

TIME SERIES DATA MINING

Ing. Petr Novák
novak@infoway.cz

Outline of talk

- Time series
- Similarity measures
- Time series representation
- Motif discovery
- TSMiner

Time series

What is time series...

- Popular statistical definition

„The time series is a sequence of values of a statistical character (indicators) arranged in terms of time away from the past to the present. The indicator changes over time, the overall development can be divided into three components - trend, periodic and random.“

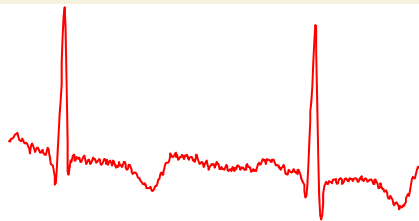
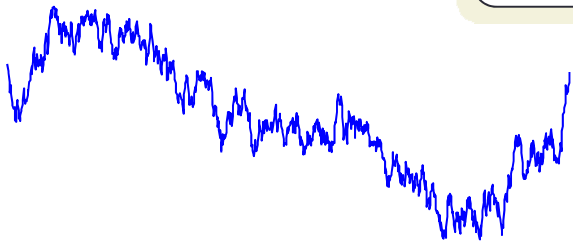
- **Is it the only option?**

Time series

People measure things...

- *Blood pressure*
- *Popularity of politicians*
- *Annual rainfall*
- *Stock value*

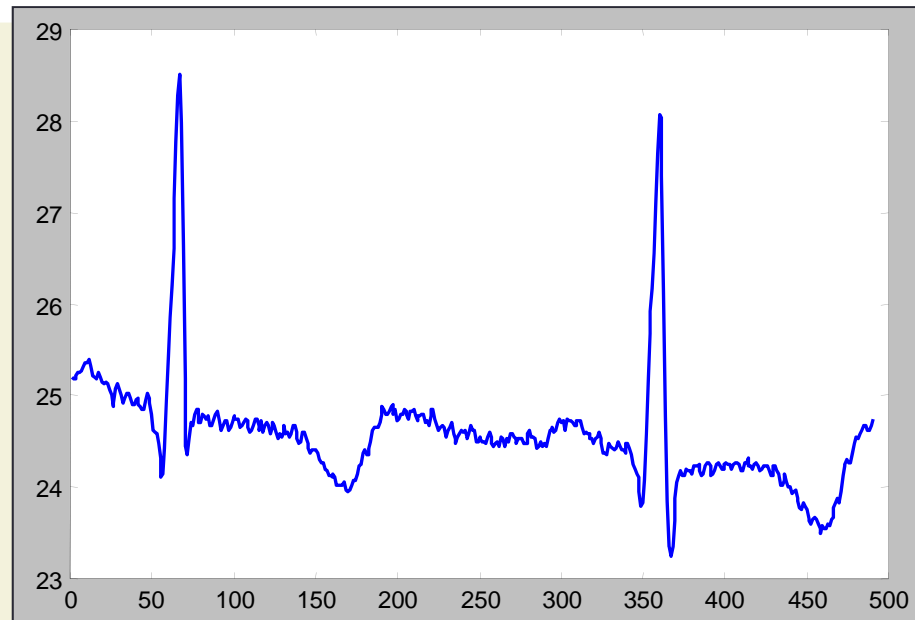
...and things change over time.



Time series

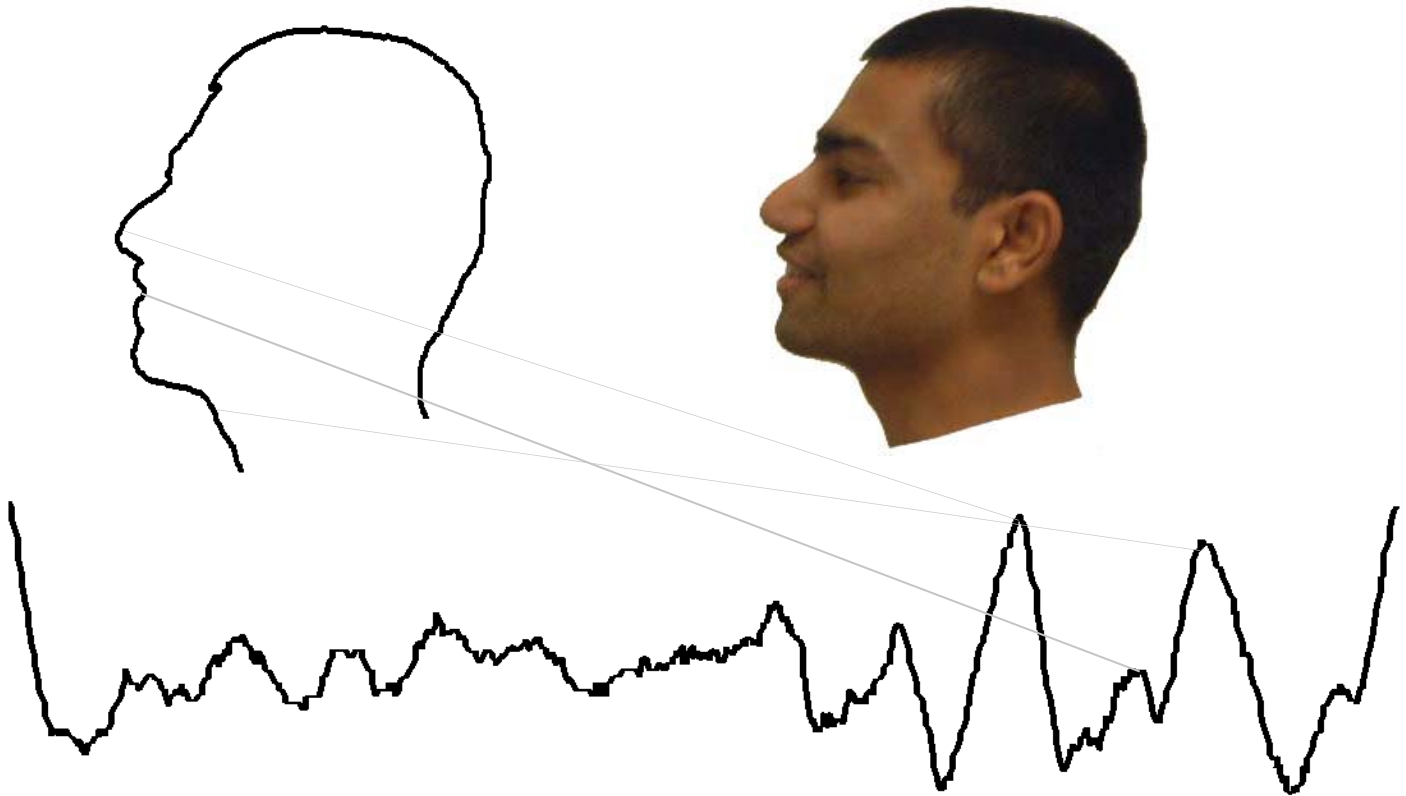
What is time series...

- „Time series is a collection of observations made sequentially in time.“



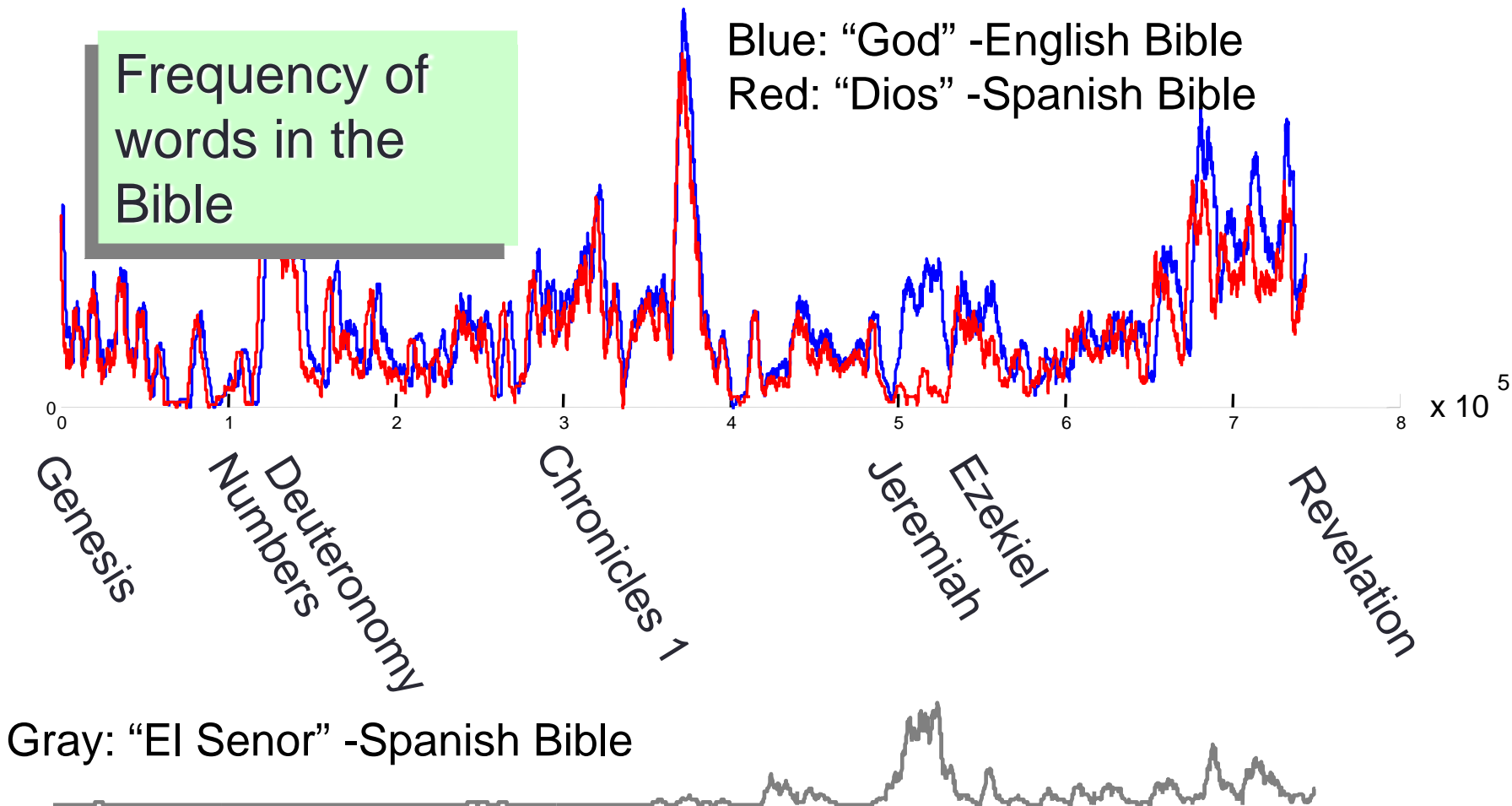
Time series

- Image data, may be thought of as time series



Time series

- Text data, may best be thought of as time series...



