

Forecasting Big Time Series: Theory and Practice



Yuyang (Bernie) Wang, Christos Faloutsos, Valentin Flunkert, Jan Gasthaus and Tim Januschowski



<https://lovvge.github.io/Forecasting-Tutorial-WWW-2020/>

Overall Outline



- Part 1: Fundamentals
- Part 2:
 - NN – from old to new
 - CNNs / Wavenet / Transformer-based models
 - Deep Probabilistic Models
 - System, practice, and the AWS Forecasting Stack



Part 1 – Fundamentals - Outline

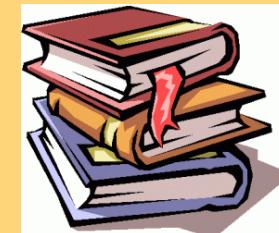


- Motivation
- P1.1. Similarity Search and Indexing
- P1.2. DSP (Digital Signal Processing)
- P1.3. Linear Forecasting
- P1.4. Non-linear forecasting
- P1.5. Tensors
- Conclusions



‘Recipe’ Structure:

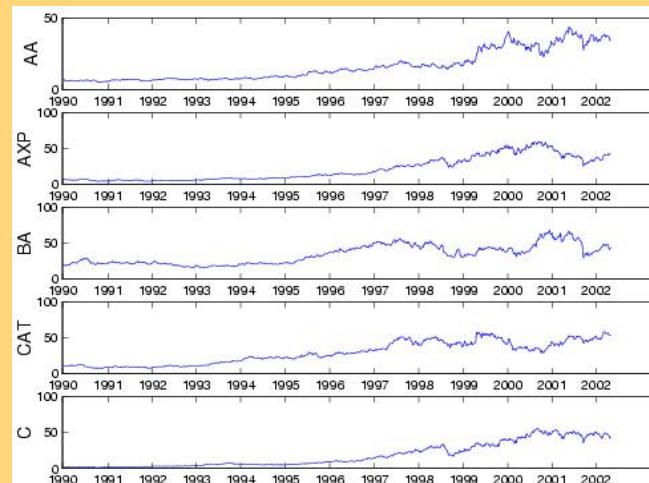
- Problem definition
- Short answer/solution
- LONG answer – details
- Conclusion/short-answer





Problem:

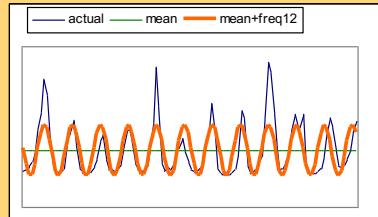
Q: mine/forecast (one, or more) time sequences



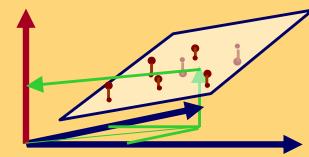
Answers



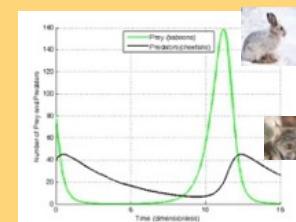
- P1.1. Similarity search: **Euclidean/time-warping; feature extraction and SAMs**
- P1.2. Periodicities: **DFT/DWT**
- P1.3. Linear Forecasting: **AR** (Box-Jenkins)
- P1.4 Non-linear forecasting: **lag-plots**
 - Gray-box modeling: **Lotka-Volterra**
- P1.5. Tensors: **PARAFAC**



WWW 2020



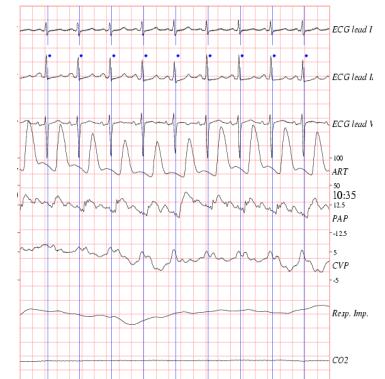
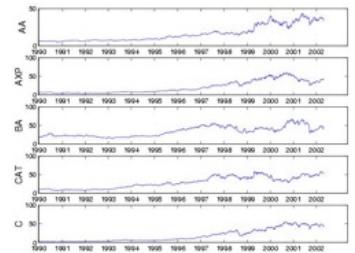
Wang, Faloutsos, Flunkert, Gasthaus,
Januschowski



6

Motivation - Applications

- Financial, sales, economic series
- Medical
 - reactions to new drugs
 - elderly care

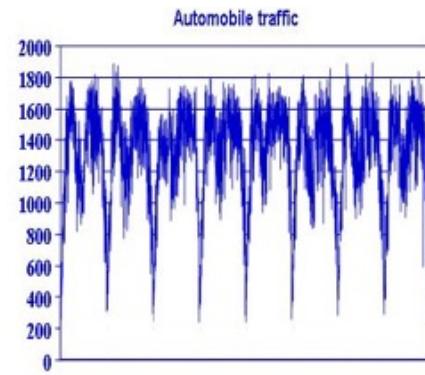


EEG - epilepsy



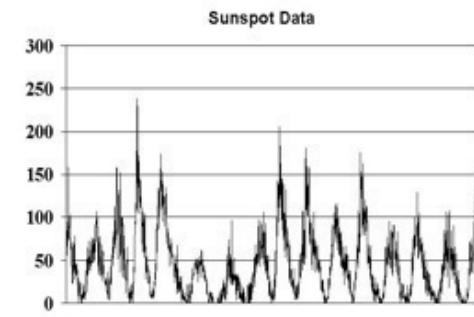
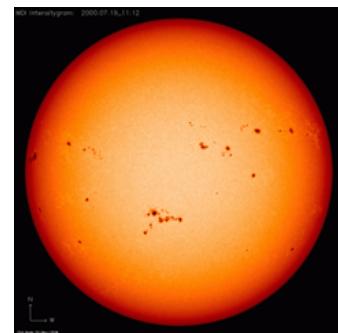
Motivation - Applications (cont'd)

- civil/automobile infrastructure
 - bridge vibrations [Oppenheim+02]
 - road conditions / traffic monitoring



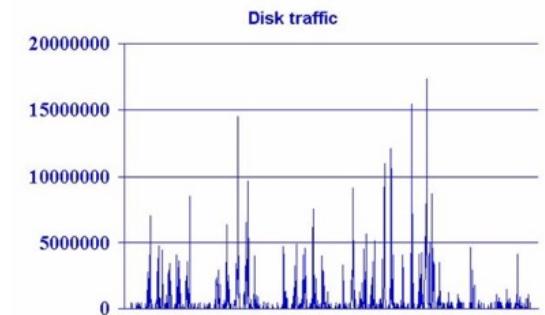
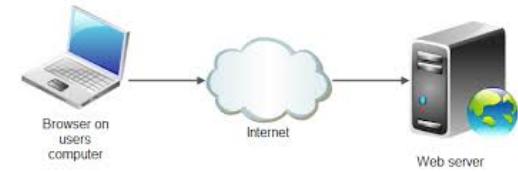
Motivation - Applications (cont'd)

- Weather, environment/anti-pollution
 - volcano monitoring
 - air/water pollutant monitoring
 - sunspots



Motivation - Applications (cont'd)

- Computer systems
 - web servers (caching, prefetching)
 - network traffic monitoring
 - ...

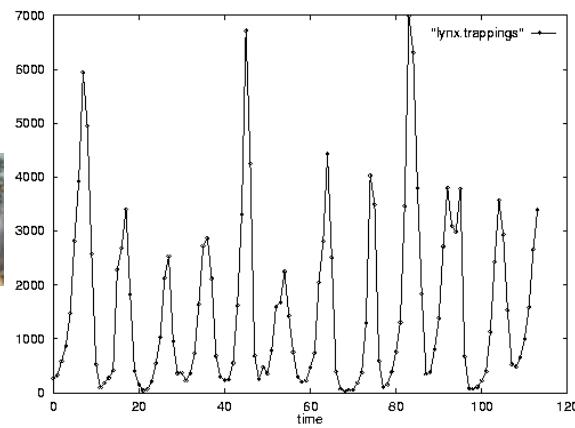


Problem #1:

Goal: given a signal (eg., #packets over time)

Find: patterns, periodicities, and/or compress

count

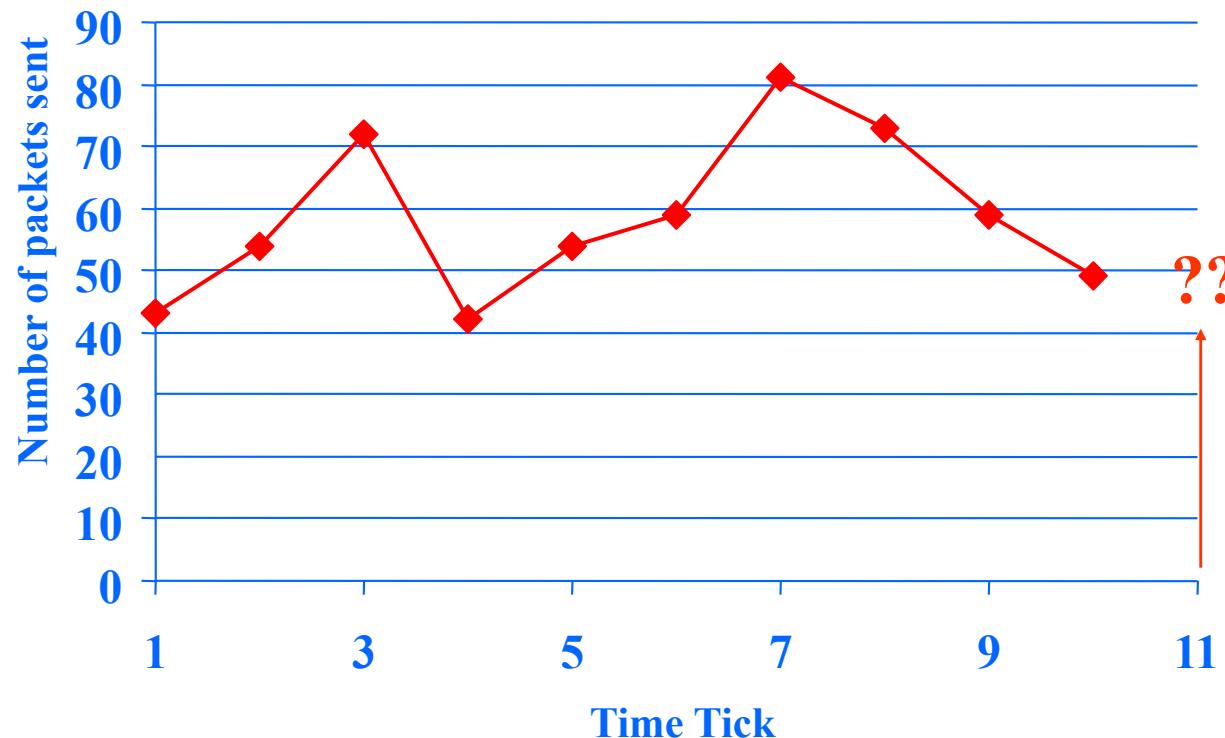


lynx caught per year
(packets per day;
temperature per day)

year

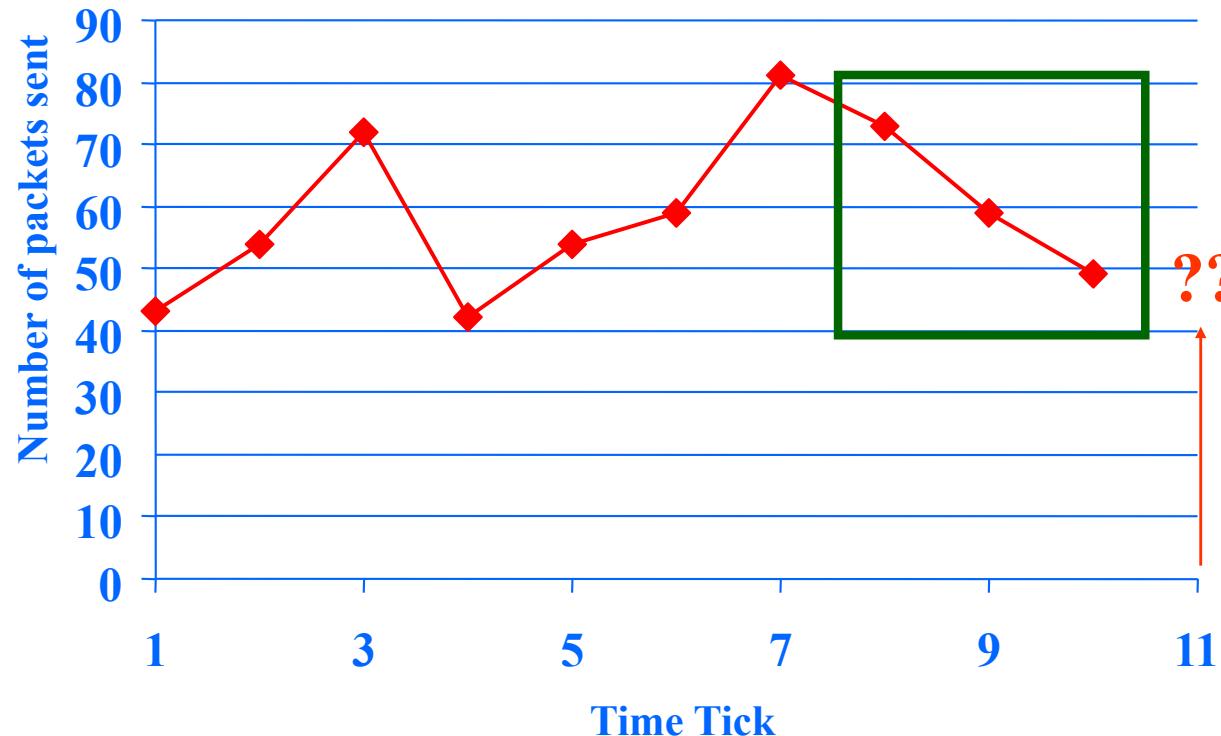
Problem#2: Forecast

Given x_t, x_{t-1}, \dots , forecast x_{t+1}



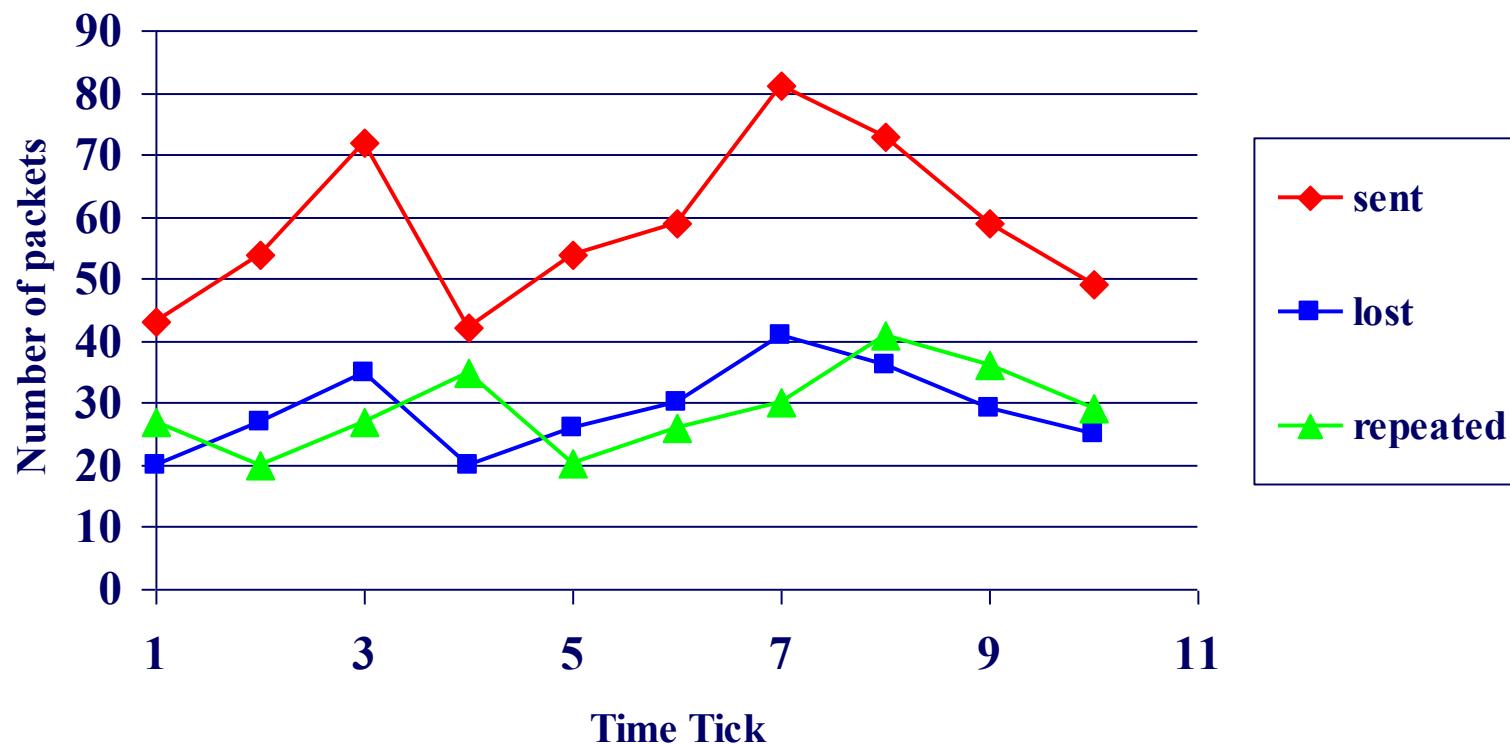
Problem#2': Similarity search

Eg., Find a 3-tick pattern, similar to the last one



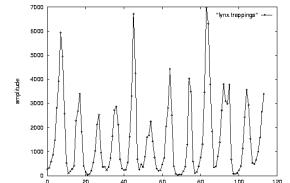
Problem #3:

- Given: A set of **correlated** time sequences
- Forecast ‘**Sent(t)**’



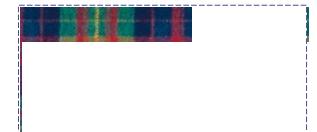


Important observations



Patterns, rules, forecasting and similarity indexing are closely related:

- To do **forecasting**, we need
 - to find **patterns/rules**
 - compress
 - to find **similar** settings in the past
- to find outliers, we need to have forecasts
 - (outlier = too far away from our forecast)



Part 1 - Outline

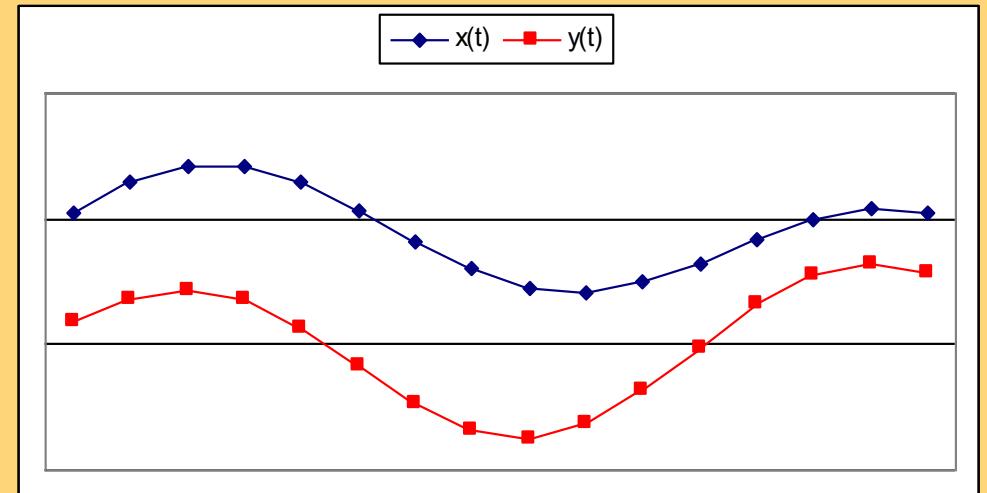
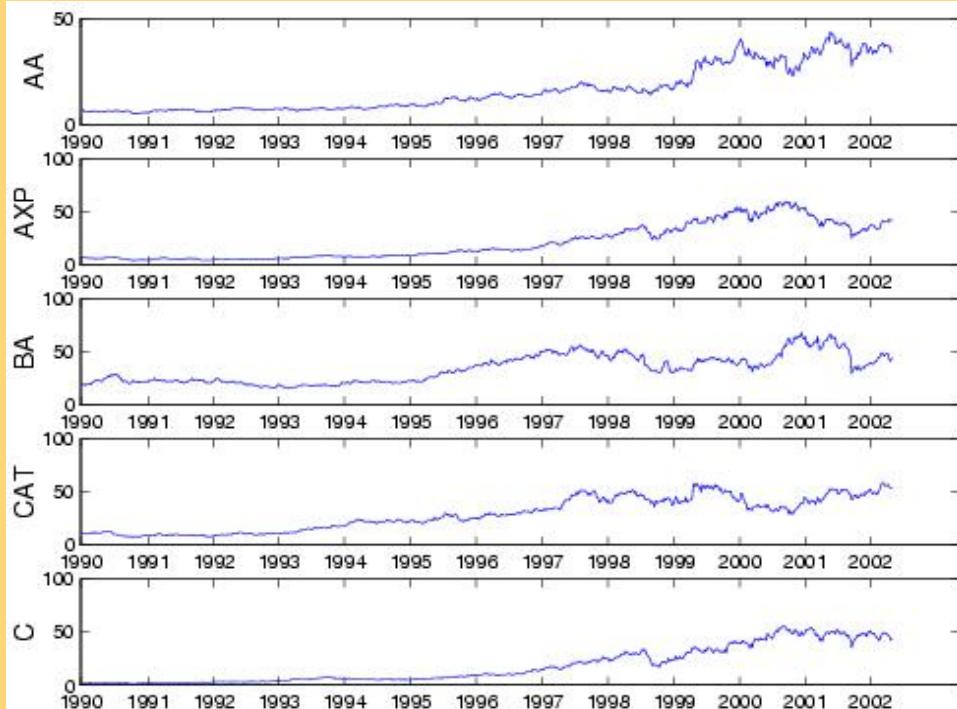


- Motivation
- P1.1. Similarity Search and Indexing
 - – distance functions: Euclidean; Time-warping
 - indexing
 - feature extraction
- P1.2. DSP
- P1.3. Linear forecasting
- ...



P1.1 - Problem:

Q: How similar are two sequences?

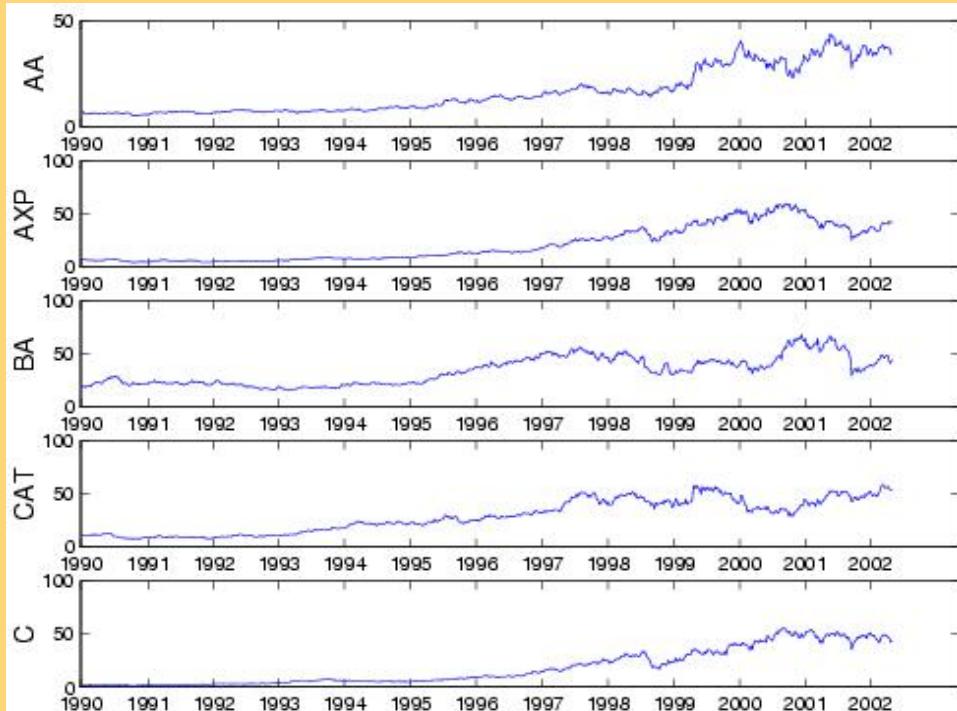


P1.1 - Answer:

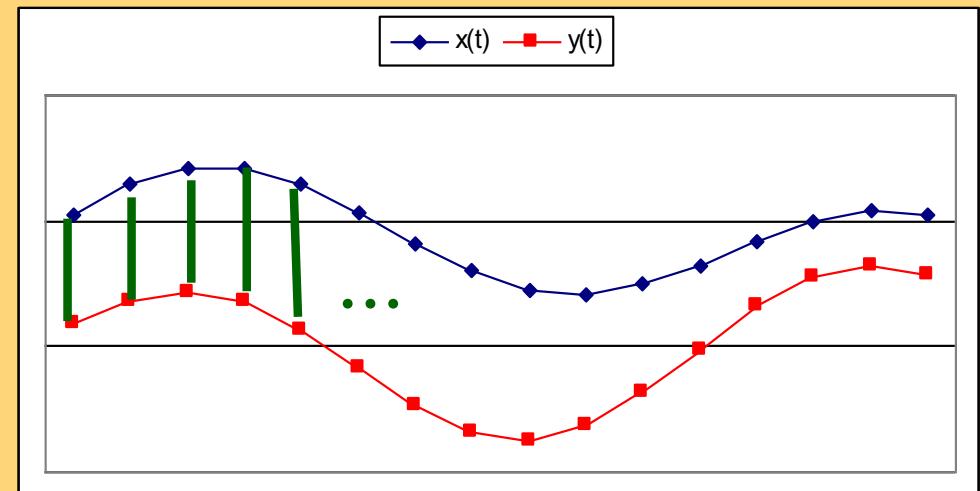


Q: How similar are two sequences?

A: Euclidean distance (<-> cosine similarity)



$$D(\vec{x}, \vec{y}) = \sum_{i=1}^n (x_i - y_i)^2$$



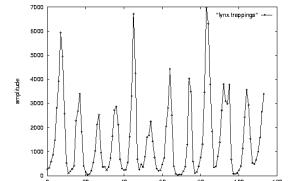
Part 1 - Outline



- Motivation
- P1.1. Similarity Search and Indexing
 - distance functions
 - indexing
 - feature extraction
-
- P1.2. DSP
- P1.3. Linear forecasting
- ...

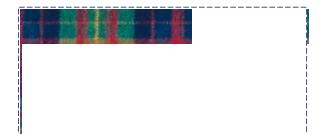


Important observations



Patterns, rules, forecasting and similarity indexing are closely related:

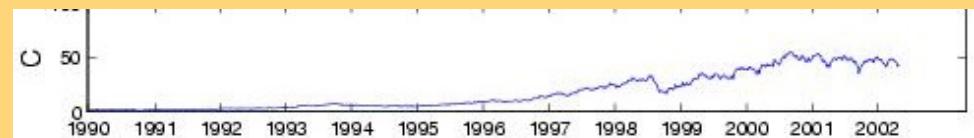
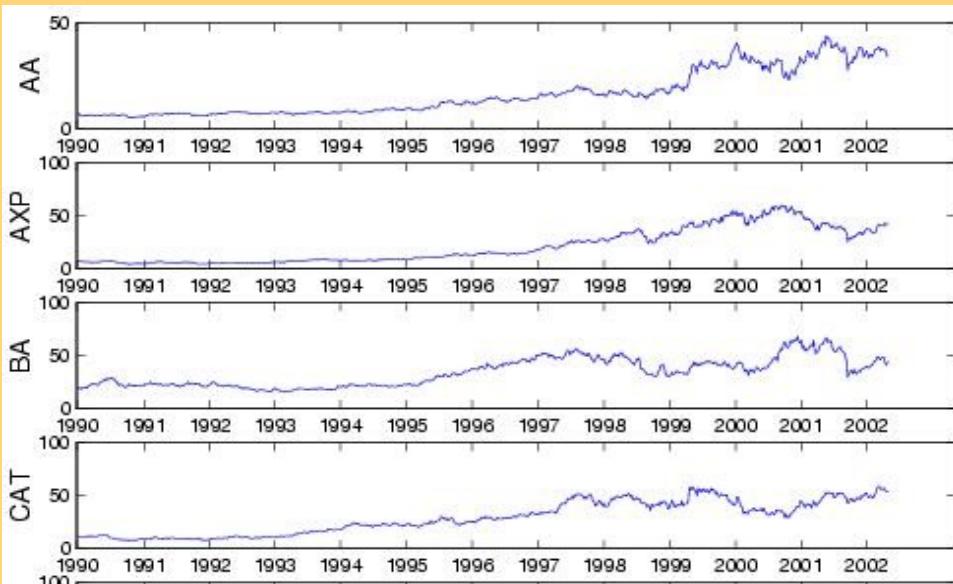
- To do forecasting, we need
 - to find patterns/rules
 - compress
 - **to find similar settings in the past**
- to find outliers, we need to have forecasts
 - (outlier = too far away from our forecast)





Problem:

Q: find quickly stocks like ‘C’ (or customers like ‘smith’)

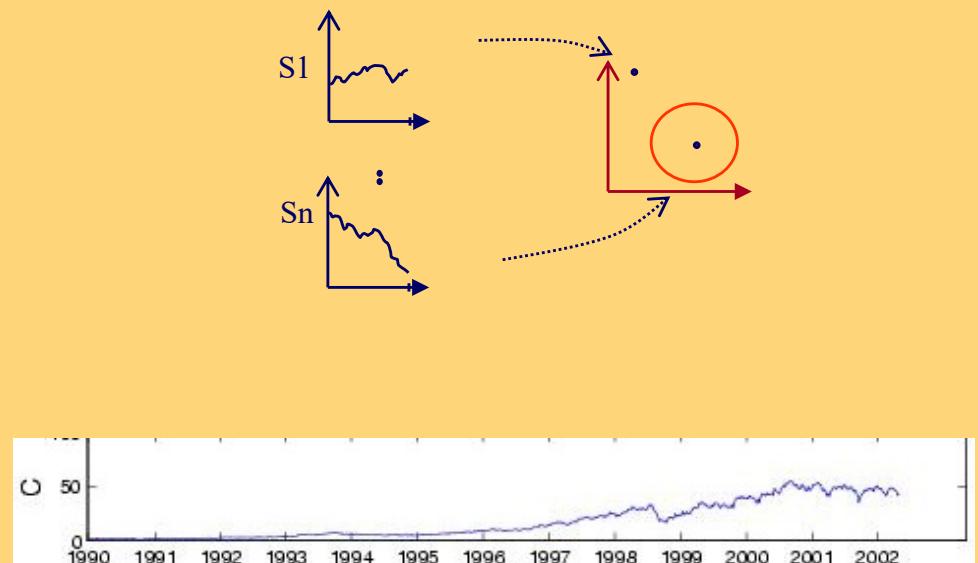
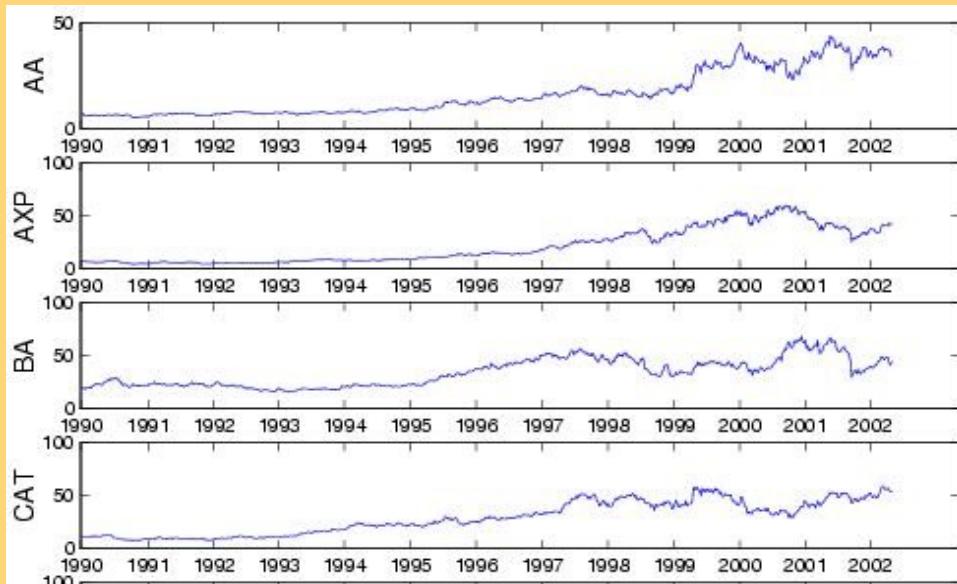


Answer:



Q: find quickly stocks like ‘C’ (or customers like ‘smith’)

A: summarize seq. to a few numbers/features (eg., avg, stdv, Fourier coeff.)



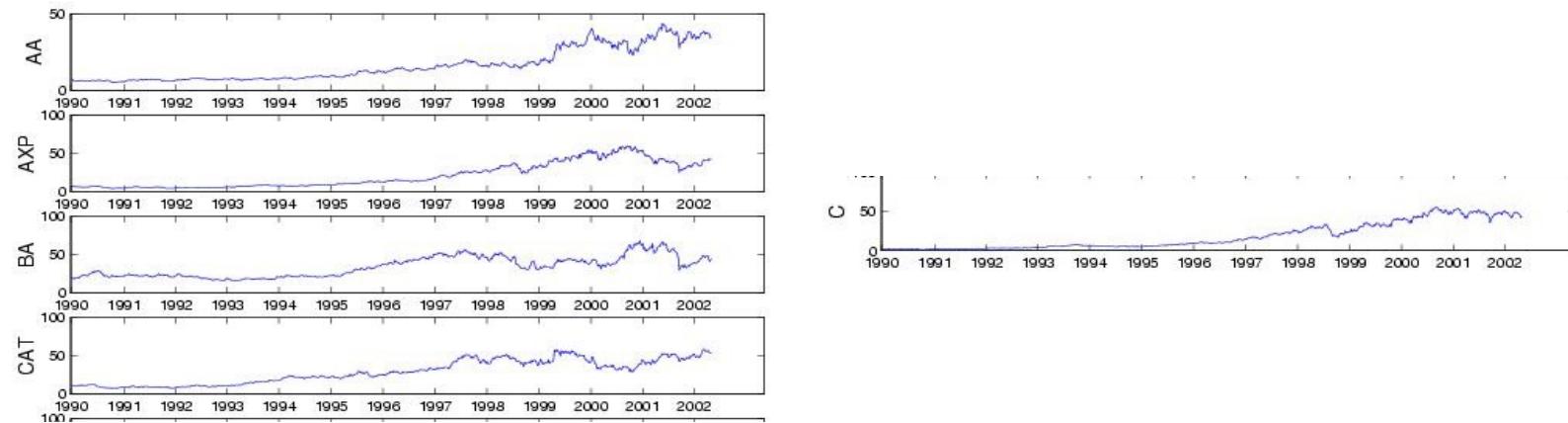
Idea: ‘GEMINI’

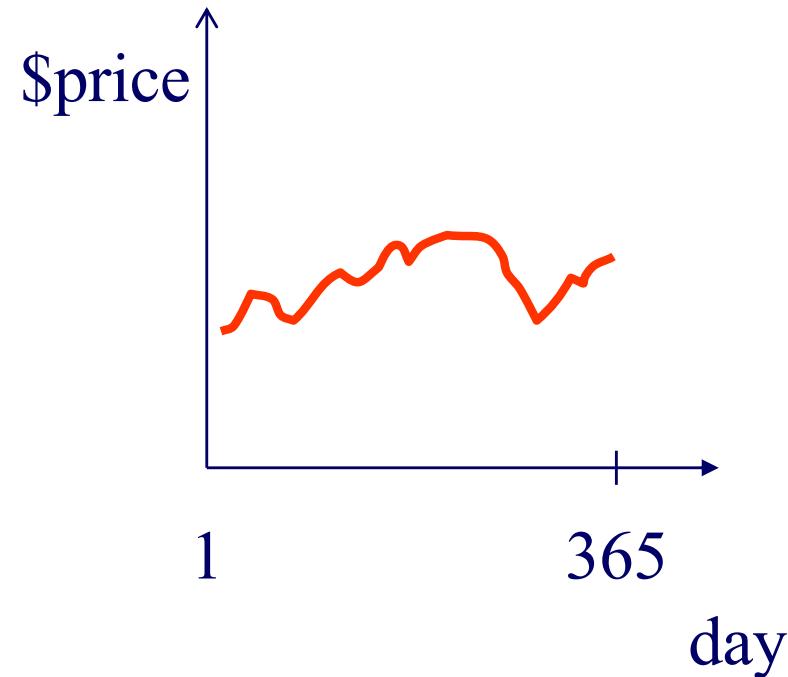
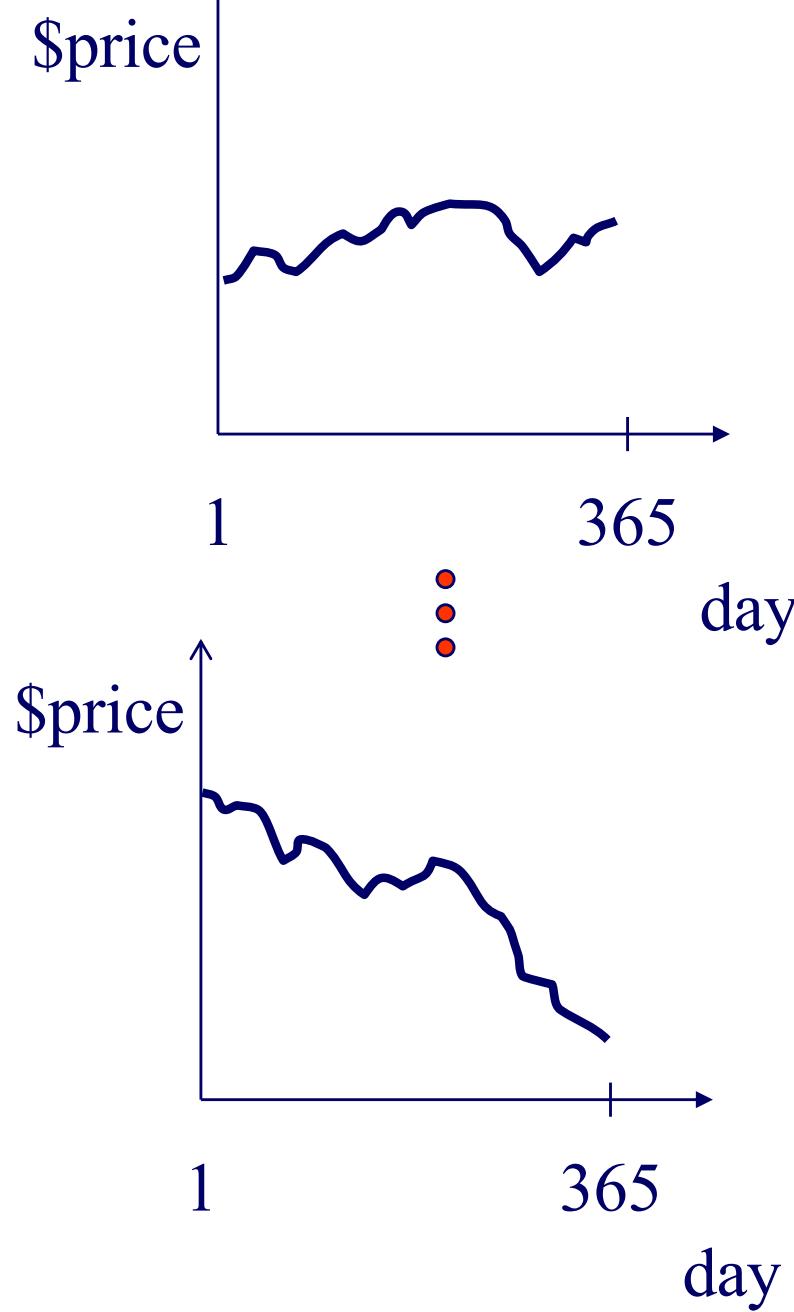
Eg., ‘*find stocks similar to MSFT*’

Seq. scanning: too slow

How to accelerate the search?

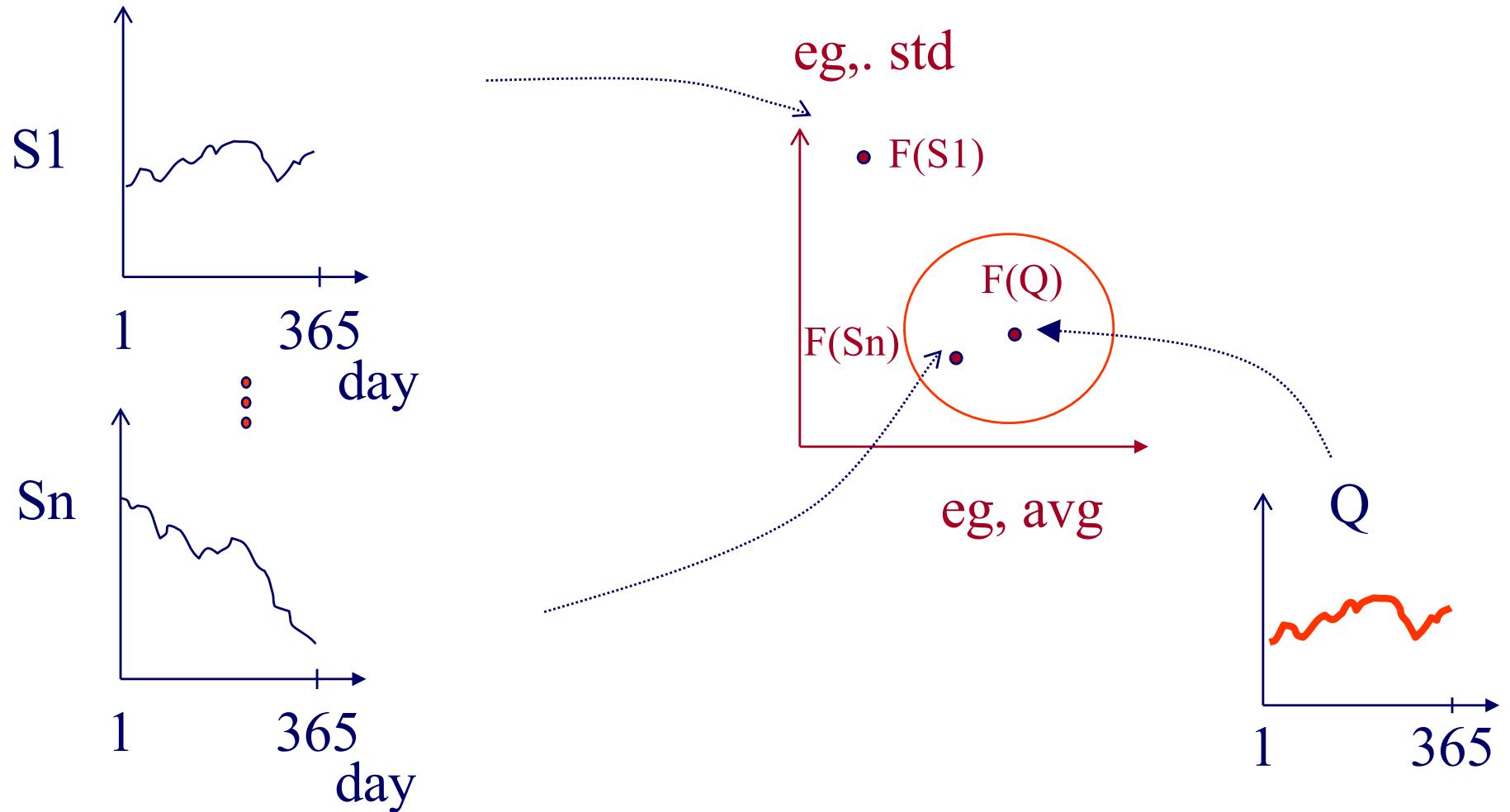
[Faloutsos96]





distance function: by **expert**

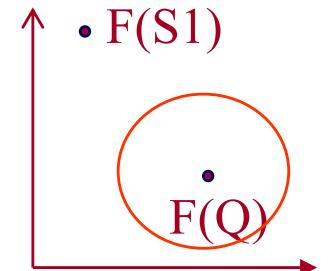
‘GEMINI’ - Pictorially



GEMINI

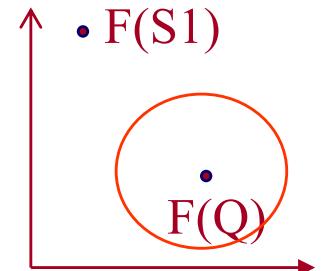
Solution: Quick-and-dirty' filter:

- extract n features (numbers, eg., avg., etc.)
- map into a point in n -d feature space
- organize points with off-the-shelf spatial access method ('SAM')
- discard false alarms



Conclusions

- Fast indexing: through GEMINI
 - feature extraction and
 - (off the shelf) Spatial Access Methods
[Gaede+98] (eg, Rtree, kdtree, even in python)



But: features?

