultural streams
- Streams might produce „Weak Signals“
- „Weak Signals“ might indicate trends
- Trends are based on pattern
 - From such pattern we might learn the future direction of individual action
- **Task**
 - Pattern Identification and Classification
(which is an native AI (Artificial Intelligence) task)

Identification of trends

- Trend occurs, if
 - Social coincidences **turn** into collective phenomena
 - Weak Signals **turn** into Tipping Points


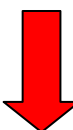
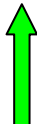




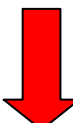




Effective computer based detection methods for future trends have **not** been found yet [Lehmann 2005]



This research project aims at solving above detection problems

Market for Automatic Trend analyses

Portal	Advantages	Disadvantages
technorati	<ul style="list-style-type: none"> • statistics provided • blog search 	<ul style="list-style-type: none"> • no trend analysis 
blogpulse	<ul style="list-style-type: none"> • multilingual • statistical monitoring based analysis (featured and individual) 	<ul style="list-style-type: none"> • no prediction capabilities 
ArgYou	<ul style="list-style-type: none"> • Offer-Demand Comparison • Measurement of potential demand • Supervised statistical monitoring on existing trends 	<ul style="list-style-type: none"> • Not specialized on trend prediction • Part-manual process 
blogscout	<ul style="list-style-type: none"> • Statistical analysis of impressions and visits 	<ul style="list-style-type: none"> • no trend analysis and prediction 
gridpatrol	<ul style="list-style-type: none"> • Automatic monitoring and result analyses • Based on sophisticated grid technology 	<ul style="list-style-type: none"> • Trend analyses for financial applications 

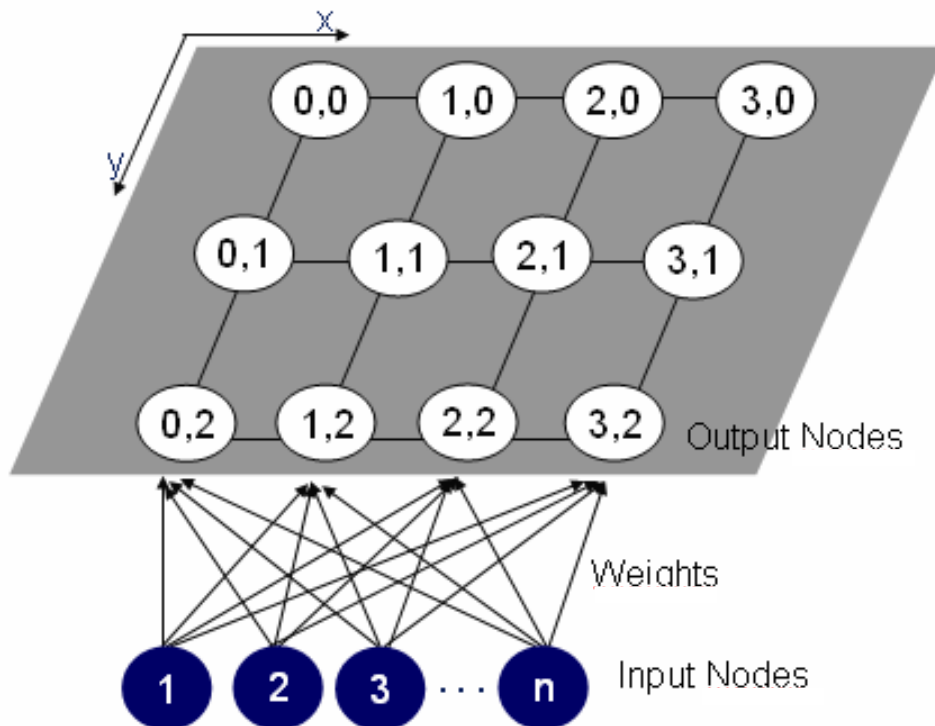
Initial Situation

- Disadvantages of „classical“ statistical methods
 - Supervised processing
 - Not sensitive to weak signals
 - Redundancies , imprecision
- Therefore:
 - Unsupervised Analysis of Weblogs based on Self Organizing Maps (SOM)
 - Investigation of capabilities of SOM concerning trend analyses