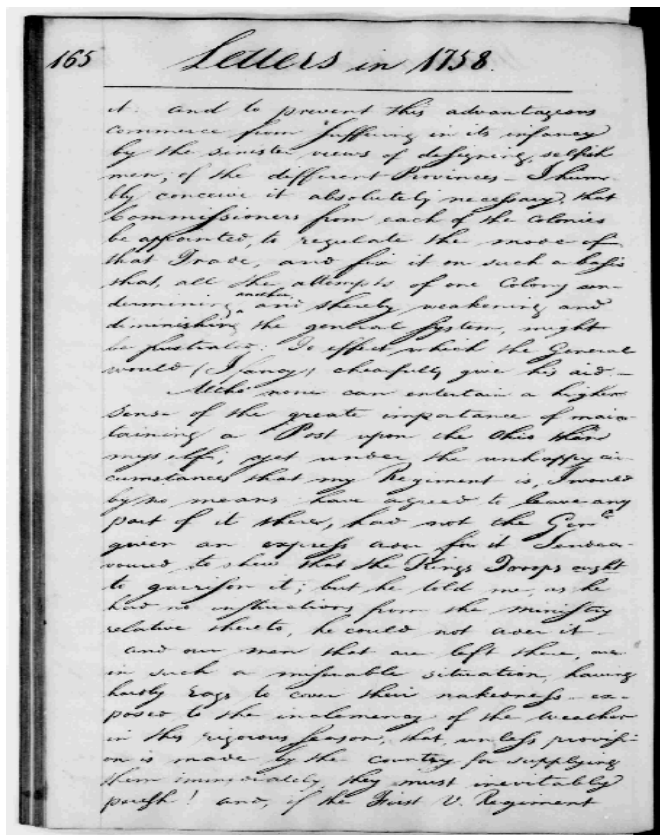
Time series

- Video data, may be thought of as time series

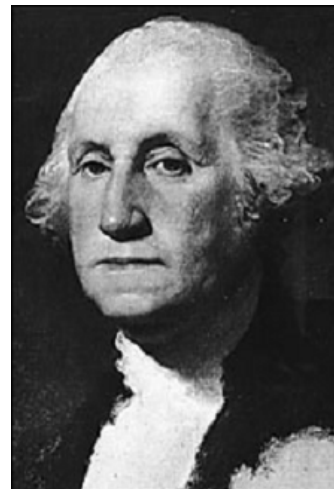


Time series

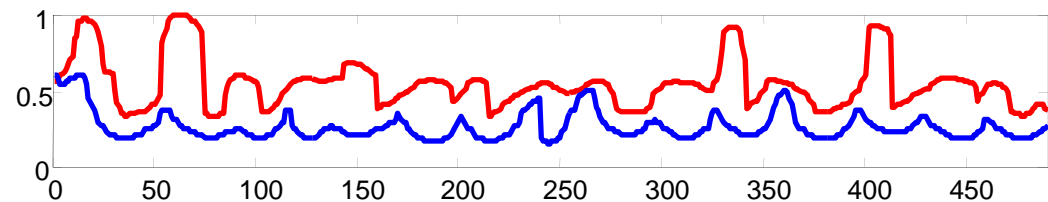
- Handwriting data, may be thought of as time series



G. Washington manuscript



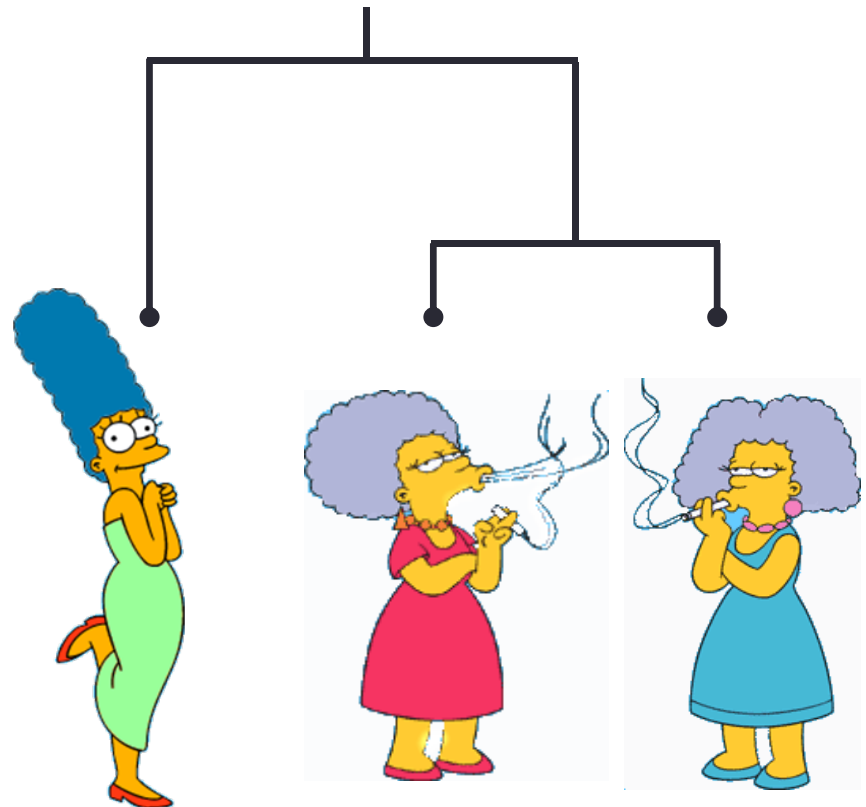
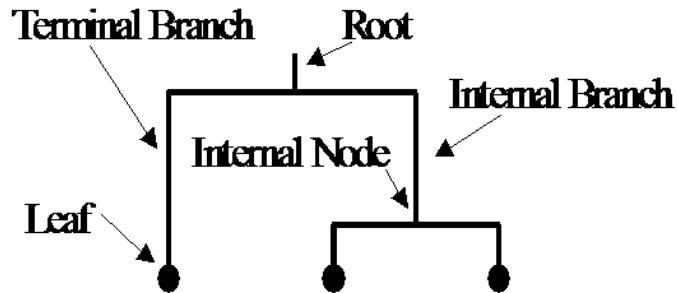
George Washington
1732-1799



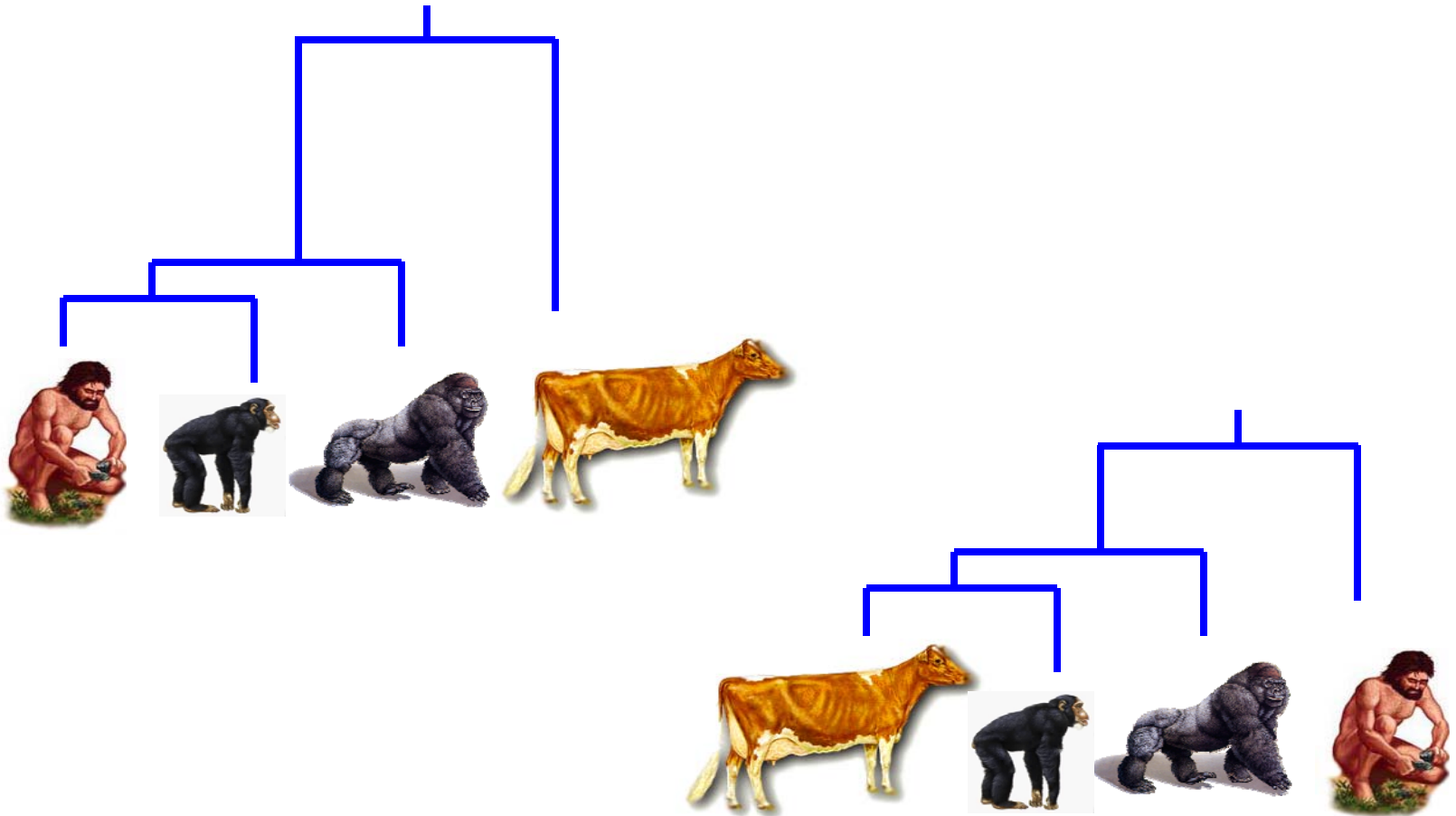
Small digression - dendrogram

- A Useful Tool for Summarizing Similarity Measurements

The similarity between two objects in a dendrogram is represented as the height of the lowest internal node they share.



Small digression - dendrogram

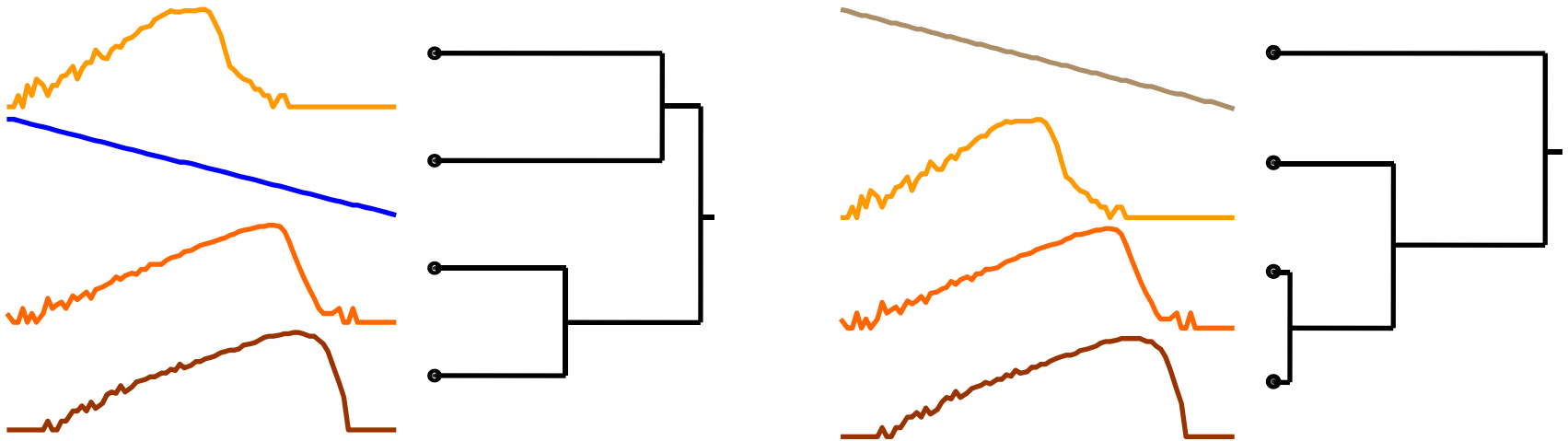


Time series

- Why is working with time series so difficult?
 - 1 hour of EKG – 1 gigabyte.
 - Space shuttle telemetry – 20 000 sensors send data every second, hundreds of GB per regular mission.
 - Database Macho – astronomical observations 20 millions stars, 3 terabytes per day.
 - 300 millions phone calls in AT&T network every day between 100 millions of customers.
 - 50 000 stock titles in USA, 100 000 trades per second.
- **We need a data representation for efficient processing**

Time series

- Why is working with time series so difficult?
 - We are dealing with subjectivity