- How to extract ‘good features’?
- A1:
 - **SVD** (Singular Value Decomposition); or
 - **ICA** (Independent Component Analysis)
- A2:
 - **DFT** (Discrete Fourier Transform)
 - **DWT** (Discrete Wavelet Transform)



Important observations:

Feature extraction =

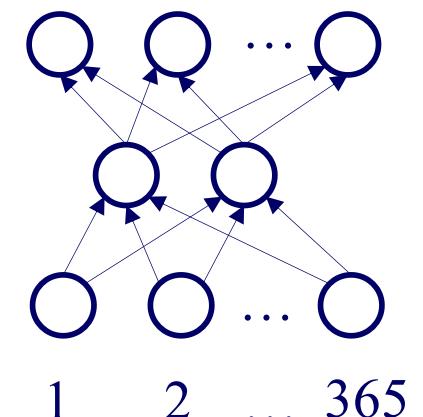
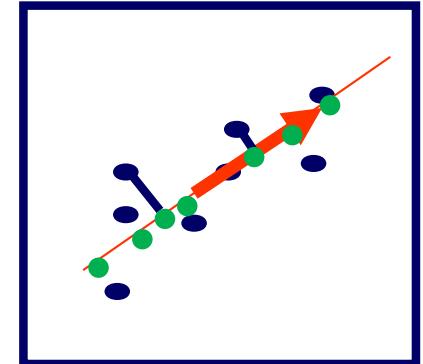
= *(lossy) compression*

= *Dimensionality reduction*

= *Embedding*

= *Latent/hidden variable detection*

= *AutoEncoding*

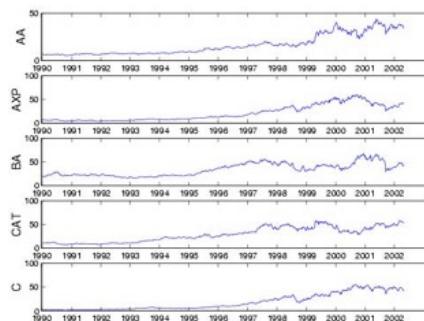


Part 1 - Outline

- Motivation
- P1.1. Similarity Search and Indexing
 - distance functions
 - indexing
 - feature extraction
 - SVD etc (data dependent)
 - (DFT, DWT, DCT (data independent))
- P1.2. DSP (DFT, DWT)



SVD



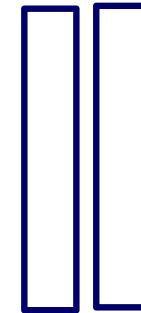
AA
AXP

t_1

t_n



\sim



$u_1 \quad u_2$

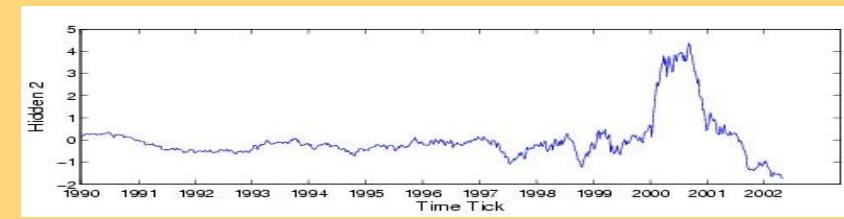
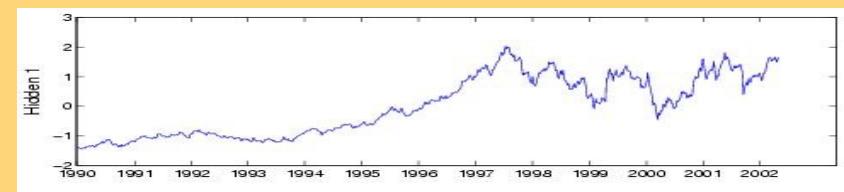
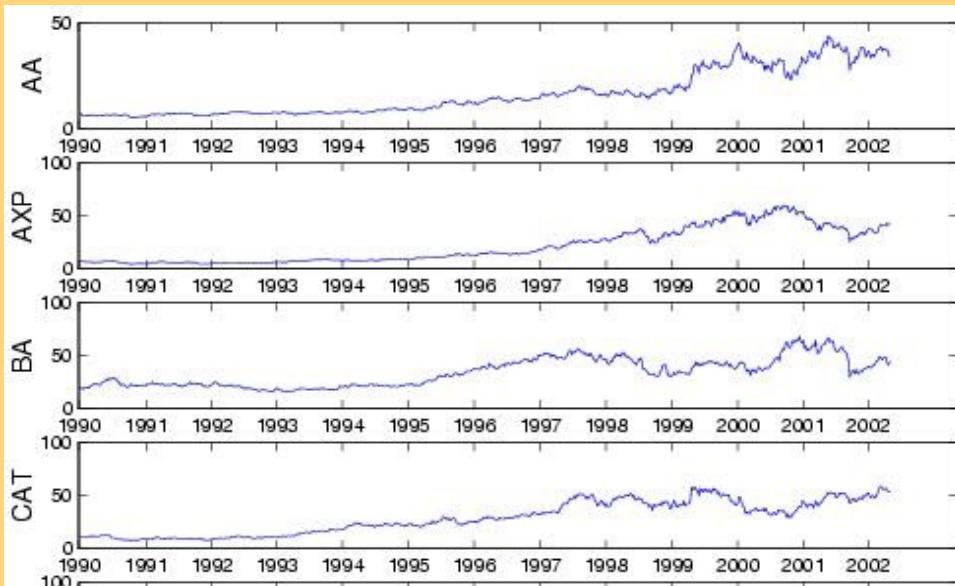
```
# svd code
import numpy as np
u, s, vh = np.linalg.svd(b)
```

$\mathbf{U} \quad \Sigma \quad \mathbf{V}^T$

Answer



Q: how to extract features (commonalities)? (given the data)
A: SVD, ICA



Citations

- SVD** • Jolliffe, I. T. (1986). *Principal Component Analysis*, Springer Verlag.
- ICA** • *AutoSplit: Fast and Scalable Discovery of Hidden Variables in Stream and Multimedia Databases*,
Jia-Yu Pan, Hiroyuki Kitagawa, Christos Faloutsos and Masafumi Hamamoto PAKDD 2004, Sydney, Australia



Books + lecture notes

- William H. Press, Saul A. Teukolsky, William T. Vetterling and Brian P. Flannery: *Numerical Recipes in C*, Cambridge University Press, 1992, 2nd Edition. (Great description, intuition and code for SVD)
- C. Faloutsos: *Searching Multimedia Databases by Content*, Kluwer Academic Press, 1996 (introduction to SVD, and GEMINI)
- CMU class slides: [SVD](#), and [ICA](#)



References

- Foltz, P. W. and S. T. Dumais (Dec. 1992). “*Personalized Information Delivery: An Analysis of Information Filtering Methods.*” Comm. of ACM (CACM) 35(12): 51-60.
- Guttman, A. (June 1984). *R-Trees: A Dynamic Index Structure for Spatial Searching.* Proc. ACM SIGMOD, Boston, Mass.

References

- Gaede, V. and O. Guenther (1998). “*Multidimensional Access Methods.*” Computing Surveys 30(2): 170-231.
- Gunopulos, D. and G. Das (2001). *Time Series Similarity Measures and Time Series Indexing*. SIGMOD Conference, Santa Barbara, CA.
- Eamonn J. Keogh, Themis Palpanas, Victor B. Zordan, Dimitrios Gunopulos, Marc Cardle: *Indexing Large Human-Motion Databases*. VLDB 2004: 780-791

Part 1.2: DSP (Digital Signal Processing)

Part 1 - Outline



- Motivation
- P1.1. Similarity Search and Indexing
- P1.2. DSP (DFT, DWT)
 - – Fourier (DFT)
 - Wavelets (DWT)
- P1.3. Linear Forecasting
- P1.4. Non-linear forecasting
- P1.5. Tensors
- Conclusions

Part 1.2 - Outline

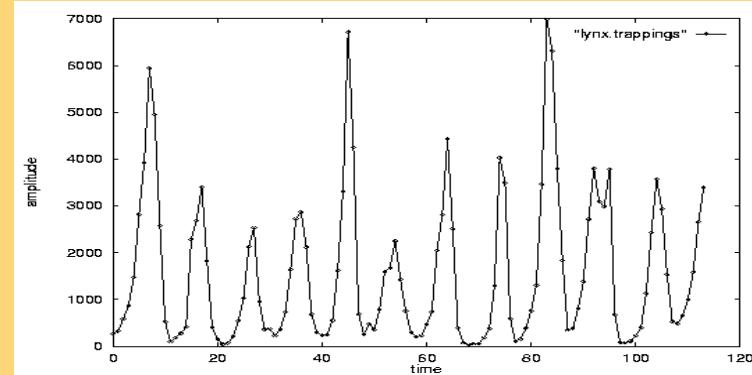
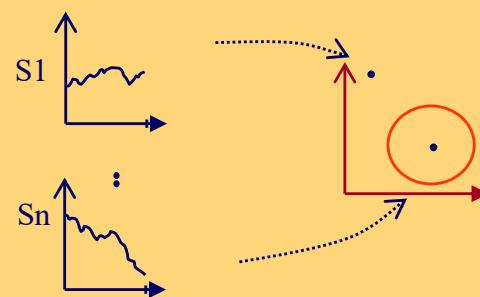


- DFT
 - Definition of DFT and properties
 - how to read the DFT spectrum
- DWT
 - Definition of DWT and properties
 - how to read the DWT scalogram



P1.2 - Problem

Q: How to summarize / extract few features



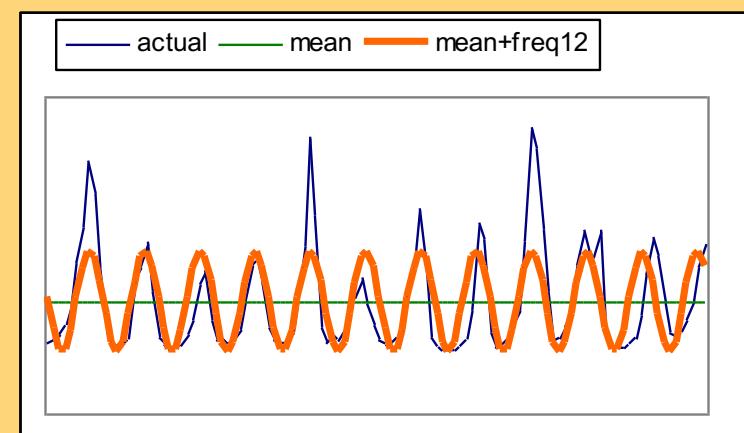
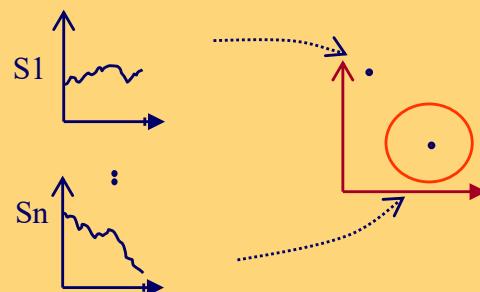
P1.2 - Answer:



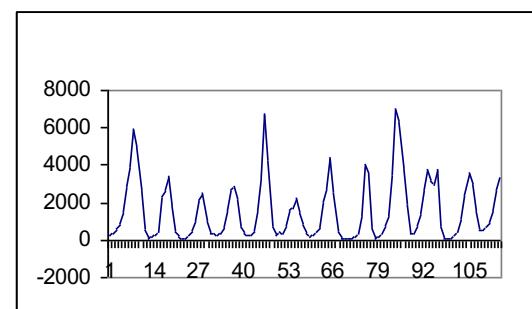
Q: How to summarize / extract few features

A1: Data dep.: SVD, ICA

A2: Data indep.: Fourier; Wavelets

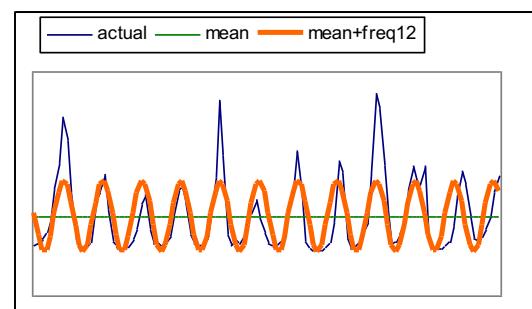


What does DFT do?



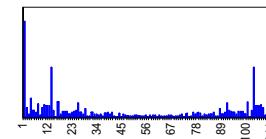
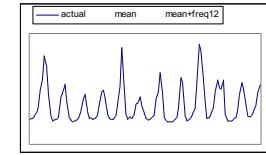
What does DFT do?

A: highlights the periodicities



DFT: definition

- For a sequence x_0, x_1, \dots, x_{n-1}
- the (n-point) Discrete Fourier Transform is
- X_0, X_1, \dots, X_{n-1} :



$$X_f = 1/\sqrt{n} \sum_{t=0}^{n-1} x_t * \exp(-j2\pi tf / n) \quad f = 0, \dots, n-1$$

($j = \sqrt{-1}$)

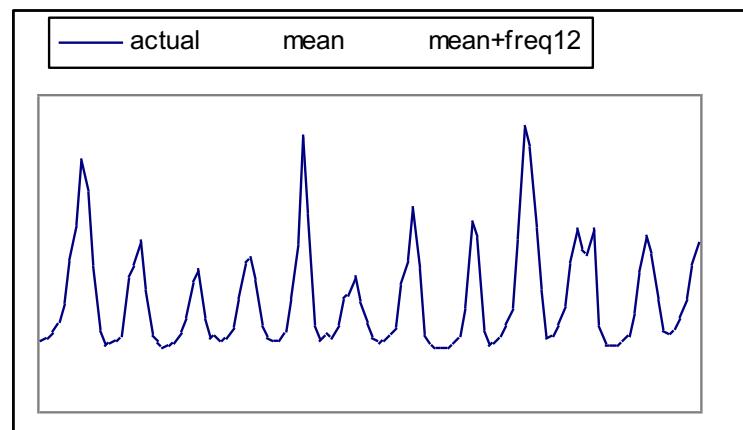
$$x_t = 1/\sqrt{n} \sum_{f=0}^{n-1} X_f * \exp(+j2\pi tf / n)$$

inverse DFT

DFT: Amplitude spectrum

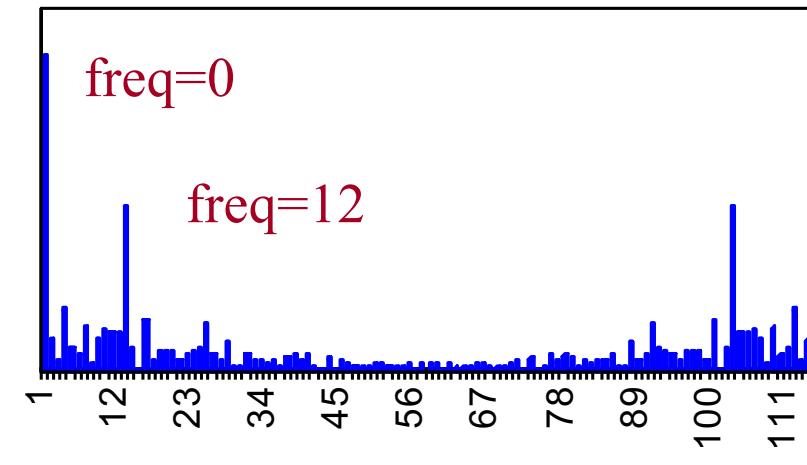
$$\text{Amplitude: } A_f^2 = \text{Re}^2(X_f) + \text{Im}^2(X_f)$$

count



year

Ampl.

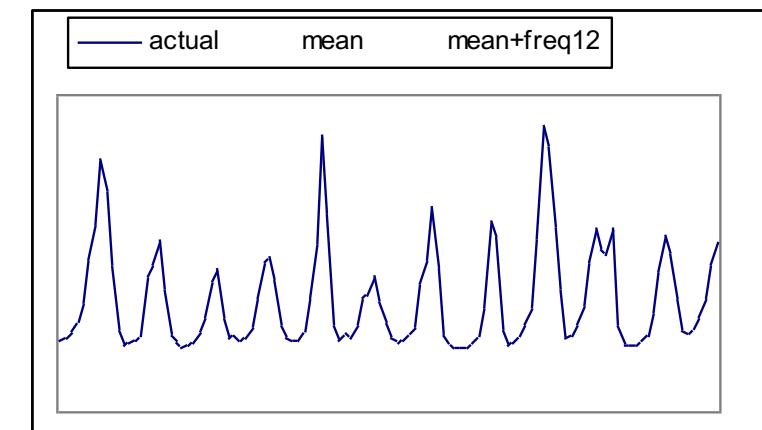


Freq.

DFT: Amplitude spectrum

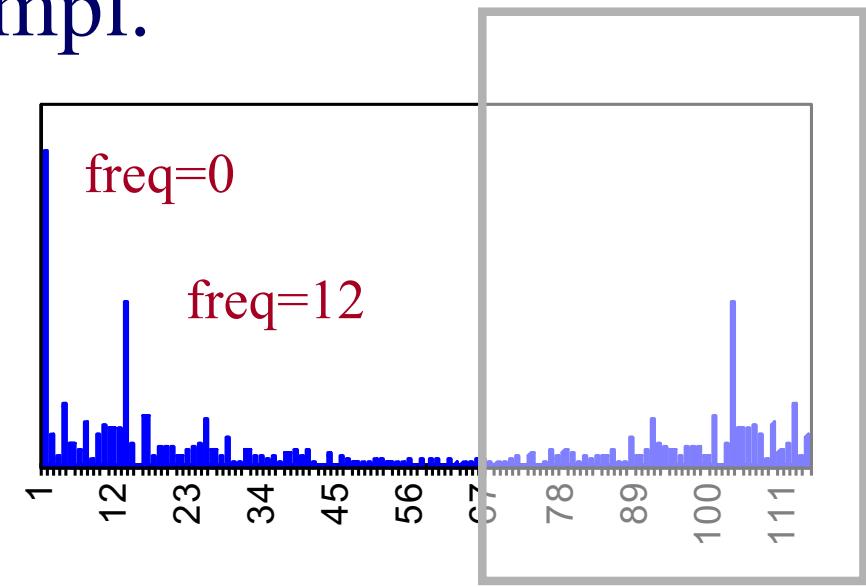
$$\text{Amplitude: } A_f^2 = \text{Re}^2(X_f) + \text{Im}^2(X_f)$$

count



year

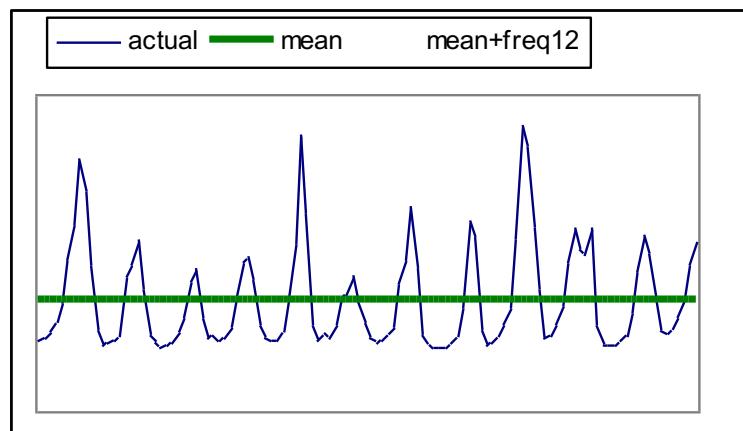
Ampl.



Freq.

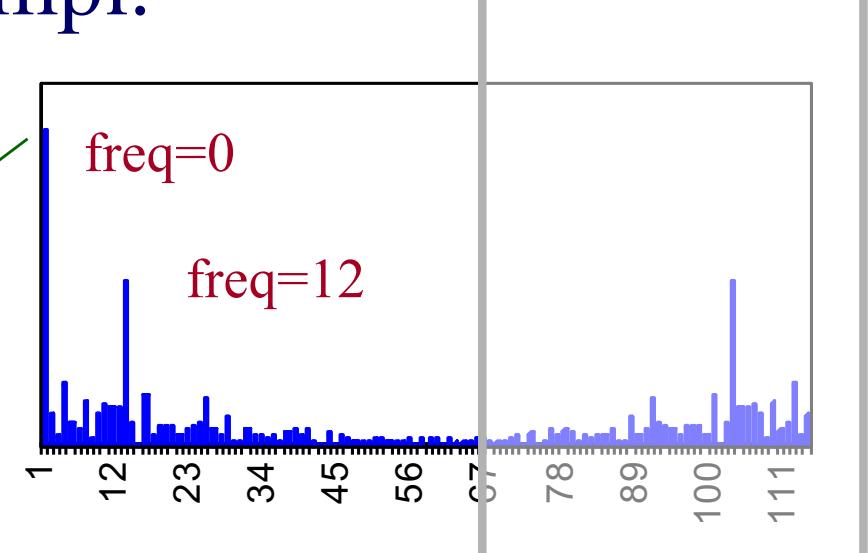
DFT: Amplitude spectrum

count



year

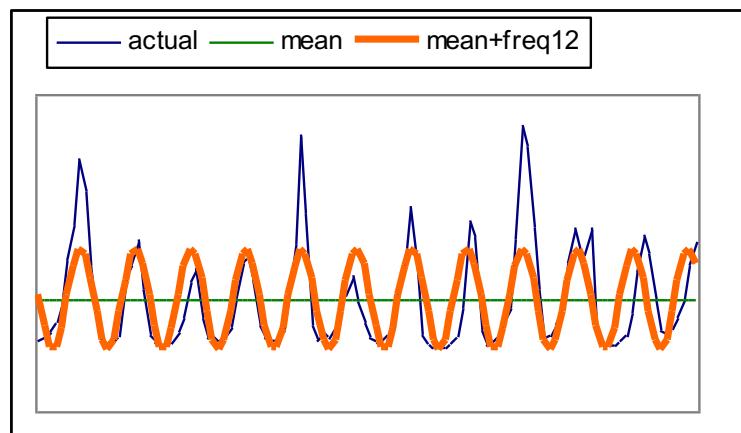
Ampl.



Freq.

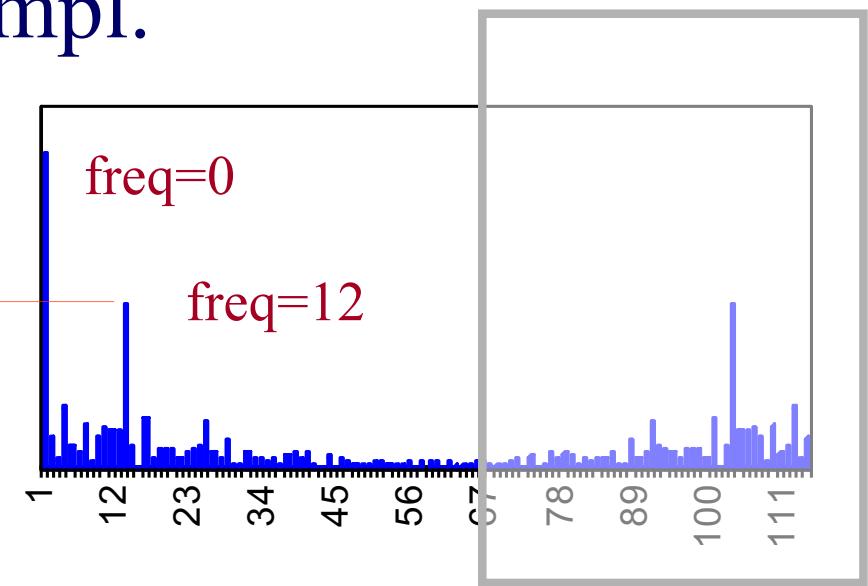
DFT: Amplitude spectrum

count



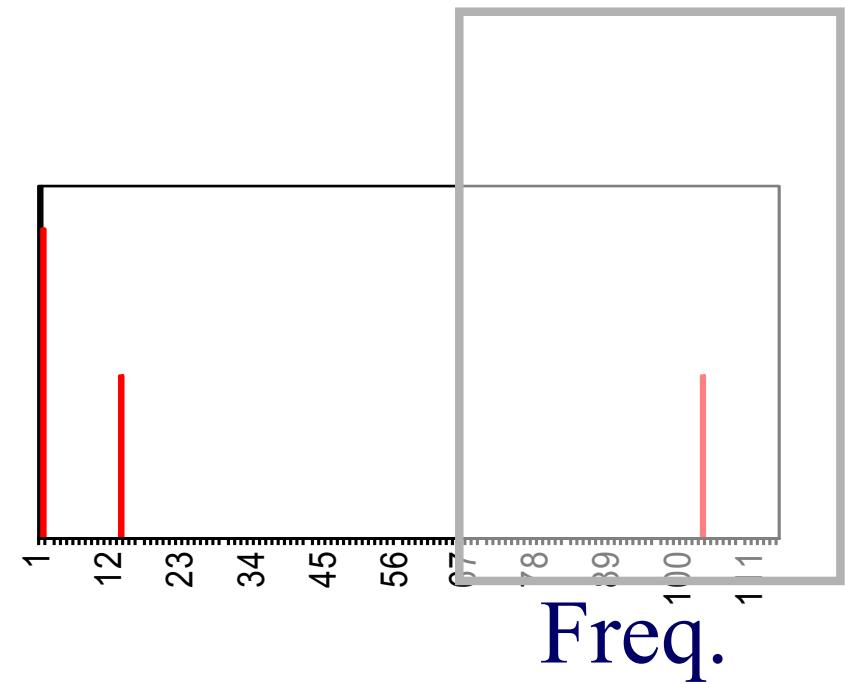
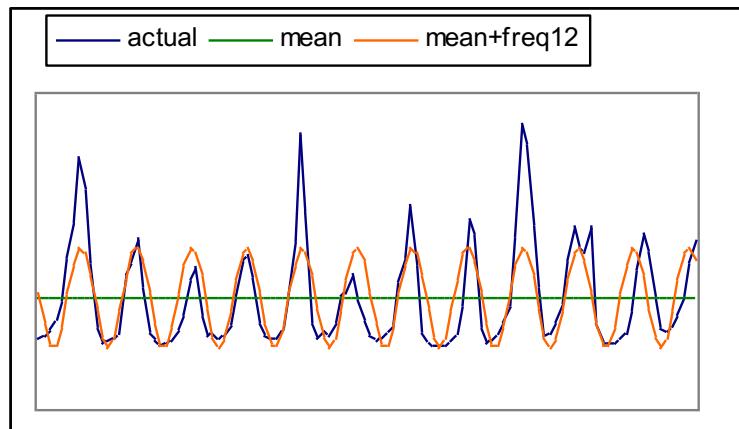
year

Ampl.



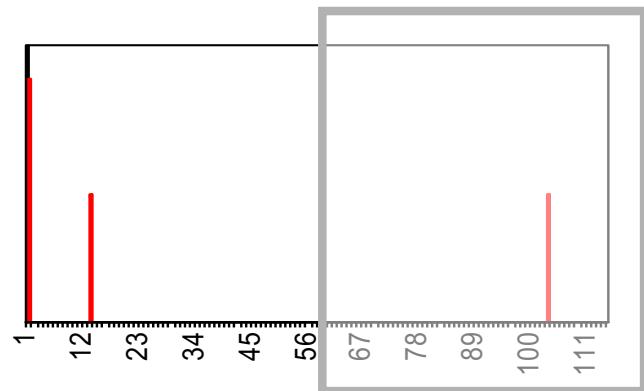
DFT: Amplitude spectrum

- excellent approximation, with only 2 frequencies! (2 int, 3 float)
- so what?



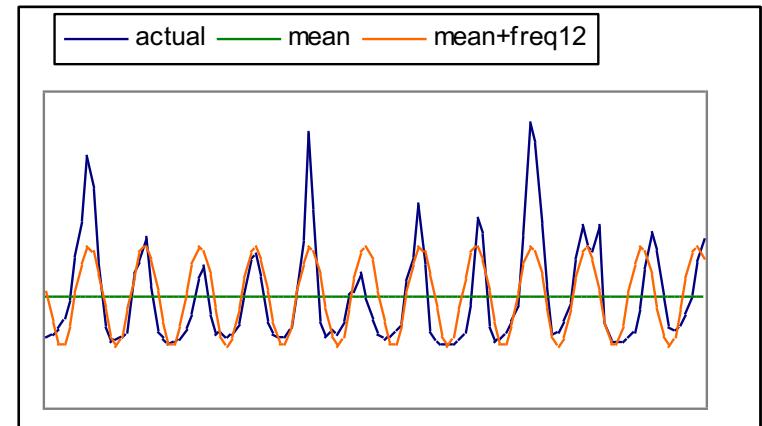
DFT: Amplitude spectrum

- excellent approximation, with only 2 frequencies!
- so what?
- A1: **(lossy) compression**
- A2: pattern discovery



DFT: Amplitude spectrum

- excellent approximation, with only 2 frequencies!
- so what?
- A1: (lossy) compression
- A2: **pattern discovery**



DFT - Conclusions

- It spots periodicities (with the '**amplitude spectrum**')
- can be quickly computed ($O(n \log n)$), thanks to the FFT algorithm.
- **standard** tool in signal processing (speech, image etc signals)
- (closely related to DCT and JPEG)

Outline



- Motivation
- P1.1. Similarity Search and Indexing
- P1.2. DSP
 - DFT
 - DWT
 - Definition of DWT and properties
 - how to read the DWT scalogram
- P1.3. Linear Forecasting
-

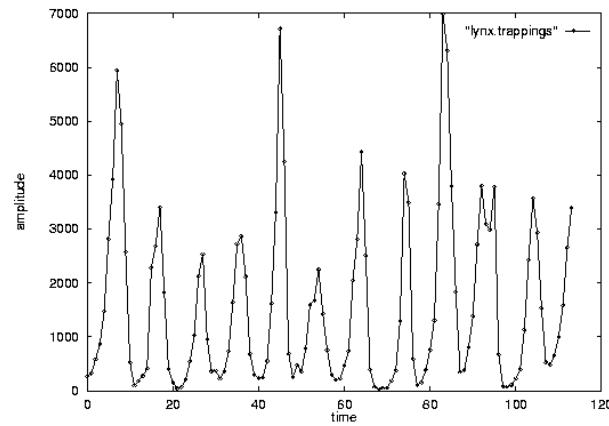


Problem #1:

Goal: given a signal (eg., #packets over time)

Find: patterns, periodicities, and/or compress

count



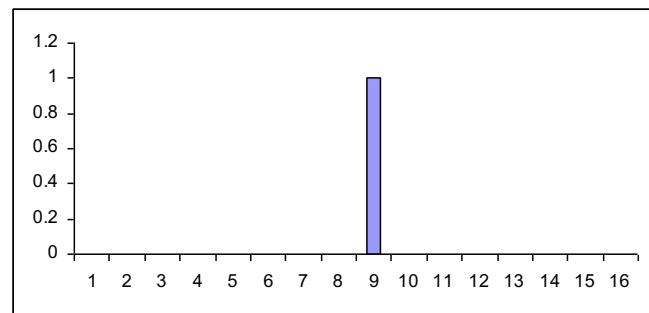
lynx caught per year
(packets per day;
virus infections per month)

year

Wavelets - DWT

- DFT is great - but, how about compressing a spike?

value

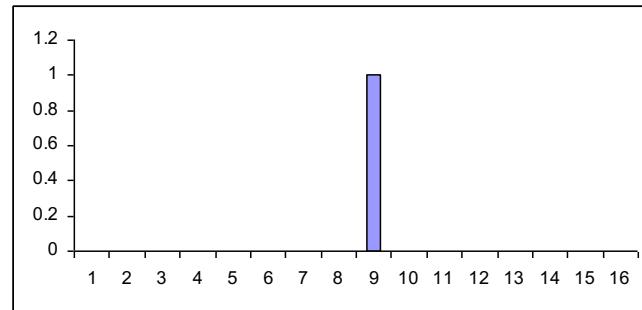


time

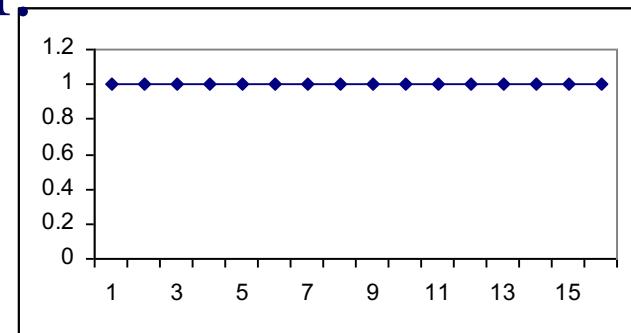
Wavelets - DWT

- DFT is great - but, how about compressing a spike?
- A: Terrible - all DFT coefficients needed!

value



Ampl.

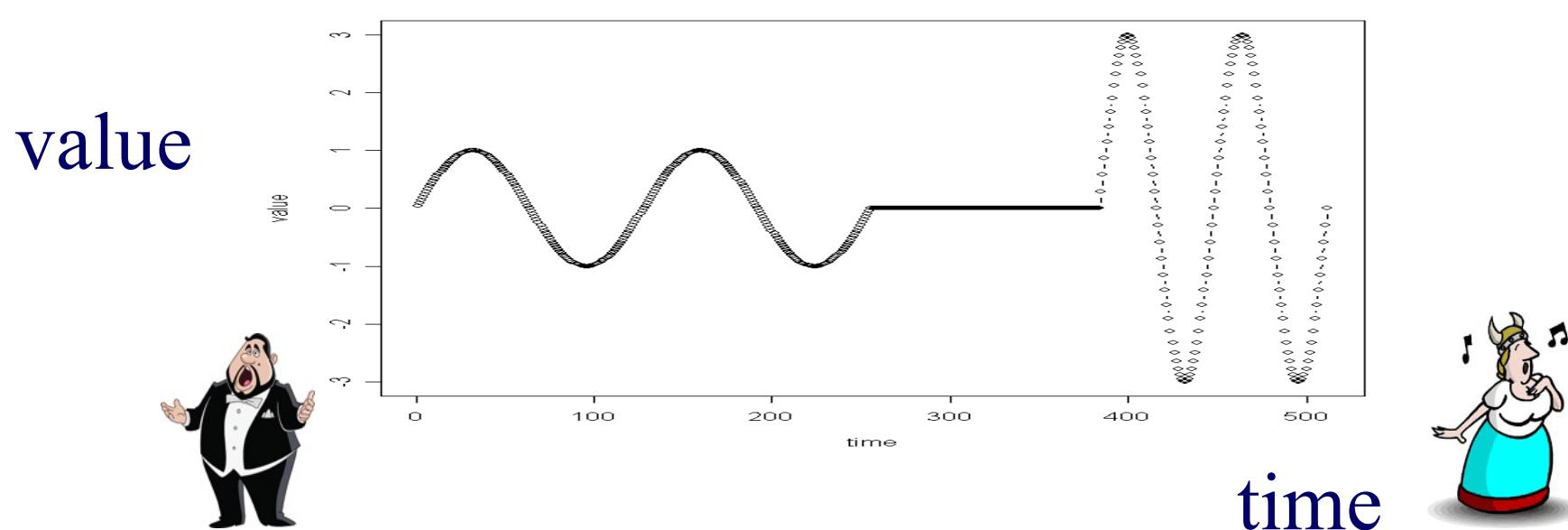


time

Freq.₅₇

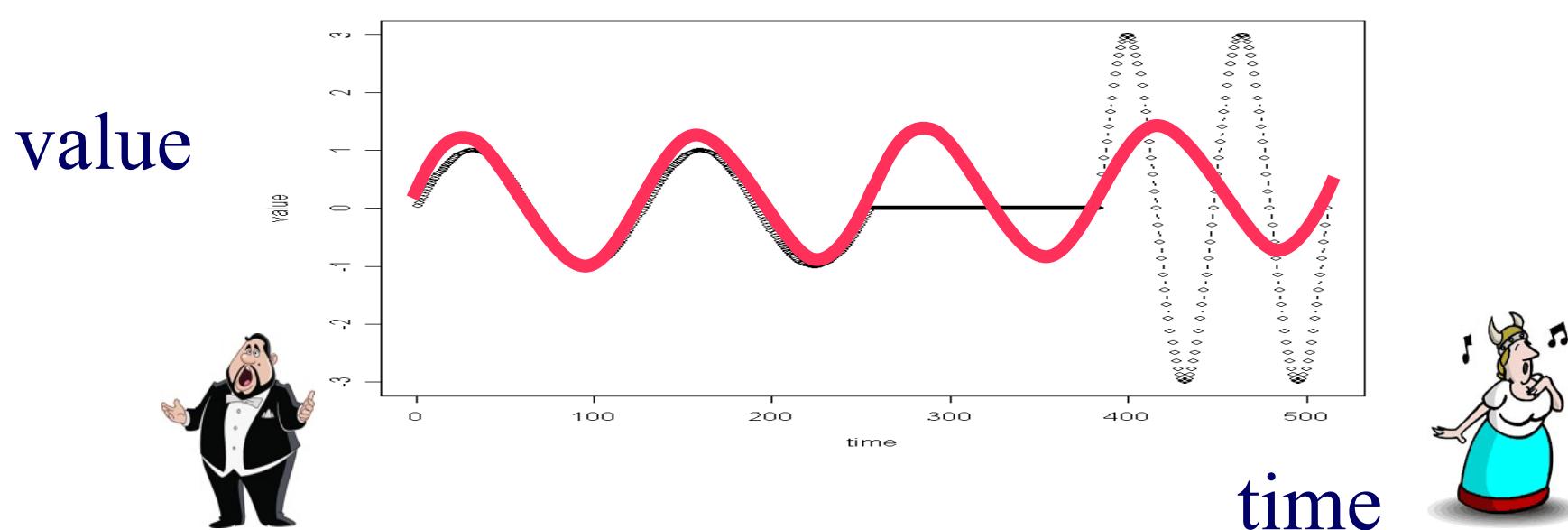
Wavelets - DWT

- Similarly, DFT suffers on short-duration waves (eg., baritone, silence, soprano)



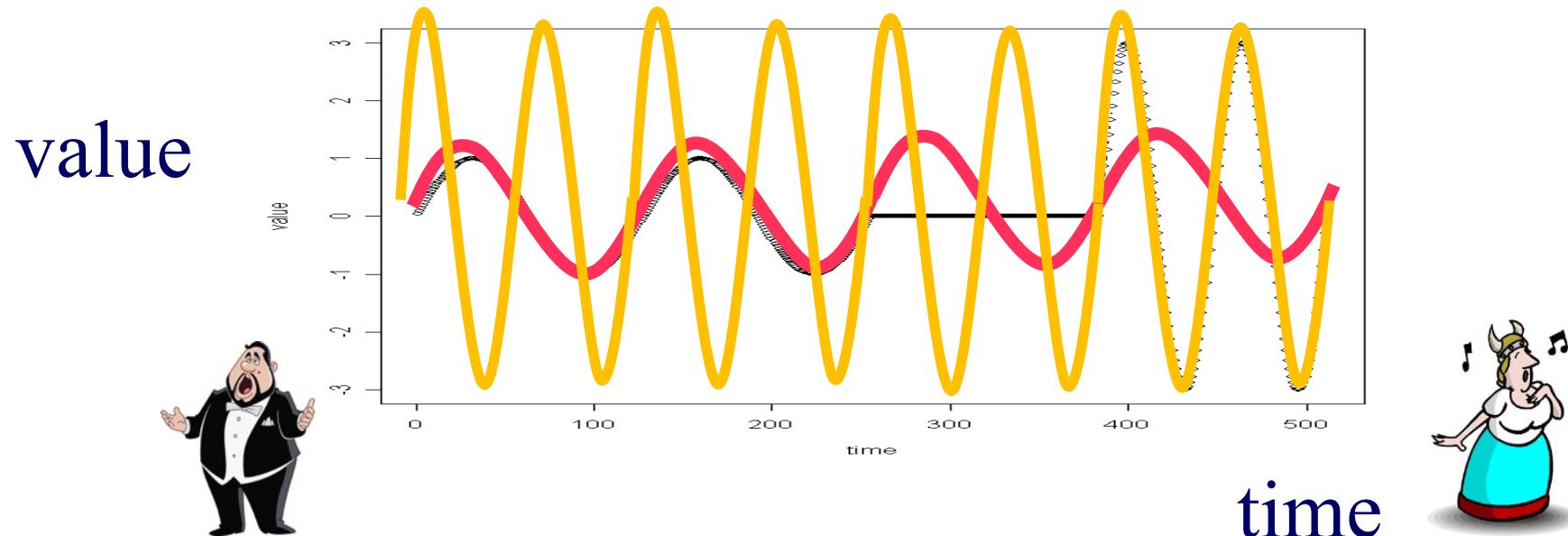
Wavelets - DWT

- Similarly, DFT suffers on short-duration waves (eg., baritone, silence, soprano)



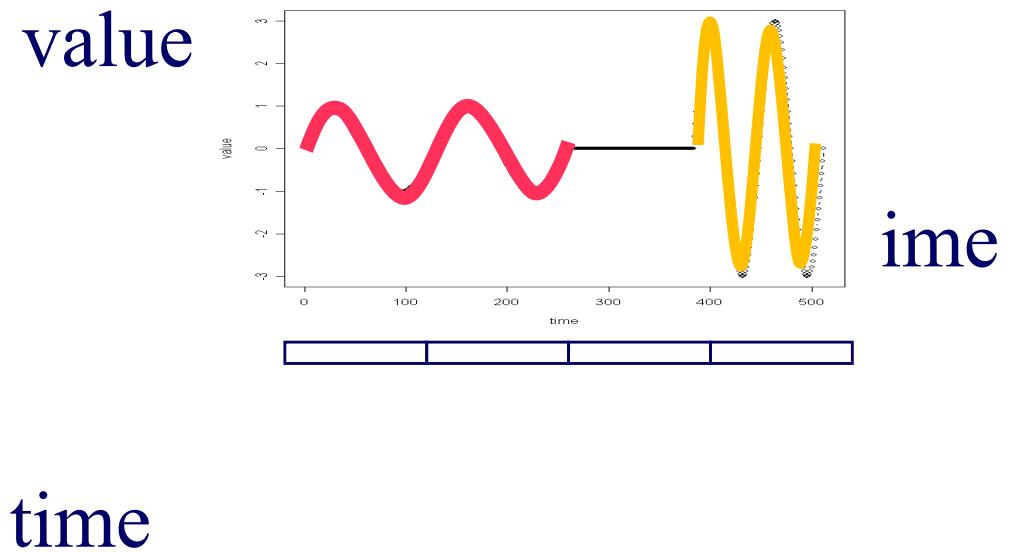
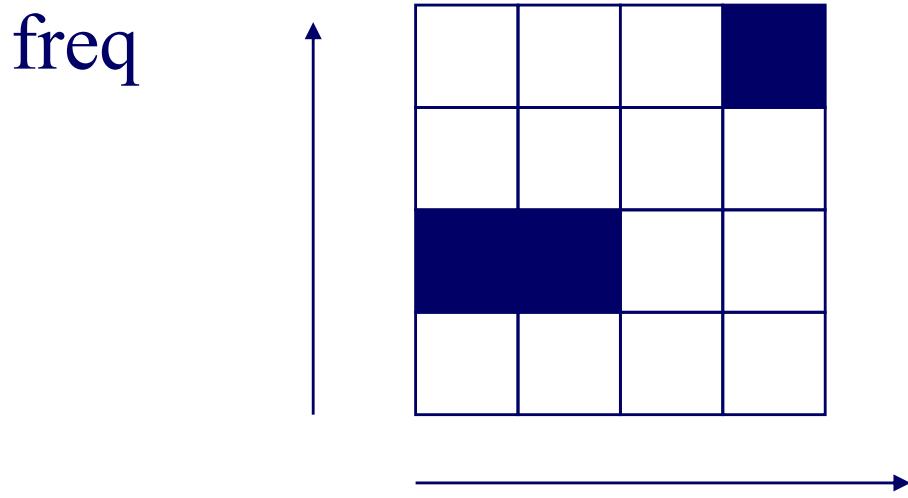
Wavelets - DWT

- Similarly, DFT suffers on short-duration waves (eg., baritone, silence, soprano)



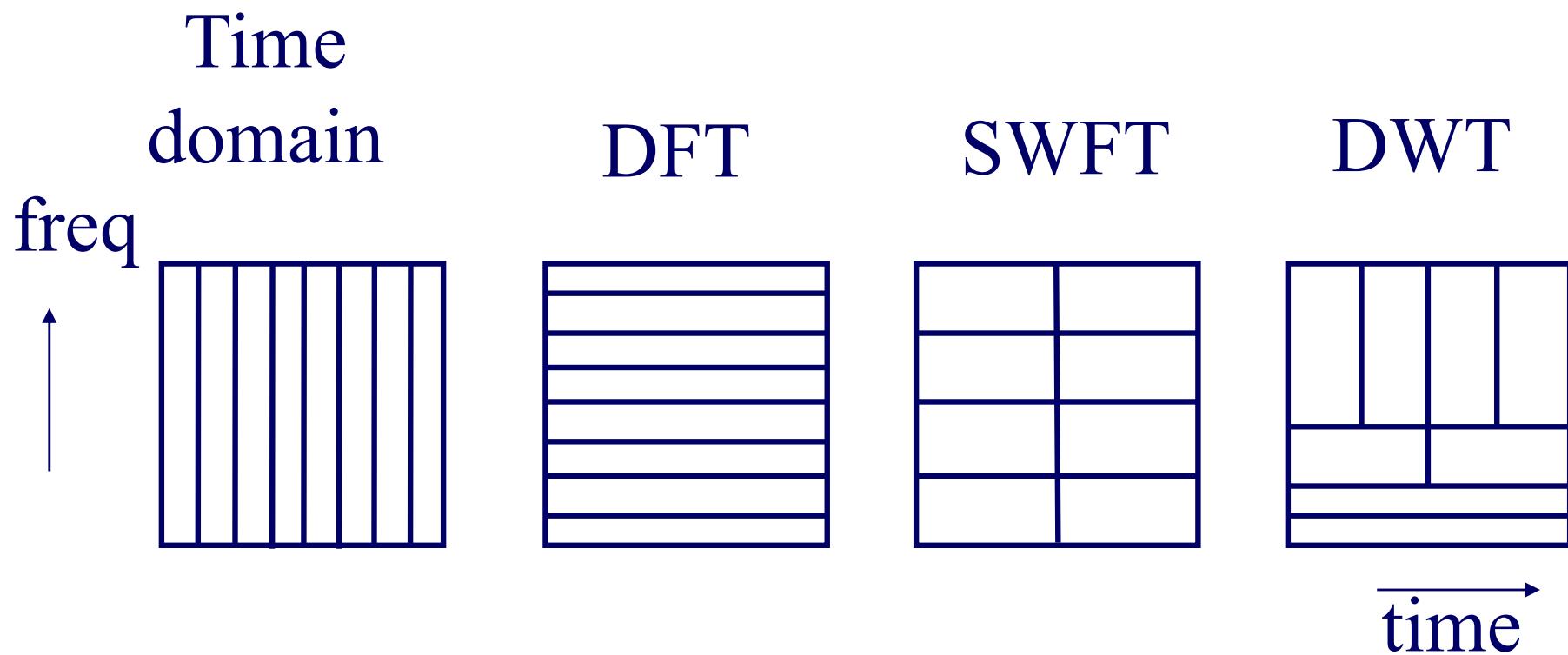
Wavelets - DWT

- Solution#1: Short window Fourier transform (SWFT)
- But: how short should be the window?