Trend analyses based on SOM

- Self Organizing Maps (SOM)
 - Unsupervised learning based classifier
 - Classes are not known prior to the training process
 - Capable of identifying clusters/ classes in a multidimensional feature vector space automatically
 - Data for training and test are described by multidimensional feature vectors
 - Classification tendencies might be interpreted as s „Weak Signals“ or „Tipping Points“ in terms of trend analyses

SOM architecture



$$X_i = \{x_{1i}, x_{2i}, \dots, x_{ki}, \dots, x_{ni}\}$$

- Similar feature combinations are represented by close regions on the map

Trend analyses based on SOM

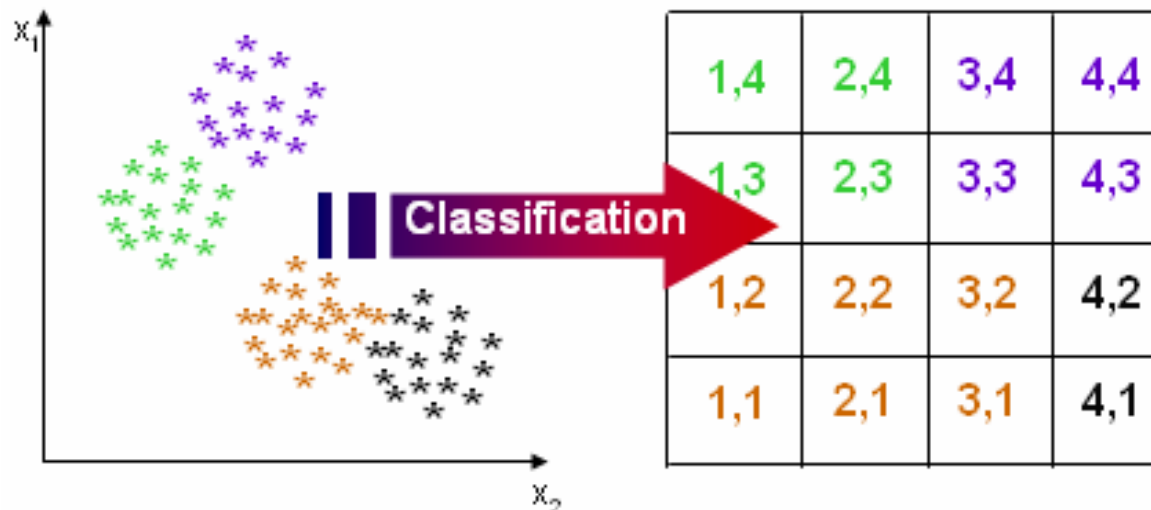
- Features are
 - Multidimensional quantitative descriptions of sequences of words
 - Represented by feature vectors

$$X_i = \{x_{1i}, x_{2i}, \dots, x_{ki}, \dots, x_{ni}\}$$

- Example
 - Context based description of „Items“
 - X_1 =Item code, x_2 =Owner, x_3 =Price, X_4 =Colour, ...
 - Sequence oriented description
 - X_1 =Wordcode1, x_2 =Wordcode2, x_3 =Wordcode3, ...

SOM training process

- SOM capable of finding clusters in the input feature vector space
- Clusters might be labelled
- SOM represents distribution of feature vectors X_i



for presentation purposes:
2- dimensional input feature vector space

How to find trends by SOM?

- Output nodes represent Weblog clusters
- Weak Signal
 - Indicated by deviation from the centre of class during test mode
- Trend, „Tipping Point“:
 - statistically **significant** deviation from the centre of class during test mode
 - Threshold to be defined

SOM based trend analyses

Vector for training

$$X_i = \{x_{1i}, x_{2i}, \dots, x_{ki}, \dots, x_{ni}\}$$

Trained SOM

$t=t_0$

1,4	2,4	3,4	4,4
1,3	2,3	3,3	4,3
1,2	2,2	3,2	4,2
1,1	2,1	3,1	4,1

Centre of class

(Strong after- training representation of a particular feature)