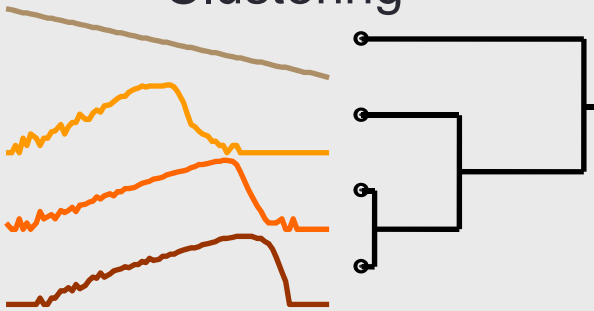
- We need to measure the similarity of time series regardless of the subjective feeling

Time series

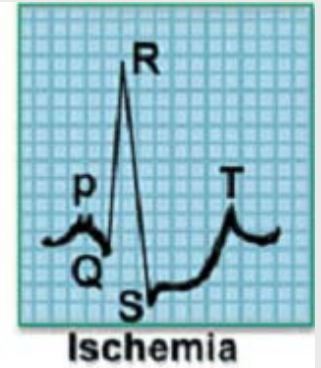
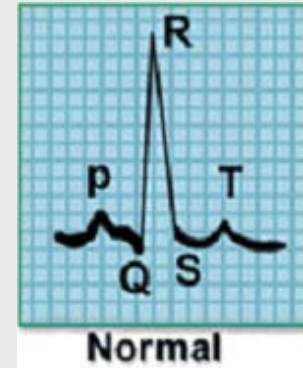
- Why is working with time series so difficult?
 - Different scales
 - Differing sampling rates
 - Noise, missing values

Time series

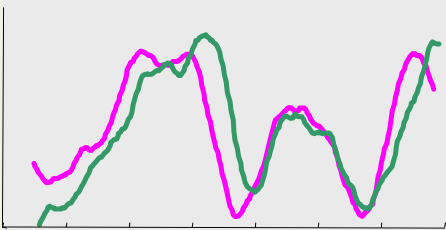
Clustering



Classifications

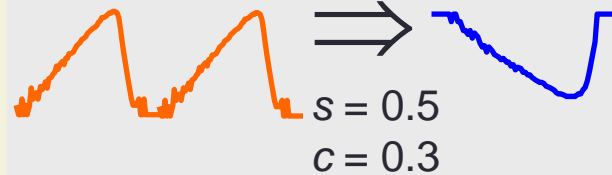


Motif discovery

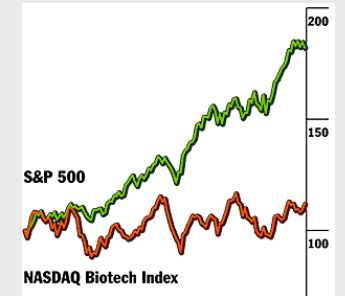


Rule discovery

10



Query by Content



Visualization



Novelty detection



Time series

- For all these activities we need to determine the **similarity** between time series.

Time series similarity

- Similarity is hard to define, but we know it when we see it



Time series similarity

- We have to define **distance measure**
- **Definition:** Let O_1 and O_2 are two objects from the universe of possible objects. The distance (dissimilarity) is denoted by $D(O_1, O_2)$
- Properties of distance measure
 - $D(A, B) = D(B, A)$, symmetry
 - $D(A, A) = 0$, identity
 - $D(A, B) = 0$ iff $A = B$ positivity
 - $D(A, B) \leq D(A, C) + D(B, C)$ triangular inequality

Time series similarity

$$D(A,B) = D(B,A)$$

Symmetry

$$D(\text{Patty}, \text{Selma}) = D(\text{Selma}, \text{Patty})$$

Otherwise you could say:



Patty looks like
Selma, but Selma
does not look like
Patty!

Time series similarity

$$D(A,A) = 0$$

Identity

$$D(\text{Patty}, \text{Patty}) = 0$$



Otherwise you could say:



Marge looks more like Patty than Patty does!

Time series similarity

$D(A,B) = 0 \text{ IIF } A = B$ *Positivity*

$D(\text{Patty}, \text{Patty}) = 0, \text{ IIF } \text{Patty} = \text{Patty}$

Otherwise you could say:

I know Patty and Marge are somehow different, but I can't tell them apart!



Time series similarity

$D(A,B) \leq D(A,C) + D(B,C)$ *Triangular inequality*

$$D(\text{Marge}, \text{Patty}) \leq D(\text{Marge}, \text{Selma}) + D(\text{Patty}, \text{Selma})$$

Otherwise you could say:

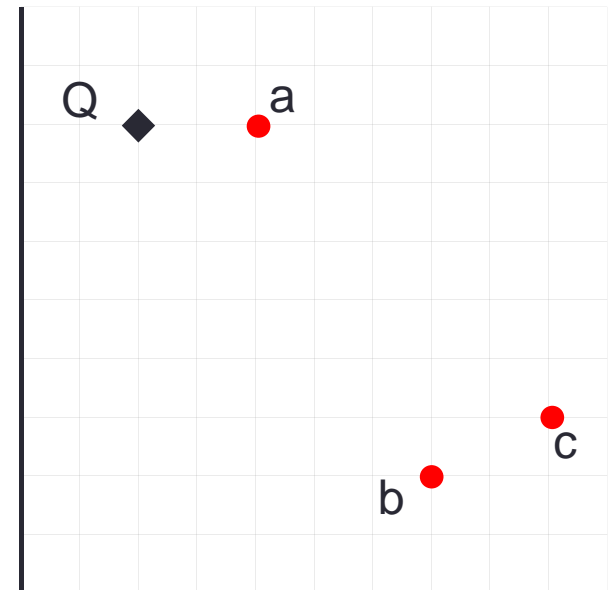


Patty looks like Marge,
Selma also looks like
Marge, but Patty looks
nothing like Selma!

Time series similarity

- Importance of triangular inequality:
- We are looking for the closest point to Q, in a database of 3 objects. Suppose that the triangular inequality holds, and that we have precompiled a table of distance between all the items in the database.

	a	b	c
a		6.70	7.07
b			2.30
c			



Time series similarity

- Importance of triangular inequality:
- We calculate that a is 2 units from Q, b is 7,81 from Q. We don't have to calculate distance from c to Q!

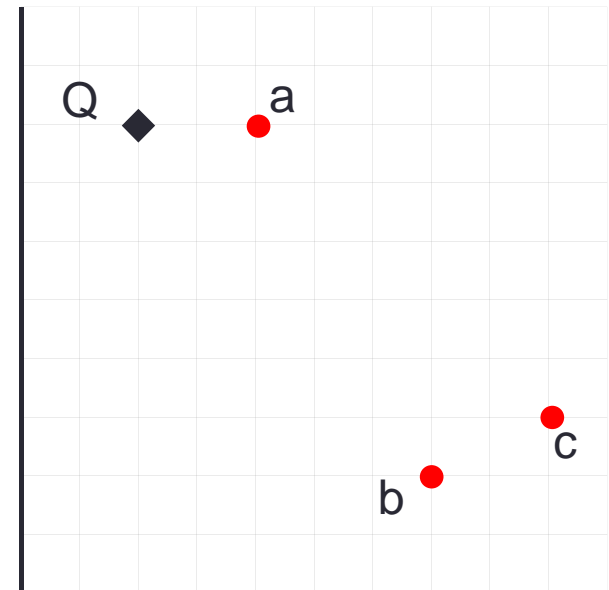
$$D(Q, \mathbf{b}) \leq D(Q, \mathbf{c}) + D(\mathbf{b}, \mathbf{c})$$

$$D(Q, \mathbf{b}) - D(\mathbf{b}, \mathbf{c}) \leq D(Q, \mathbf{c})$$

$$7.81 - 2.30 \leq D(Q, \mathbf{c})$$

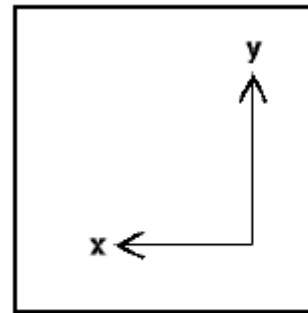
$$5.51 \leq D(Q, \mathbf{c})$$

	a	b	c
a		6.70	7.07
b			2.30
c			

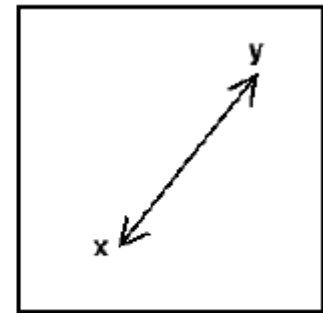


Time series similarity

- So we are looking for a suitable distance measure between two series
- Frequently used measures of similarity are based on a comparison of the overall shape of time-series
- **Minkowski metrics**



Manhattan



Euclidean

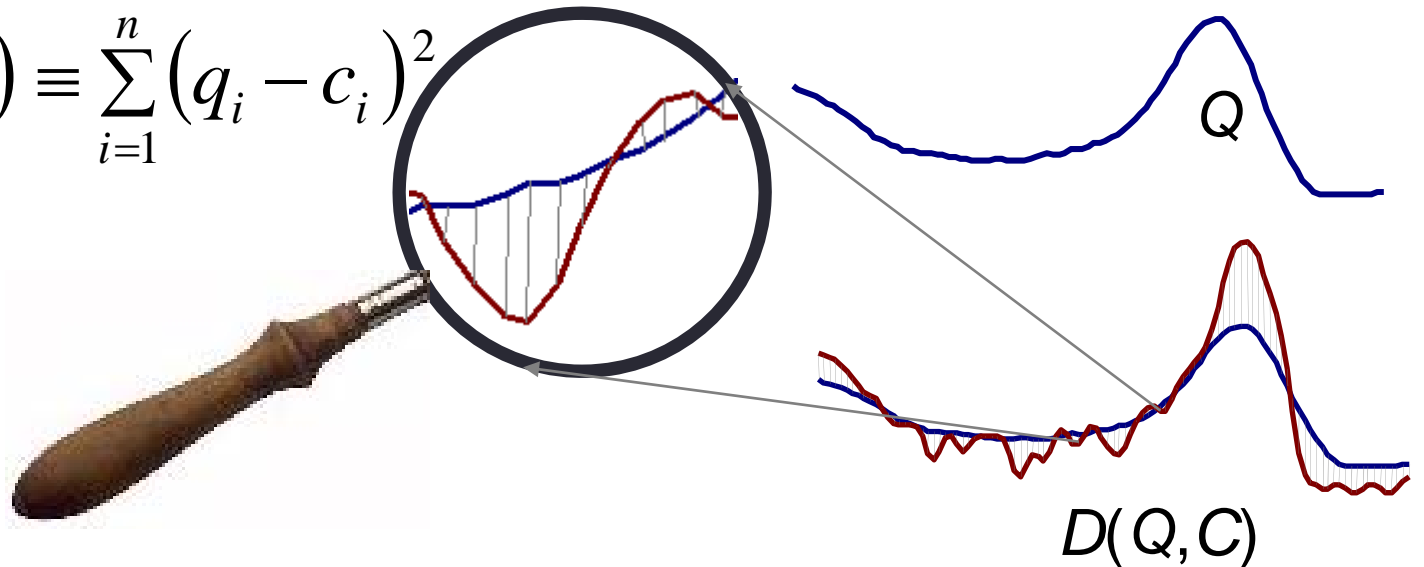
- Frequently used is not always the best