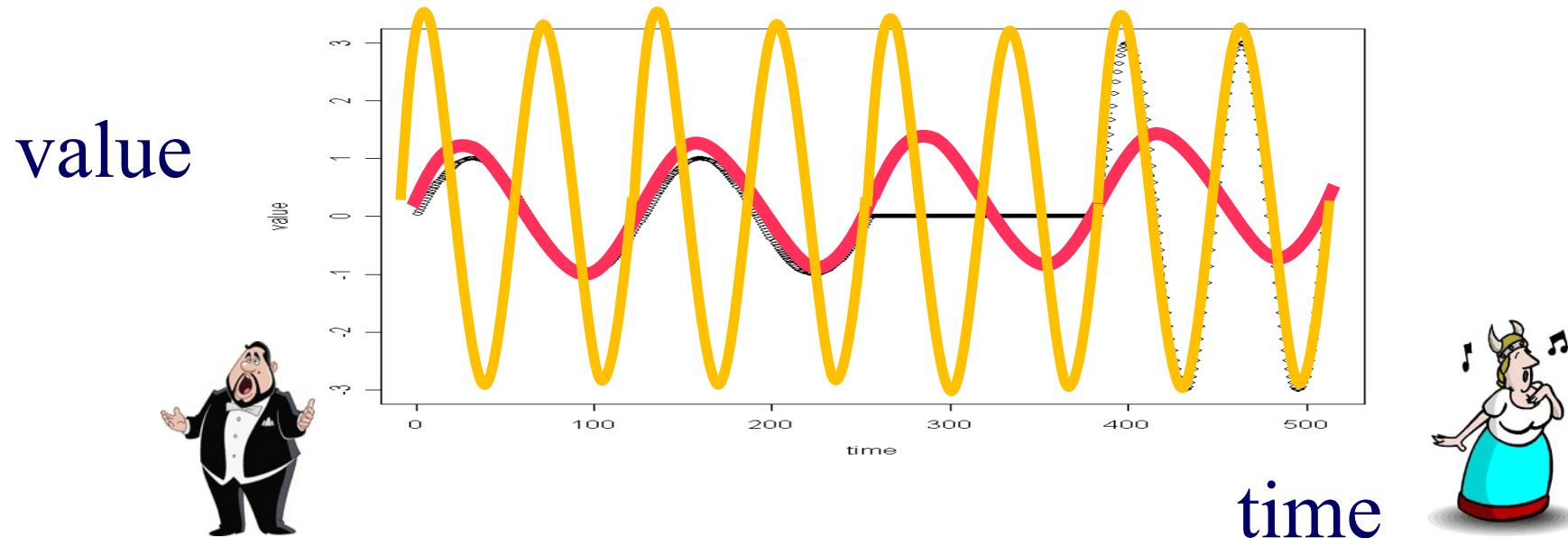
Wavelets - DWT

- Answer: **multiple** window sizes! \rightarrow DWT



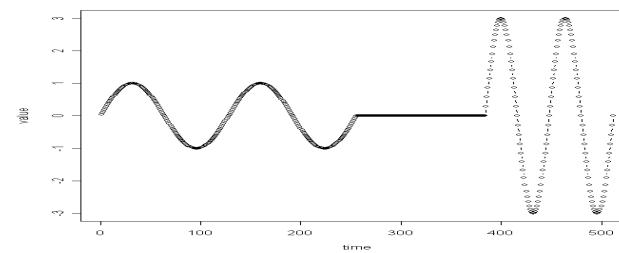
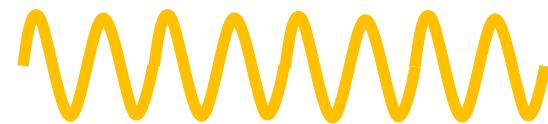
(Basis functions of DFT)

- Sine waves of several frequencies (1, 2, 3, ...)



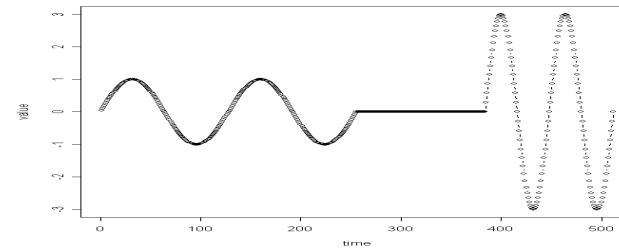
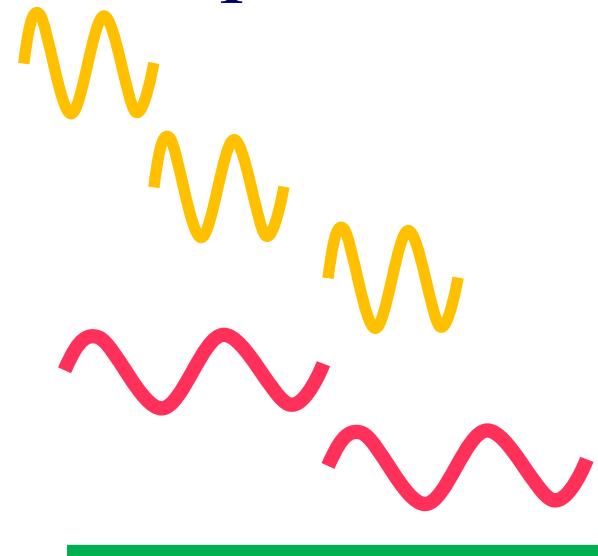
(Basis functions of DFT)

- Sine waves of several frequencies (1, 2, 3, ...)



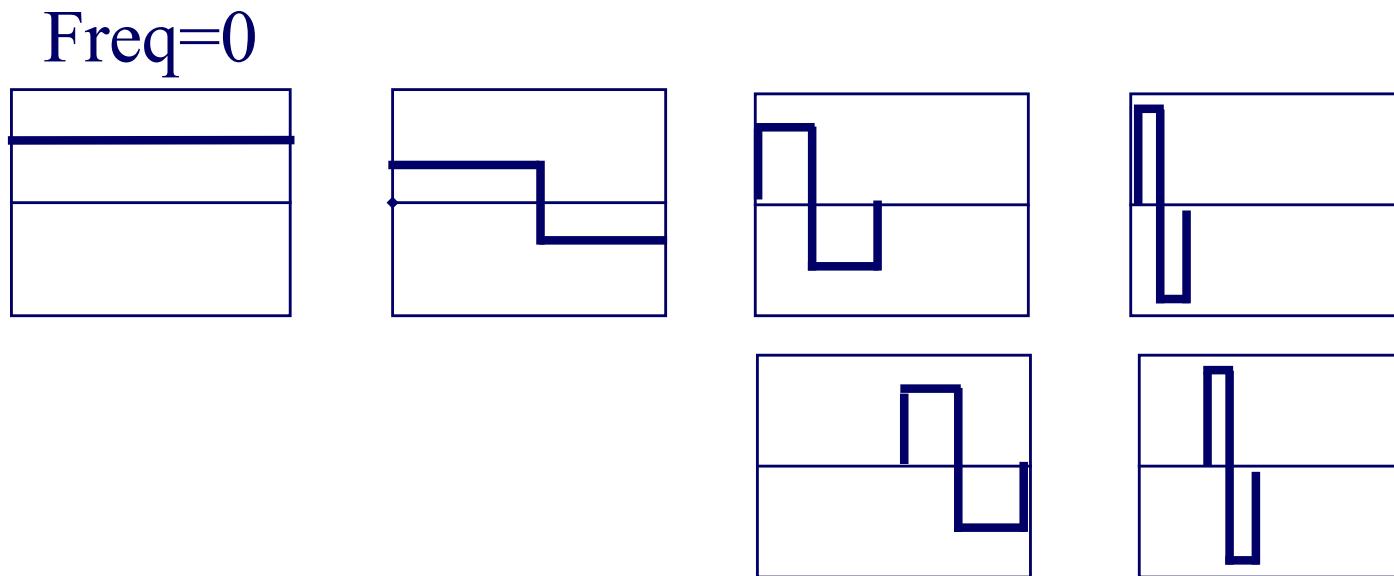
(Basis functions of DWT)

- (\sim Sine) waves of several frequencies (1,2,4,8,...)
- and LIMITED timespan: ‘little waves’



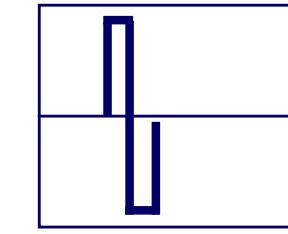
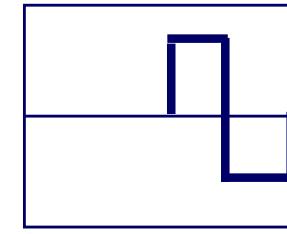
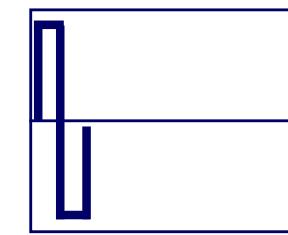
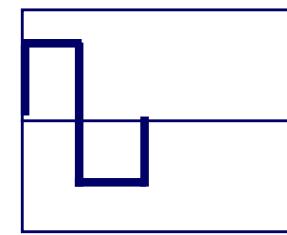
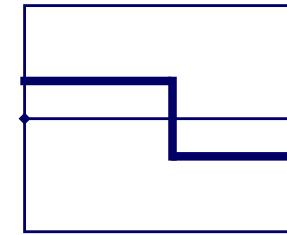
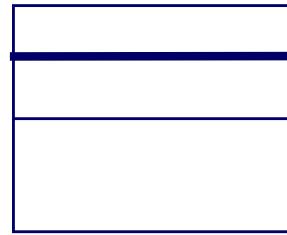
Specifically, Haar Wavelets:

- subtract sum of left half from right half
- repeat recursively for quarters, eighths, ...



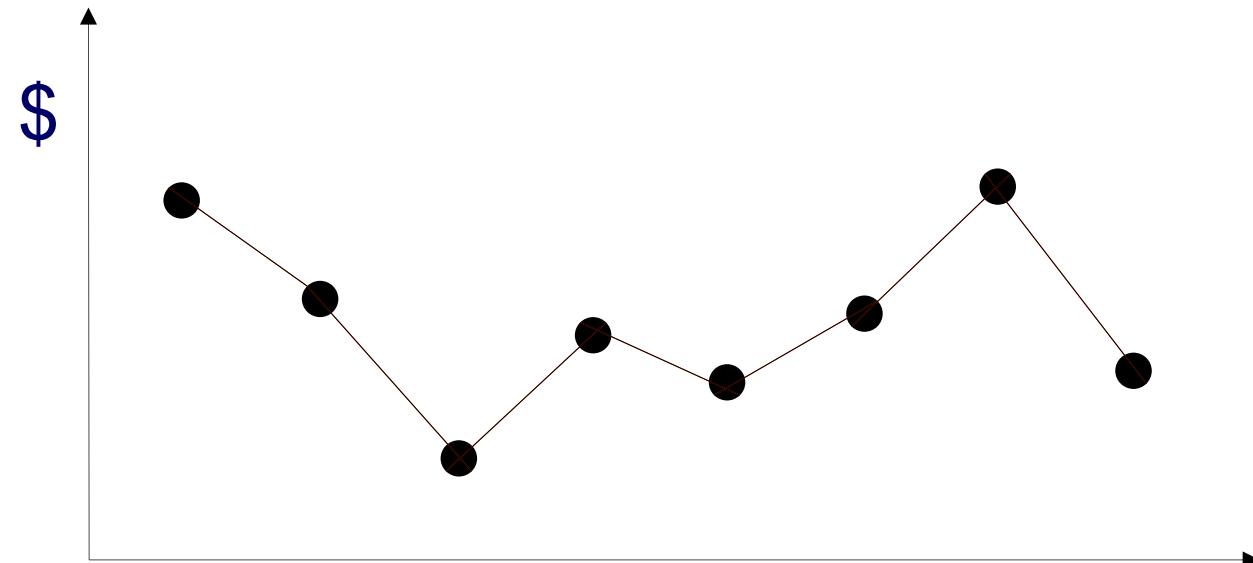
Haar Wavelets

- subtract sum of left half from right half
- repeat recursively for quarters, eight-ths, ...



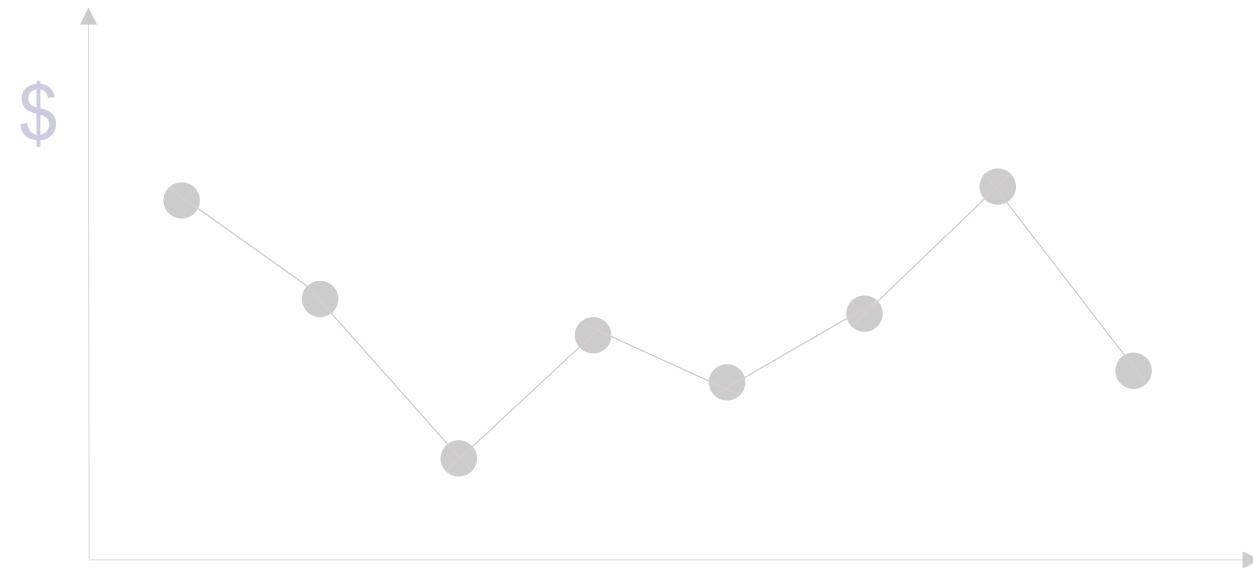
Similar to conv. N.N.

Wavelets - construction



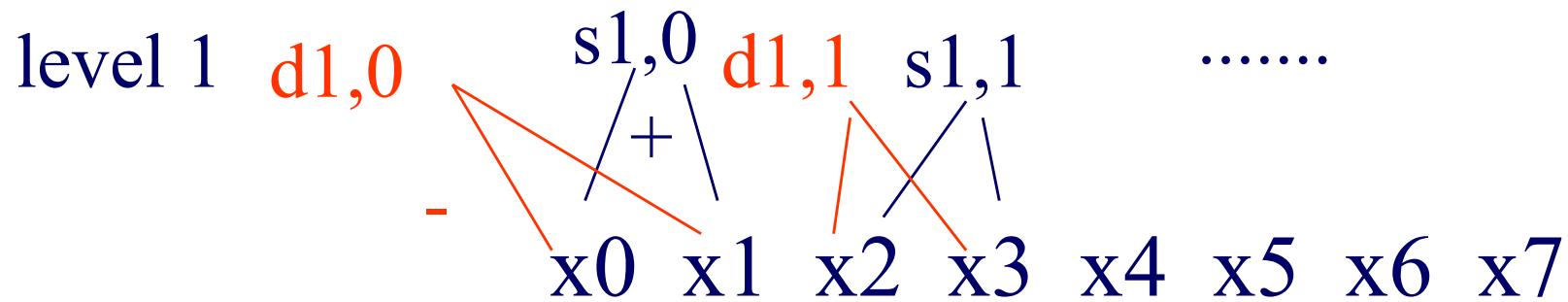
$x_0 \ x_1 \ x_2 \ x_3 \ x_4 \ x_5 \ x_6 \ x_7$

Wavelets - construction

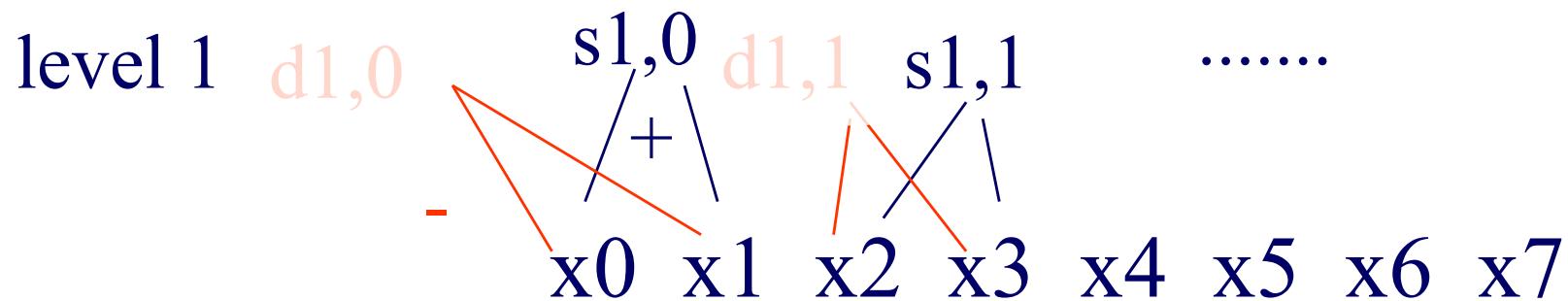


x0 x1 x2 x3 x4 x5 x6 x7

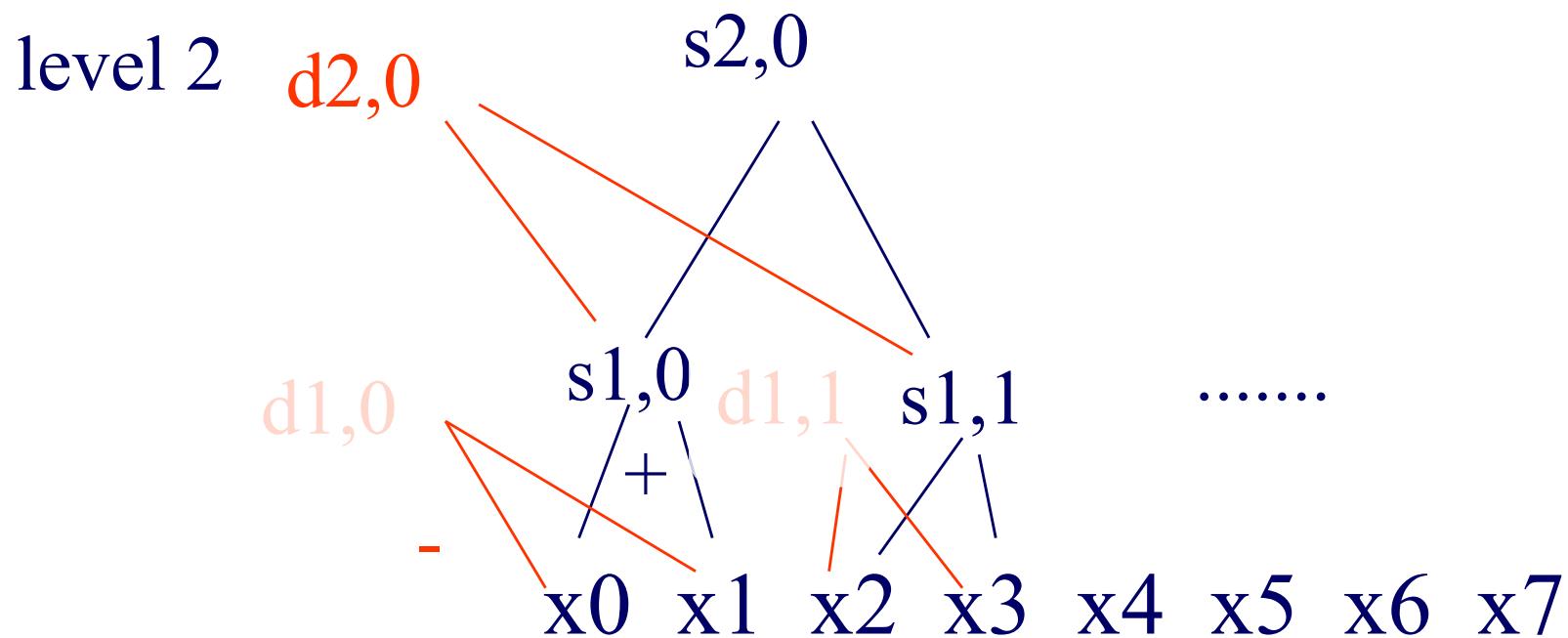
Wavelets - construction



Wavelets - construction



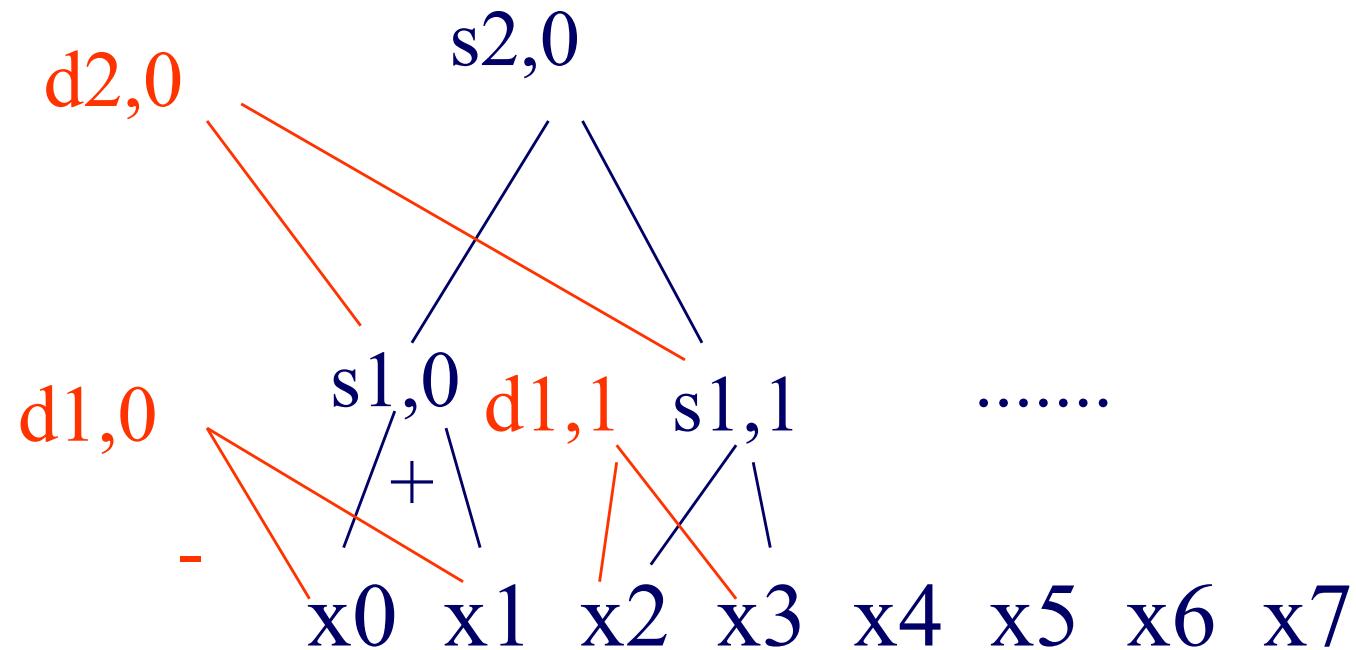
Wavelets - construction



Wavelets - construction



etc ...

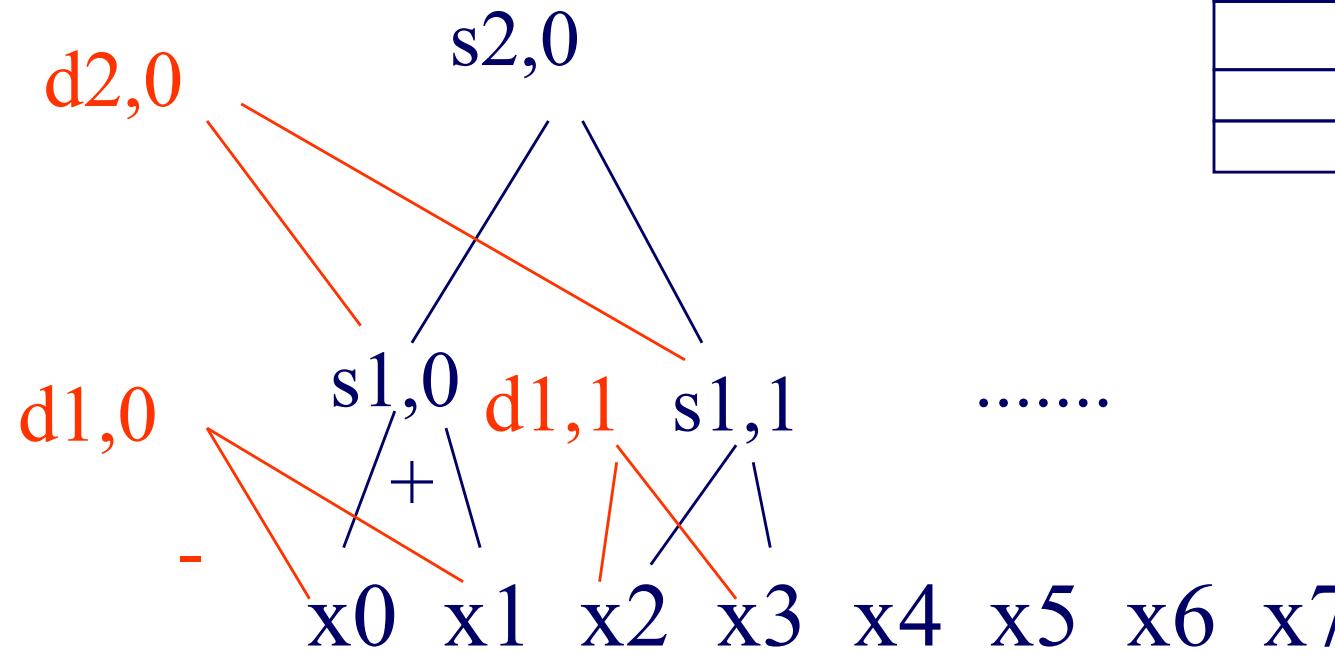


Wavelets - construction



Q: map each coefficient

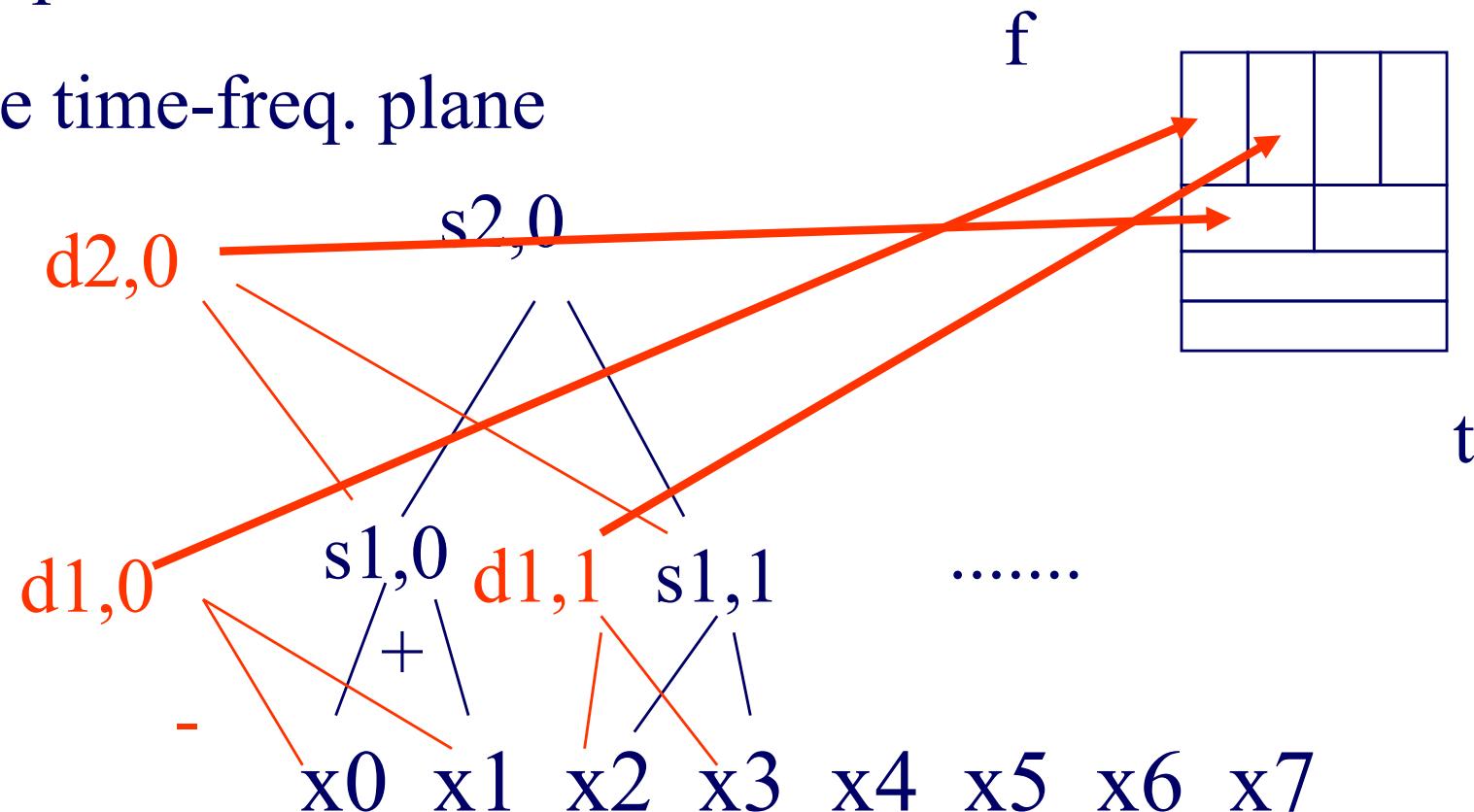
on the time-freq. plane



Wavelets - construction



Q: map each coefficient
on the time-freq. plane



Haar wavelets - code

```

#!/usr/bin/perl5
# expects a file with numbers
# and prints the dwt transform
# The number of time-ticks should be a power of 2
# USAGE
#  haar.pl <fname>

my @vals=();
my @smooth; # the smooth component of the signal
my @diff; # the high-freq. component

# collect the values into the array @val
while(<>){
    @vals = ( @vals , split );
}

```

```

my $len = scalar(@vals);
my $half = int($len/2);
while($half >= 1 ){
    for(my $i=0; $i< $half; $i++){
        $diff[$i] = ($vals[2*$i] - $vals[2*$i + 1]) / sqrt(2);
        print "\t", $diff[$i];
        $smooth[$i] = ($vals[2*$i] + $vals[2*$i + 1]) / sqrt(2);
    }
    print "\n";
    @vals = @smooth;
    $half = int($half/2);
}
print "\t", $vals[0], "\n" ; # the final, smooth component

```

Wavelets - construction

Observation1:

‘+’ can be some weighted addition

‘-’ is the corresponding weighted difference
('Quadrature mirror filters')

Observation2: unlike DFT/DCT,

there are *many* wavelet bases: Haar, Daubechies-4, Daubechies-6, Coifman, Morlet, Gabor, ...

Part 1 - Outline

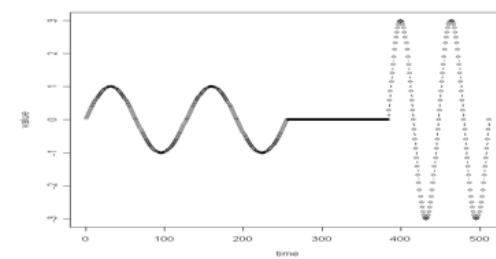
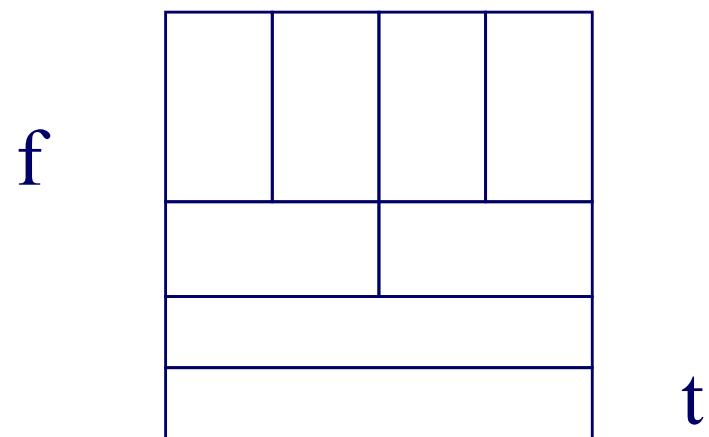


- Motivation
- P1.1. Similarity Search and Indexing
- P1.2. DSP
 - DFT
 - DWT
 - Definition of DWT and properties
 - how to read the DWT scalogram
- P1.3. Linear Forecasting
- ...



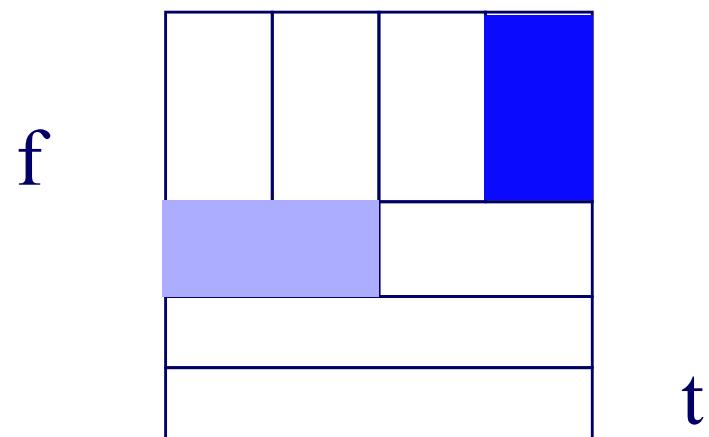
Wavelets - Drill#1:

- Q: baritone/silence/soprano - DWT?