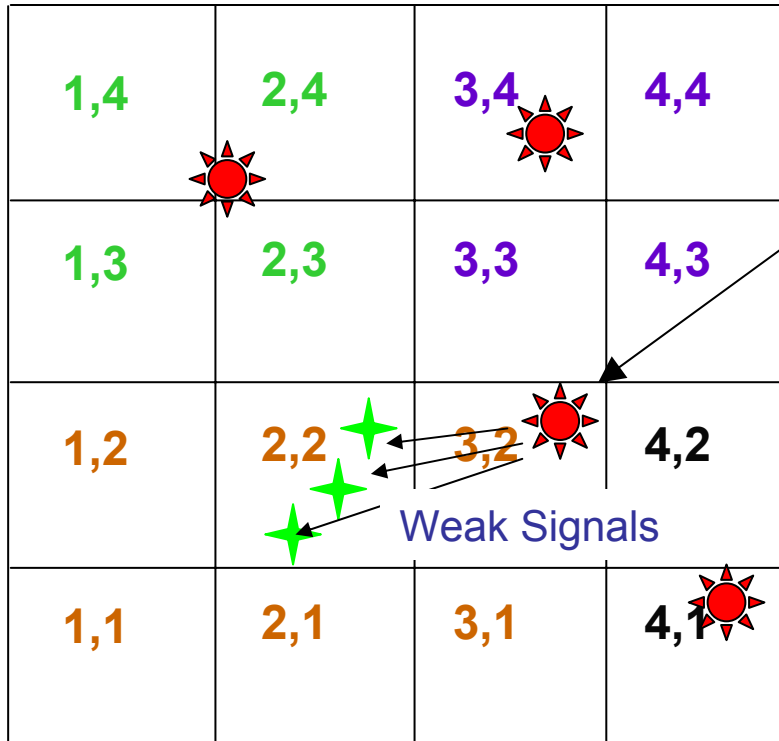
SOM based trend analyses

Vector for training

$$X_i = \{x_{1i}, x_{2i}, \dots, x_{ki}, \dots, x_{ni}\}$$

SOM in **Testmode**

$$t=t_0+t_{t1}$$



Centre of class

(Strong after- training representation of a particular feature)

Weak Signals

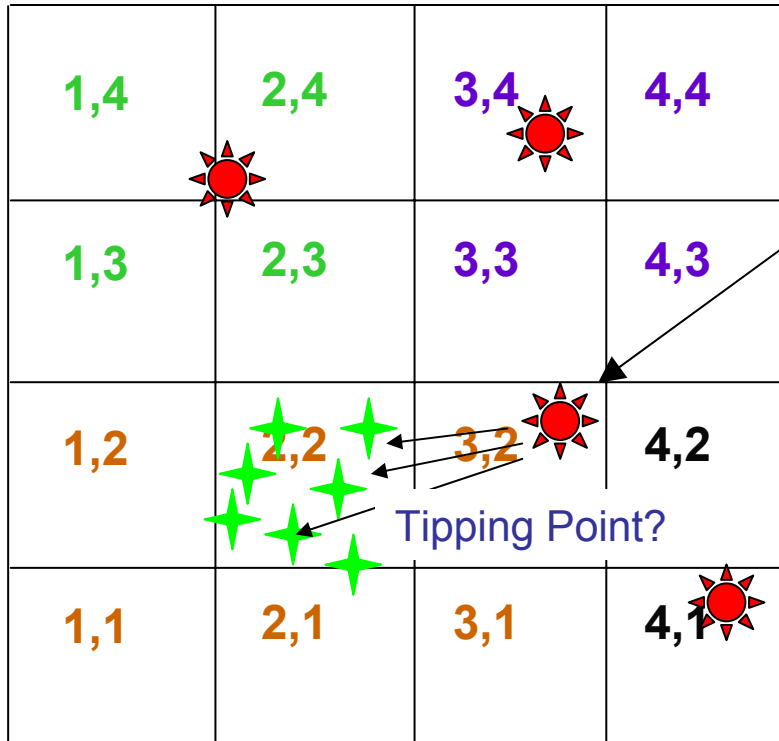
SOM based trend analyses

Vector for training

$$X_i = \{x_{1i}, x_{2i}, \dots, x_{ki}, \dots, x_{mi}\}$$

SOM in **Testmode**

$$t = t_0 + t_1 + t_2$$



Centre of class

(Strong after-training representation of a particular feature)