Time series similarity

- **Euclidean Distance Metric**
- Let Q and C are time series

$$D(Q, C) \equiv \sqrt{\sum_{i=1}^n (q_i - c_i)^2}$$

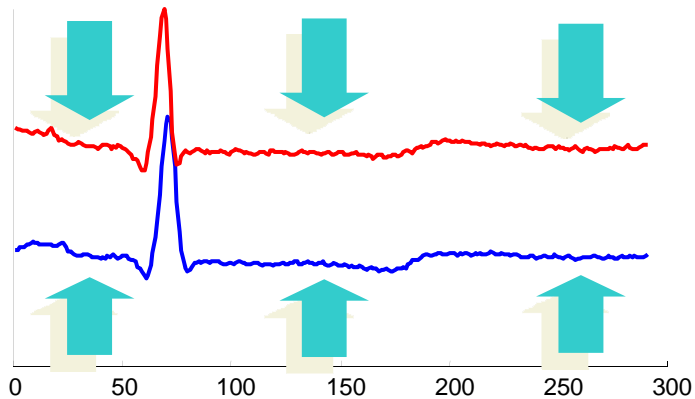
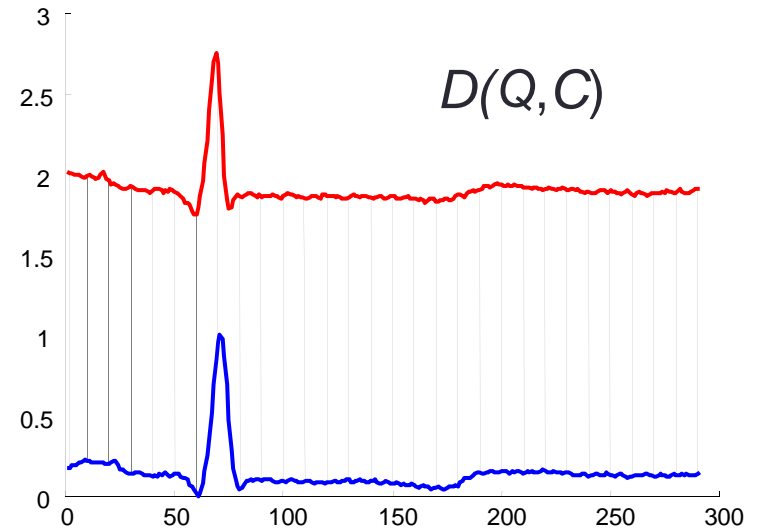
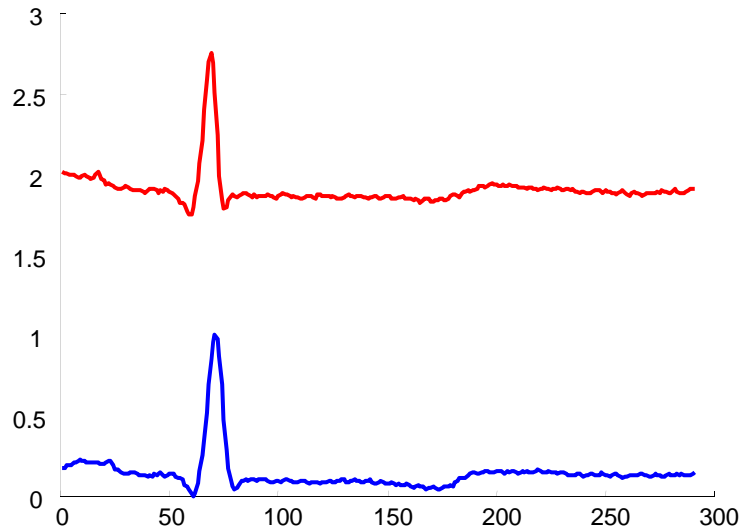
$$D_{squared}(Q, C) \equiv \sum_{i=1}^n (q_i - c_i)^2$$



Time series similarity

- In most cases, the Euclidean metric does not give suitable results, due to "misrepresentation" in the source data.
 - Offset translation
 - Amplitude scaling
 - Linear trend
 - Noise

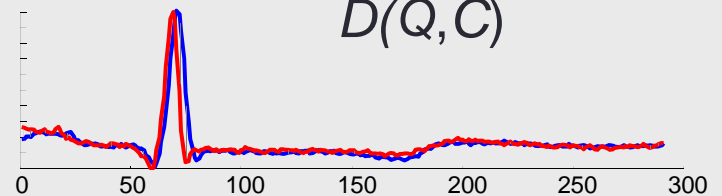
Time series similarity



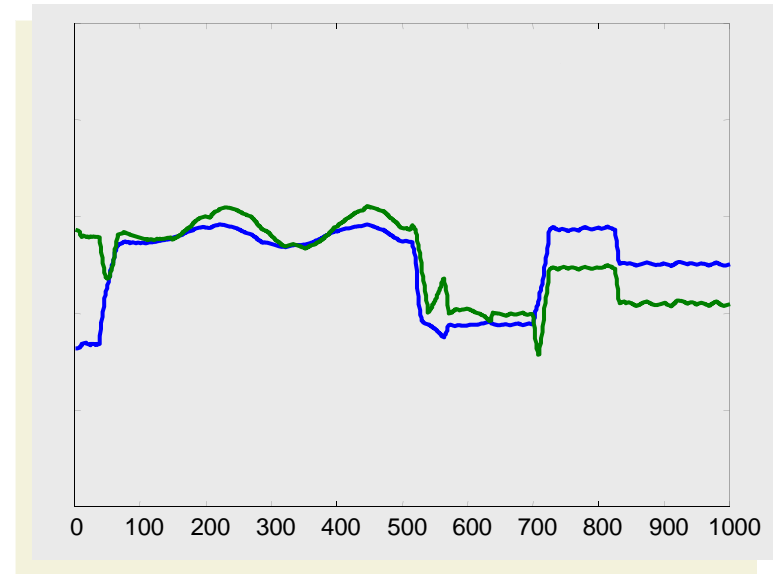
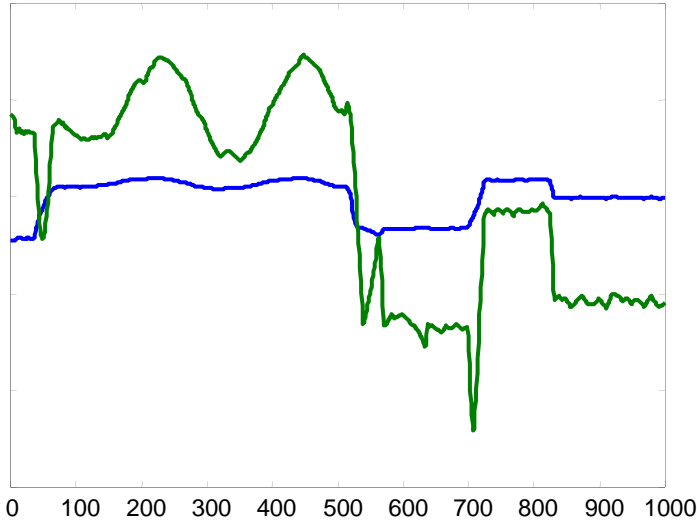
$$Q = Q - \text{mean}(Q)$$

$$C = C - \text{mean}(C)$$

$$D(Q, C)$$



Time series similarity

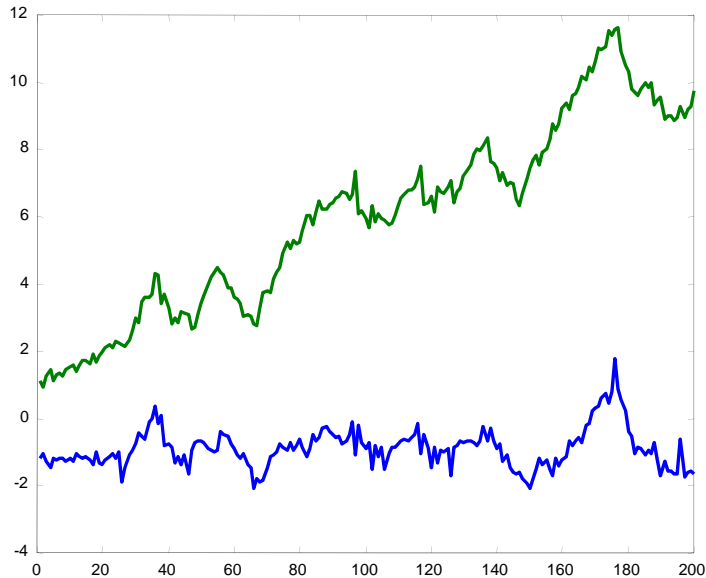


$$Q = (Q - \text{mean}(Q)) / \text{std}(Q)$$

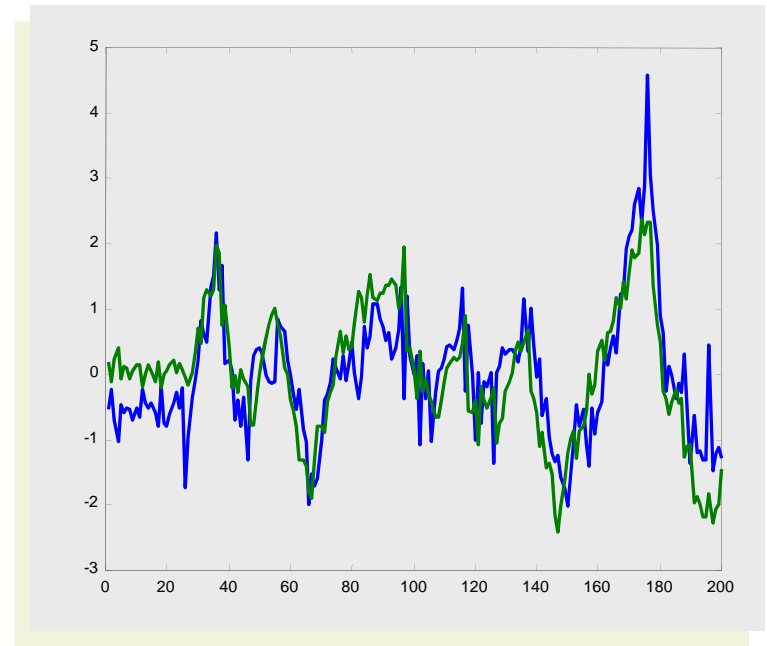
$$C = (C - \text{mean}(C)) / \text{std}(C)$$

$$D(Q, C)$$

Time series similarity

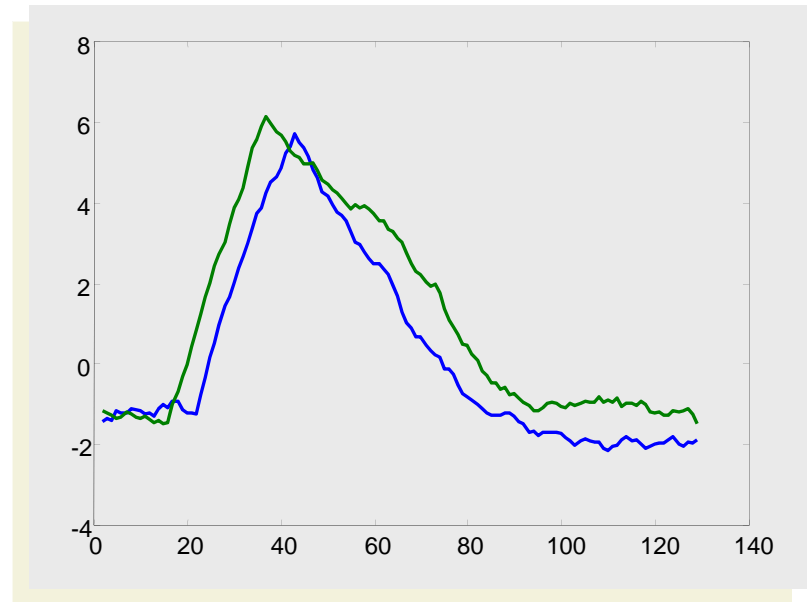
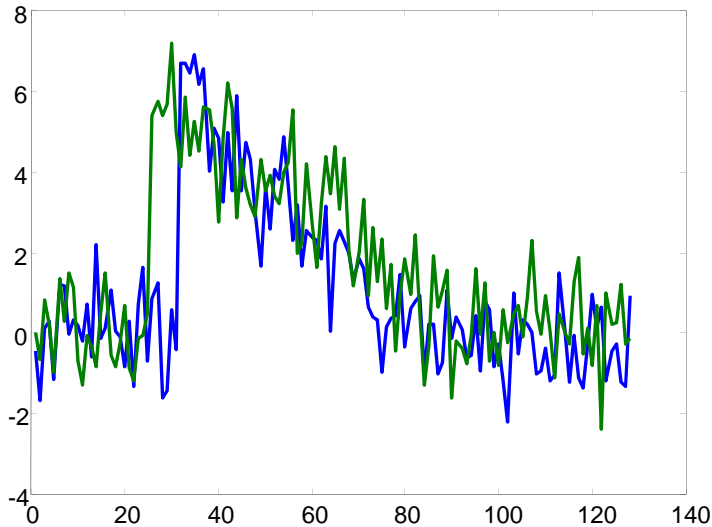


Find best linear approximation, then subtract that line from time series.