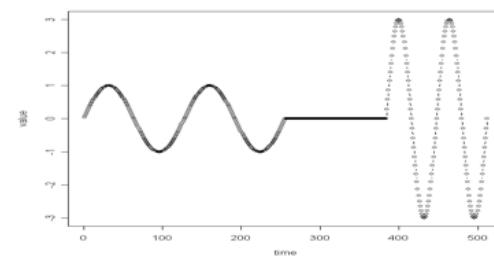


Wavelets - Drill#1:

- Q: baritone/silence/soprano - DWT?



value



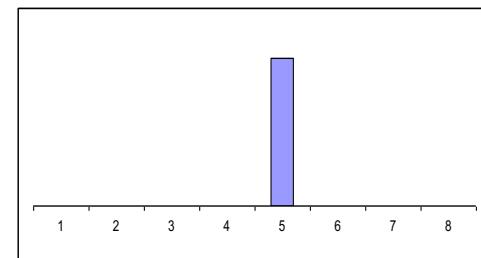
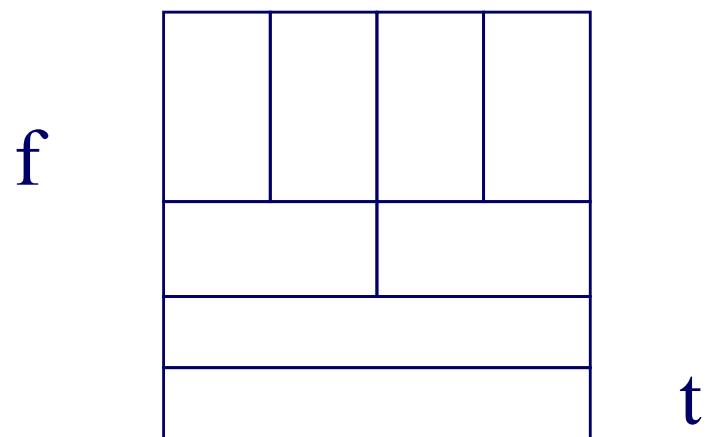
t

time



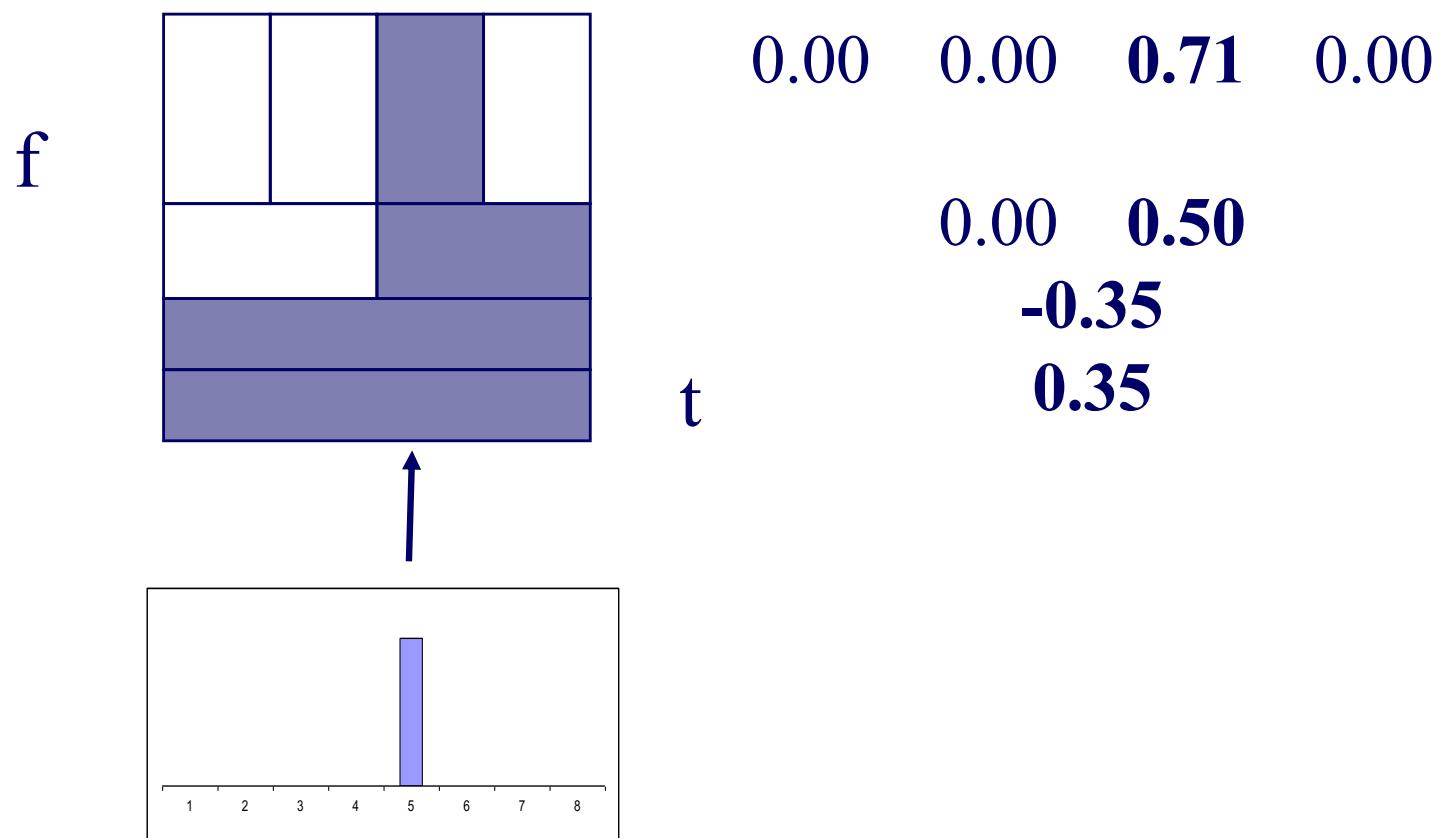
Wavelets - Drill#2:

- Q: spike - DWT?



Wavelets - Drill#2:

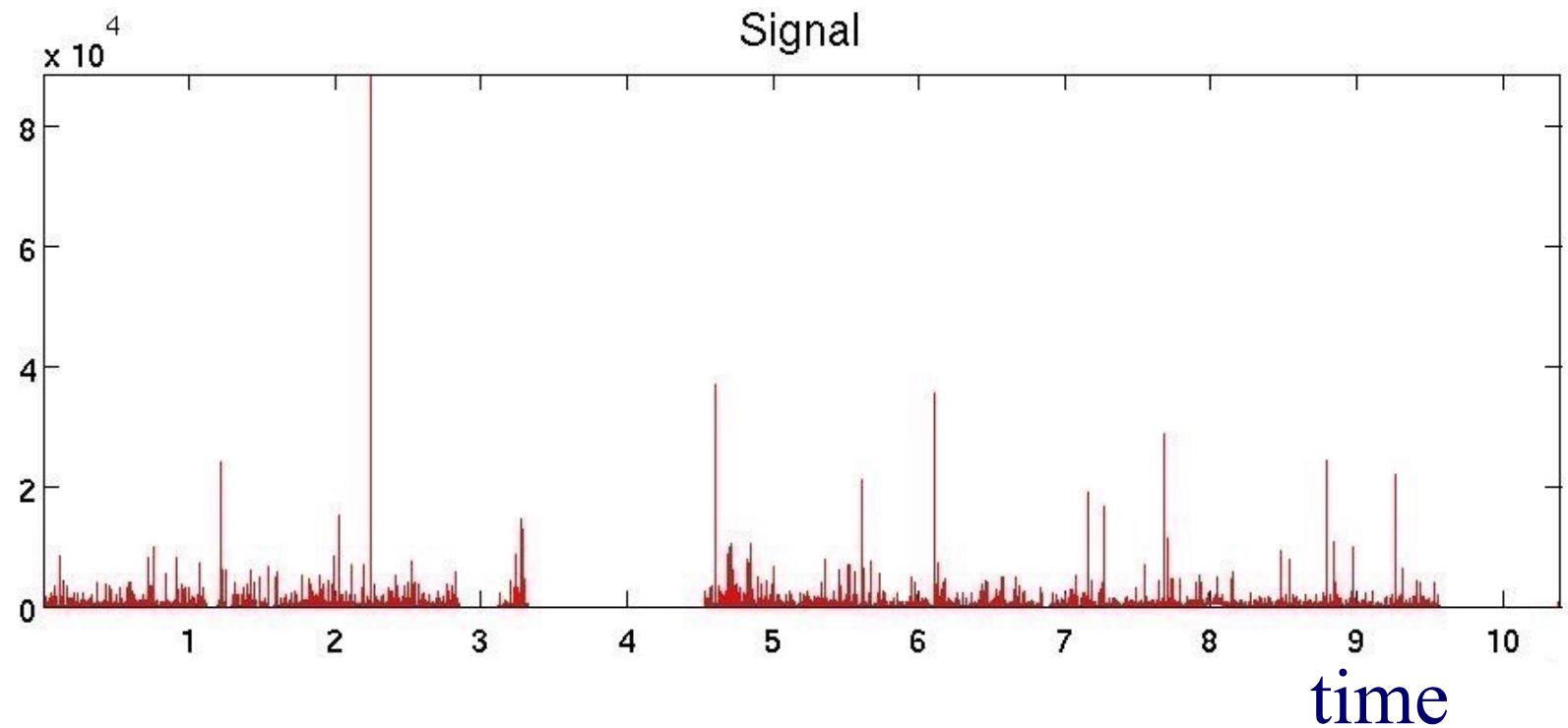
- Q: spike - DWT?



Wavelets in action

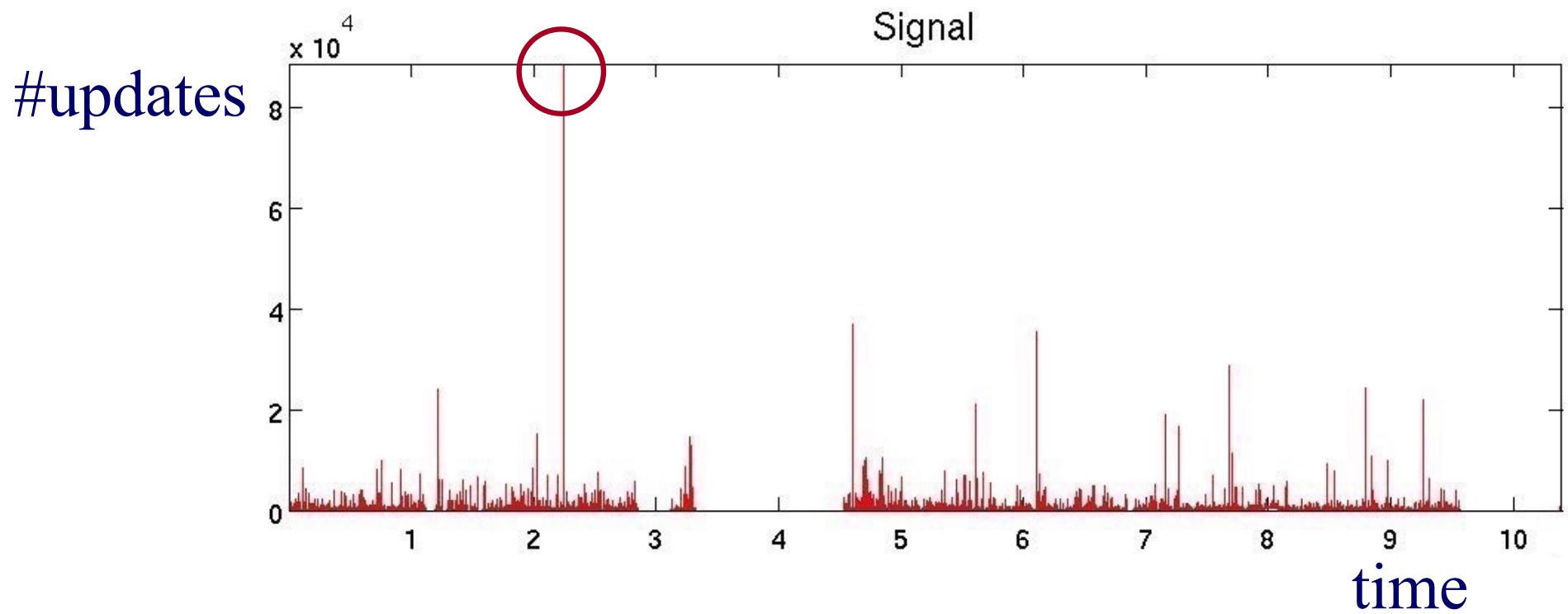
More examples (BGP updates)

#updates



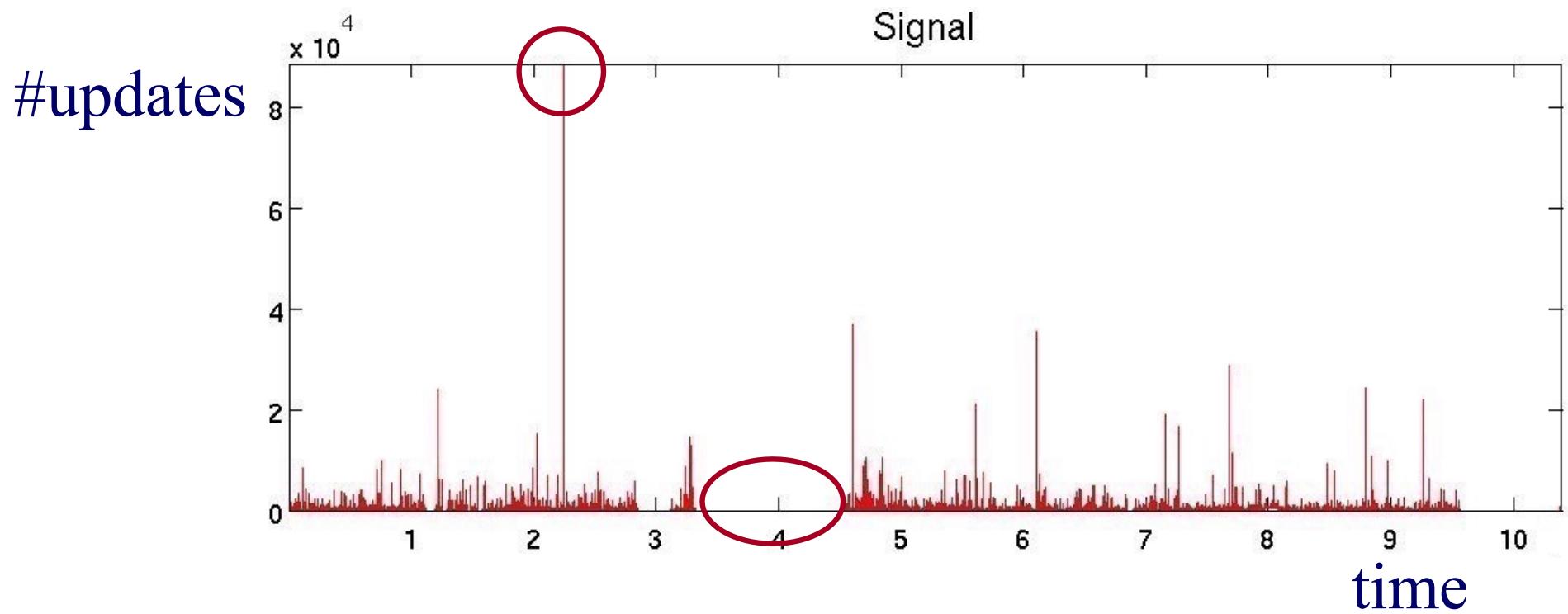
*BGP-lens: Patterns and Anomalies in Internet Routing
Updates* B. Aditya Prakash et al, SIGKDD 2009

More examples (BGP updates)



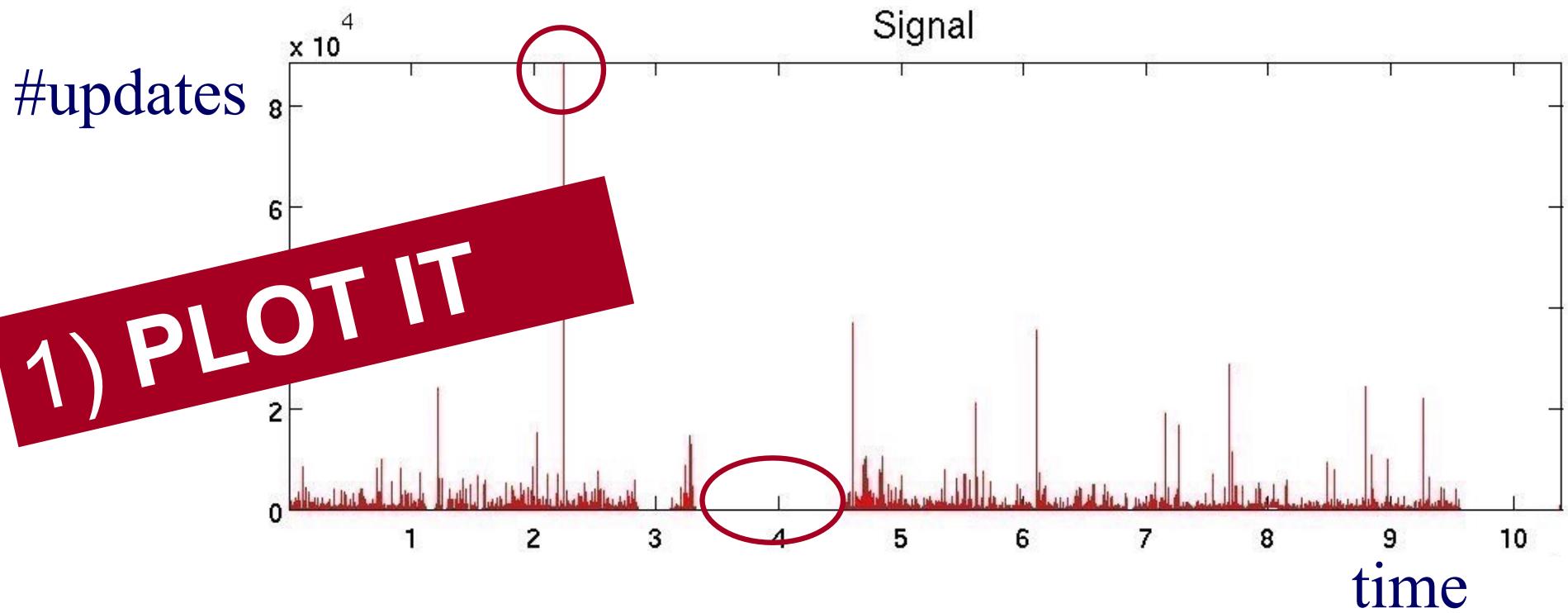
*BGP-lens: Patterns and Anomalies in Internet Routing
Updates* B. Aditya Prakash et al, SIGKDD 2009

More examples (BGP updates)



*BGP-lens: Patterns and Anomalies in Internet Routing
Updates* B. Aditya Prakash et al, SIGKDD 2009

More examples (BGP updates)



*BGP-lens: Patterns and Anomalies in Internet Routing
Updates* B. Aditya Prakash et al, SIGKDD 2009

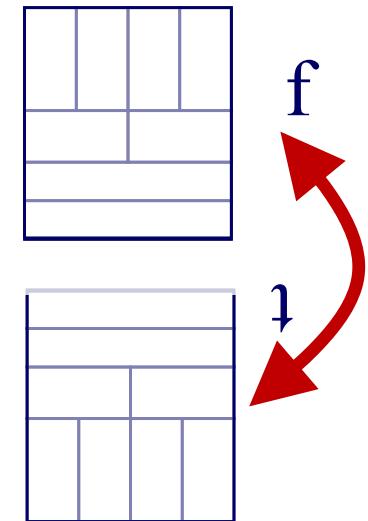
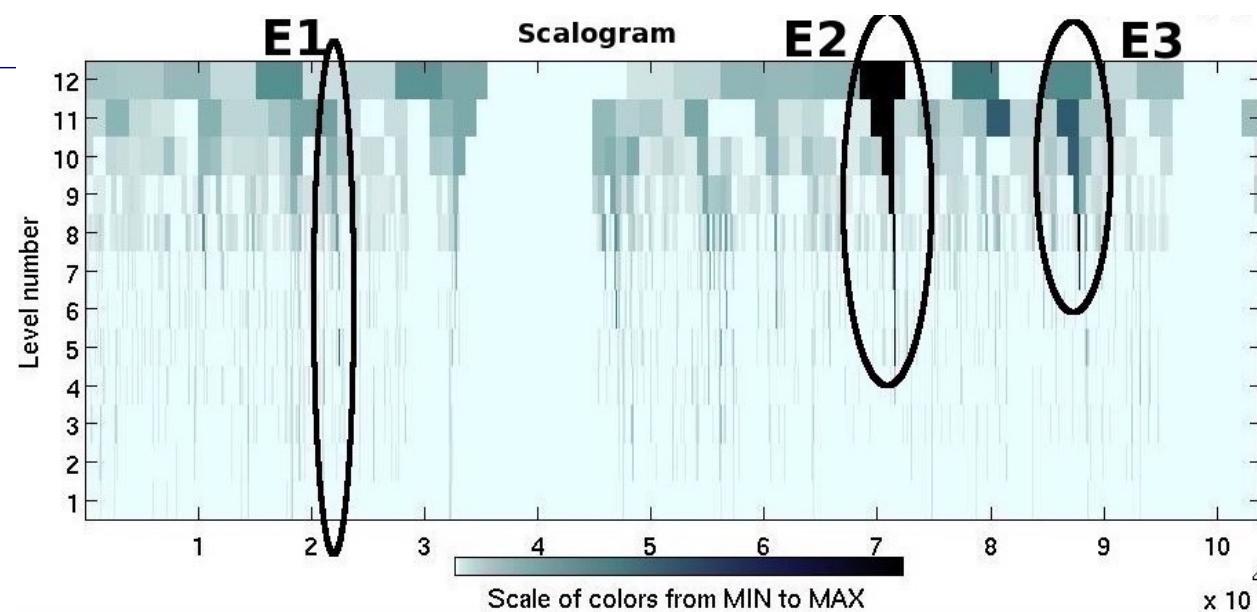
More examples (BGP updates)

- 
- 1) PLOT IT
 - 2) DFT / DWT

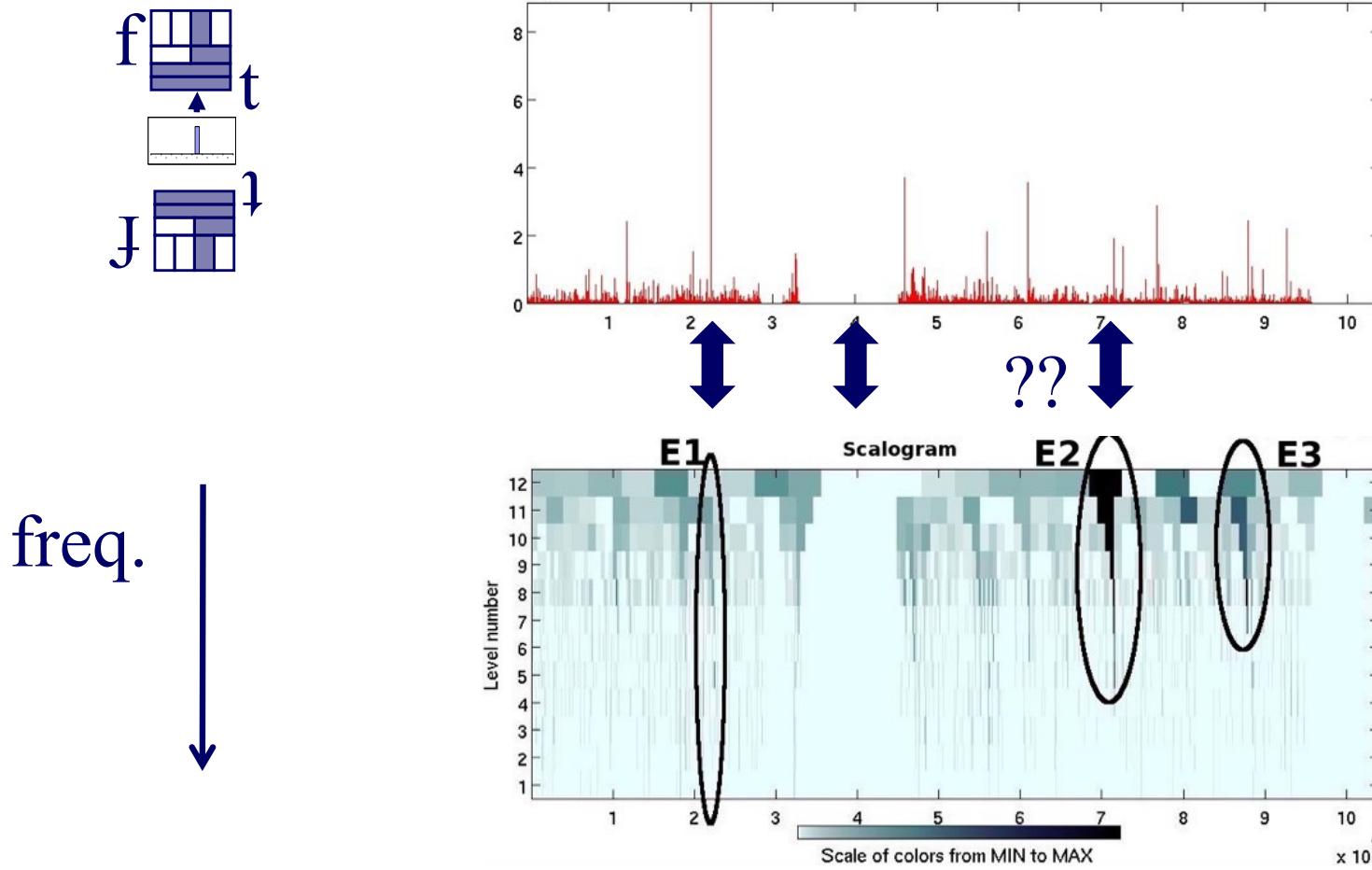
More examples (BGP updates)

Low freq.:
omitted

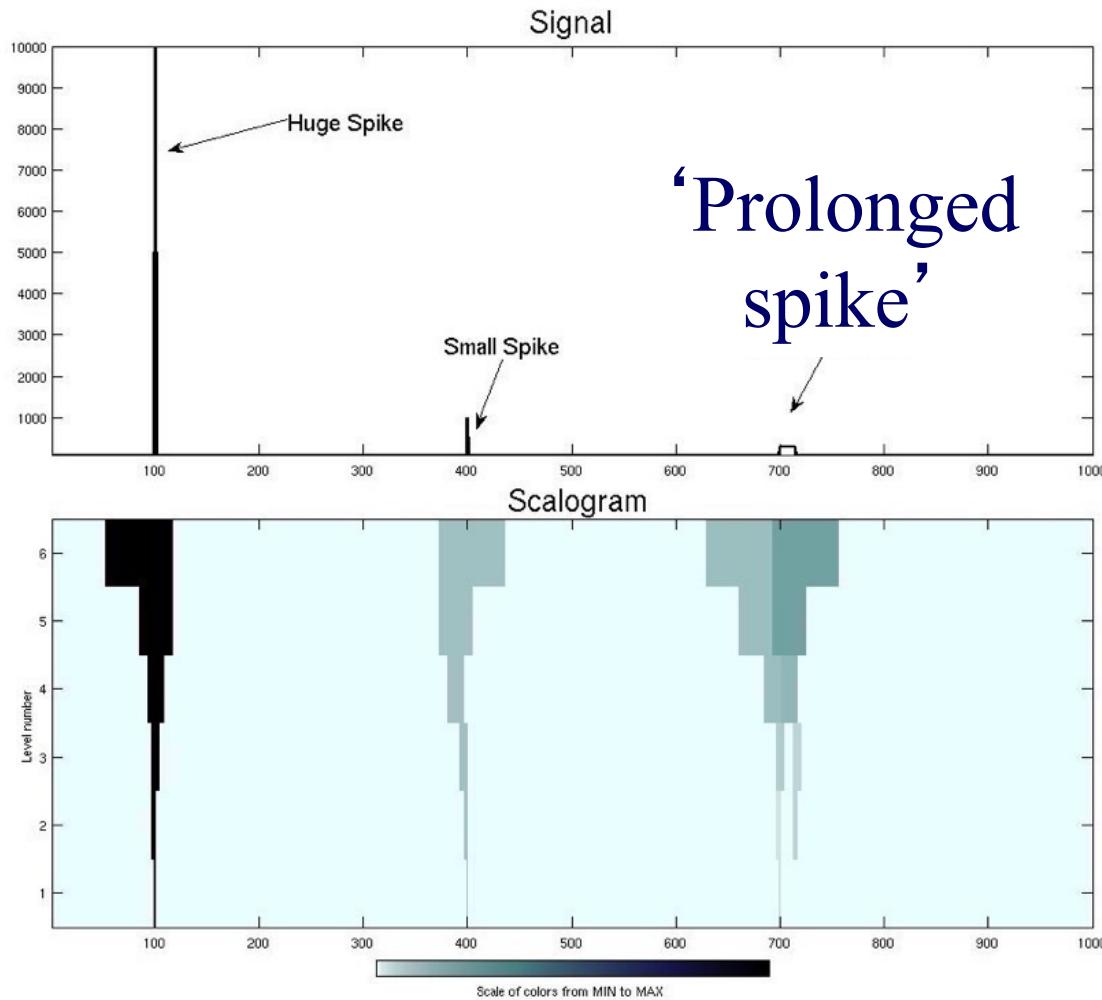
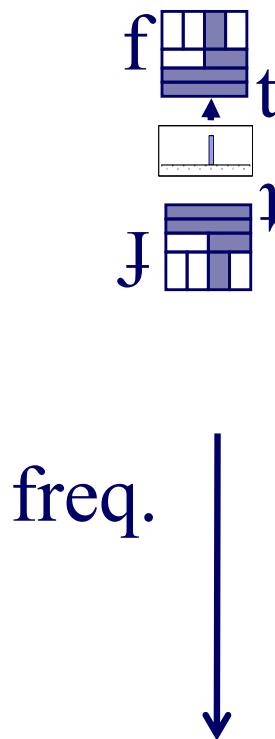
freq.
↓



More examples (BGP updates)



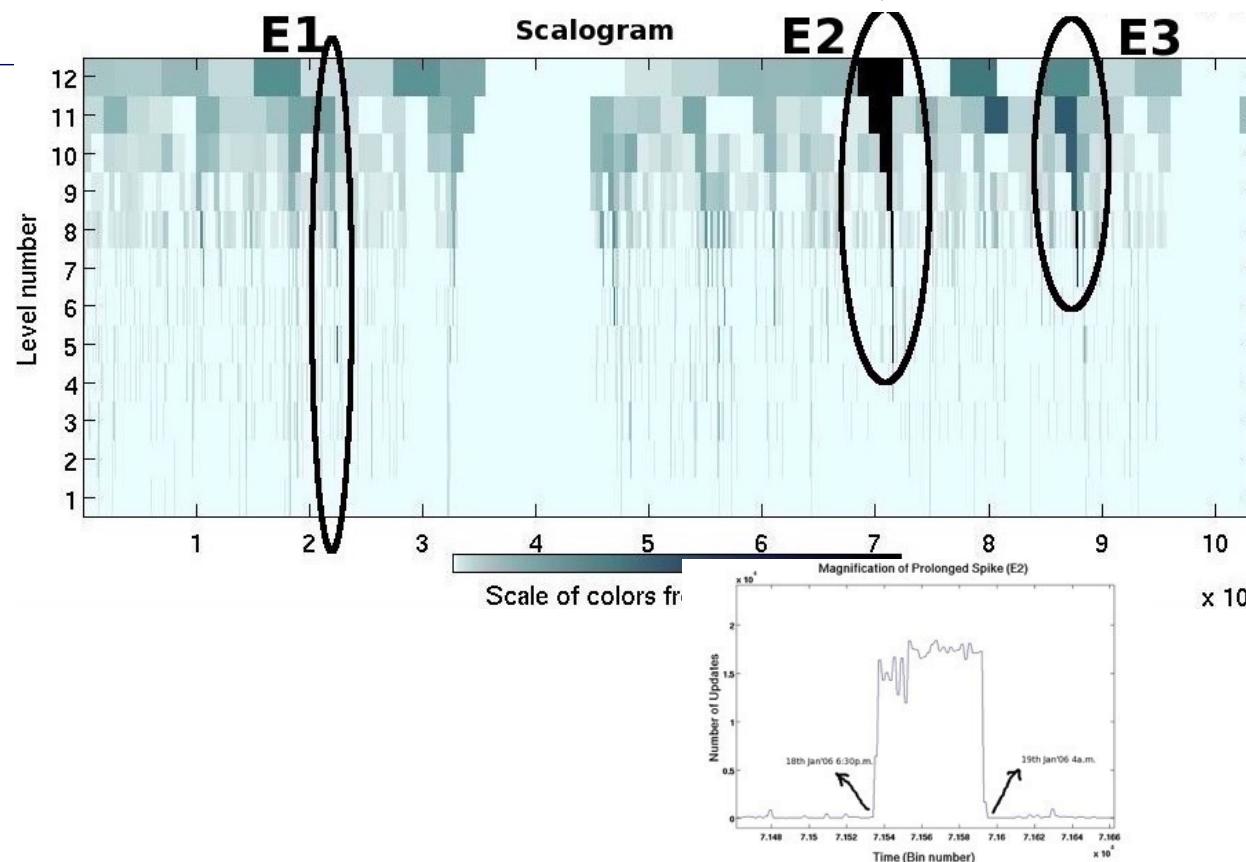
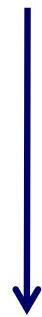
More examples (BGP updates)



More examples (BGP updates)

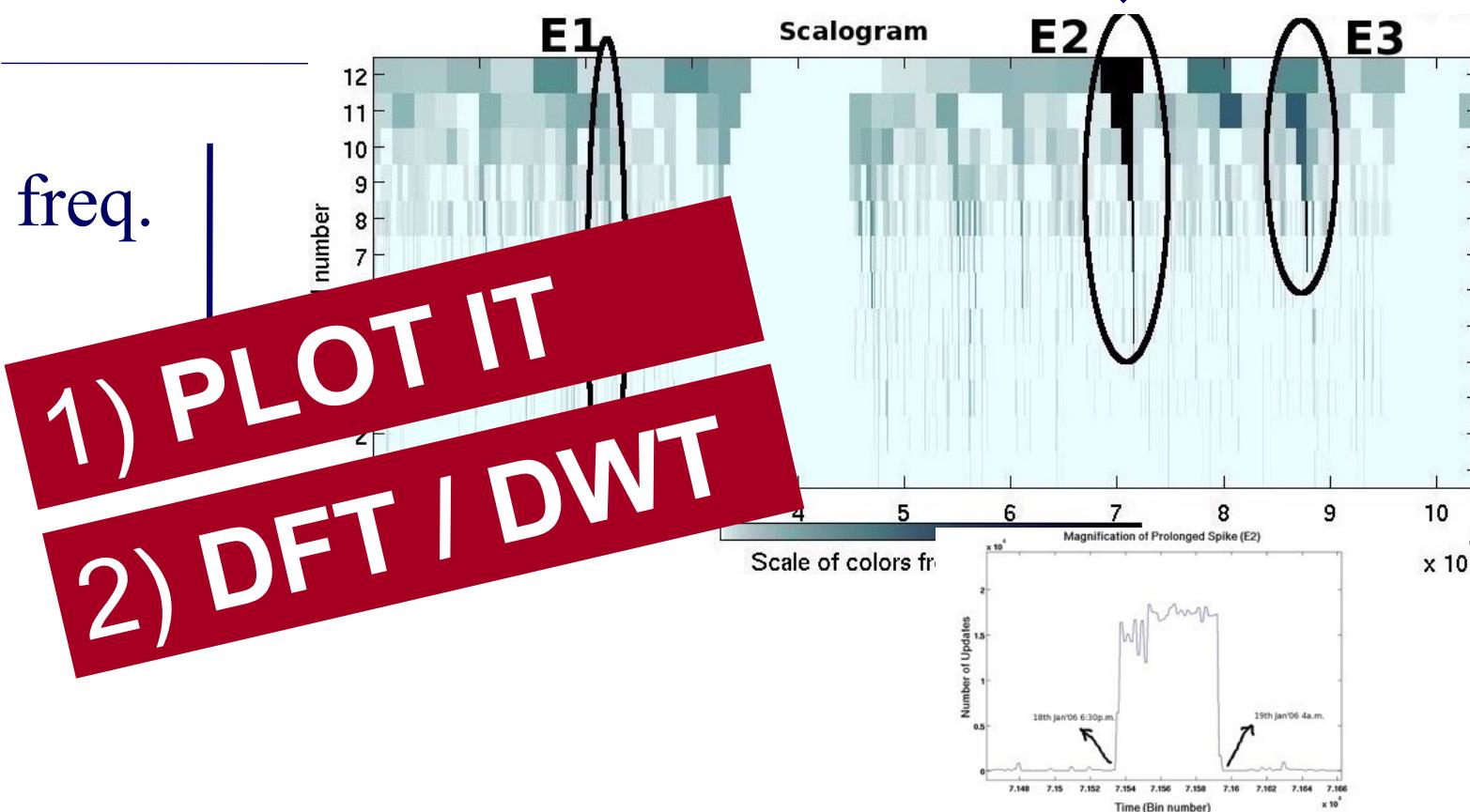
15K msgs, for several hours: 6pm-4am

freq.



More examples (BGP updates)

15K msgs, for several hours: 6pm-4am



Advantages of Wavelets

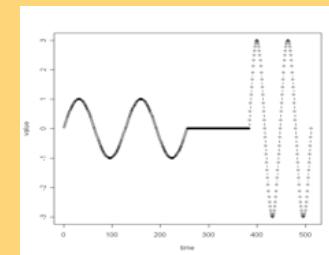
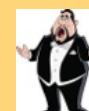
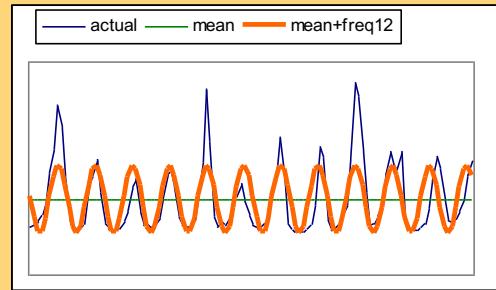
- Better compression (better RMSE with same number of coefficients - used in JPEG-2000)
- fast to compute (usually: $O(n)$)
- good for ‘spikes’
- good for de-noising
- mammalian eye and ear: Gabor wavelets



Part 1.2: Conclusions



- DFT spots periodicities
- DWT : multi-resolution - matches processing of mammalian ear/eye better
- Both: powerful tools for **compression, pattern detection** in real signals

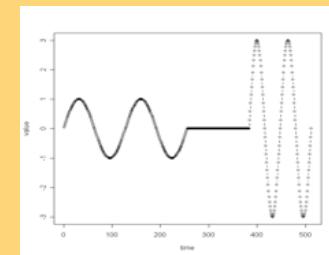
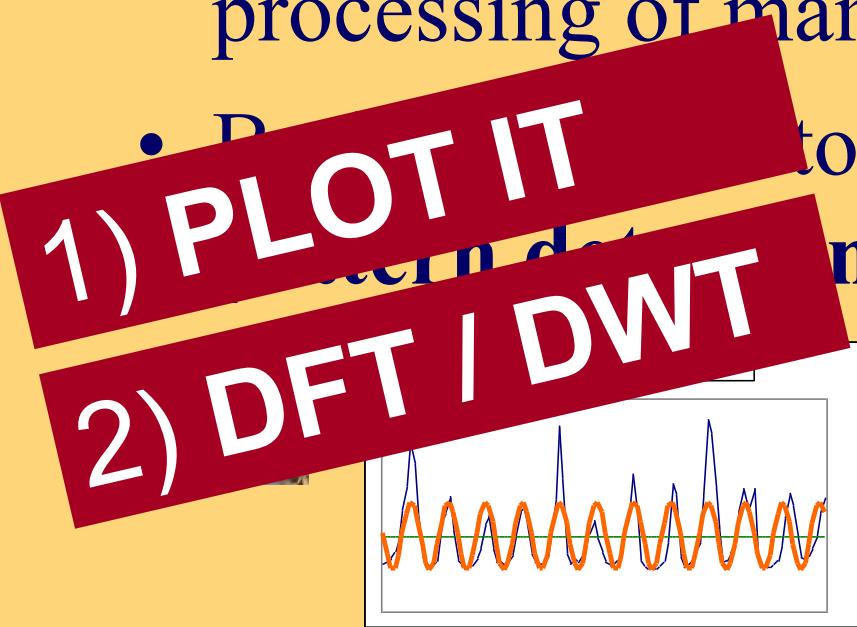


Part 1.2: Conclusions



- DFT spots periodicities
- DWT : multi-resolution - matches processing of mammalian ear/eye better

- DFT / DWT are tools for **compression**,
in real signals



Resources: software

- *xwpl*: open source wavelet package from Yale, with excellent GUI
- PyWavelets
- R, ‘wavelets’
- Mathematica, matlab, etc



Books

- ★ • William H. Press, Saul A. Teukolsky, William T. Vetterling and Brian P. Flannery: *Numerical Recipes in C*, Cambridge University Press, 1992, 2nd Edition. (Great description, intuition and code for DFT, DWT)
- C. Faloutsos: *Searching Multimedia Databases by Content*, Kluwer Academic Press, 1996 (introduction to DFT, DWT)
- Rabiner, L. and B.-H. Juang (1993). *Fundamentals of Speech Recognition*, Prentice Hall.

Part 1 - Outline



- Motivation
- P1.1. Similarity Search and Indexing
- P1.2. DSP
- ➡ • P1.3. Linear Forecasting
- P1.4. Non-linear forecasting
- P1.5. Tensors
- Conclusions

Part 1.3.