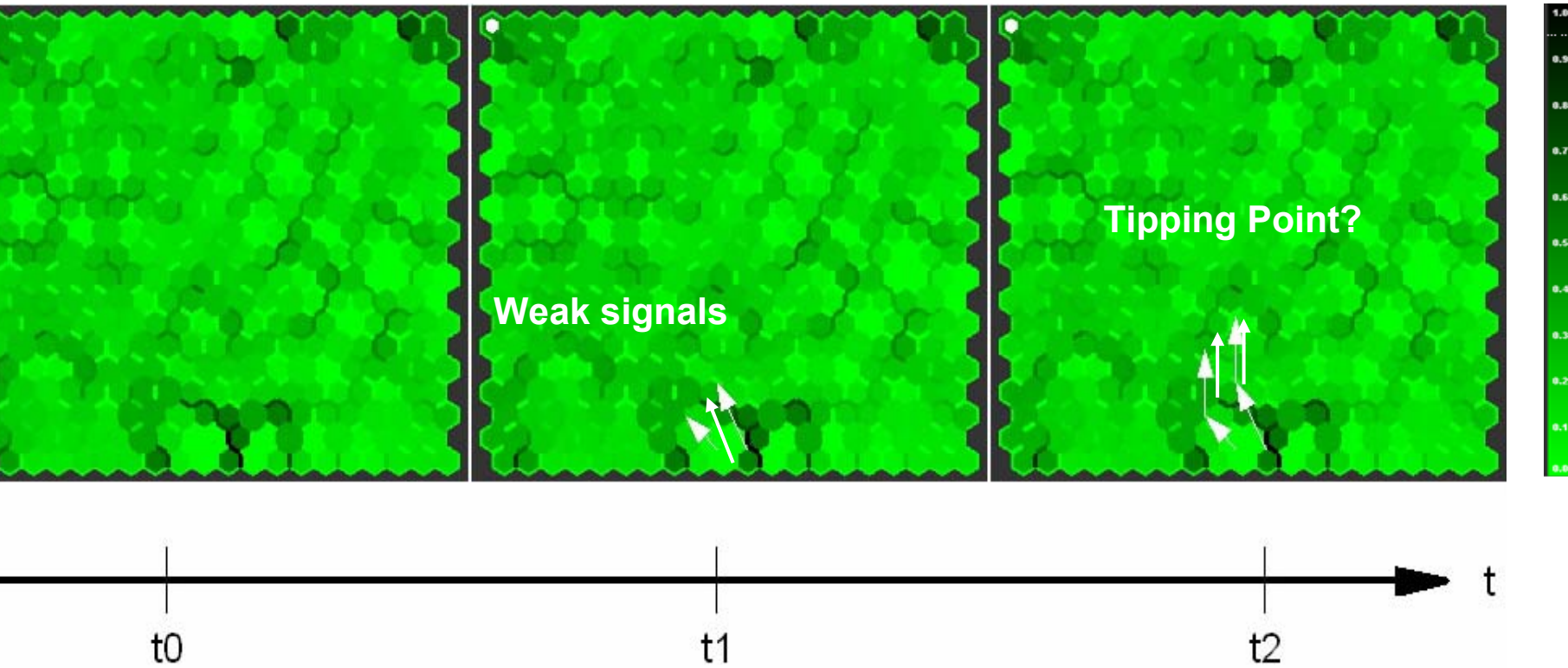
Trend?!

Tipping Point?

Experimental Result

- Procedure
 1. Setup corpus for training and test
 2. Identify and extract features
 3. Define SOM parameters
 4. Training of SOM
 5. Testing of SOM
 6. Identify Weak Signals, Tipping Points

Experimental Result



10 different source texts
4*5=20 features (Sequence oriented description)
SOM 20*20

Summary

- Weblogs reflect socio- cultural streams -> trends
- SOM based identification of trends is a promising method
- Most important feature of the method
 - Unsupervised processing
 - Classes do not have to be known prior to training
- First experimental results support technical idea
- Problem: effective feature extraction