Removed linear trend
Removed offset translation
Removed amplitude scaling

Time series similarity



Replacing points
with average of
their neighbors

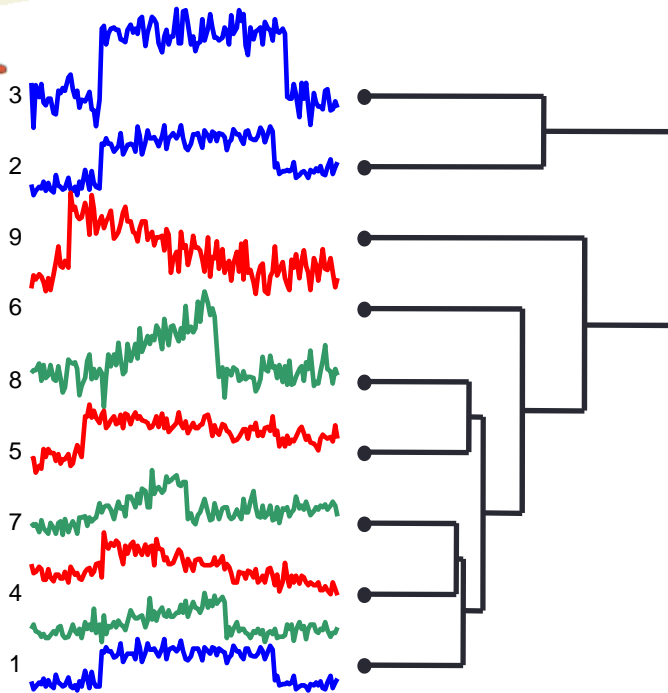
$$Q = \text{smooth}(Q)$$

$$C = \text{smooth}(C)$$

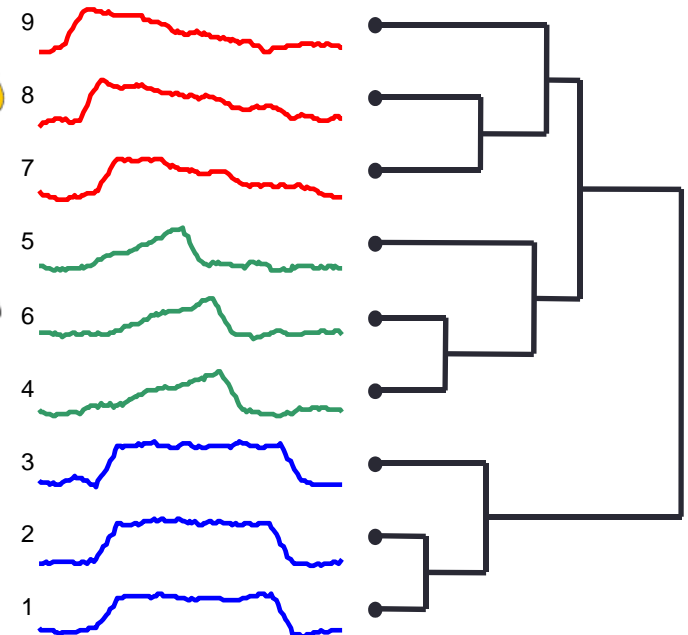
$$D(Q, C)$$

Time series similarity

Clustered using Euclidean distance on the raw data.



Clustered using Euclidean distance, after removing noise, linear trend, offset translation and amplitude scaling.



Time series similarity

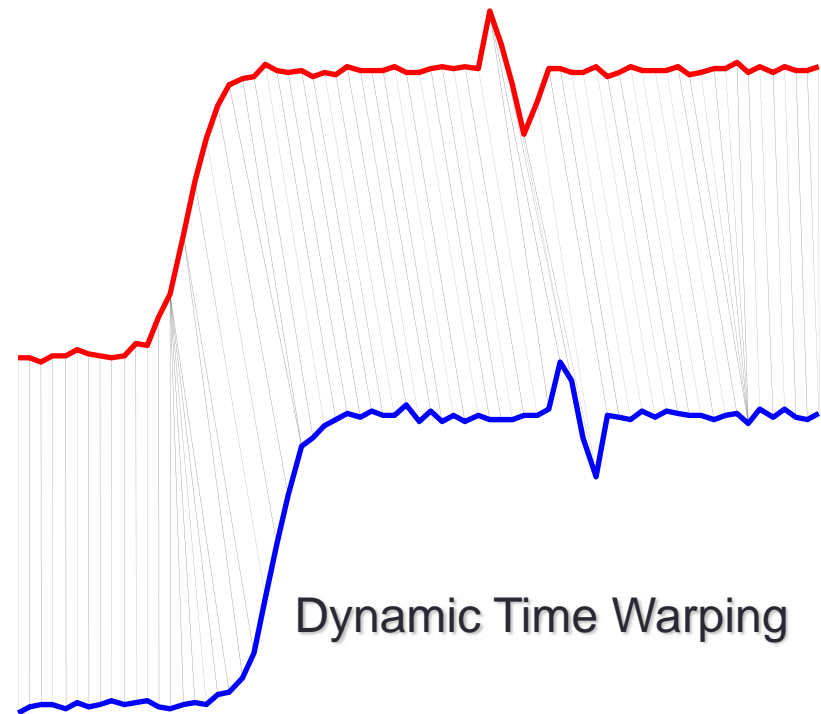
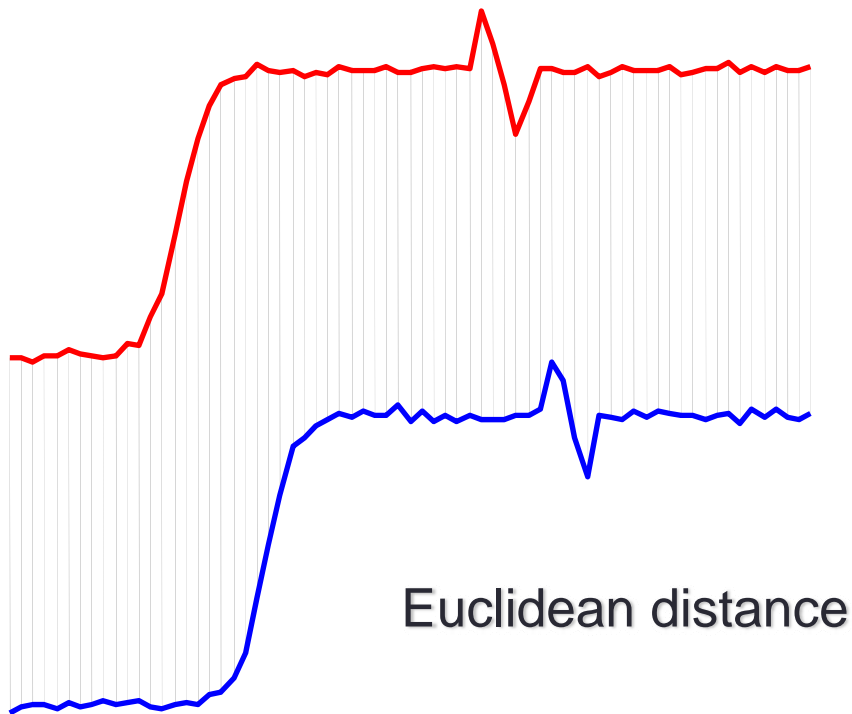
- To remember – the “raw” time series may have distortions which we should remove before clustering, classification etc.

BUT

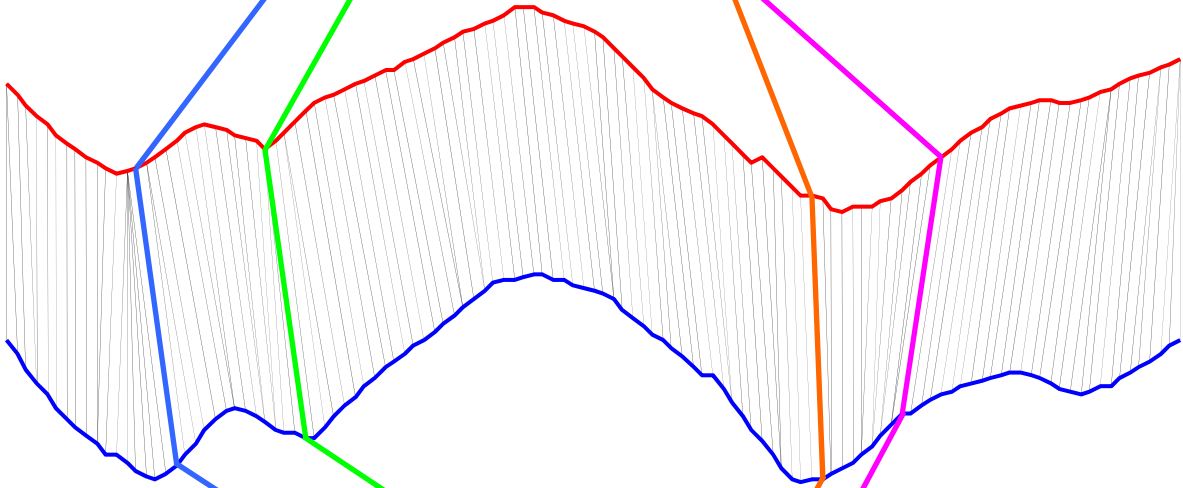
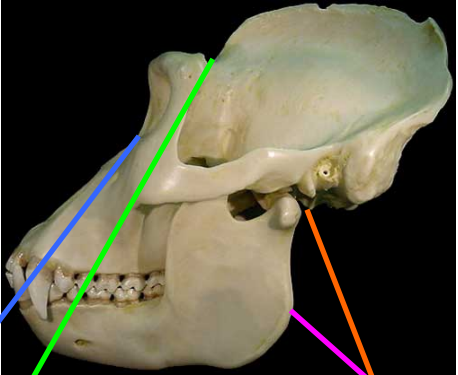
- Sometimes the distortions are the most interesting thing about the data!