Linear Forecasting

Forecasting

"Prediction is very difficult, especially about the future." - Nils Bohr

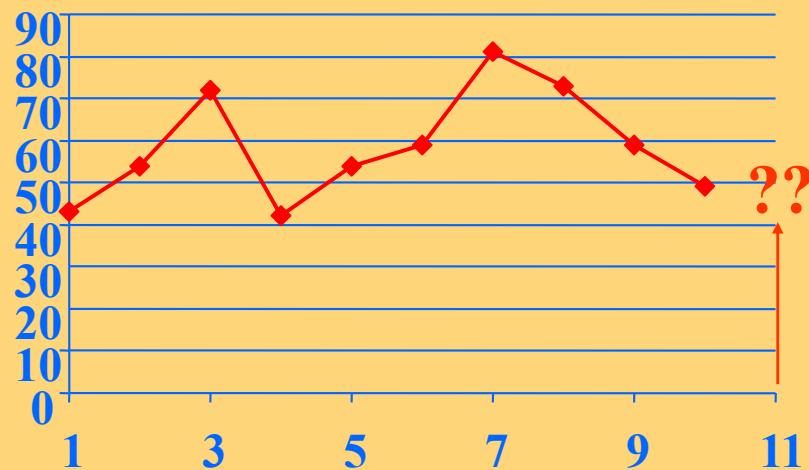
<http://www.hfac.uh.edu/MediaFutures/thoughts.html>



Problem#2: Forecast



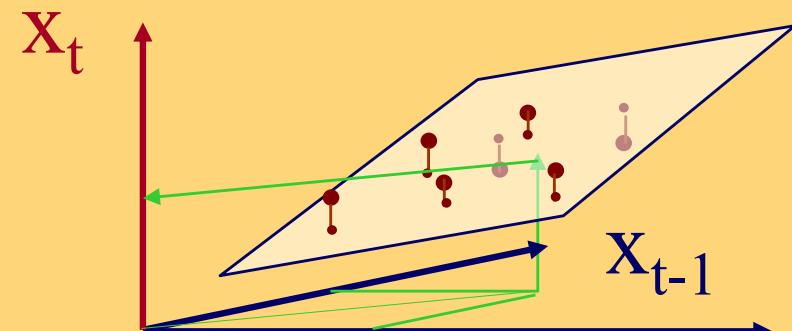
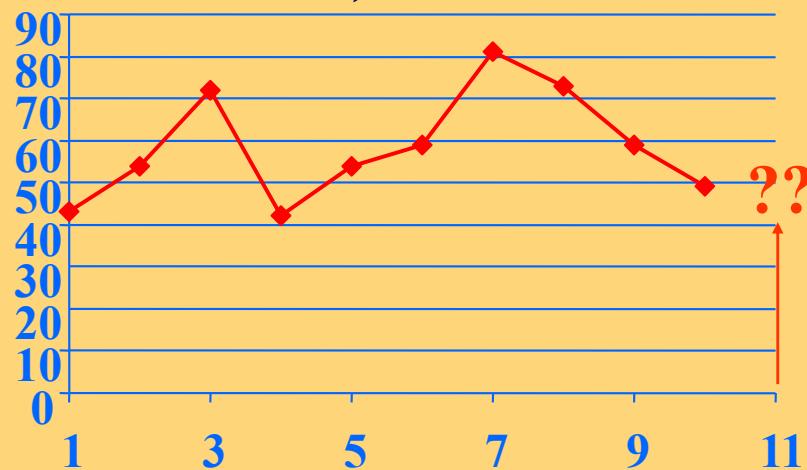
- given $x_{t-1}, x_{t-2}, \dots,$
- Q: forecast x_t



Solution: AR(IMA)



- given $x_{t-1}, x_{t-2}, \dots,$
- Q: forecast x_t
- A: AR(IMA) = Box-Jenkins (< Holt-Winters, Kalman)



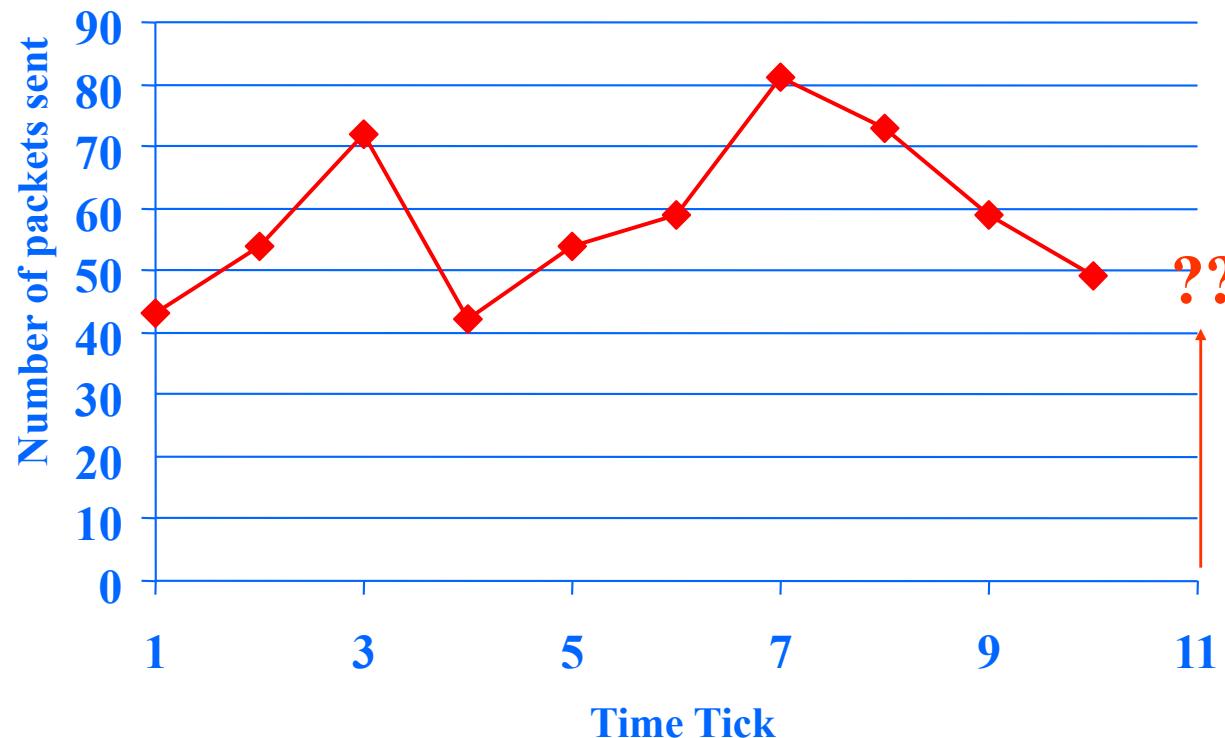
Part 1 - Outline



- Motivation
- ...
- P1.3. Linear Forecasting
 - – Auto-regression: Least Squares; RLS
 - Co-evolving time sequences
 - Examples
 - Conclusions

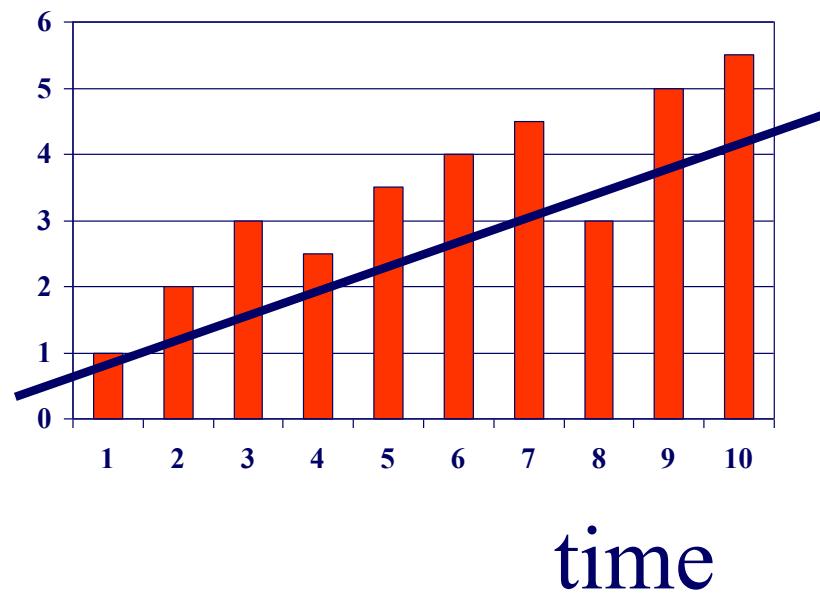
Problem#2: Forecast

- Example: give x_{t-1}, x_{t-2}, \dots , forecast x_t

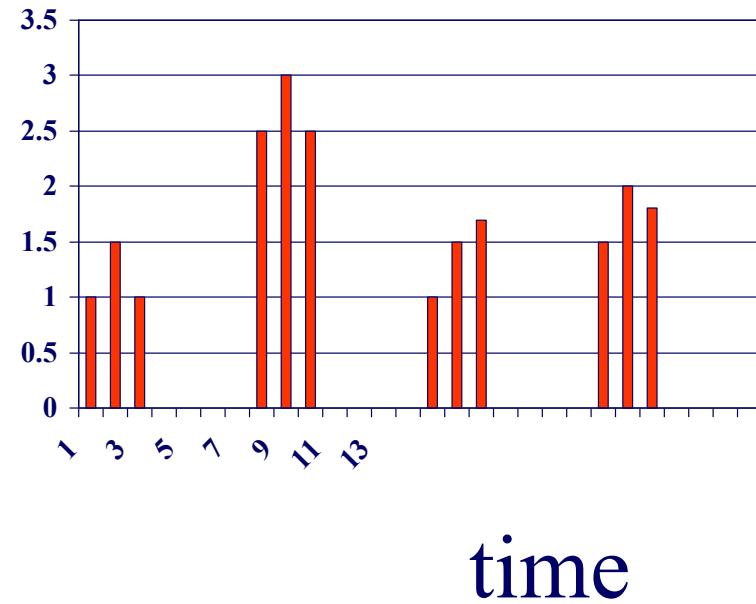


Forecasting: Preprocessing

MANUALLY:
remove trends



spot periodicities
7 days



Problem#2: Forecast

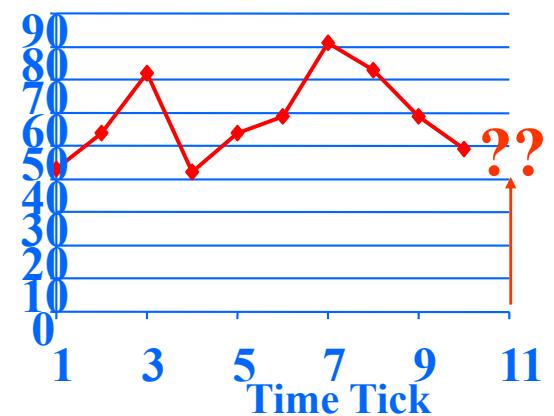
- Solution: try to express

x_t

as a linear function of the past: $x_{t-2}, x_{t-3}, \dots,$
(up to a window of w)

Formally:

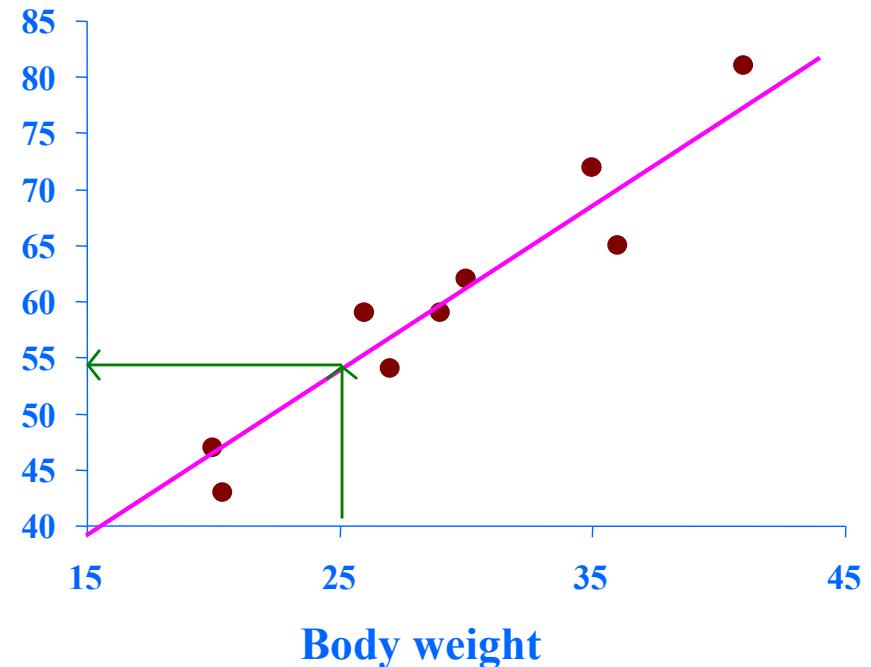
$$x_t \approx a_1 x_{t-1} + \dots + a_w x_{t-w} + \text{noise}$$



Linear Regression: idea

patient	weight	height
1	27	43
2	43	54
3	54	72
...
N	25	??

Body height



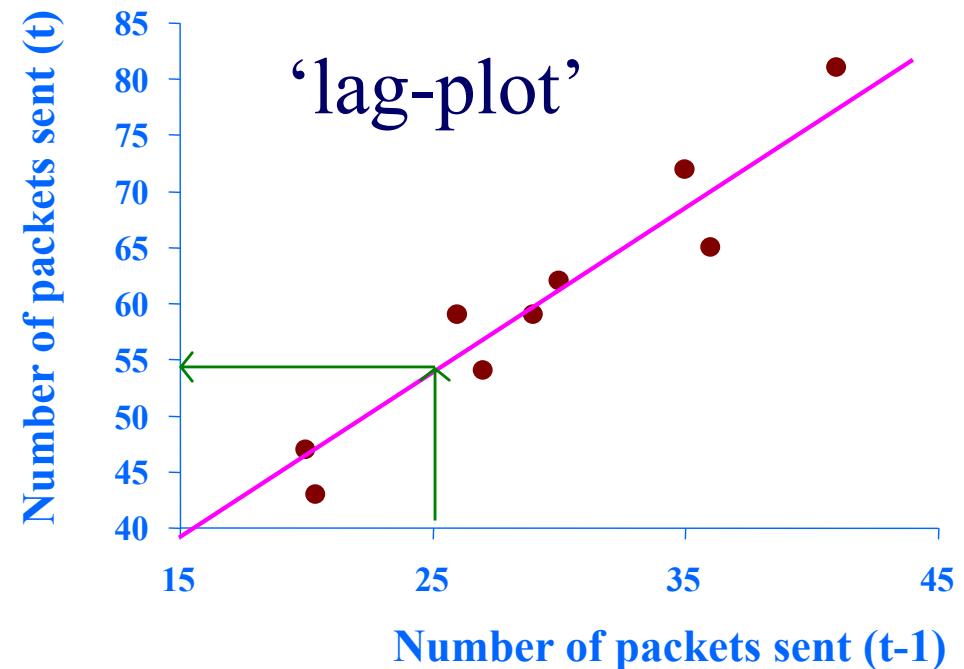
- express what we don't know (= 'dependent variable')
- as a linear function of what we know (= 'indep. variable(s)')

Linear Auto Regression:

Time	<i>Packets Sent(t)</i>
1	43
2	54
3	72
...	...
N	??

Linear Auto Regression:

Time	<i>Packets Sent (t-1)</i>	<i>Packets Sent(t)</i>
1	-	43
2	43	54
3	54	72
...
N	25	??



- lag $w=1$
- Dependent variable = # of packets sent ($S[t]$)
- Independent variable = # of packets sent ($S[t-1]$)

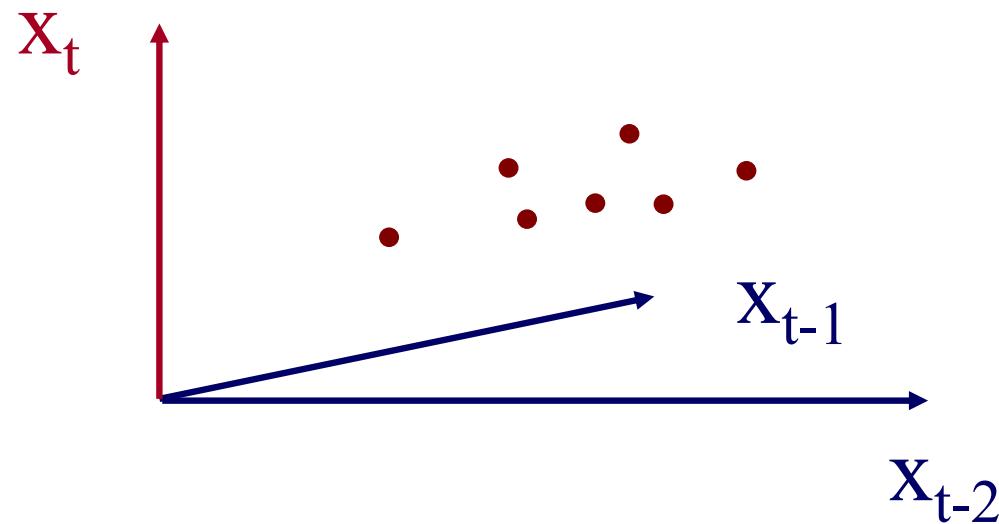
Part 1 - Outline



- Motivation
- ...
- P1.3. Linear Forecasting
 - – Auto-regression: Least Squares; RLS
 - Co-evolving time sequences
 - Conclusions

More details:

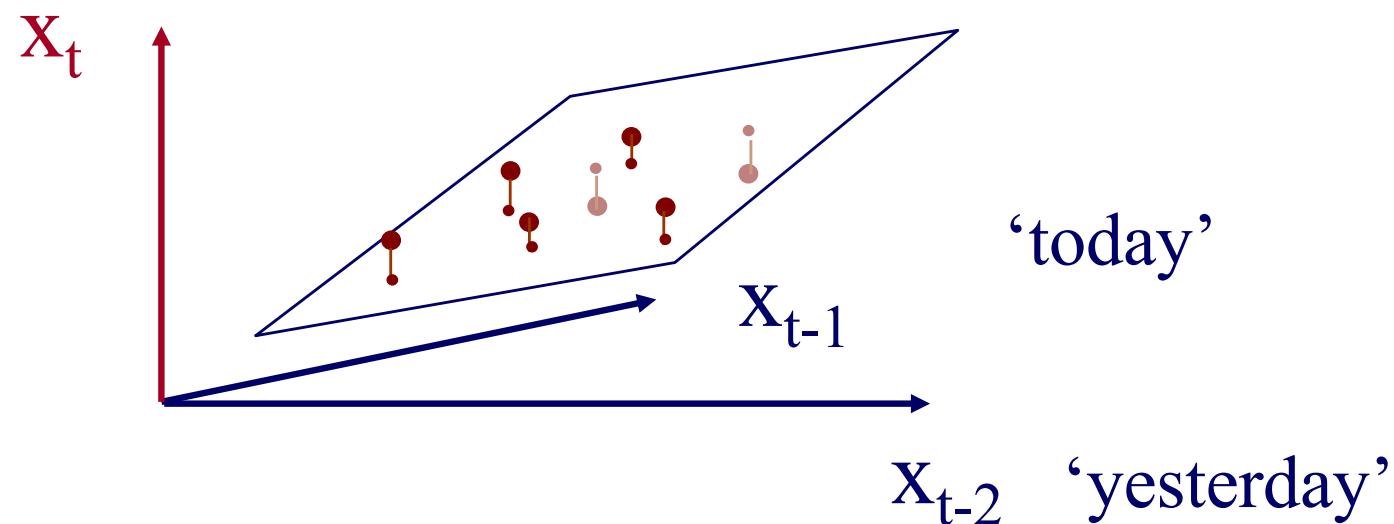
- Q1: Can it work with window $w>1$?
- A1: YES!



More details:

- Q1: Can it work with window $w>1$?
- A1: YES! (we'll fit a hyper-plane, then!)

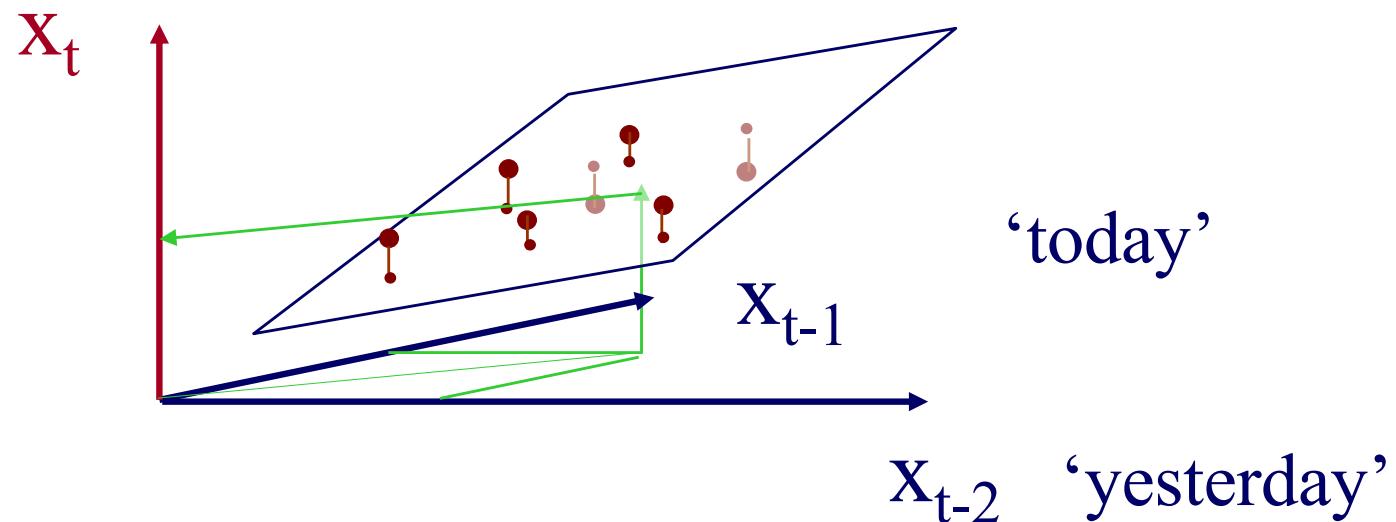
‘tomorrow’



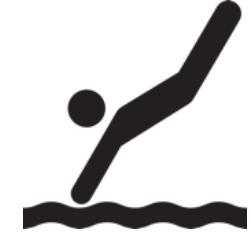
More details:

- Q1: Can it work with window $w>1$?
- A1: YES! (we'll fit a hyper-plane, then!)

‘tomorrow’



More details:



- Q1: Can it work with window $w > 1$?
- A1: YES! The problem becomes:

$$\mathbf{X}_{[N \times w]} \times \mathbf{a}_{[w \times 1]} = \mathbf{y}_{[N \times 1]}$$

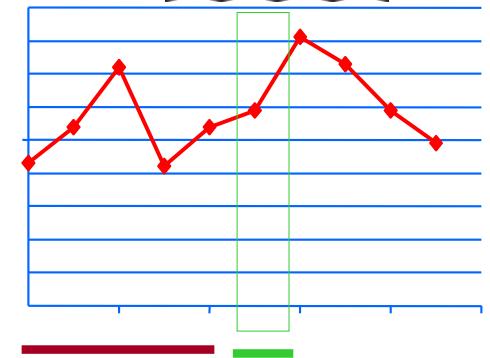
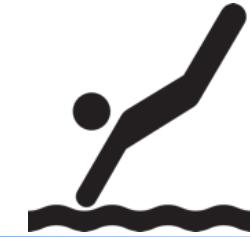
- OVER-CONSTRAINED
 - \mathbf{a} is the vector of the regression coefficients
 - \mathbf{X} has the N values of the w indep. variables
 - \mathbf{y} has the N values of the dependent variable

More details:

- $\mathbf{X}_{[N \times w]} \times \mathbf{a}_{[w \times 1]} = \mathbf{y}_{[N \times 1]}$

Ind-var1 Ind-var-w

$$\begin{array}{c}
 \text{time} \\
 \downarrow \\
 \left[\begin{array}{c}
 \overset{\downarrow}{X_{11}, X_{12}, \dots, X_{1w}} \\
 \hline
 \overset{\swarrow}{X_{21}, X_{22}, \dots, X_{2w}} \\
 \vdots \\
 \vdots \\
 \vdots \\
 \vdots \\
 X_{N1}, X_{N2}, \dots, X_{Nw}
 \end{array} \right] \times \left[\begin{array}{c}
 a_1 \\
 a_2 \\
 \vdots \\
 a_w
 \end{array} \right] = \left[\begin{array}{c}
 y_1 \\
 \hline
 y_2 \\
 \vdots \\
 \vdots \\
 \vdots \\
 y_N
 \end{array} \right]
 \end{array}$$



More details:

- $\mathbf{X}_{[N \times w]} \times \mathbf{a}_{[w \times 1]} = \mathbf{y}_{[N \times 1]}$

Ind-var1 Ind-var-w

$$\begin{array}{c}
 \text{time} \\
 \downarrow \\
 \left[\begin{array}{c} X_{11}, X_{12}, \dots, X_{1w} \\ X_{21}, X_{22}, \dots, X_{2w} \\ \hline \vdots \\ \vdots \\ \vdots \\ X_{N1}, X_{N2}, \dots, X_{Nw} \end{array} \right] \times \left[\begin{array}{c} a_1 \\ a_2 \\ \vdots \\ a_w \end{array} \right] = \left[\begin{array}{c} y_1 \\ y_2 \\ \vdots \\ \vdots \\ y_N \end{array} \right]
 \end{array}$$





More details

- Q2: How to estimate $a_1, a_2, \dots, a_w = \mathbf{a}$?
- A2: with Least Squares fit
 - $\mathbf{a} = (\mathbf{X}^T \times \mathbf{X})^{-1} \times (\mathbf{X}^T \times \mathbf{y})$
 - (Moore-Penrose pseudo-inverse)
 - \mathbf{a} is the vector that minimizes the RMSE from \mathbf{y}

Even more details

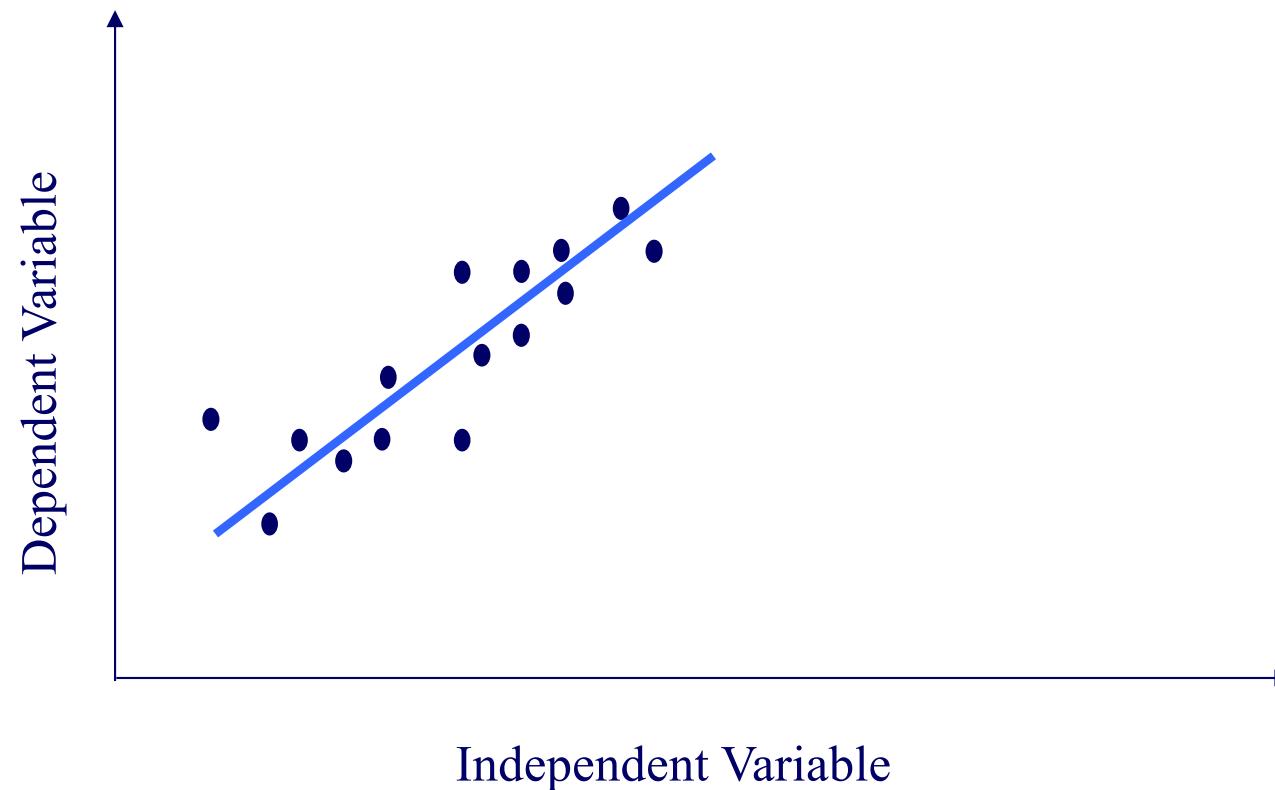


- Q3: Can we estimate \mathbf{a} incrementally?
- A3: Yes, with the brilliant, classic method of ‘Recursive Least Squares’ (RLS) (see, e.g., [Yi+00], for details) - pictorially:

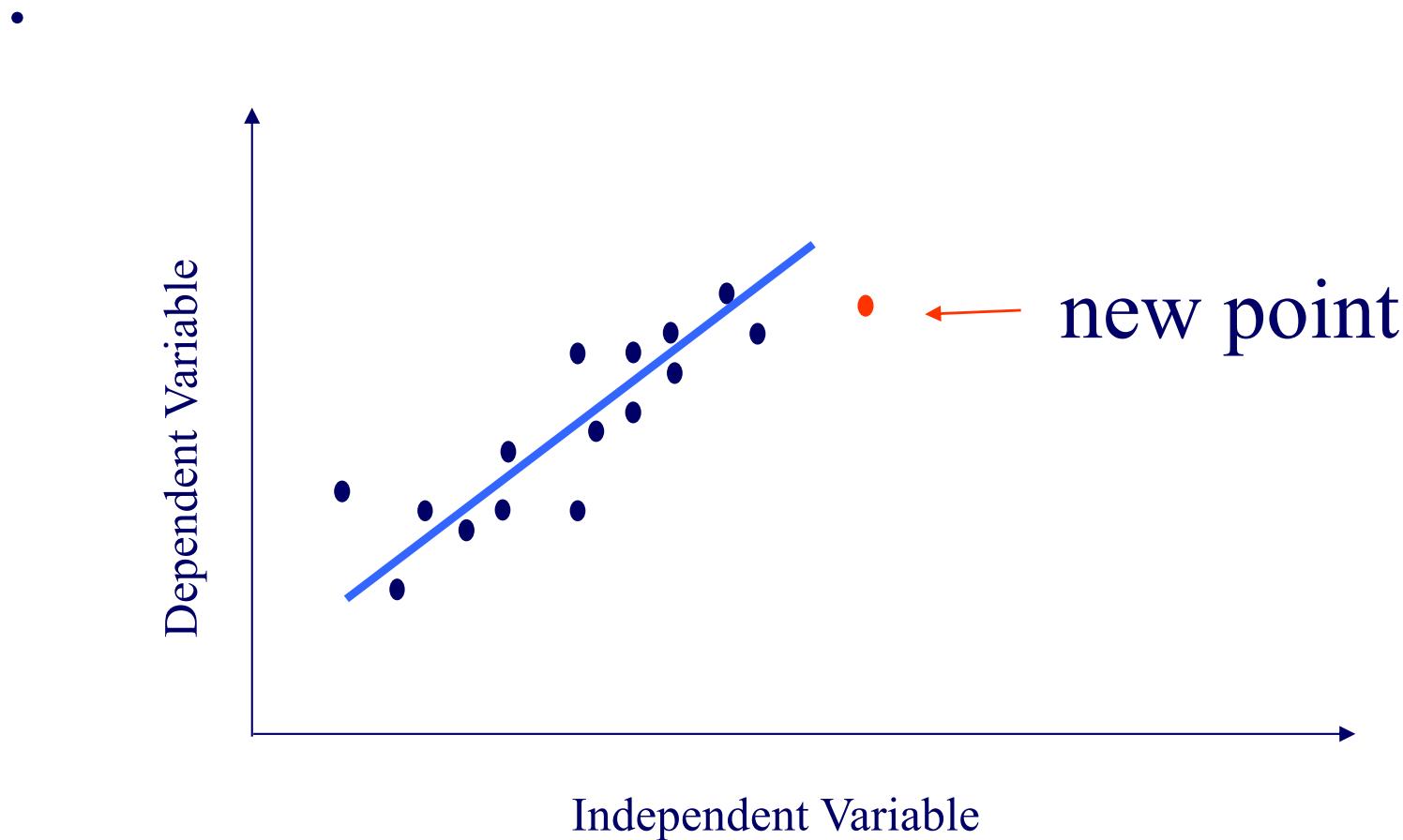
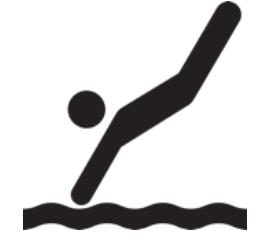
Even more details



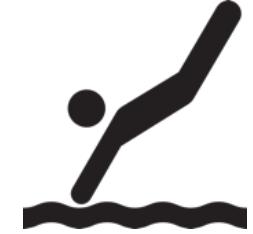
- Given:



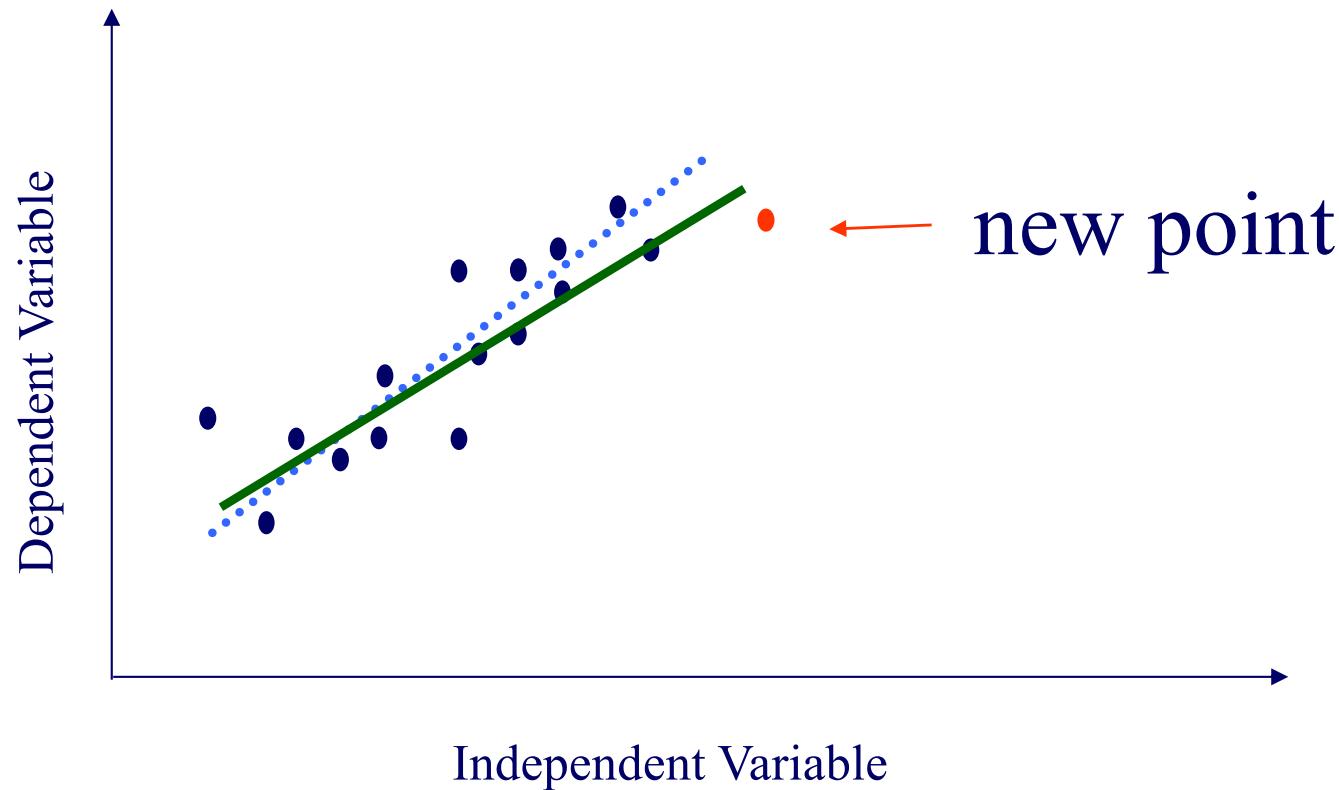
Even more details



Even more details



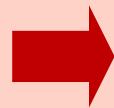
RLS: quickly compute new best fit



Part 1 - Outline

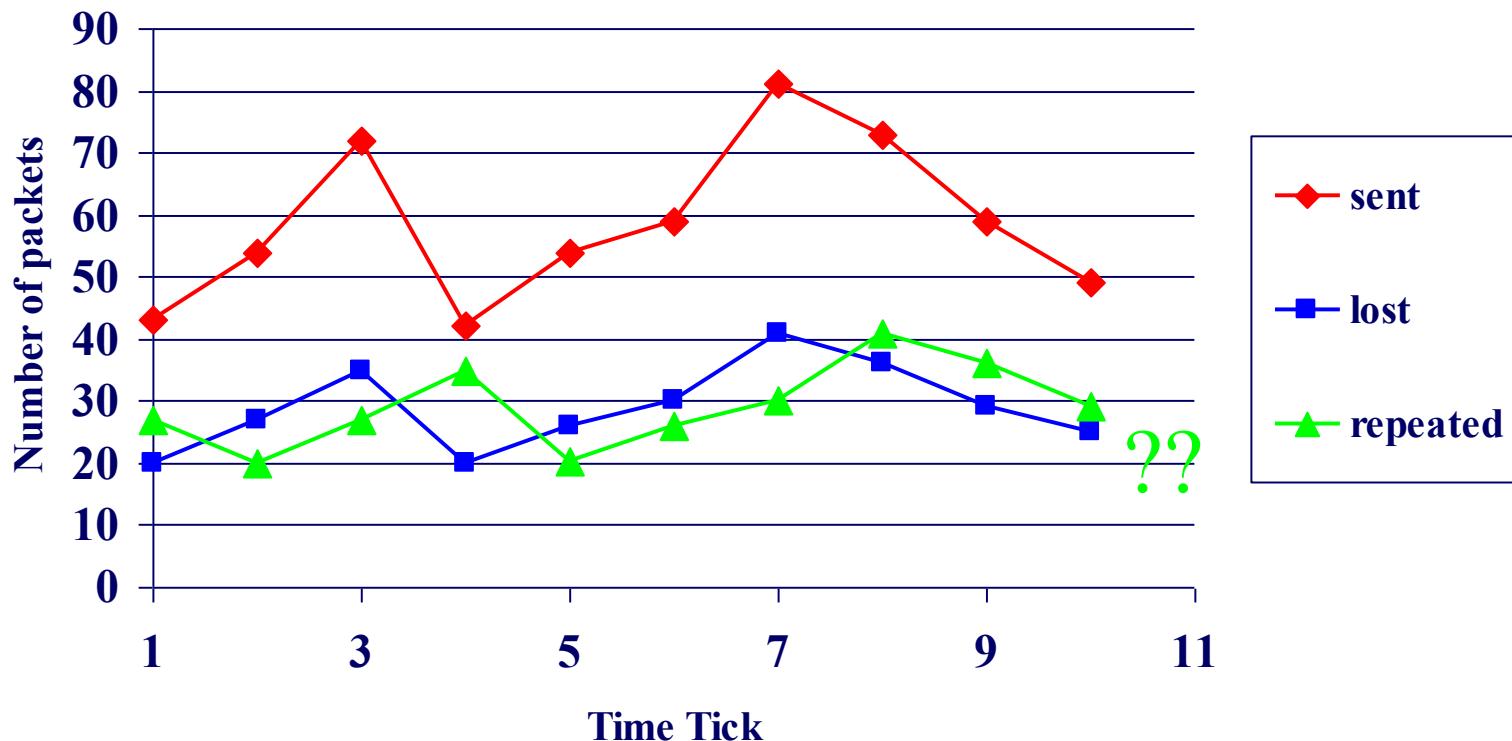


- Motivation
- ...
- P1.3. Linear Forecasting
 - Auto-regression: Least Squares; **RLS**
 - Co-evolving time sequences
 - Conclusions



Co-Evolving Time Sequences

- Given: A set of **correlated** time sequences
- Forecast ‘**Repeated(t)**’



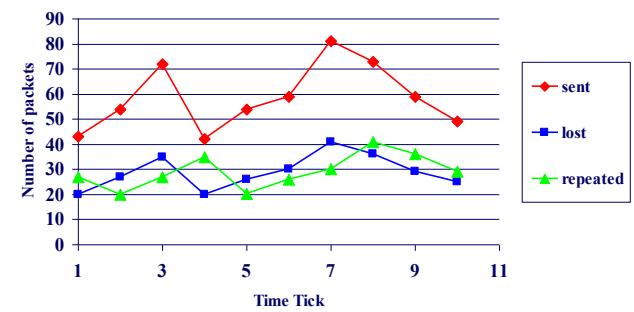
Solution:

Q: what should we do?

Solution: Vector ARIMA

Least Squares, with

- Dep. Variable: Repeated(t)
- Indep. Variables: Sent(t-1) ... Sent(t-w); Lost(t-1) ... Lost(t-w); Repeated(t-1), ...
- (named: ‘MUSCLES’ [Yi+00])



Part 1 - Outline



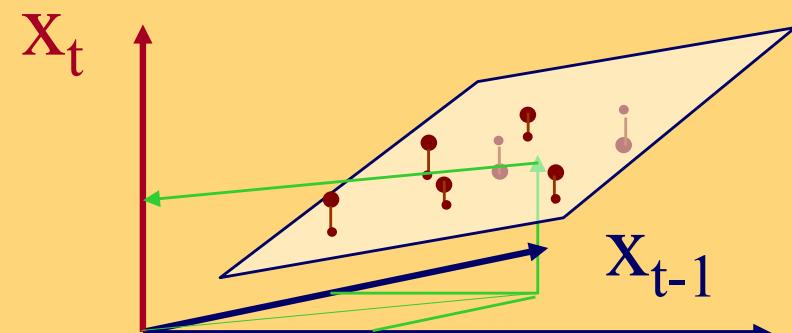
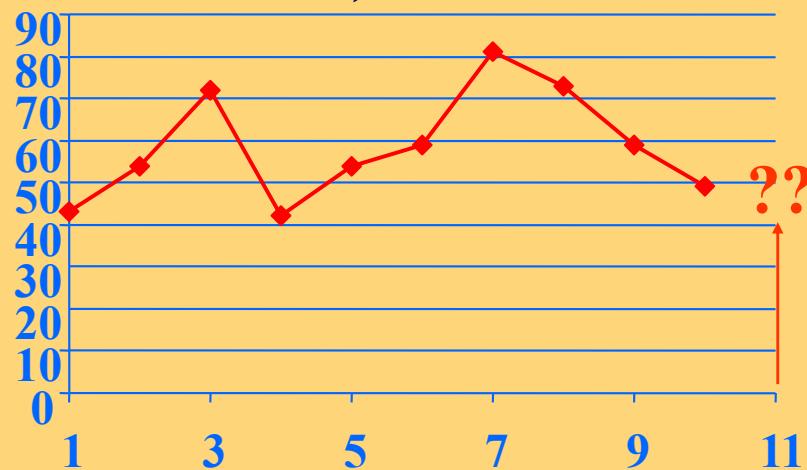
- Motivation
- ...
- P1.3. Linear Forecasting
 - Auto-regression: Least Squares; **RLS**
 - Co-evolving time sequences
 - Conclusions



Solution: AR(IMA)



- given $x_{t-1}, x_{t-2}, \dots,$
- Q: forecast x_t
- A: AR(IMA) = Box-Jenkins (< Holt-Winters, Kalman)



Books

- ★ • George E.P. Box and Gwilym M. Jenkins and Gregory C. Reinsel, *Time Series Analysis: Forecasting and Control*, Prentice Hall, 1994 (the classic book on ARIMA, 3rd ed.)
- Brockwell, P. J. and R. A. Davis (1987). Time Series: Theory and Methods. New York, Springer Verlag.