Time series similarity

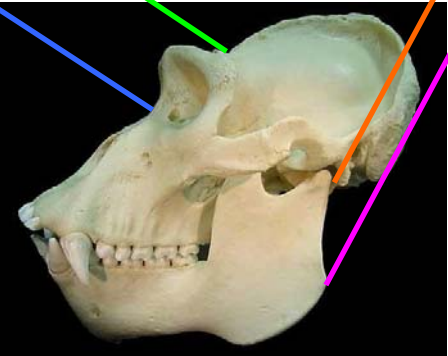
- Dynamic Time Warping (DTW)
- One method to deal with a phase shift between time series

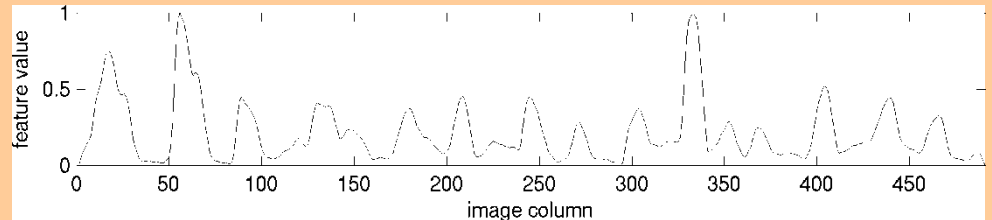
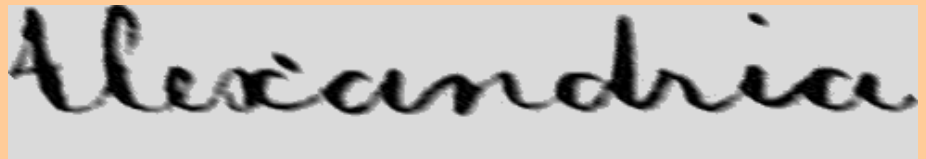
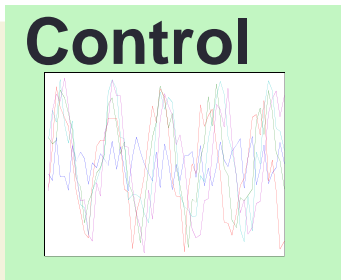
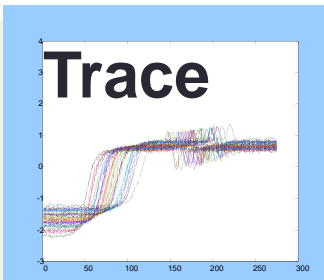
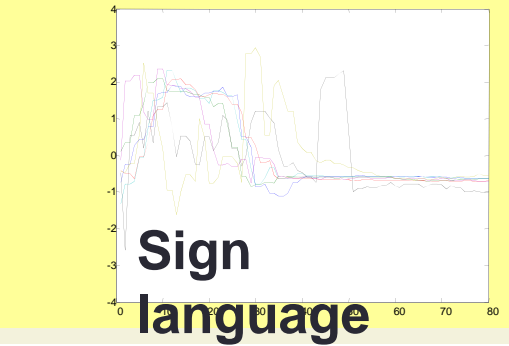
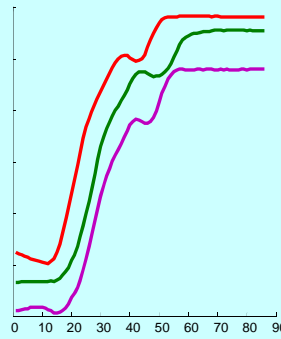
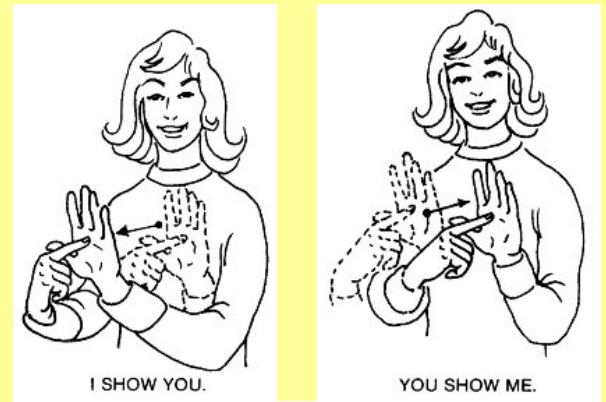
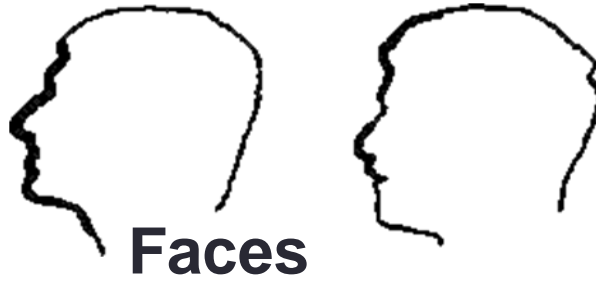


Lowland Gorilla
Gorilla gorilla graueri



Mountain Gorilla
Gorilla gorilla beringei

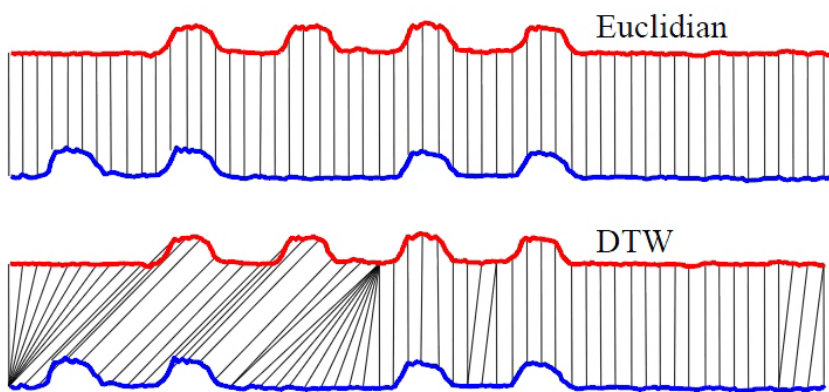




Word Spotting

Time series similarity

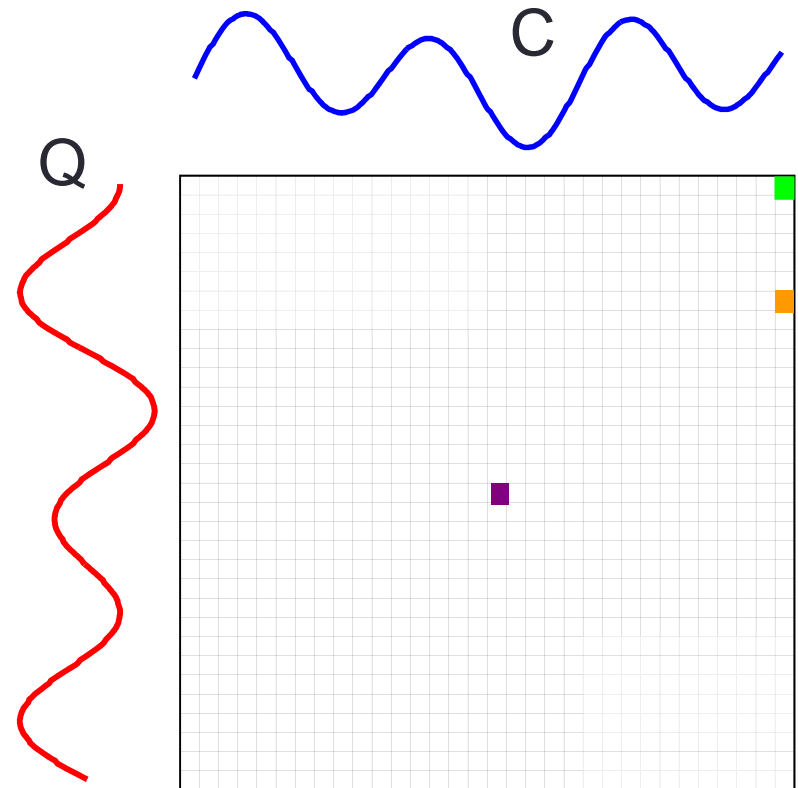
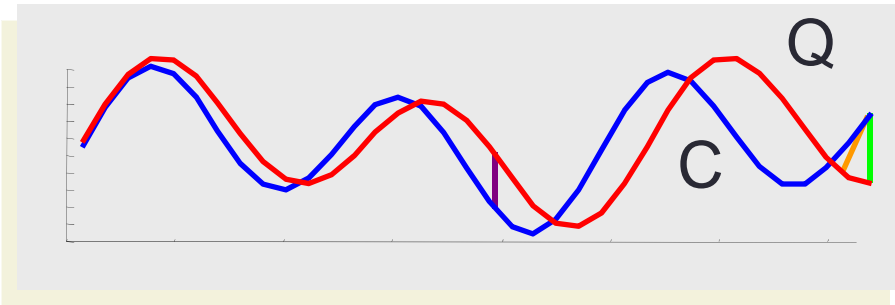
DTW is two to three orders of magnitude slower than Euclidean Distance (time in msec)



Dataset	Euclidean	DTW
Word Spotting	40	8,600
Sign language	10	1,110
GUN	60	11,820
Nuclear Trace	210	144,470
Leaves	150	51,830
(4) Faces	50	45,080
Control Chart	110	21,900
2-Patterns	16,890	545,123

Time series similarity

- We create a matrix the size of $|Q|$ by $|C|$, then fill it in with the distance between every pair of point in our two time series.

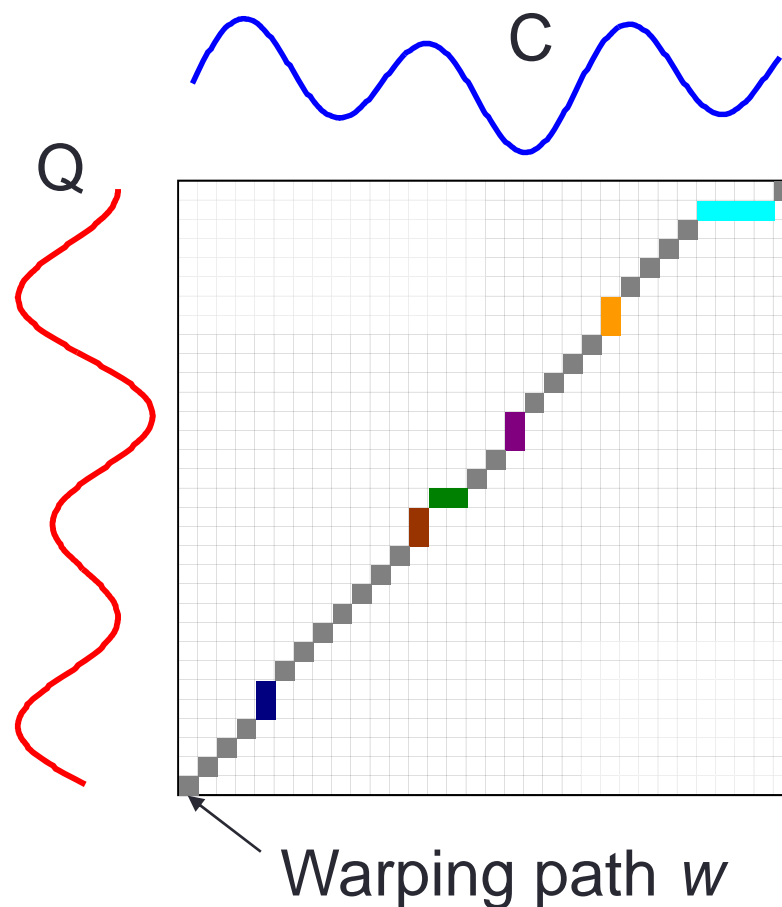


Time series similarity

- Every possible warping between two time series, is a path through the matrix. We want the best one.