Additional Reading

- [Papadimitriou+ vldb2003] Spiros Papadimitriou, Anthony Brockwell and Christos Faloutsos *Adaptive, Hands-Off Stream Mining* VLDB 2003, Berlin, Germany, Sept. 2003
- [Yi+00] Byoung-Kee Yi et al.: *Online Data Mining for Co-Evolving Time Sequences*, ICDE 2000. (Describes MUSCLES and Recursive Least Squares)
- Kalman filters: KDD'10 tutorial
<https://lileicc.github.io/TALKS/2010-KDD/kdd2010-tut-partB.pdf>

Part 14: chaos and non-linear forecasting

Part 1 - Outline

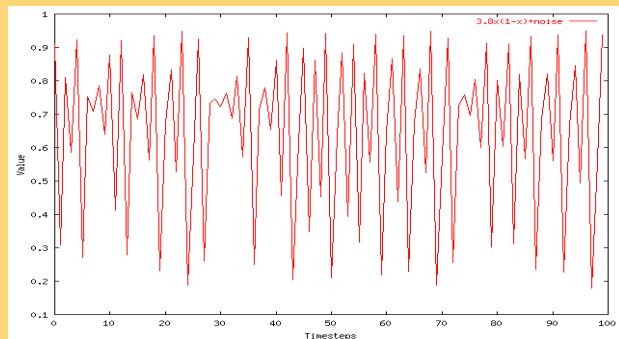


- ...
- P1.3. Linear forecasting
- P1.4. Non-linear forecasting
 - – Problem
 - Idea
 - How-to
 - Experiments
 - (Lotka-Volterra equations)
 - Conclusions
- P1.5. Tensors



Problem: Forecast

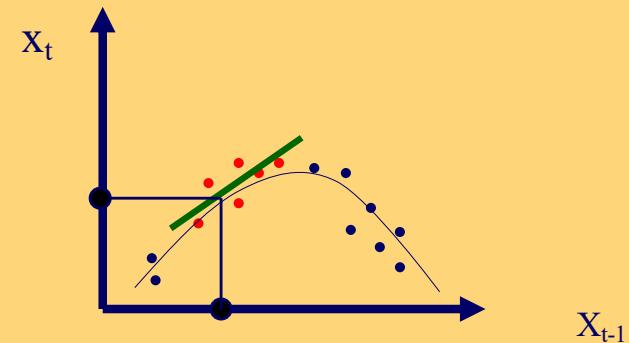
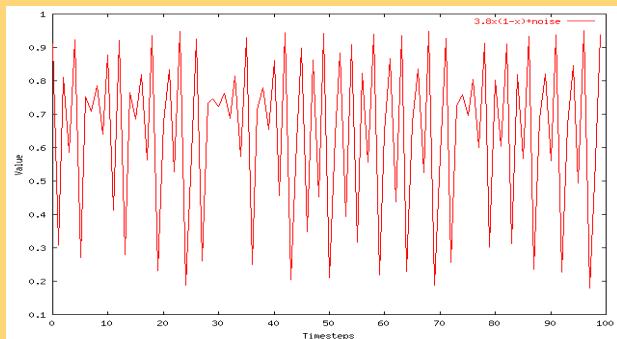
- given x_{t-1}, x_{t-2}, \dots , ('chaotic'/non-linear)
- Q: forecast x_t



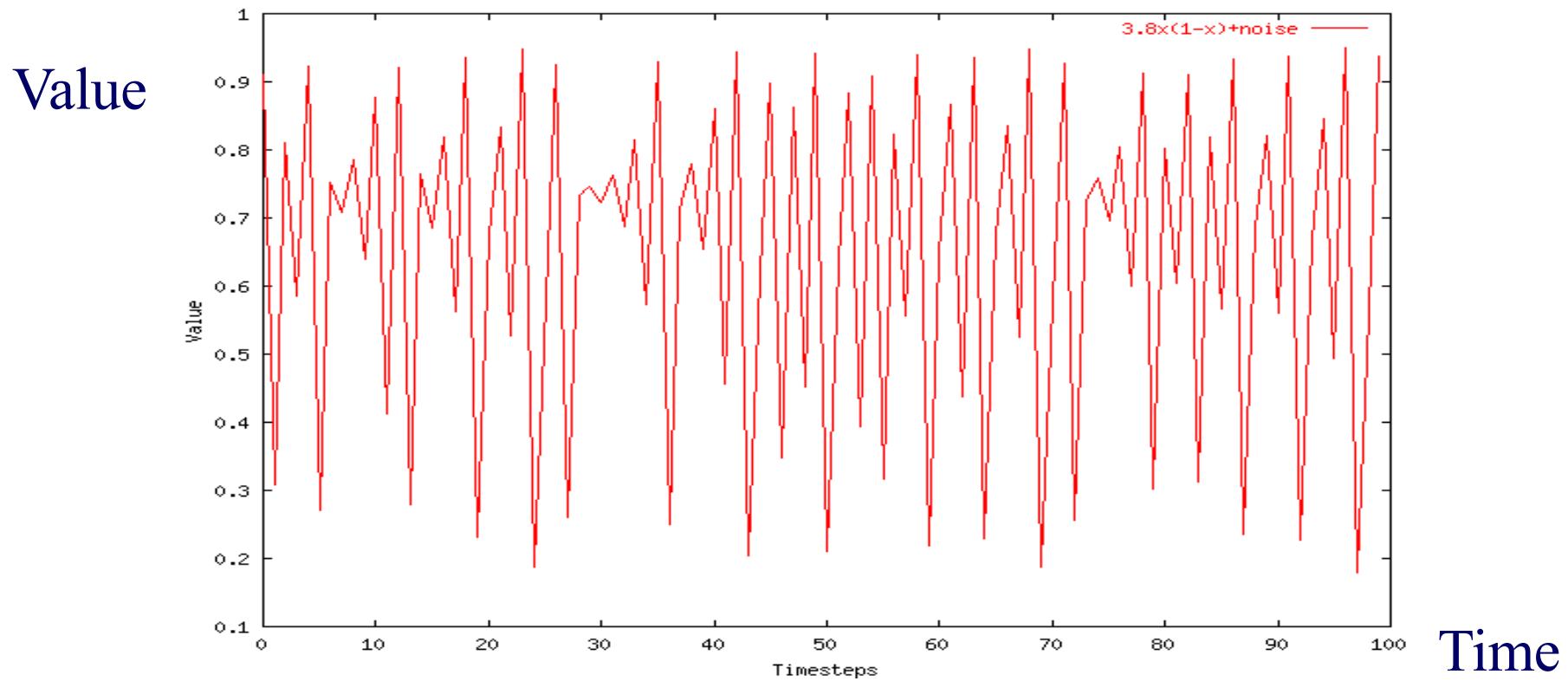
Solution



- given x_{t-1}, x_{t-2}, \dots , ('chaotic'/non-linear)
- Q: forecast x_t
- A: lag-plots + sim. search (= 'Delayed Coordinate Embedding')



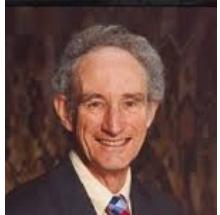
Recall: Problem #1



Given a time series $\{x_t\}$, predict its future course, that is, x_{t+1}, x_{t+2}, \dots

ARIMA pitfall

Example: logistic parabola

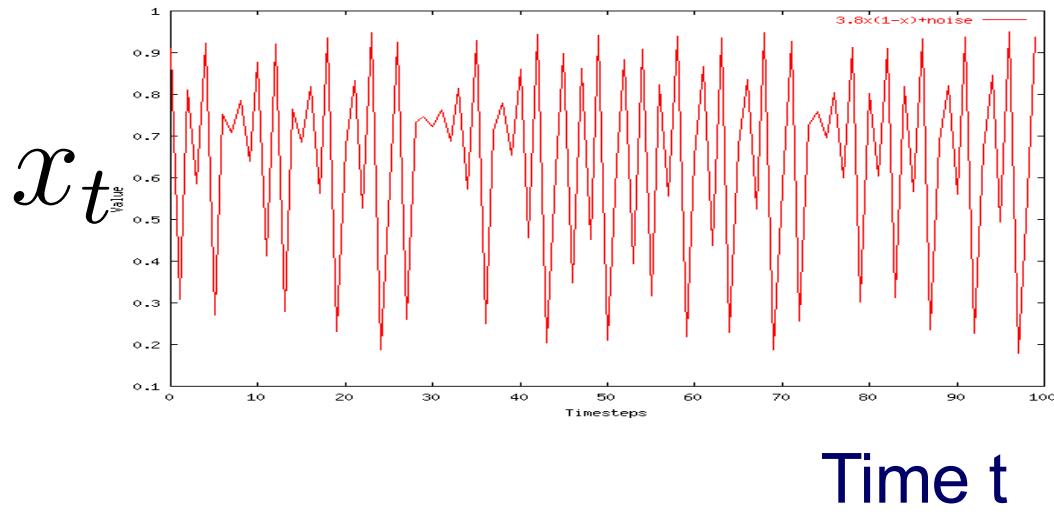


Models population of flies [R. May/1976]

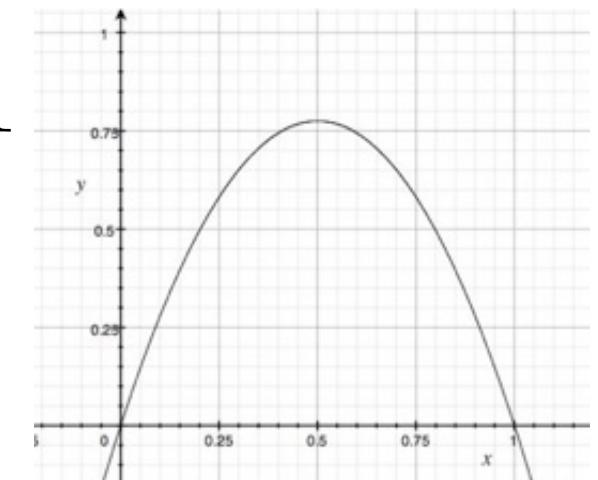


Logistic map

$$x_{t+1} = ax_t \cdot (1 - x_t)$$



Time t



x_t

ARIMA pitfall

Example: logistic parabola

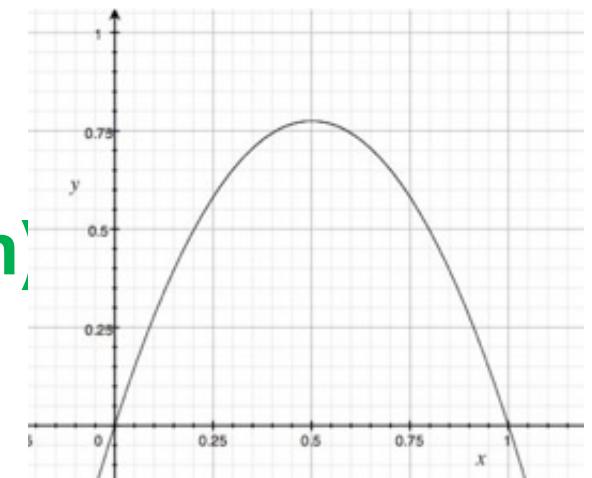
Models population of flies [R. May/1976]



$$x_{t+1} = ax_t \cdot (1 - x_t)$$

Logistic map

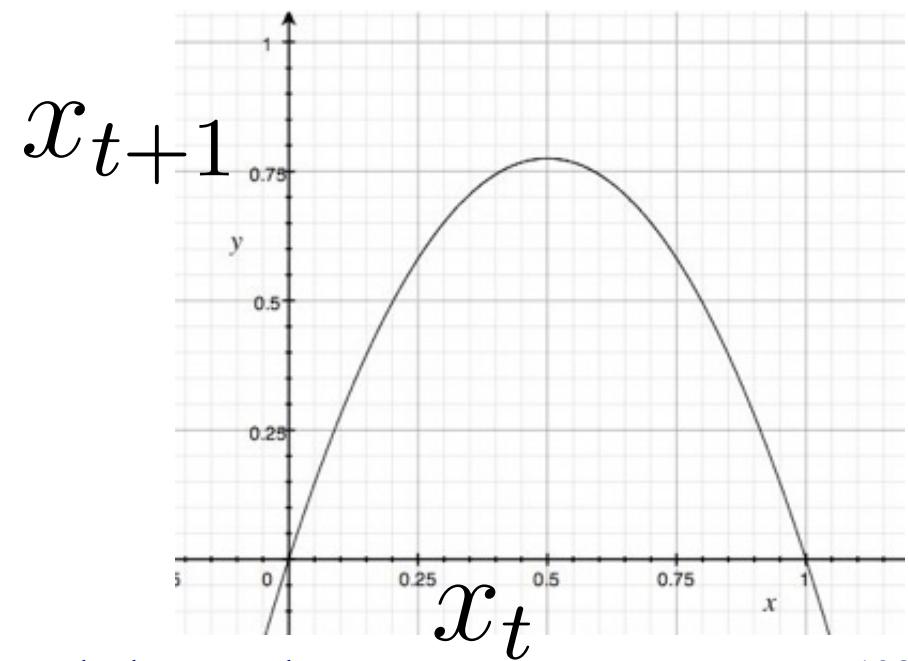
- = SI virus prop. model
- ~ Bass equation (market penetration)
- Special case of Lotka-Volterra



x_t

ARIMA pitfall

Linear equations, e.g., AR, ARIMA, ...



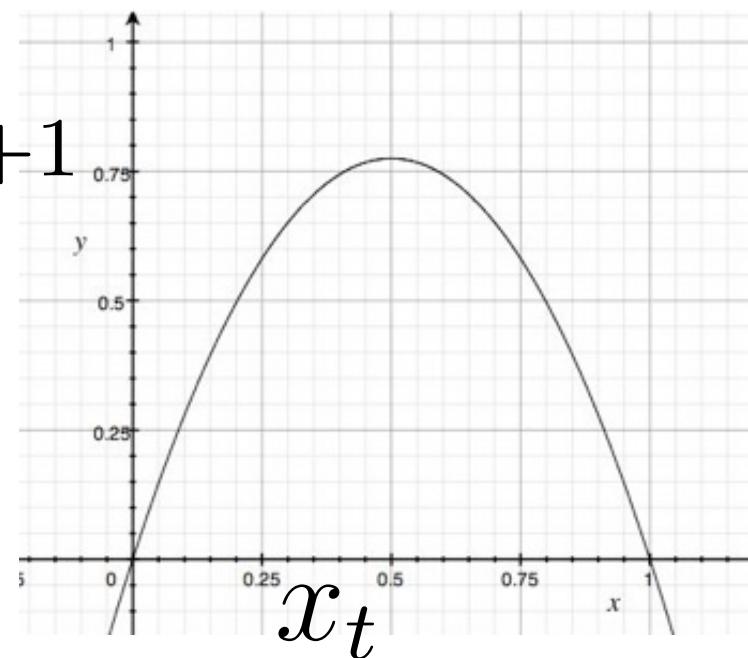
ARIMA pitfall

Linear equations, e.g., AR, ARIMA, ...

e.g., AR(1)

$$x_{t+1} = ax_t + \epsilon$$

x_{t+1}



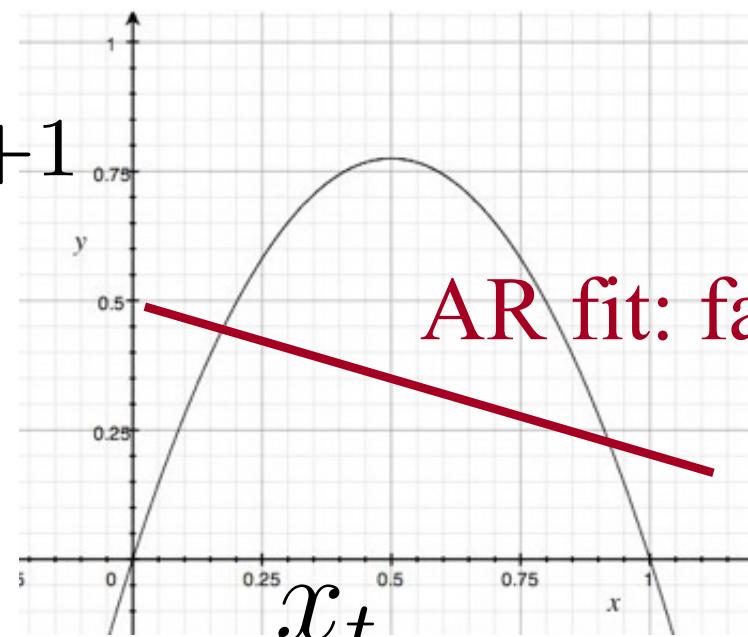
ARIMA pitfall

Linear equations, e.g., AR, ARIMA, ...

e.g., AR(1)

$$x_{t+1} = ax_t + \epsilon$$

x_{t+1}



AR fit: fails

Solution?

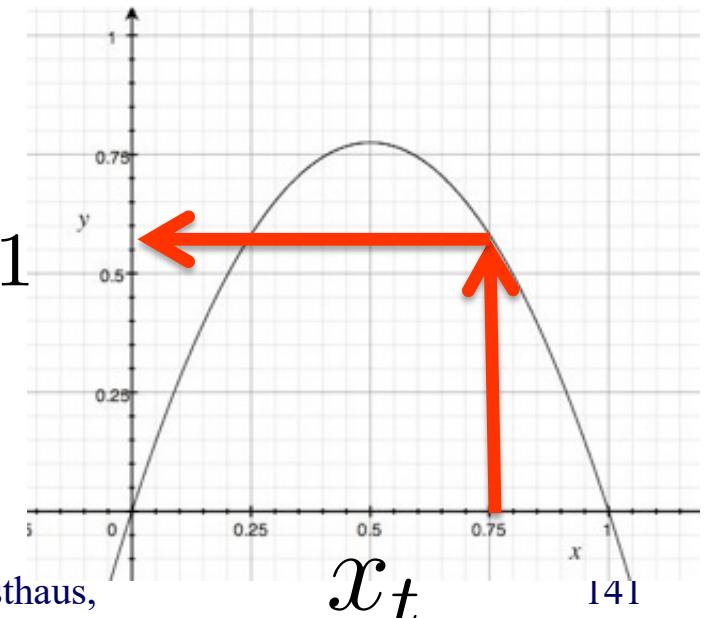
“Delayed Coordinate Embedding”

= Lag Plots

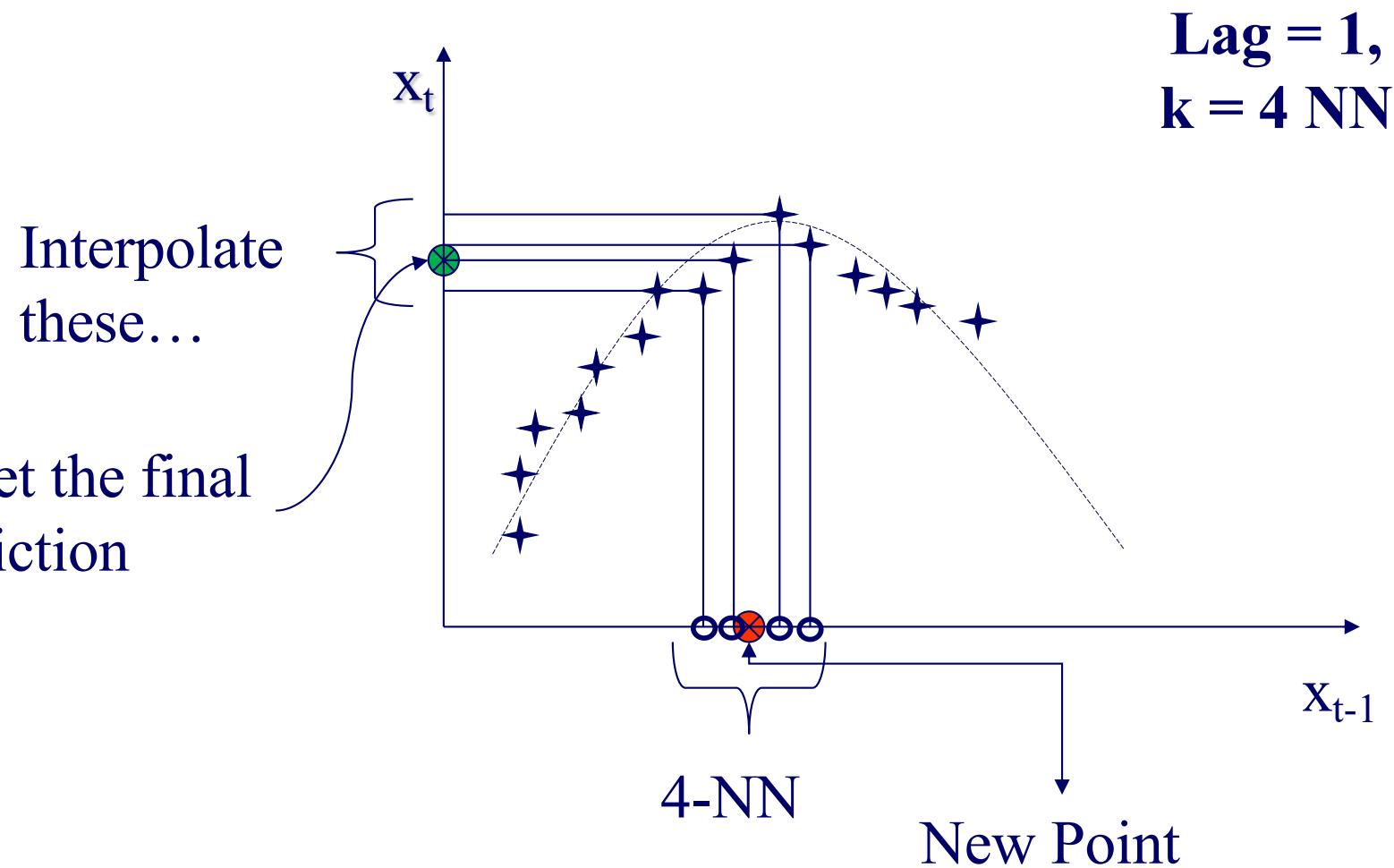
[Sauer94]

k-nearest neighbor search

x_{t+1}



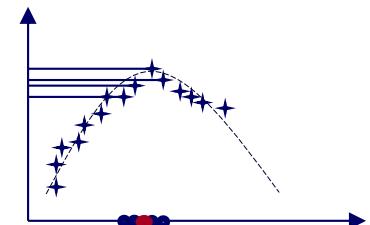
General Intuition (Lag Plot)



Q: How to interpolate?



How do we interpolate between the k nearest neighbors?



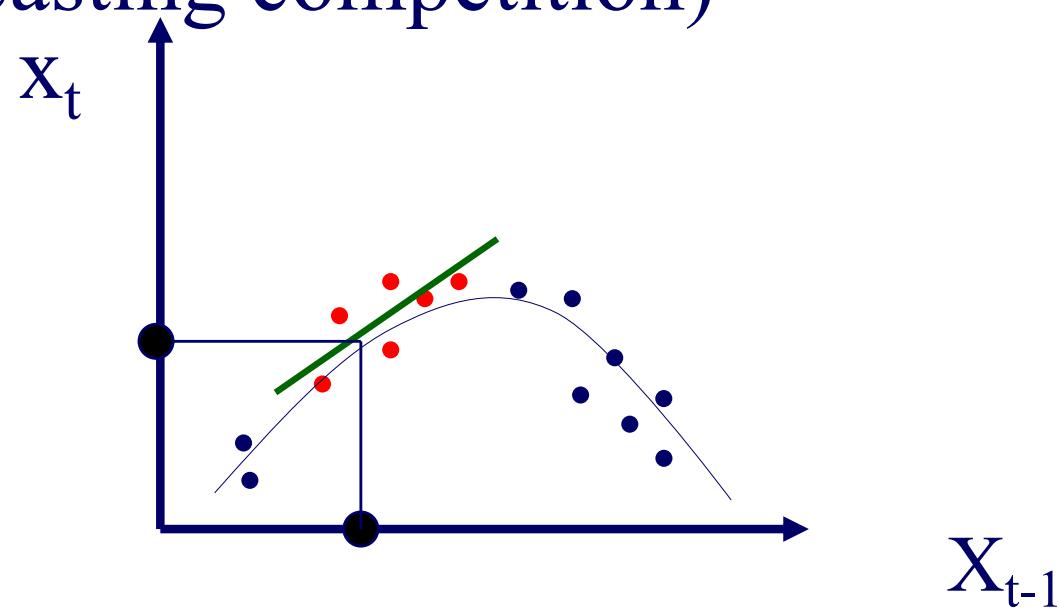
A1: Average

A2: Weighted average (weights drop with distance - how?)

Q: How to interpolate?



A3: Using SVD - seems to perform best
([Sauer94] - first place in the Santa Fe forecasting competition)



Solution?

“Delayed Coordinate Embedding”

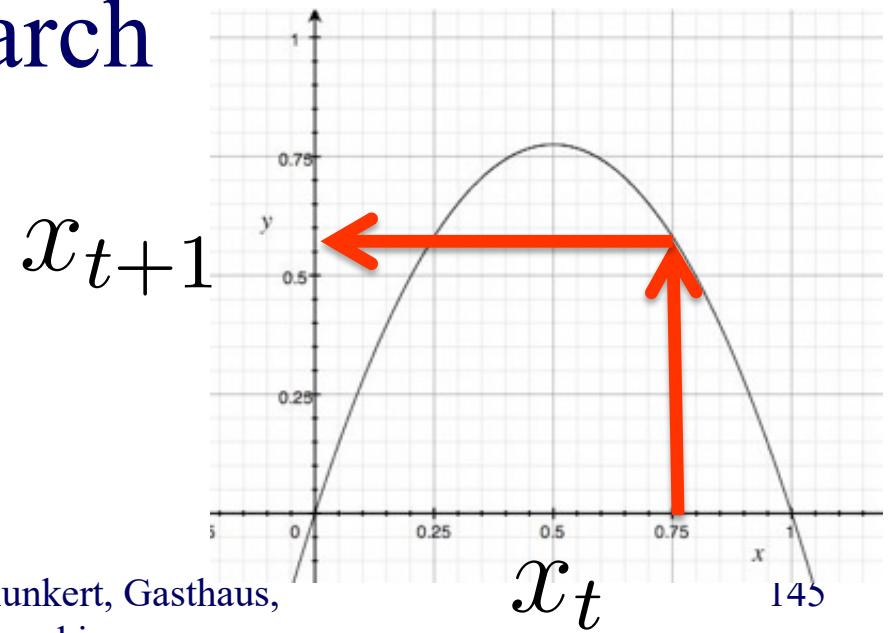
= Lag Plots

[Sauer94]

k-nearest neighbor search

$k = ??$

Lag = ??



Solution?

“Delayed Coordinate Embedding”

= Lag Plots

[Sauer94]

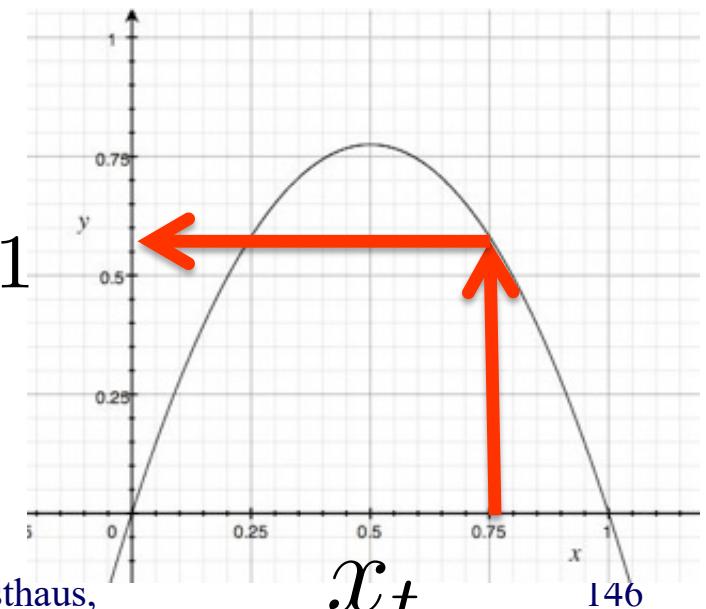
k-nearest neighbor search

$k = \sim 5$

Lag ~ 16

x_{t+1}

x_t



Q: Any theory behind it?



A: YES!

Theoretical foundation



- Based on the “Takens’ Theorem” [Takens81]
- which says that long enough delay vectors can do prediction, even if there are unobserved variables in the dynamical system (= diff. equations)

Theoretical foundation



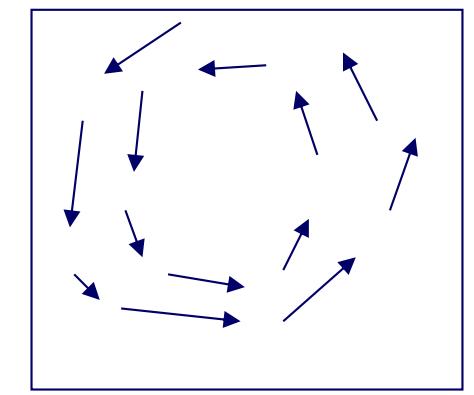
Example: Lotka-Volterra equations



$$\frac{dH}{dt} = r H - a H \cdot P$$

$$\frac{dP}{dt} = b H \cdot P - m P$$

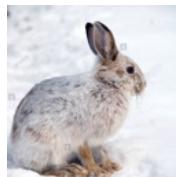
P



H is count of prey (e.g., hare)

P is count of predators (e.g., lynx)

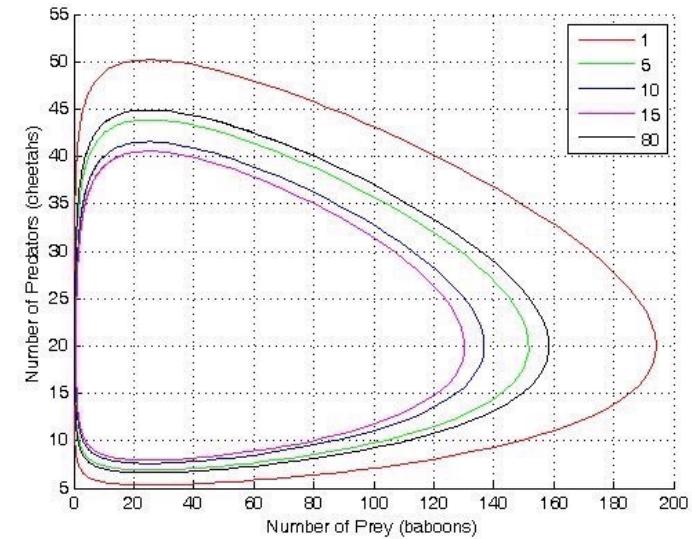
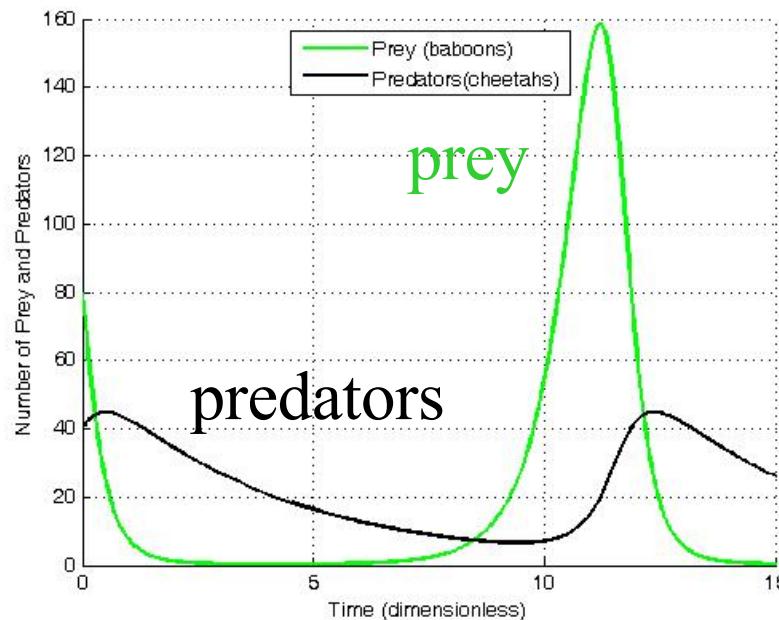
Suppose only $P(t)$ is observed ($t=1, 2, \dots$).



Solution to Volterra-Lotka eq.



predators

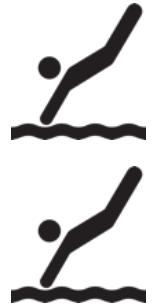


time

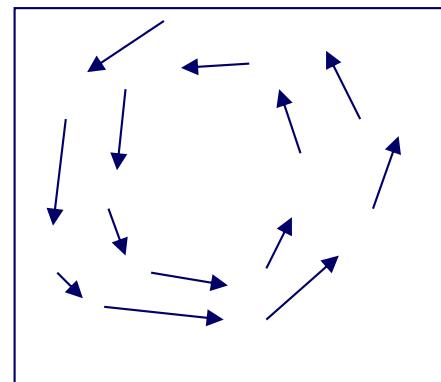
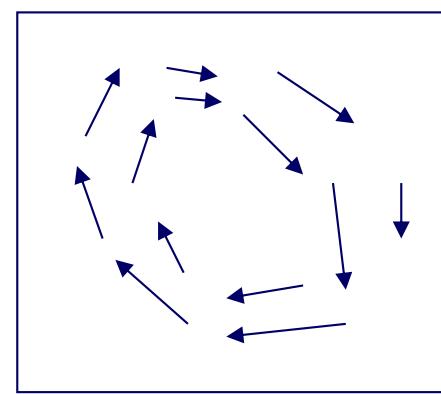
from wikipedia



Theoretical foundation



- But the delay vector space is a faithful reconstruction of the internal system state
- So prediction in **delay vector space** is as good as prediction in **state space**

 P  $P(t)$  $P(t-1)$

Theoretical foundation

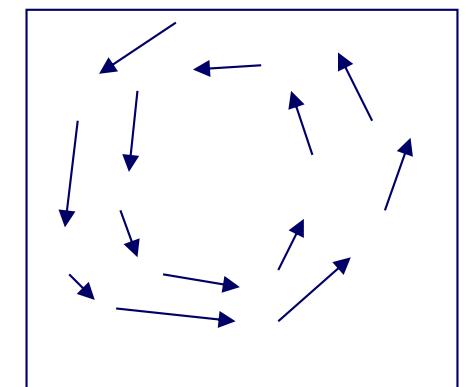


Example: Lotka-Volterra equations

$$\frac{dH}{dt} = r H - a H * P$$

$$\frac{dP}{dt} = b H * P - m P$$

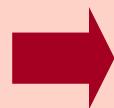
P



More on LV: later

Part 1 - Outline

- ...
- P1.3. Linear forecasting
- P1.4. Non-linear forecasting
 - Problem
 - Idea
 - How-to
 - Experiments
 - (Lotka-Volterra equations)
 - Conclusions
- P1.5. Tensors

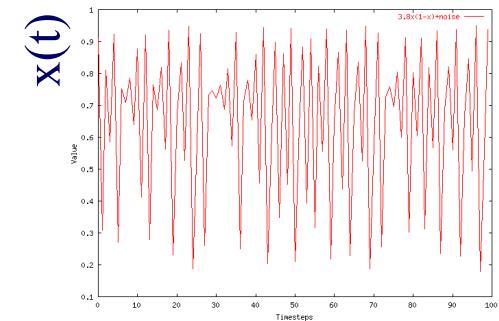


Datasets

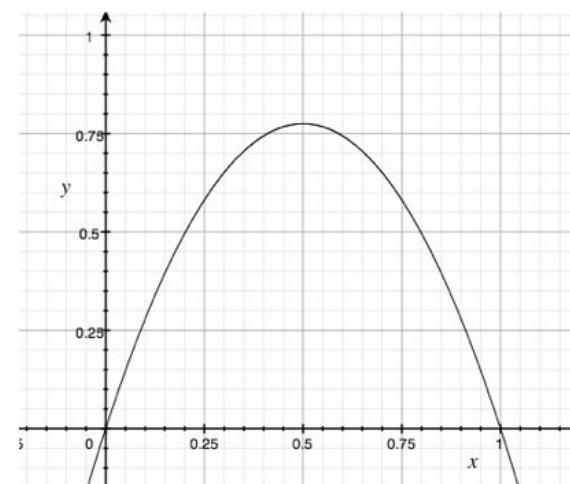
Logistic Parabola:

$$x_t = ax_{t-1}(1-x_{t-1}) + \text{noise}$$

Models population of flies [R. May/1976]



time



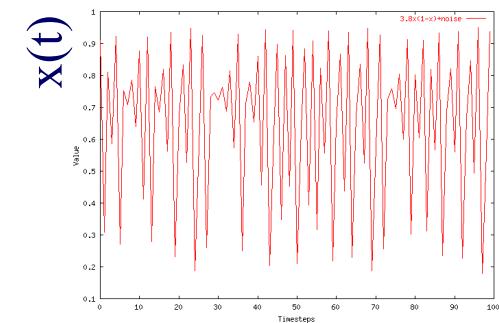
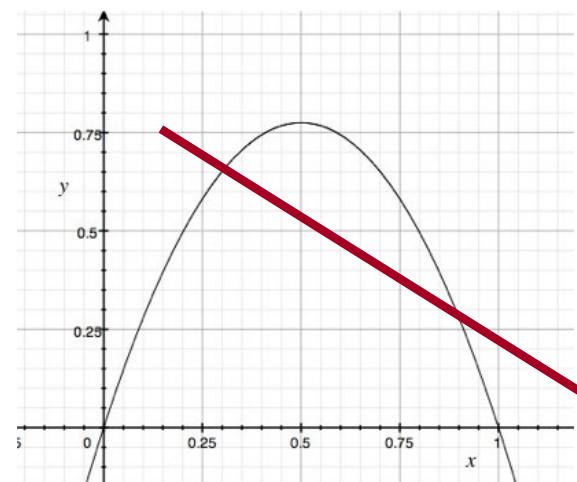
Lag-plot

Datasets

Logistic Parabola:

$$x_t = ax_{t-1}(1-x_{t-1}) + \text{noise}$$

Models population of flies [R. May/1976]

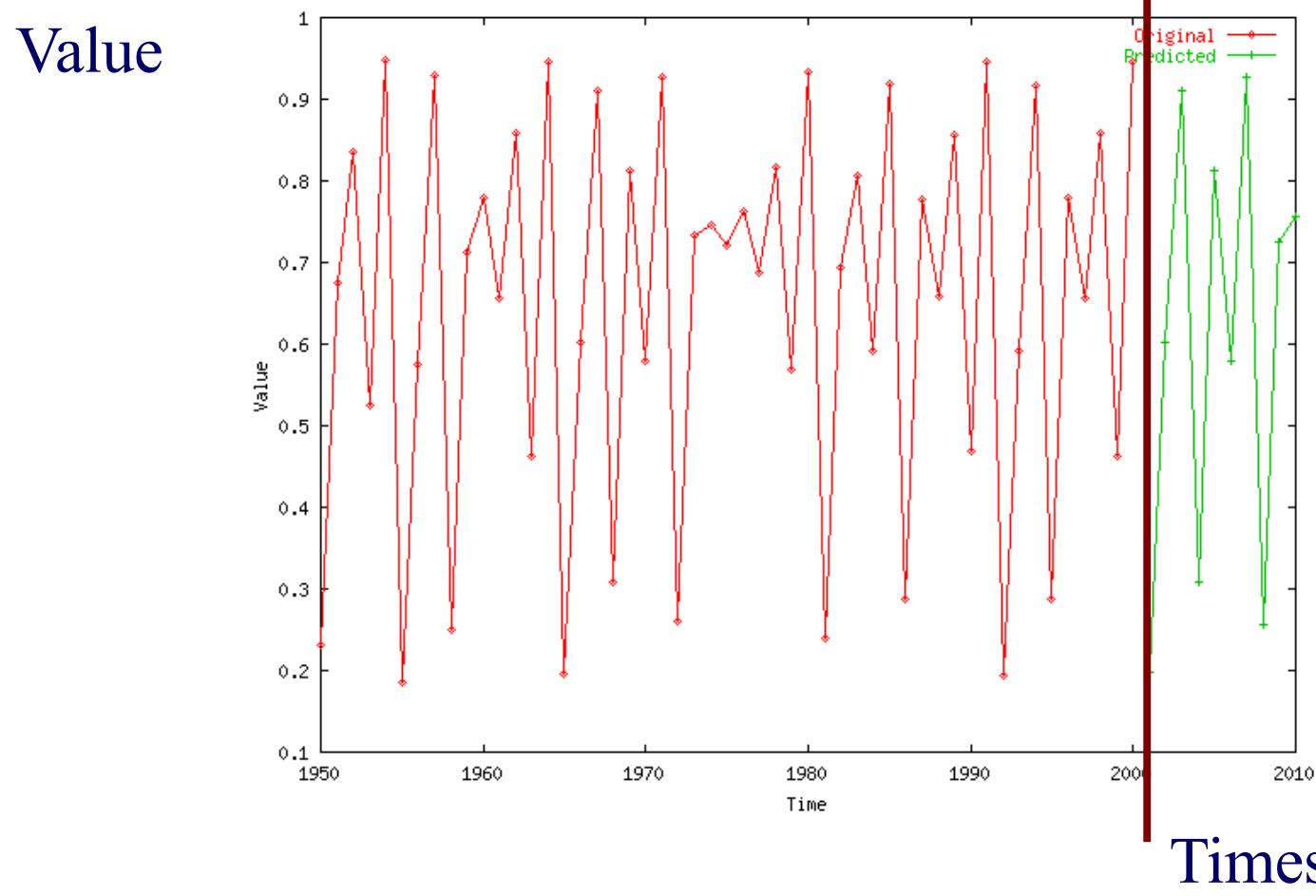


time

Lag-plot
ARIMA: fails

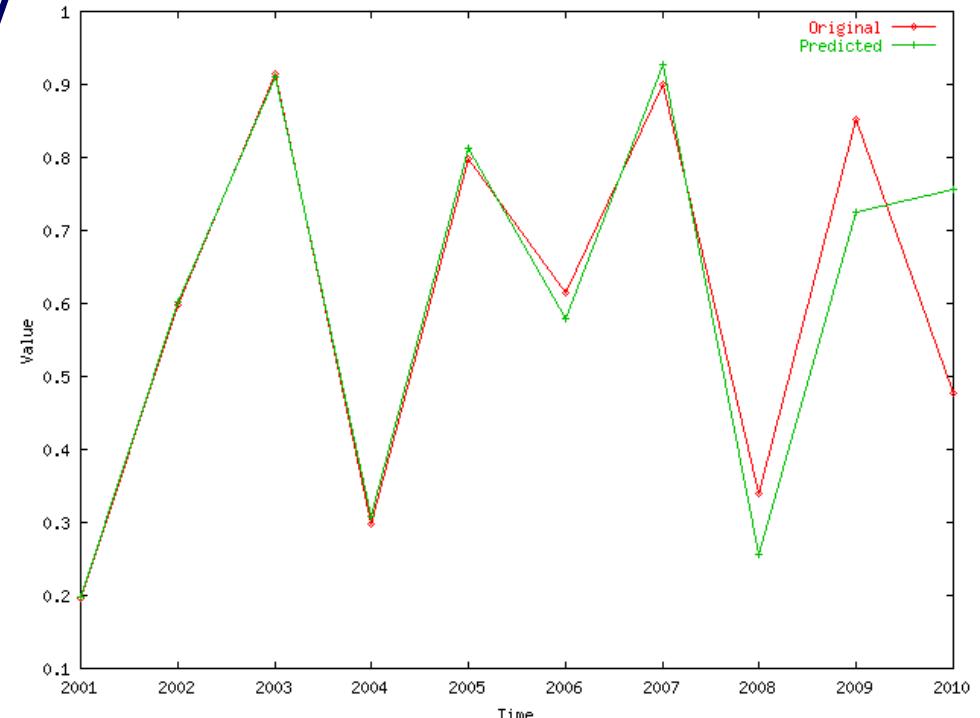
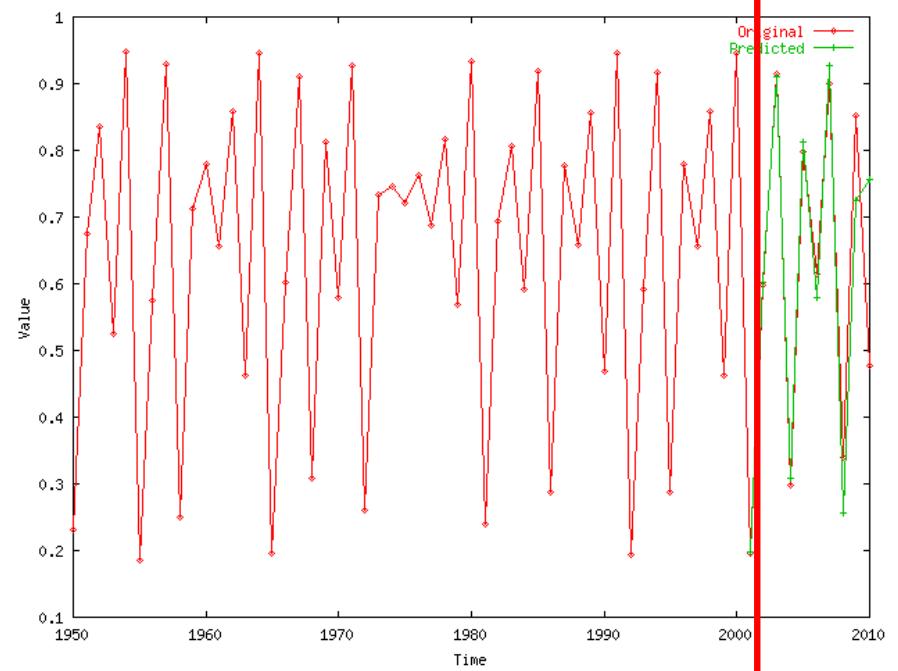
Logistic Parabola