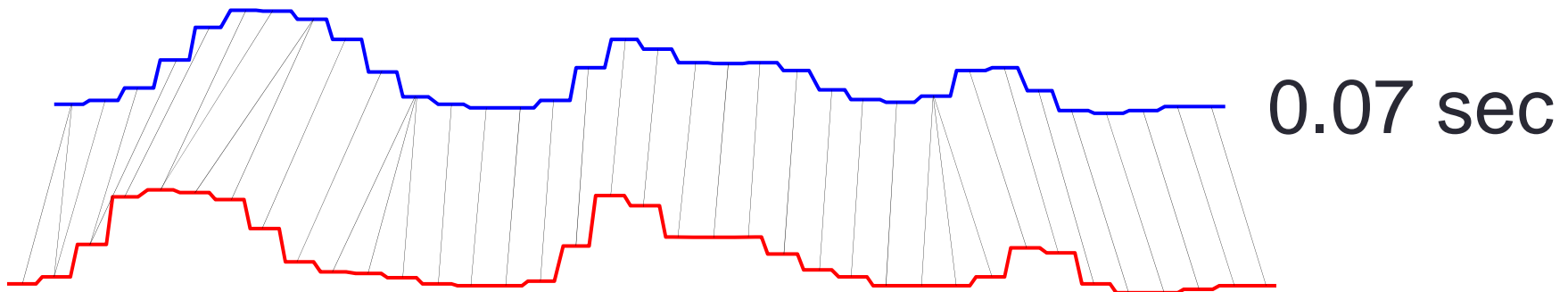
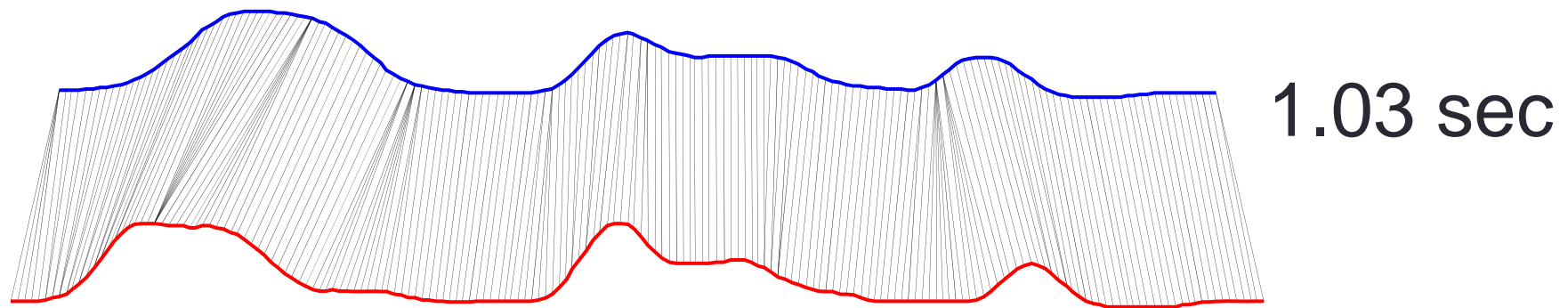
$$DTW(Q, C) = \min \left\{ \sqrt{\sum_{k=1}^K w_k} / K \right.$$

- Dynamic Time Warping gives much better results than Euclidean distance, but DTW is very slow to calculate!



Time series similarity

- Simple idea - approximate the time series with some compressed or downsampled representation, and do DTW on the new representation.



Time series similarity

- In general, it's hard to speed up a single DTW calculation
- However, if we have to make many DTW calculations (which is almost always the case), we can potentially speed up the whole process by lowerbounding.

- $DTW(A,B)$

The true DTW function is very slow...

- $lower_bound_distance(A,B)$

The lower bound function is very fast...

$$lower_bound_distance(A,B) \leq DTW(A,B)$$

Time series similarity

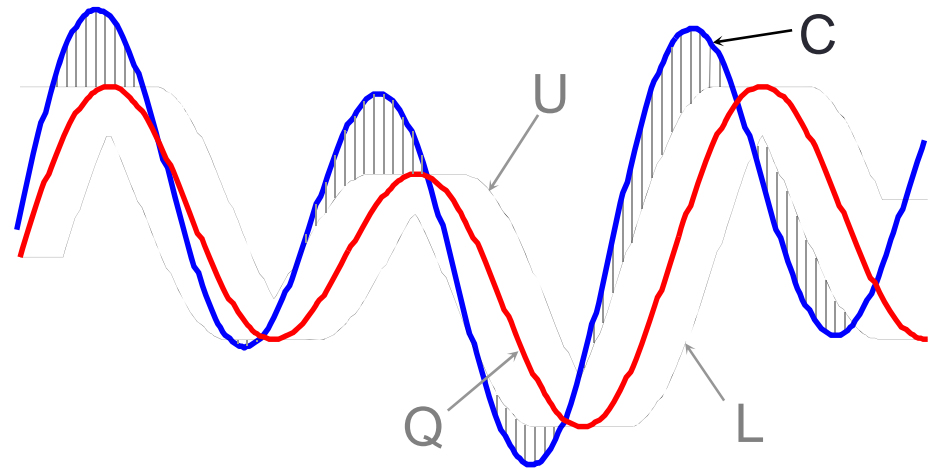
- Lowerbounding – distance of points in space is less than or equal to the actual distance

Algorithm Lower_Bounding_Sequential_Scan(Q)

```
1.  best_so_far ← infinity;
2.  for all sequences in database
3.    LB_dist = lower_bound_distance( $C_i$ , Q);
4.    if LB_dist < best_so_far
5.      true_dist = DTW( $C_i$ , Q);
6.      if true_dist < best_so_far
7.        best_so_far ← true_dist;
8.        index_of_best_match ← i;
9.      endif
10.   endif
11. endfor
```

Time series similarity

Envelope - Based Lower Bound



 **LB_Keogh**

$$LB_Keogh(Q, C) = \sum_{i=1}^n \begin{cases} (q_i - U_i)^2 & \text{if } q_i > U_i \\ (q_i - L_i)^2 & \text{if } q_i < L_i \\ 0 & \text{otherwise} \end{cases}$$

Time series – data preparation

Data cleansing

- Missing values
 - Skip values
 - Value estimation
 - Linear interpolation
- Noise reduction
 - Binning
 - Moving average

Data normalization

- Transforming data into the same range

- Min-max

$$x_{norm} = \frac{x - x_{min}}{x_{max} - x_{min}}$$

- Z-score

$$x_{norm} = \frac{x - \mu}{\sigma}$$

Time Series Representations

- We have already told how to define similarity, but how find it quickly?
- Generic Data Mining Algorithm:
 - Create an approximation of the data, which will fit in main memory, yet retains the essential features of interest.
 - Approximately solve the problem at hand in main memory
 - Make few accesses to the original data on disk to confirm the solution obtained in step 2
- **But which approximation should we use?**

Time Series Representations

Time Series Representations

Model Based

Hidden Markov Models
Statistical Models

Data Adaptive

Sorted Coefficients
Piecewise Polynomial
Singular Value Approximation
Symbolic
Trees

Piecewise Linear Approximation
Interpolation
Regression

Adaptive Piecewise Constant Approximation

Natural Language
Strings
Symbolic Aggregate Approximation
Non Lower Bounding
Value Based
Slope Based

Non Data Adaptive

Wavelets
Random Mappings
Spectral
Piecewise Aggregate Approximation

Orthonormal
Haar
Daubechies dbn n > 1

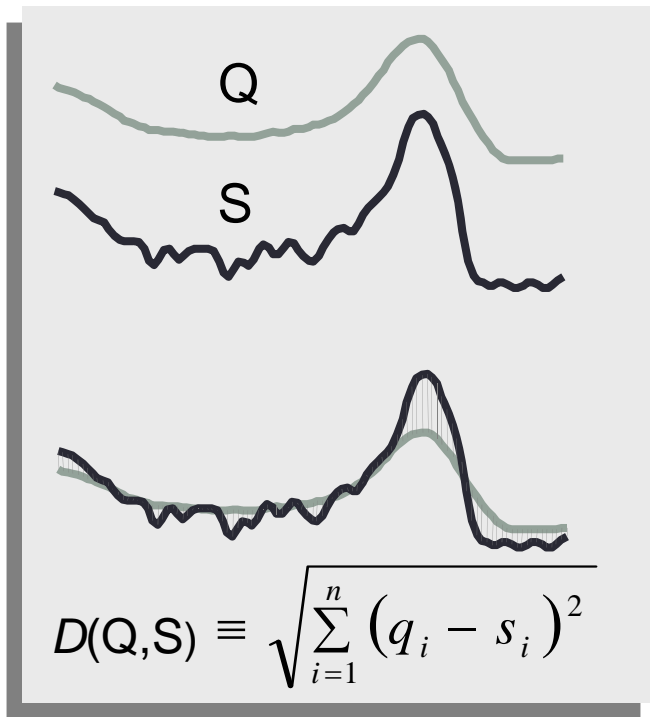
Bi-Orthonormal
Coiflets
Symlets

Discrete Fourier Transform
Discrete Cosine Transform
Chebyshev Polynomials

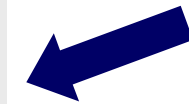
Data Dictated

Grid
Clipped Data

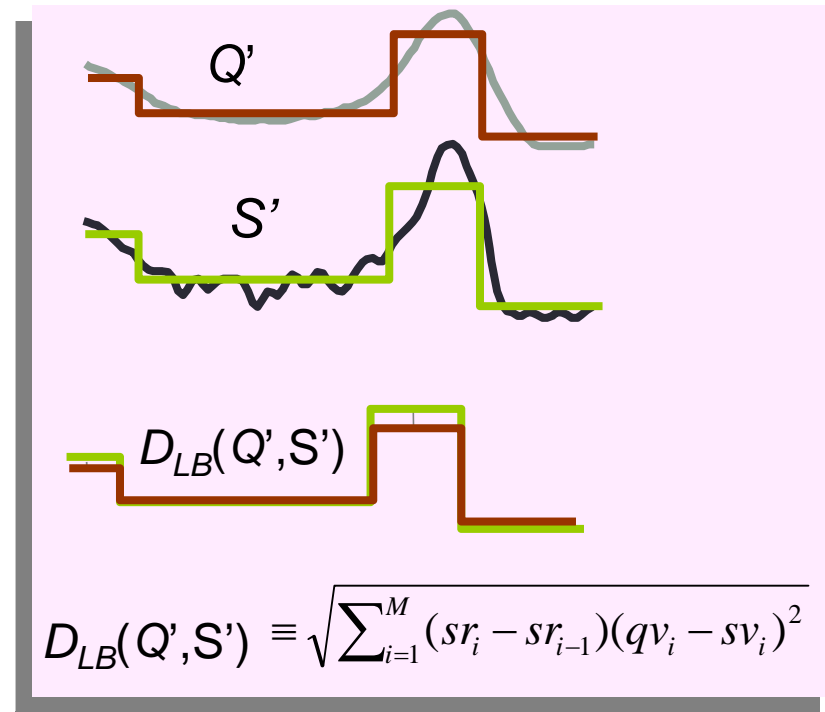
Time Series Representations



Raw data



Approximation

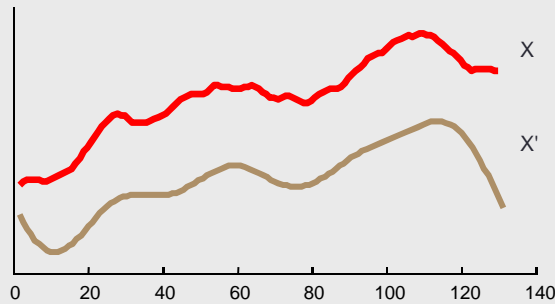


For every two time series Q and S approximation has to allow lower bounding

$$D_{LB}(Q',S') \leq D(Q,S)$$



Discrete Fourier Transform



Represent the time series as a linear combination of sines and cosines.