Our Prediction from here



Logistic Parabola

Comparison of prediction
to correct values



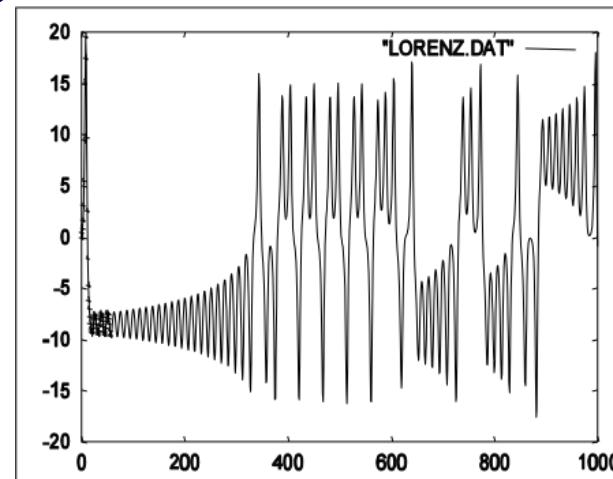
Datasets

LORENZ: Models convection currents in the air

$$\frac{dx}{dt} = a(y - x)$$

$$\frac{dy}{dt} = x(b - z) - y$$

$$\frac{dz}{dt} = xy - c z$$



Datasets

LORENZ: Models convection currents in the air

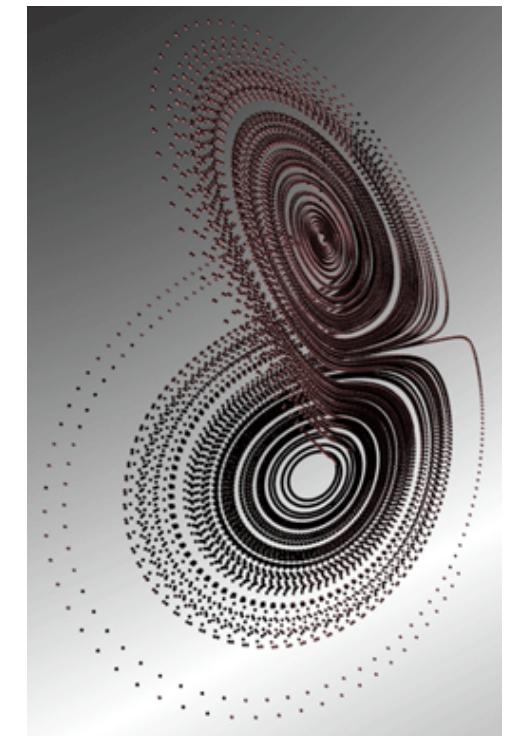
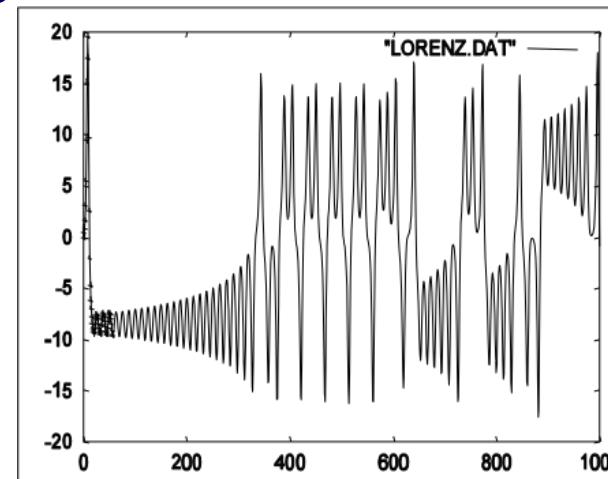
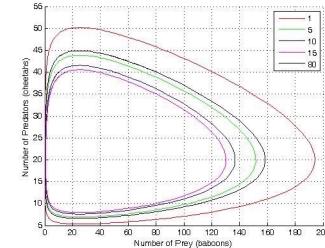
$$\frac{dx}{dt} = a(y - x)$$

$$\frac{dy}{dt} = x(b - z) - y$$

$$\frac{dz}{dt} = xy - c z$$

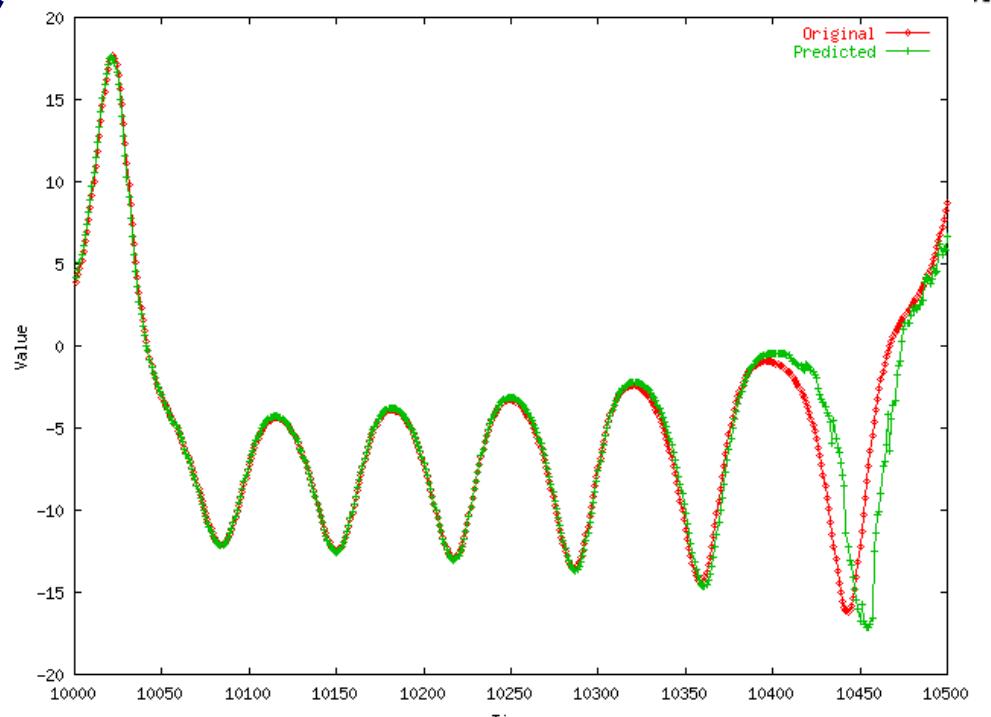
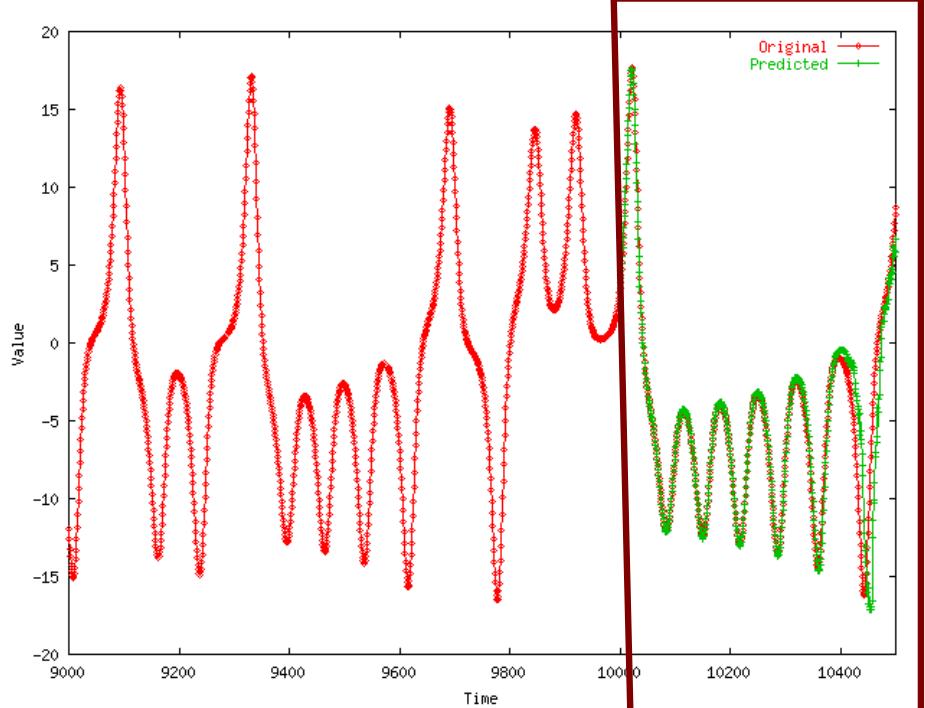


‘butterfly effect’



LORENZ

Comparison of prediction
to correct values



Timesteps



Part 1 - Outline



- ...
- P1.3. Linear forecasting
- P1.4. Non-linear forecasting
 - Problem
 - Idea
 - How-to
 - Experiments
 - (Gray box modeling - Lotka-Volterra equations)
 - Conclusions
- P1.5. Tensors

Notice: LV are vital!

Example: Lotka-Volterra equations

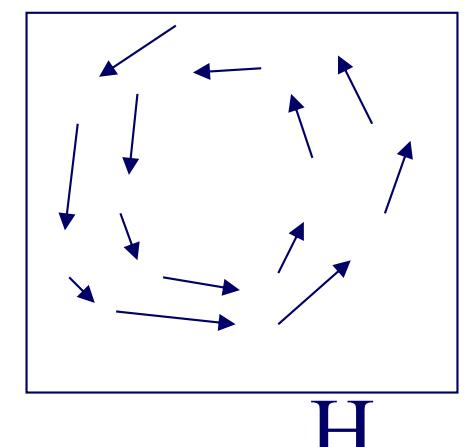
$$dH/dt = r H - a H * P$$

$$dP/dt = b H * P - m P$$

- Prey-predator
- Competing animals (rabbits/goats)
- “ products (stocks/bonds)
- Self-competition (Bass model)
- Virus propagation (SI model)
- ...



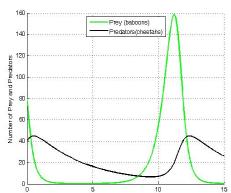
P



Notice: LV are vital!

Example: Lotka-Volterra equations

$$\begin{aligned} dH/dt &= r H - a H * P \\ dP/dt &= b H * P - m P \end{aligned}$$

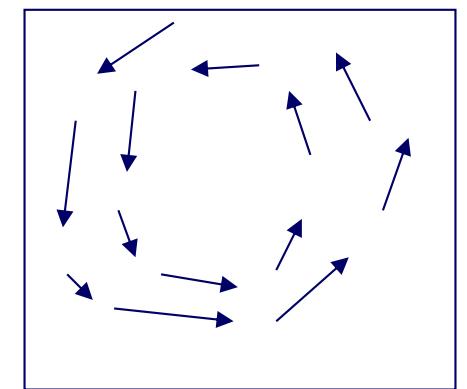


Gray-box modeling:

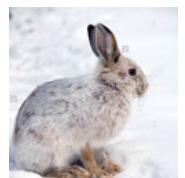
- Given historical data,
- Assume LV
- Estimate highlighted coefficients
- and forecast



P



H



The Web as a Jungle: Non-Linear Dynamical Systems for Co-evolving Online Activities



Yasuko Matsubara (Osaka University)



Yasushi Sakurai (Osaka University)

Christos Faloutsos (CMU)

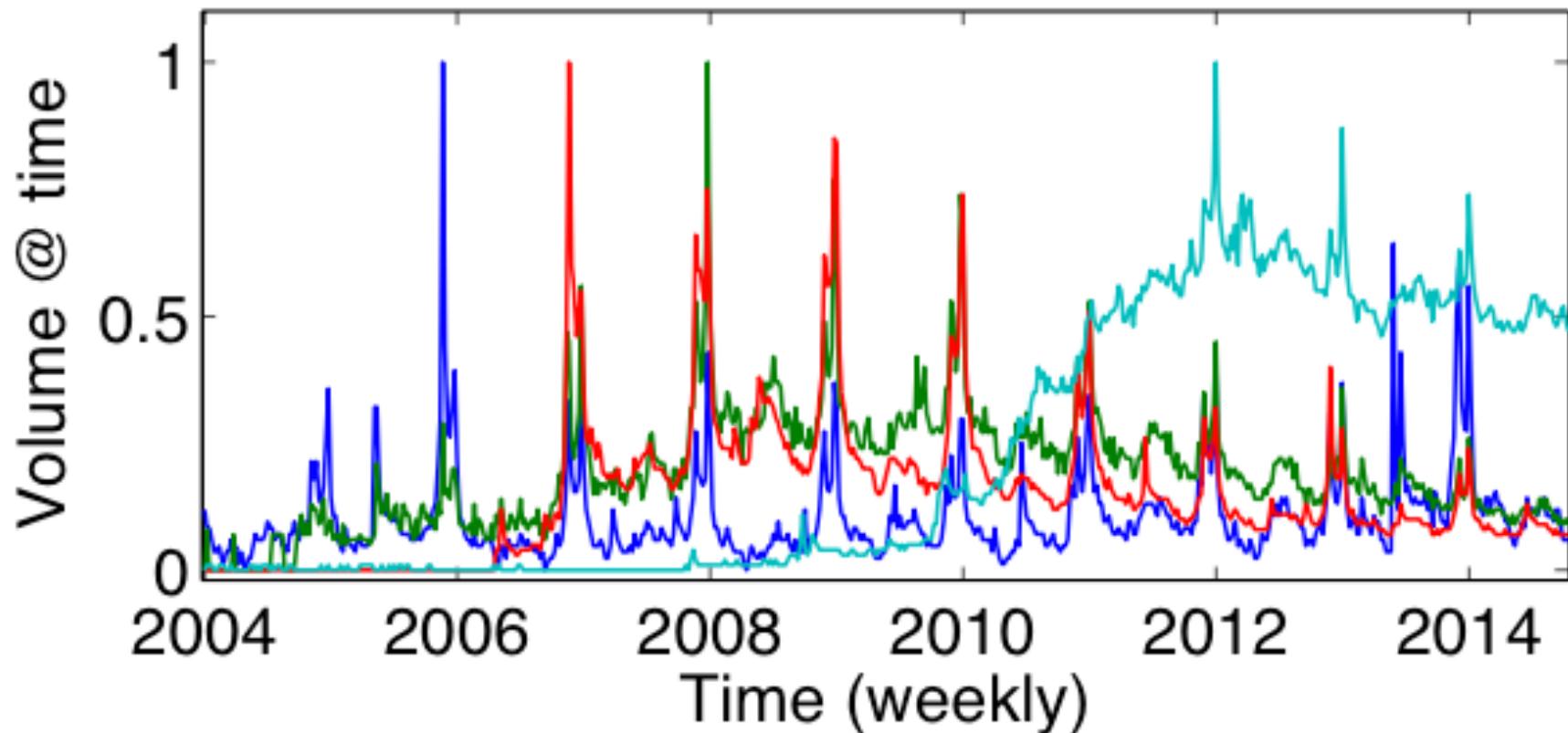
WWW 2015

Open-source code: [here](#)



Given: online user activities

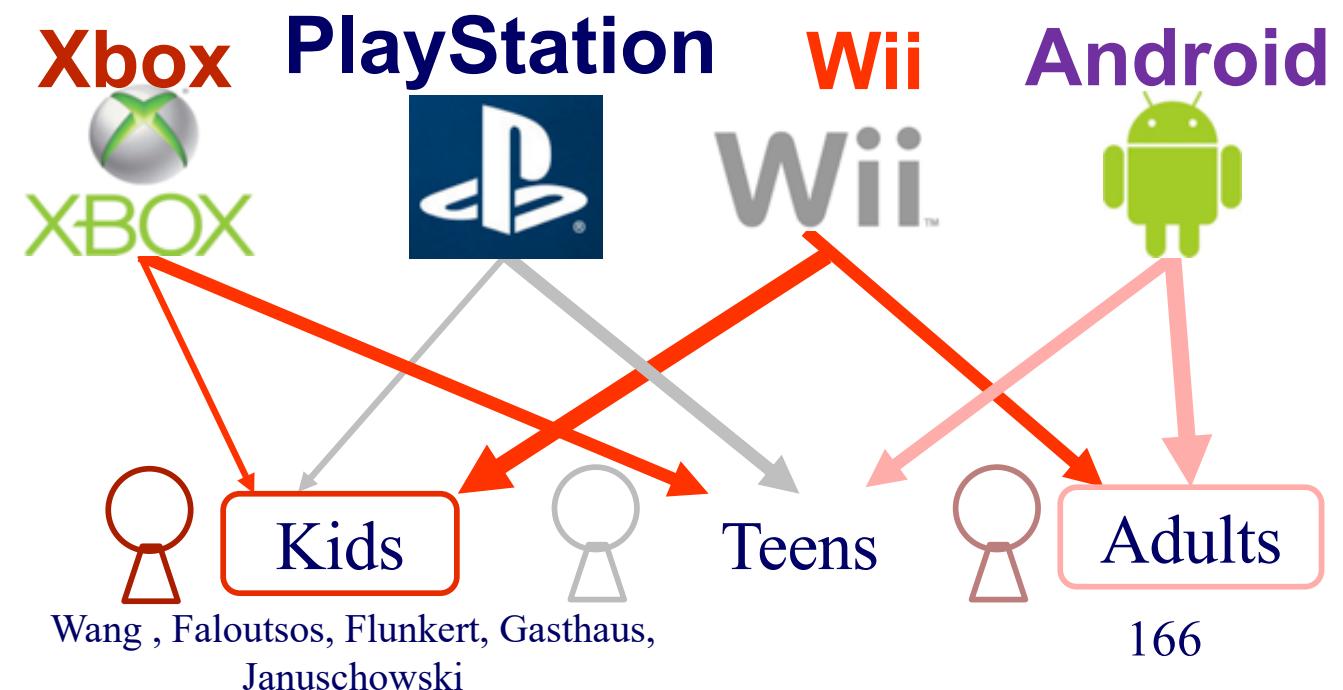
Xbox, PlayStation, Wii, Android



The Web as a jungle



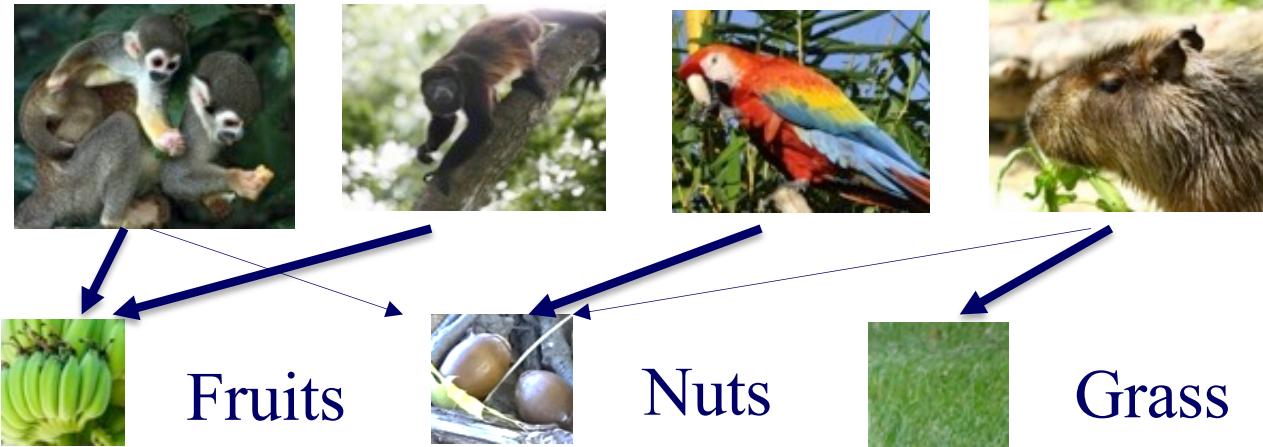
Ecosystem on the Web



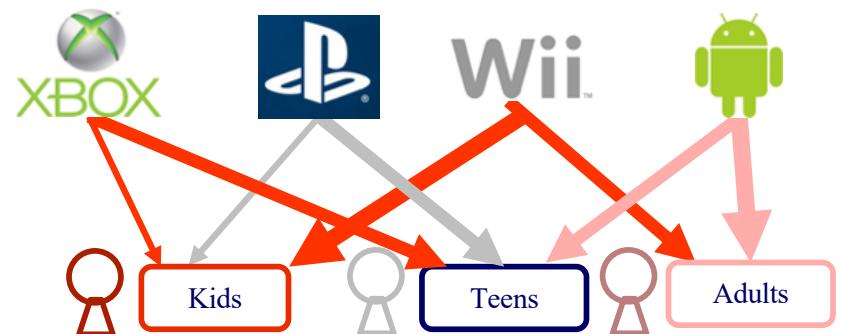


The Web as a jungle

Squirrel monkeys Spider monkeys Macaws Capybaras



Ecosystem in
the
Jungle





LV equations

Interaction between multiple (' d ') species/products/viruses



...



$$P_i(t+1) = P_i(t) \left[1 + r_i \left(1 - \frac{\sum_{j=1}^d a_{ij} P_j(t)}{K_i} \right) \right], \quad (i = 1, \dots, d),$$

$a_{i,j}$: effect (harm) of species j on species i
 (positive: hurts)

$a_{i,i}$: always positive

$a_{i,j}, a_{j,i}$: if both negative: symbiosis (eg., sharks/pilot-fish)



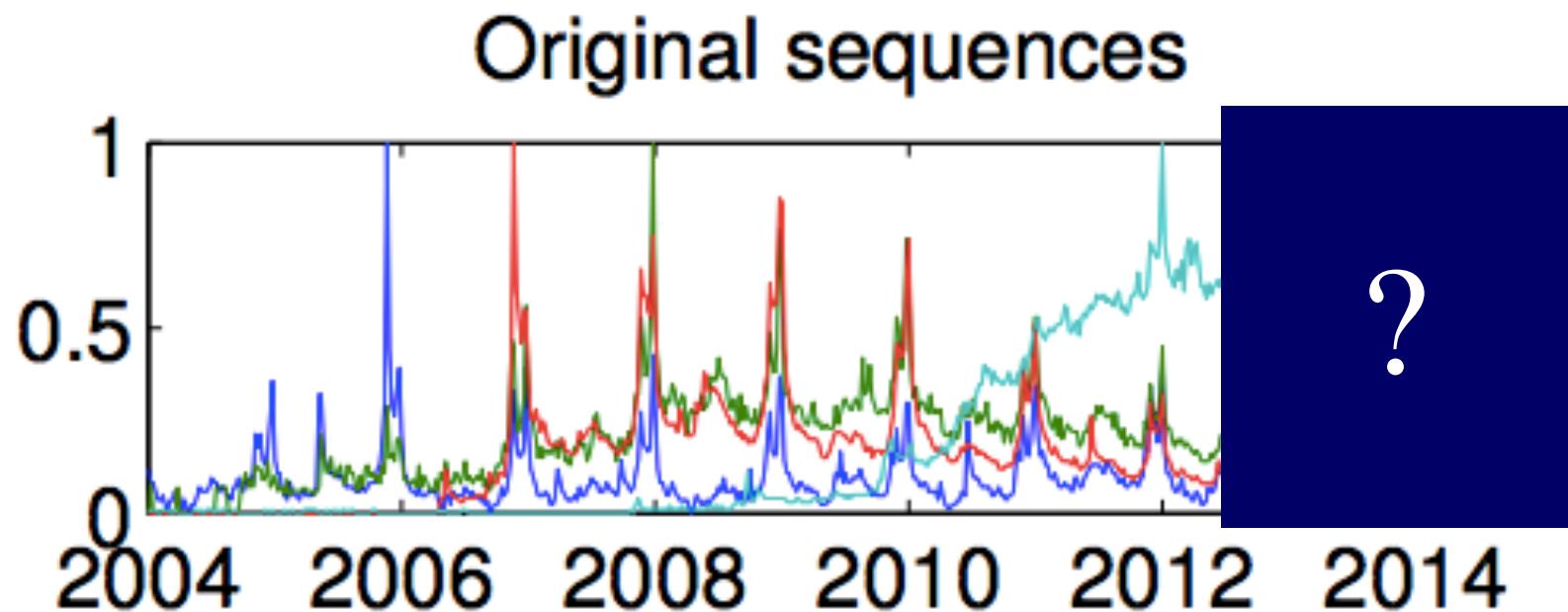
EcoWeb at work - forecasting

Train:

2/3 sequences

Forecast:

1/3 following years

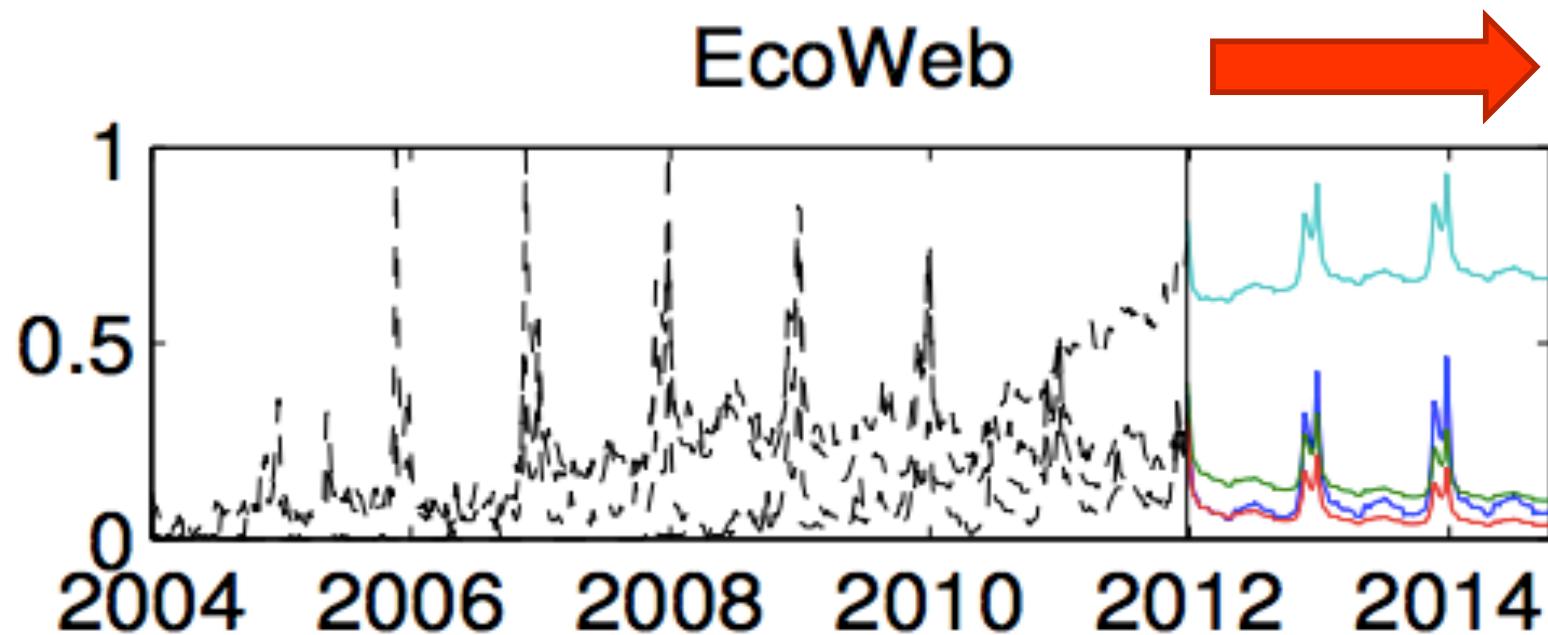




EcoWeb at work - forecasting

Train:
2/3 sequences

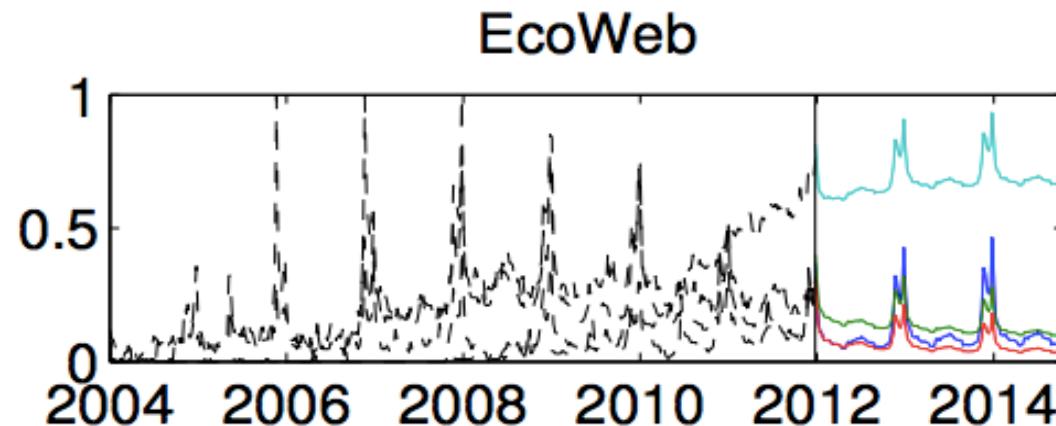
Forecast:
1/3 following years



EcoWeb can capture future patterns



EcoWeb at work - forecasting



Again: Open-source code: [here](#)

Part 1 - Outline

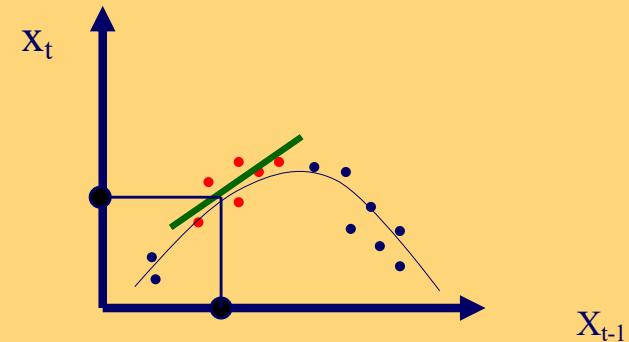
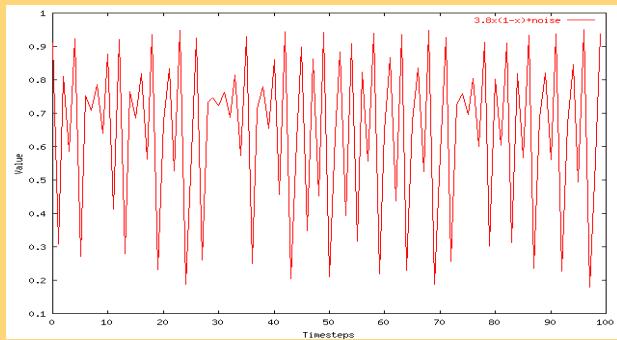


- ...
- P1.3. Linear forecasting
- P1.4. Non-linear forecasting
 - Problem
 - Idea
 - How-to
 - Experiments
 - (Gray box modeling - Lotka-Volterra equations)
 - Conclusions
- P1.5. Tensors

Solution



- given x_{t-1}, x_{t-2}, \dots , ('chaotic'/non-linear)
- Q: forecast x_t
- A: lag-plots + sim. search (= 'Delayed Coordinate Embedding')



References



- Deepay Chakrabarti and Christos Faloutsos *F4: Large-Scale Automated Forecasting using Fractals* CIKM 2002, Washington DC, Nov. 2002.
- Sauer, T. (1994). *Time series prediction using delay coordinate embedding*. (in book by Weigend and Gershenfeld, below) Addison-Wesley.
- Takens, F. (1981). *Detecting strange attractors in fluid turbulence*. Dynamical Systems and Turbulence. Berlin: Springer-Verlag.

References

- Weigend, A. S. and N. A. Gerschenfeld (1994). *Time Series Prediction: Forecasting the Future and Understanding the Past*, Addison Wesley. (Excellent collection of papers on chaotic/non-linear forecasting, describing the algorithms behind the winners of the Santa Fe competition.)

Part 1 - Outline

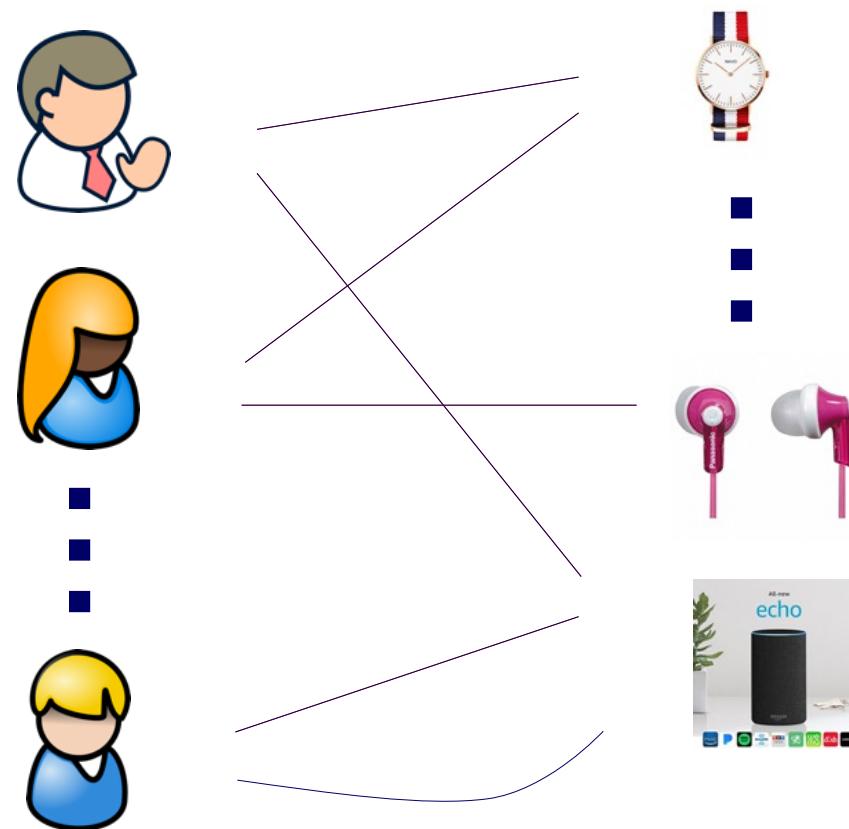


- Motivation
- P1.1. Similarity Search and Indexing
- P1.2. DSP
- P1.3. Linear Forecasting
- P1.4. Non-linear forecasting
- ➡ • P1.5. Tensors
- Conclusions

Part 1.5: Tensors - time evolving graphs

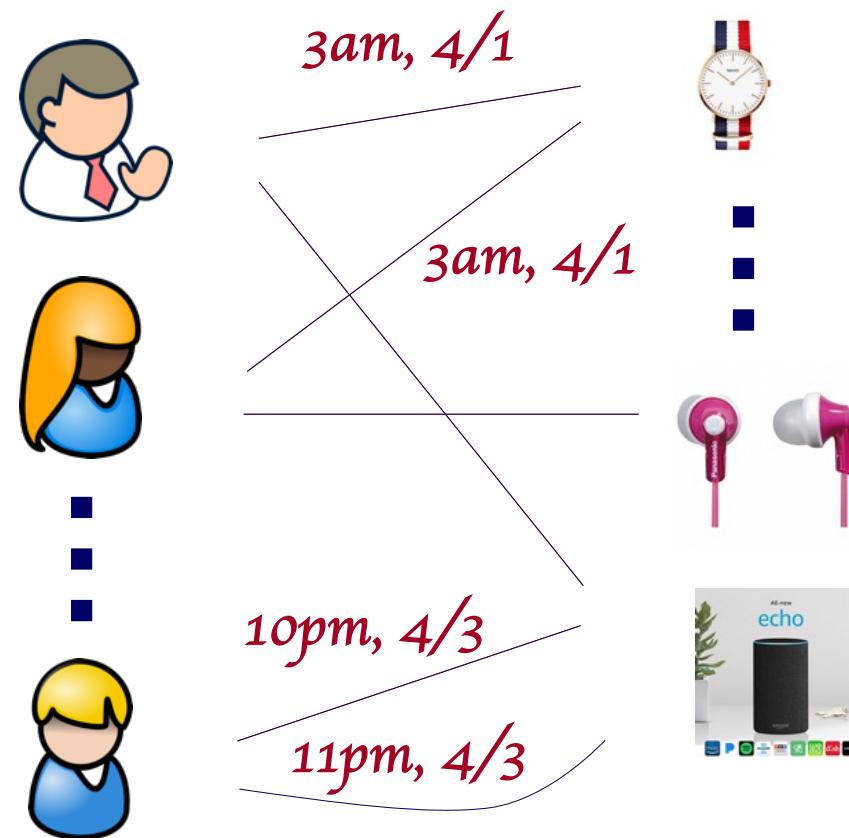
Binary relationships: graph

- Who – buys - what



Time evolving graphs - model?

- Who – buys – what - when



A: tensors

- ... for time evolving graphs (and higher-order relationships)

Tensor examples

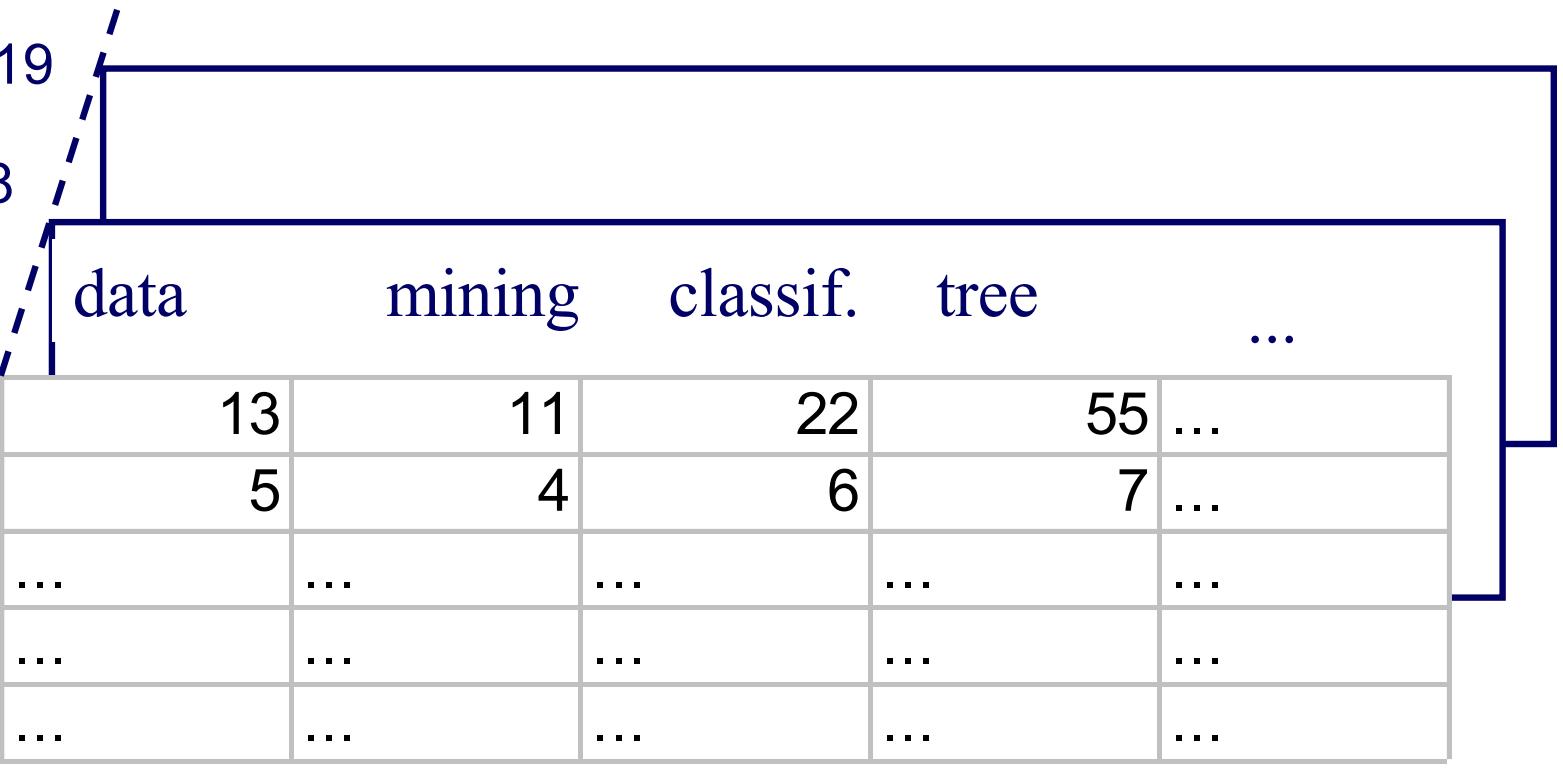
- A: N-D generalization of matrix:

WWW' 17

	data	mining	classif.	tree	...
John	13	11	22	55	...
Peter	5	4	6	7	...
Mary
Nick
...

Tensor examples

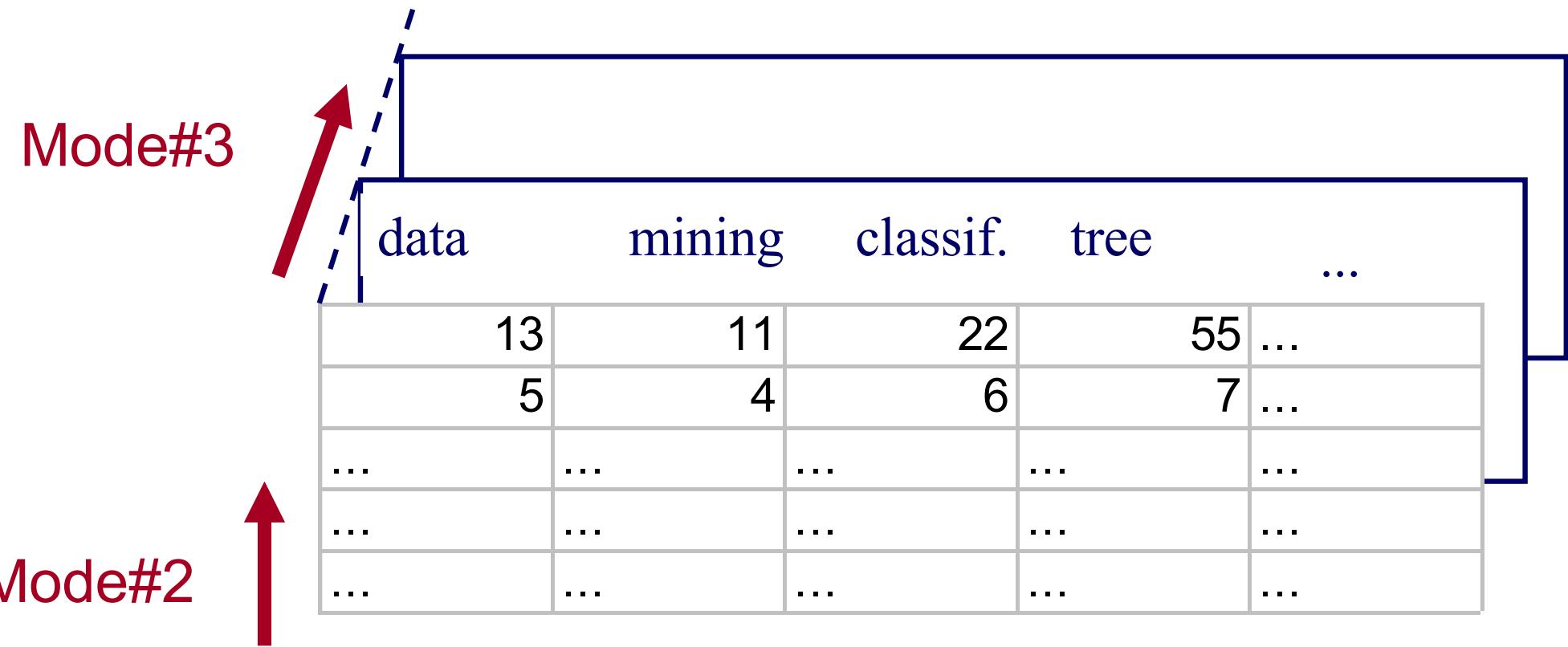
- A: N-D generalization of matrix:



John	13	11	22	55	...
Peter	5	4	6	7	...
Mary
Nick
...

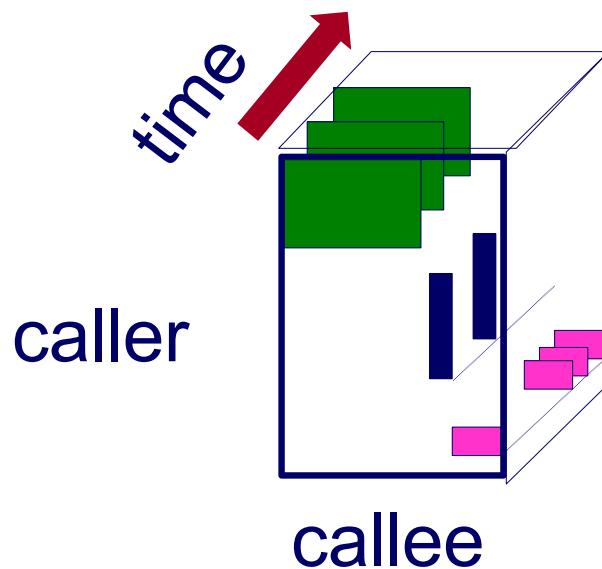
Tensors are useful for 3 or more modes

Terminology: ‘mode’ (or ‘aspect’):



Problem: co-evolving graphs

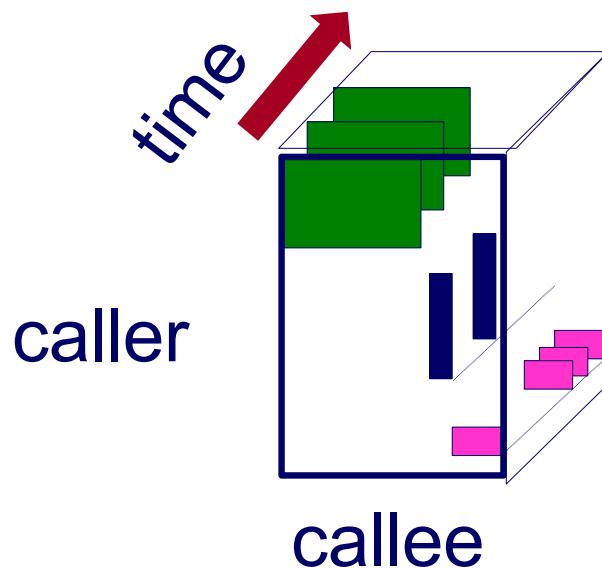
- How to forecast?
 - $4M \times 4M \times 15$ days



$4M \times 4M$ AR models?

Problem: co-evolving graphs

- How to forecast?
 - $4M \times 4M \times 15$ days



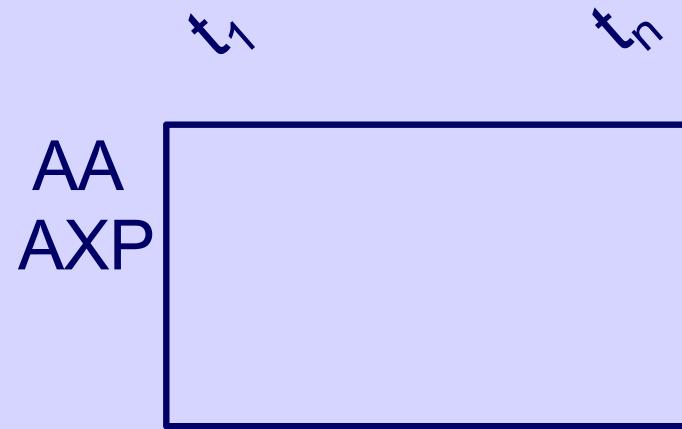
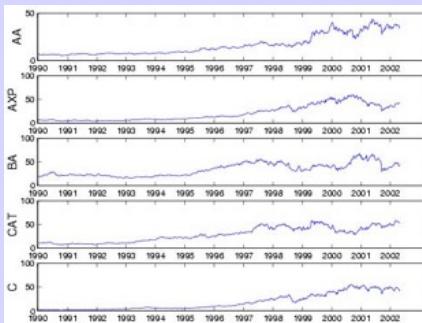
WWW 2020



Wang , Faloutsos, Flunkert, Gasthaus,
Januschowski

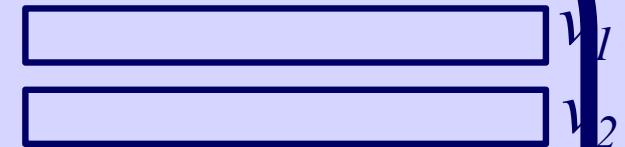
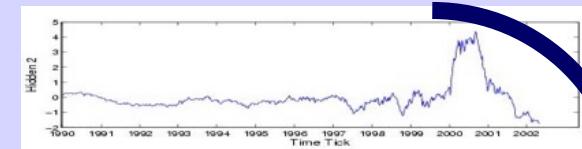
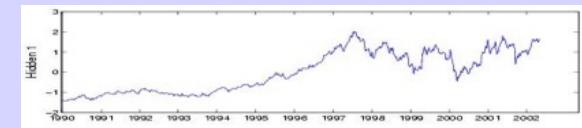
Tensor Basics: PARAFAC decomposition

SVD



$$\sim u_1 \quad u_2$$

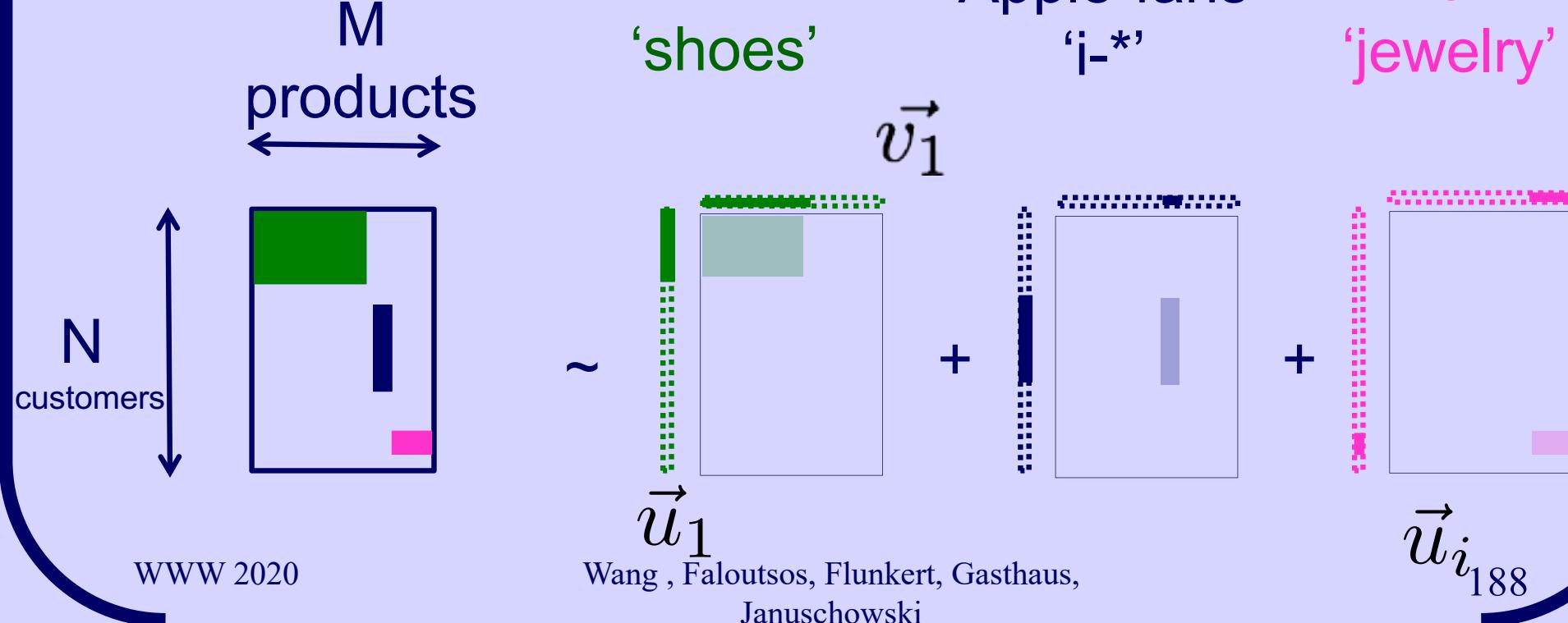
```
# svd code
import numpy as np
u, s, vh = np.linalg.svd(b)
```



$$\mathbf{U} \quad \Sigma \quad \mathbf{V}^T$$

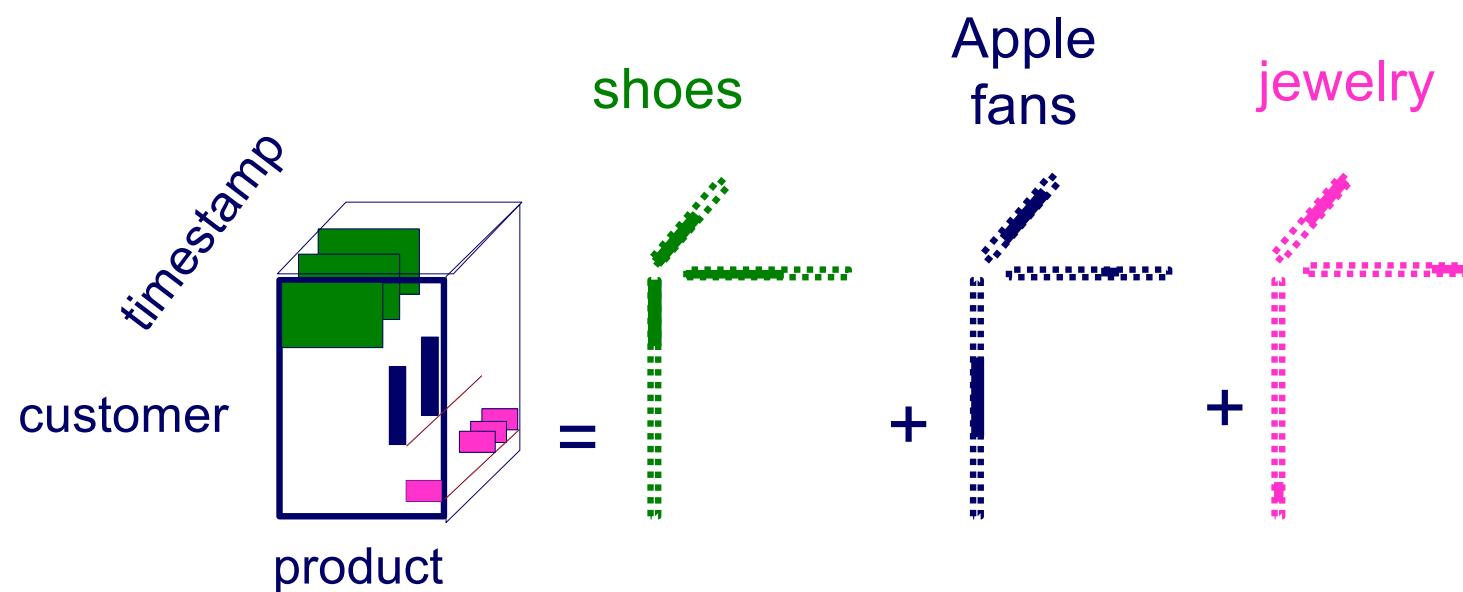
Tensor factorization

- (SVD) matrix factorization: finds blocks



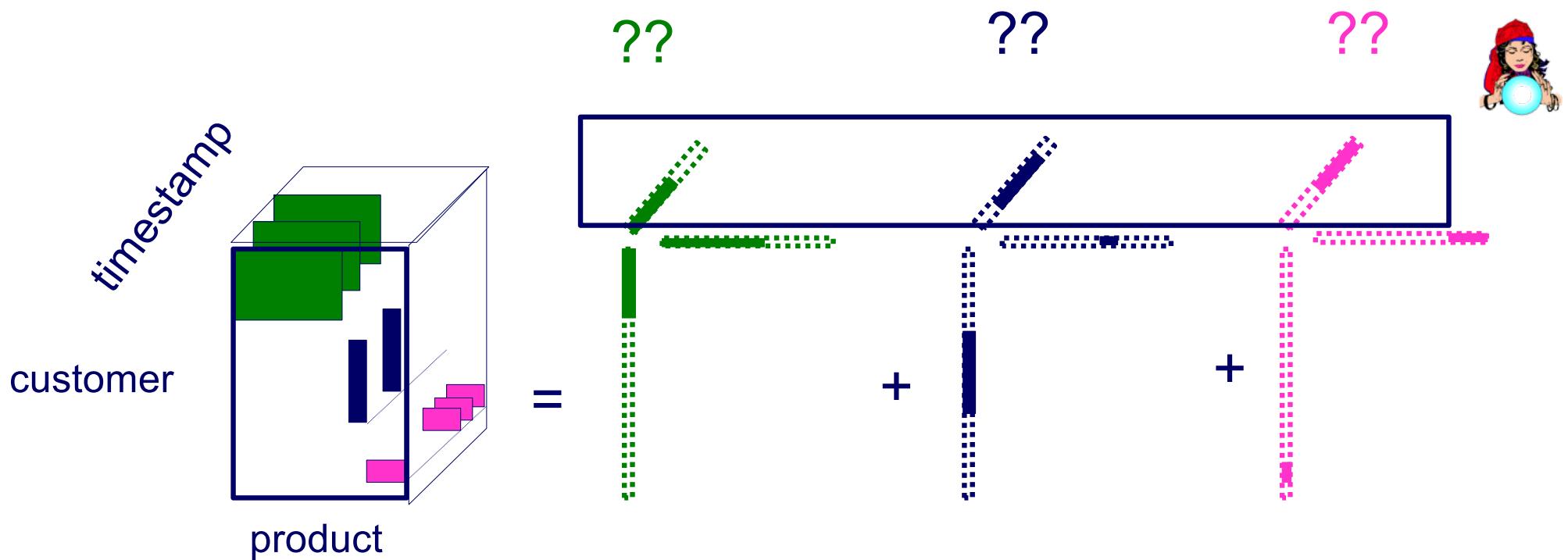
Tensor factorization

- PARAFAC decomposition



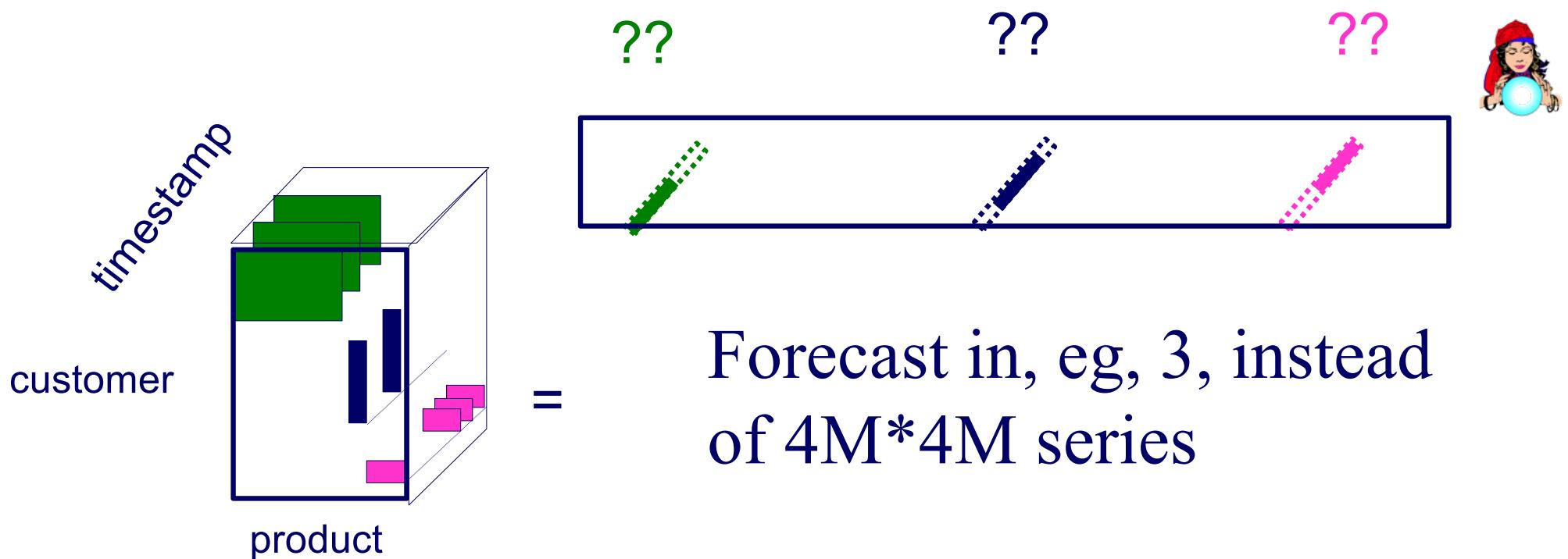
Tensor factorization

- PARAFAC decomposition



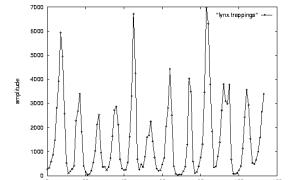
Tensor factorization

- PARAFAC decomposition



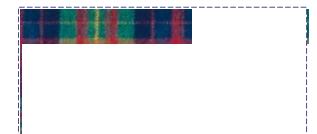


Important observations



Patterns, rules, forecasting and similarity indexing are closely related:

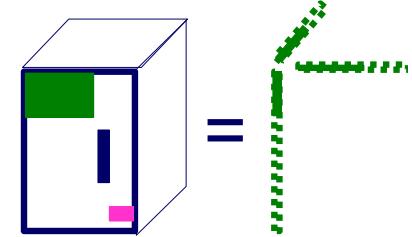
- To do forecasting, we need
 - to find **patterns**/rules
 - compress
 - to find similar settings in the past
- to find outliers, we need to have forecasts
 - (outlier = too far away from our forecast)



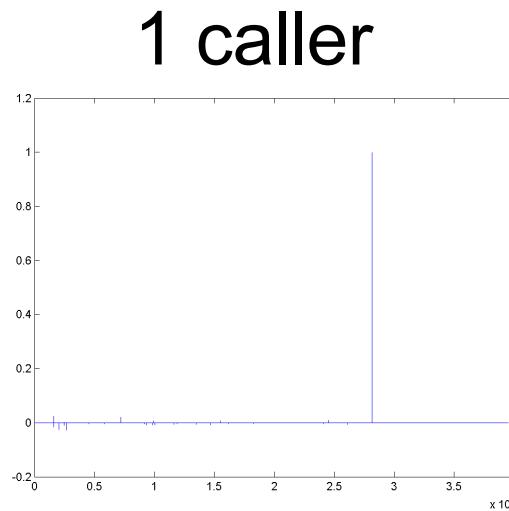
Applications

- • TA1: Phonecall
 - TA2: Network traffic
 - TA3: FaceBook

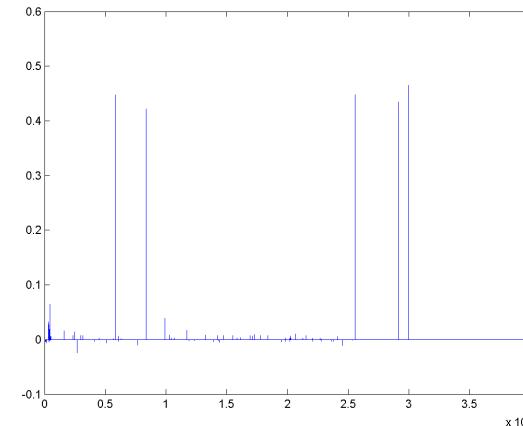
TA1: Anomaly detection in time-evolving graphs



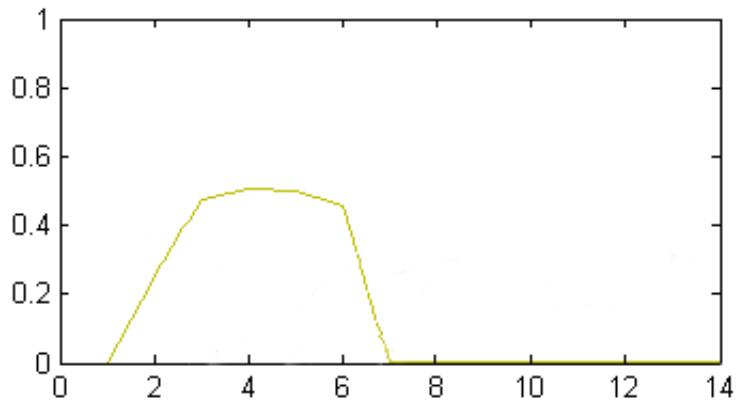
- Anomalous communities in phone call data:
 - European country, 4M clients, data over 2 weeks



5 receivers

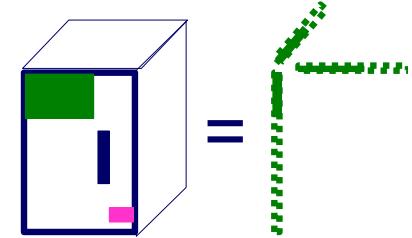


4 days of activity

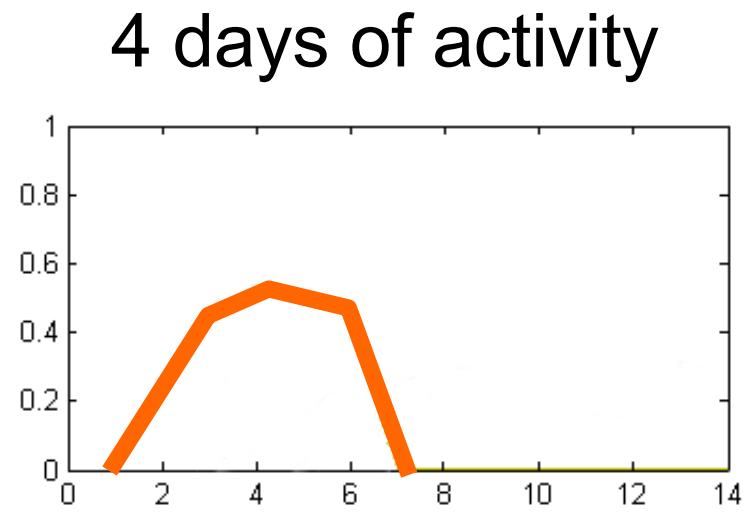
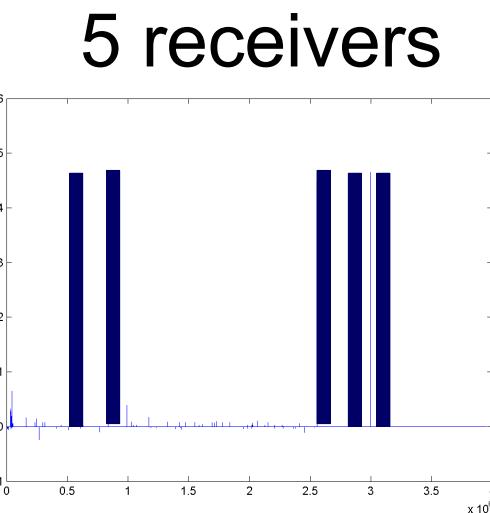
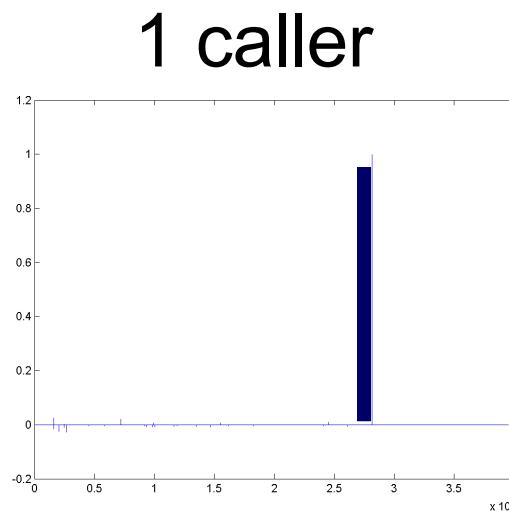


~200 calls to EACH receiver on EACH day!

TA1: Anomaly detection in time-evolving graphs

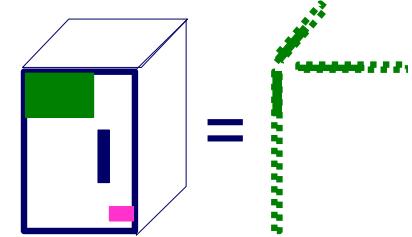


- Anomalous communities in phone call data:
 - European country, 4M clients, data over 2 weeks



~200 calls to EACH receiver on EACH day!

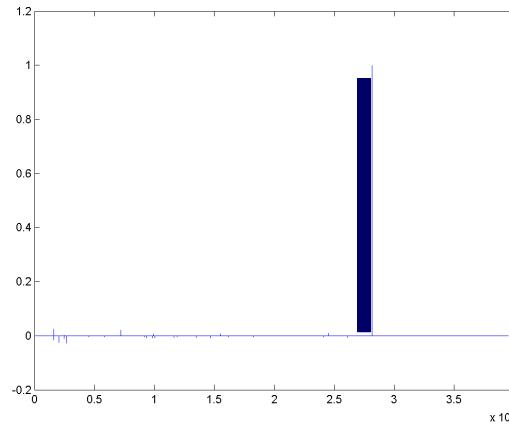
TA1: Anomaly detection in time-evolving graphs



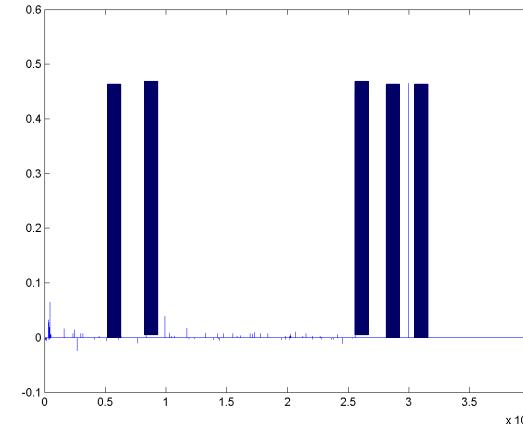
- Anomalous communities in phone call data:
 - European country, 4M clients, data over 2 weeks



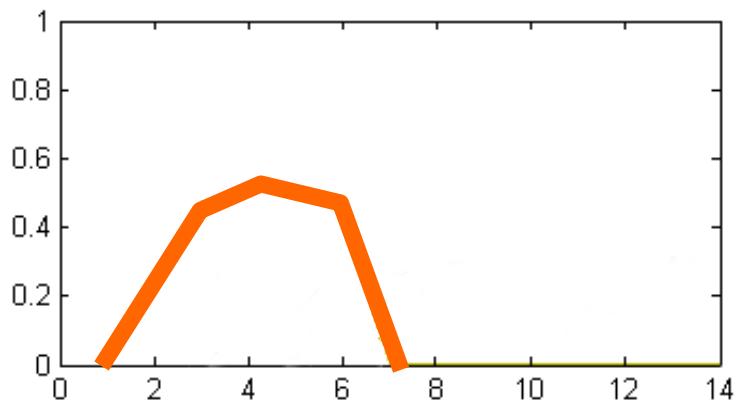
1 caller



5 receivers

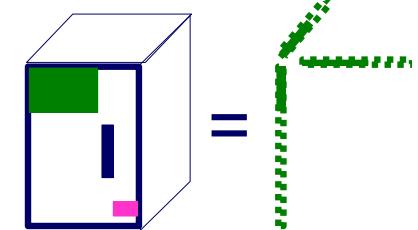


4 days of activity



~200 calls to EACH receiver on EACH day!

TA1: Anomaly detection in time-evolving graphs



- Anomalous communities in phone call data:
 - European country, 4M clients, data over 2 weeks



**Miguel Araujo, Spiros Papadimitriou, Stephan Günnemann,
Christos Faloutsos, Prithwish Basu, Ananthram Swami,
Evangelos Papalexakis, Danai Koutra. Com2: Fast
Automatic Discovery of Temporal (Comet) Communities.
PAKDD 2014, Tainan, Taiwan.**

Applications

- TA1: Phonecall
- • TA2: Network traffic
- TA3: FaceBook



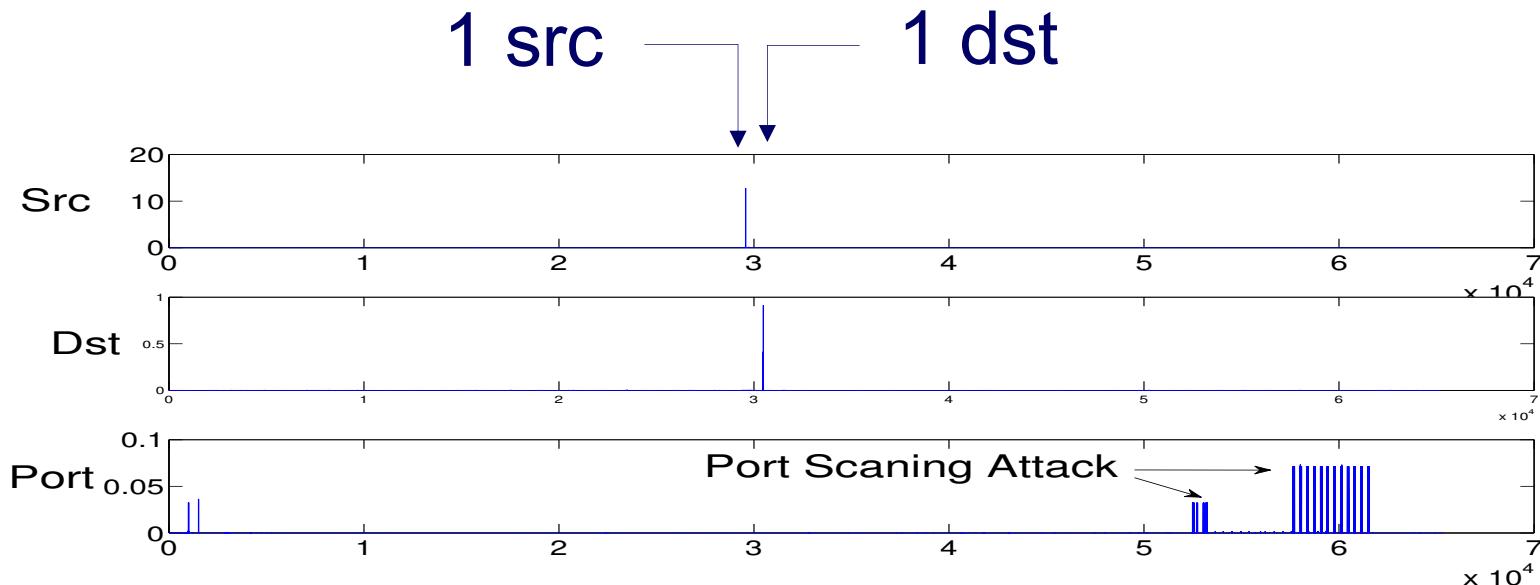
ParCube: Sparse Parallelizable Tensor Decompositions

**Evangelos E. Papalexakis, Christos Faloutsos,
Nikos Sidiropoulos, ECML/PKDD 2012**

<http://www.cs.ucr.edu/~epapalex>



TA2: LBNL Network Data

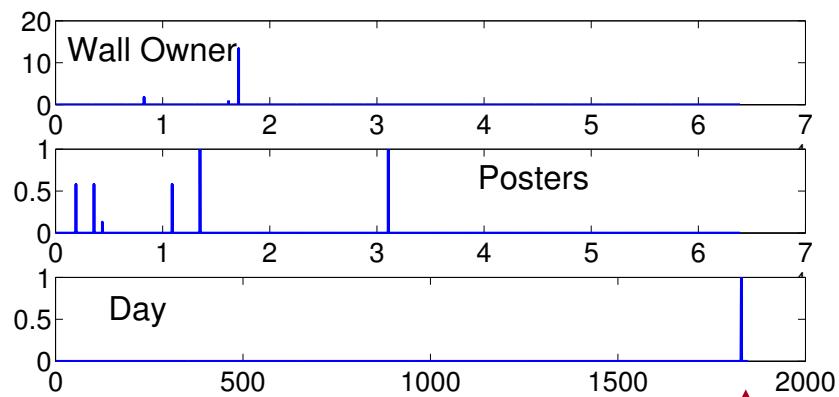


- Modes: src IP, dst IP, port #
- ~ **Port Scanning Attack**

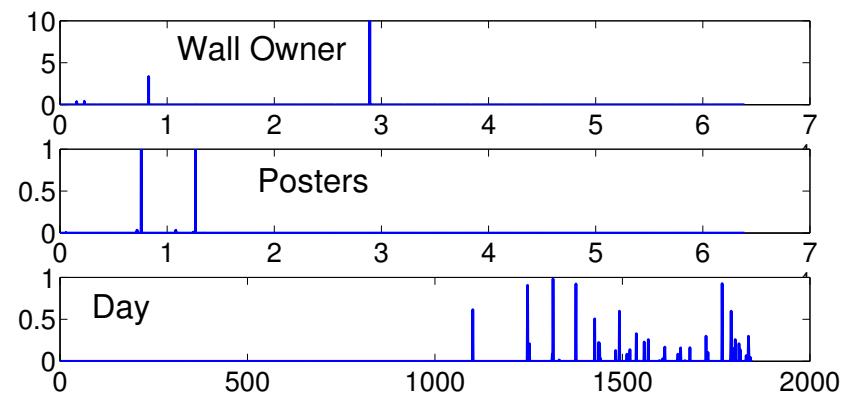


TA3: FACEBOOK Wall posts

1 Wall



(a) FACEBOOK anomaly



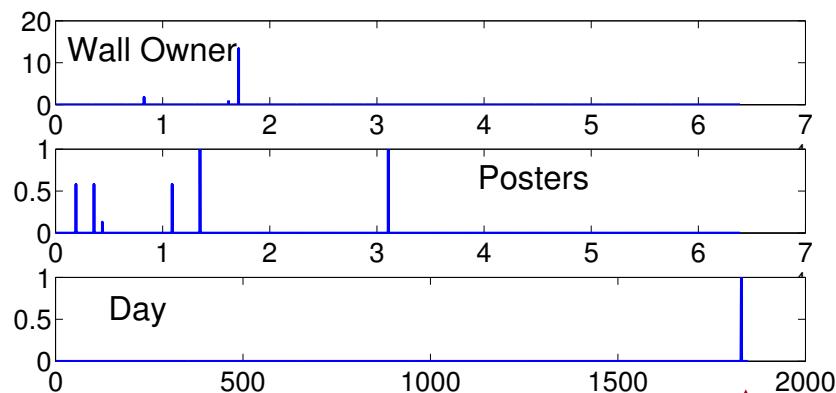
(b) FACEBOOK normal activity

- Modes: wall-owner, poster, timestamp



TA3: FACEBOOK Wall posts

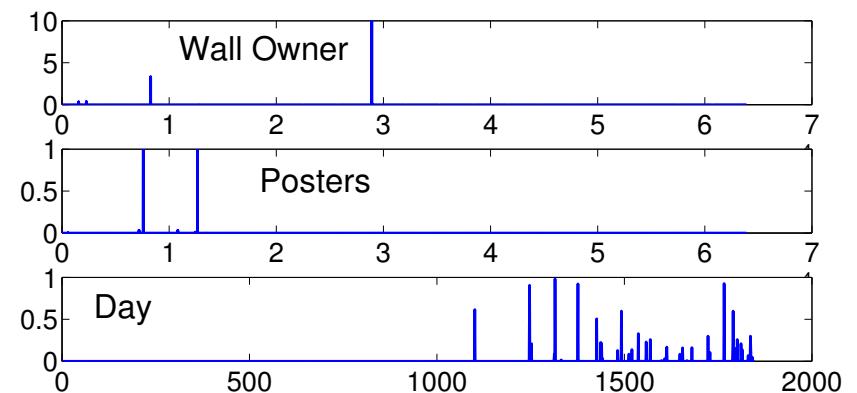
1 Wall



(a) FACEBOOK anomaly (Wall owner's birthday)



1 day



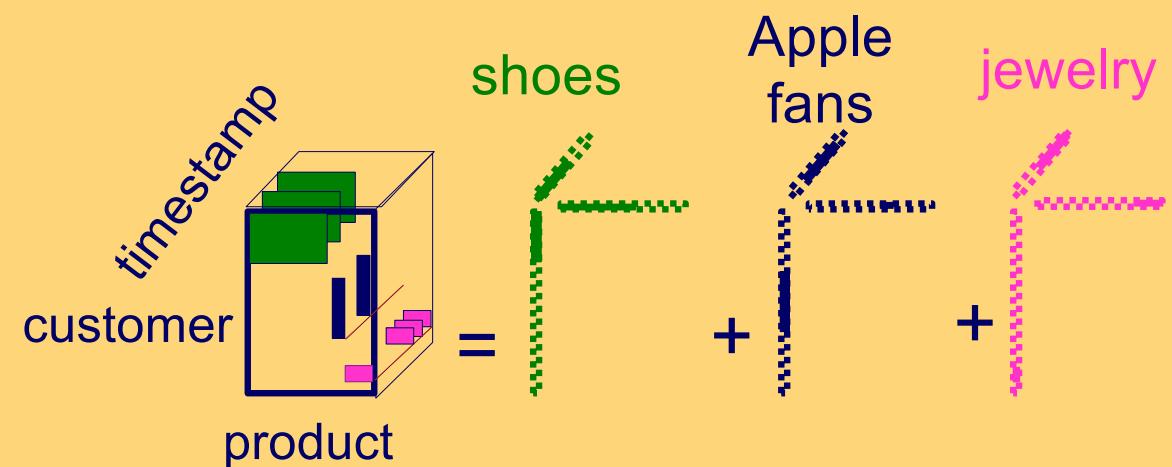
(b) FACEBOOK normal activity

- Modes: wall-owner, poster, timestamp
- Discovery: birthday-like event.



Conclusions (P1.5)

- Tensor analysis finds latent variables (market-segments, lockstep-groups, etc
 - Deviations → anomalies
- Extends SVD/factorization, to higher-modes



Tensors - More references

Tensor survey

- Tamara G. Kolda and Brett W. Bader
Tensor Decompositions and Applications
SIAM Rev., 51(3), pp 455–500, 2009.

Tensors - More references

Tensor survey #2

- Nicholas D. Sidiropoulos, Lieven De Lathauwer,, Xiao Fu,, Kejun Huang, Evangelos E. Papalexakis, and Christos Faloutsos, *Tensor Decomposition for Signal Processing and Machine Learning*, IEEE TSP, 65(13), July 1, 2017

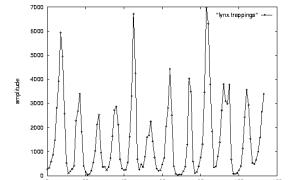
Part 1 - Outline



- Motivation
- P1.1. Similarity Search and Indexing
- P1.2. DSP
- P1.3. Linear Forecasting
- P1.4. Non-linear forecasting
- P1.5. Tensors
- • Conclusions

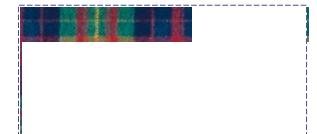


Important observations



Patterns, rules, forecasting and similarity indexing are closely related:

- To do forecasting, we need
 - to find patterns/rules
 - compress
 - to find similar settings in the past
- to find outliers, we need to have forecasts
 - (outlier = too far away from our forecast)



Forecasting Big Time Series: Theory and Practice



*Yuyang (Bernie) Wang, Christos Faloutsos, Valentin
Flunkert, Jan Gasthaus and Tim Januschowski*

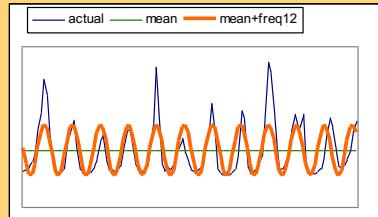


<https://lovvge.github.io/Forecasting-Tutorial-WWW-2020/>

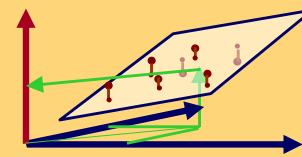
Part 1 - conclusions



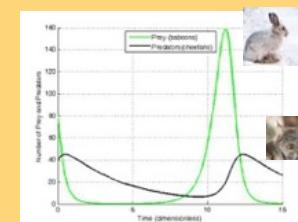
- P1.1. Similarity search: **Euclidean/time-warping; feature extraction and SAMs**
- P1.2. Periodicities: **DFT/DWT**
- P1.3. Linear Forecasting: **AR** (Box-Jenkins)
- P1.4 Non-linear forecasting: **lag-plots**
 - Gray-box modeling: **Lotka-Volterra**
- P1.5. Tensors: **PARAFAC**



WWW 2020



Wang, Faloutsos, Flunkert, Gasthaus,
Januschowski



209

Part 1 - conclusions



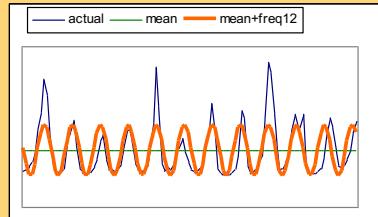
- P1.1. Similarity search: Euclidean/time-warping; **feature extraction** and SAMs

- P1.2. Periodicities: DFT/DWT

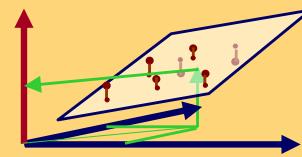
1) PLOT IT
Linear E
Forecasting: AR (Box-Jenkins)

2) DFT / DWT
Non-linear forecasting: lag-plots
Gray-box modeling: Lotka-Volterra

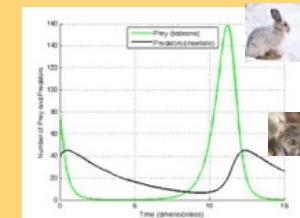
- P1.5. Tensors: PARAFAC



WWW 2020



Wang, Faloutsos, Flunkert, Gasthaus,
Januschowski



210