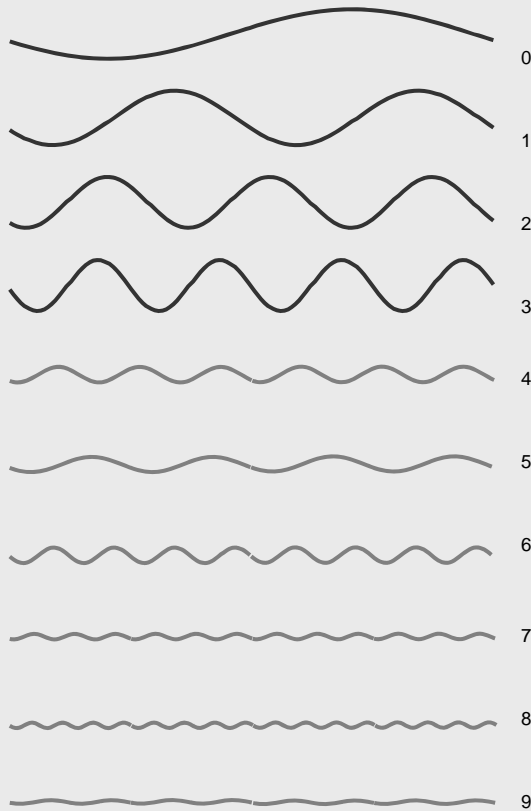


Jean Fourier

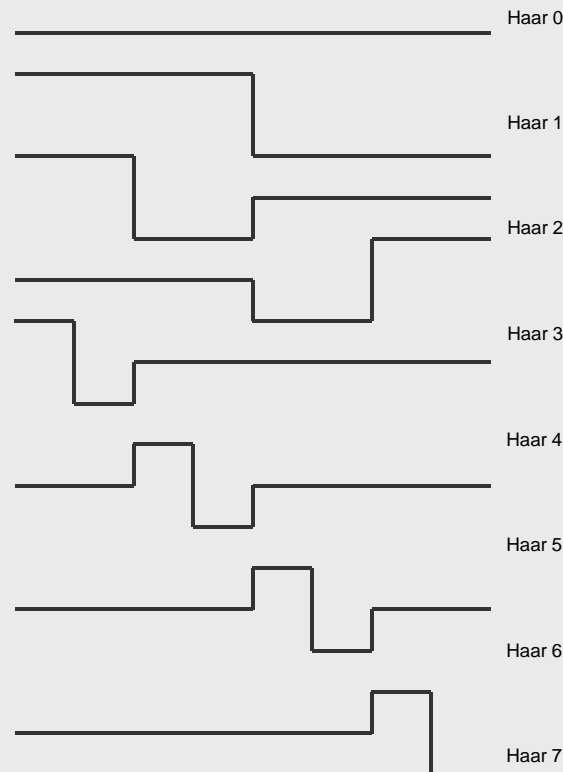
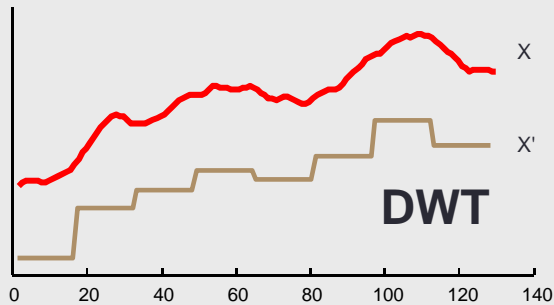
1768-1830

$$C(t) = \sum_{k=1}^n (A_k \cos(2\pi w_k t) + B_k \sin(2\pi w_k t))$$

- Good ability to compress most natural signals, $O(n \cdot \log(n))$.
- Difficult to deal with sequences of different lengths.



Discrete Wavelet Transform



Represent the time series as a linear combination of Wavelet basis functions.

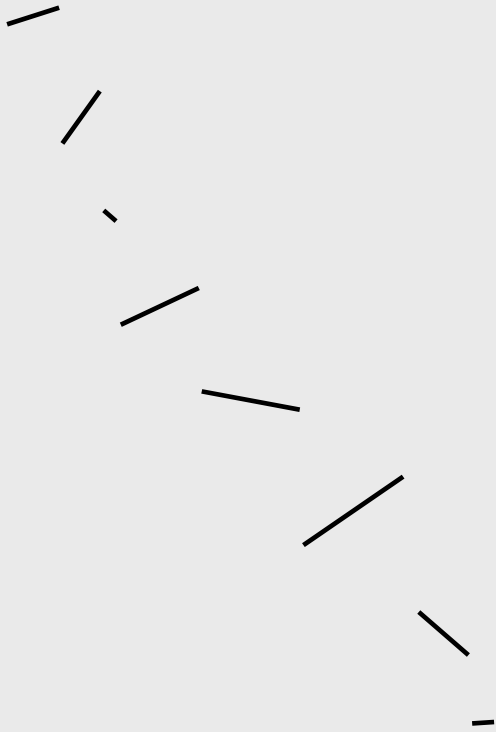
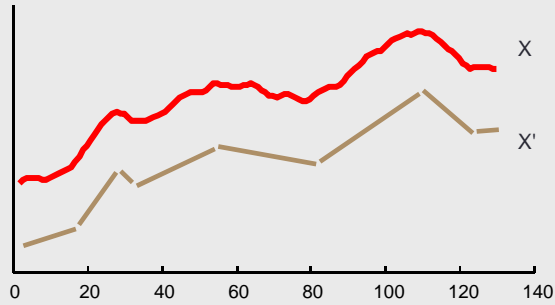


Alfred Haar

1885-1933

- Good ability to compress
- Fast linear time algorithms for DWT
- Able to support some interesting non-Euclidean similarity measures
- Signal must have a length $n = 2^{\text{int}}$

Piecewise Linear Approximation

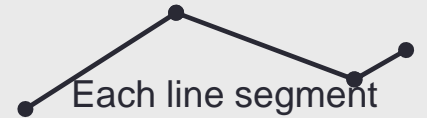


Represent the time series as a sequence of straight lines. Lines could be **connected** ($N/2$) or **disconnected** ($N/3$). Series is replaced by segments, their number is much smaller than no of points in the original series.

- $O(n^2N)$, which is too slow for data mining
- Faster heuristic solutions
 - Top-Down
 - Bottom-Up
 - Sliding Window
- Not suitable for indexing



Karl Friedrich Gauss
1777 - 1855



Each line segment has

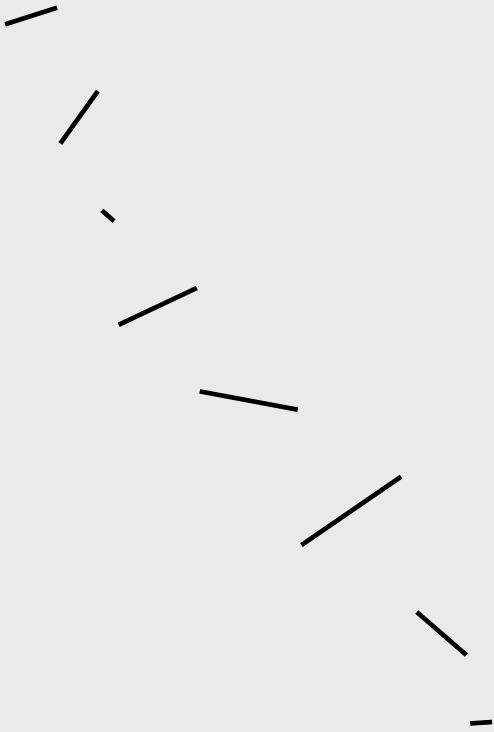
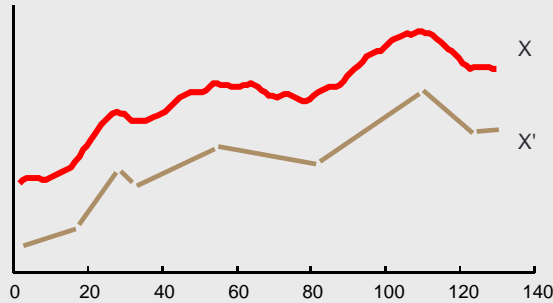
- length
- left_height (right_height can be inferred by looking at the next segment)



Each line segment has

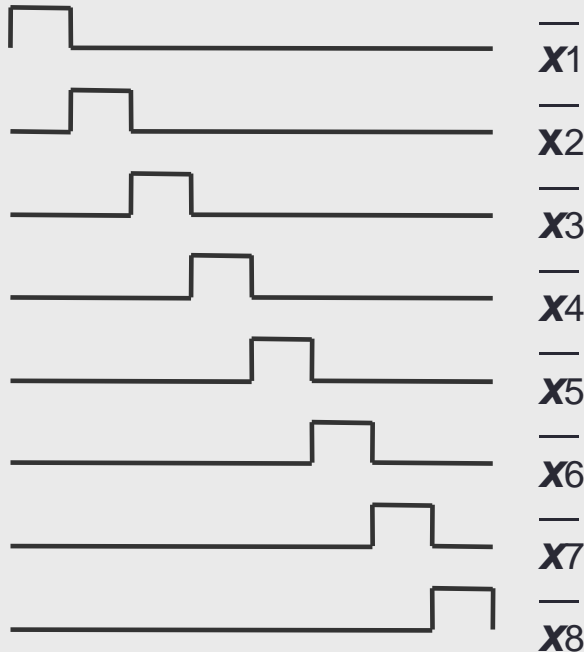
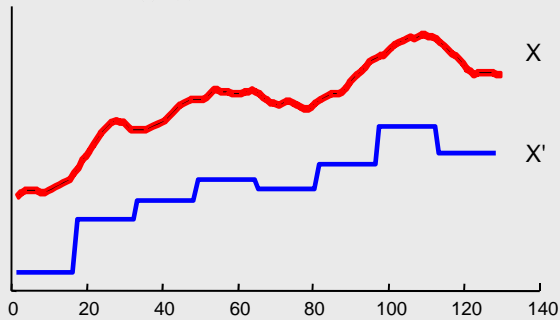
- length
- left_height
- right_height

Piecewise Linear Approximation



- **Sliding Windows**
 - Gradually from left increase the potential segment until the deviation from the original series does not exceed the limit set by the user
 - Relatively good example for stock data
 - Beware of extreme values
 - Variant with k-value adding is faster
- **Top-down**
 - Divide series at a suitable location (eg, minimum or maximum values) into two segments. Then further divide into smaller segments until the stopping criterion.
- **Bottom-up**
 - Complement previous - series is divided into $n / 2$ segments, which gradually combine.
- **SWAB** – sliding window and bottom-up
 - Buffer segments of length w is filled with sliding window, then bottom-up.

Piecewise Aggregate Approximation



Split series into n segments, calculate average for each segment and this value will replace all points in the segment. All segments have the same length.

$$\bar{x}_i = \frac{N}{n} \sum_{j=\frac{n}{N}(i-1)+1}^{\frac{n}{N}i} x_j$$

Given the reduced dimensionality representation we can calculate the approximate Euclidean distance

$$DR(\bar{X}, \bar{Y}) \equiv \sqrt{\frac{n}{N}} \sqrt{\sum_{i=1}^N (\bar{x}_i - \bar{y}_i)^2}$$

- This measure is provably lower bounding.
- Extremely fast to calculate
- Support series of arbitrary lengths
- Can support any Minkowski metric
- Supports non Euclidean measures
- Simple, intuitive

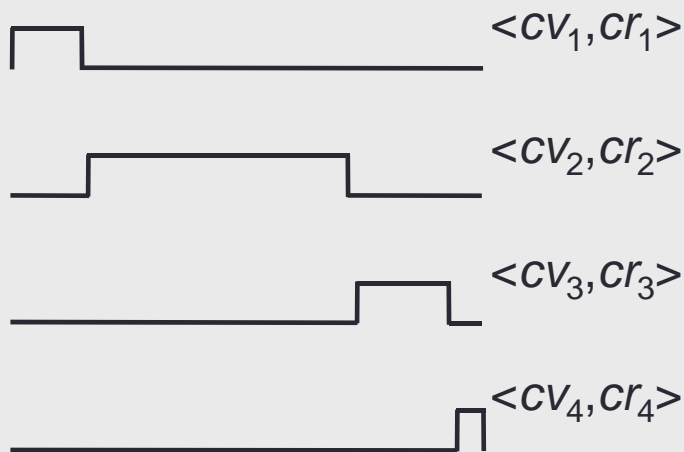
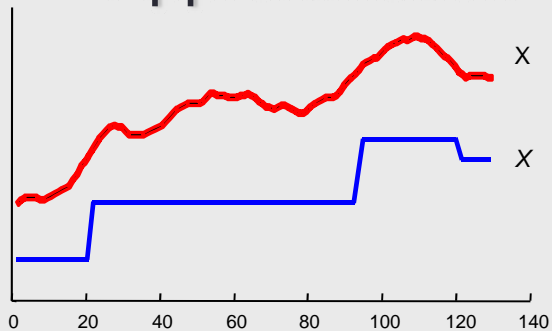
Time Series Representations



Piecewise
Aggregate
Approximation

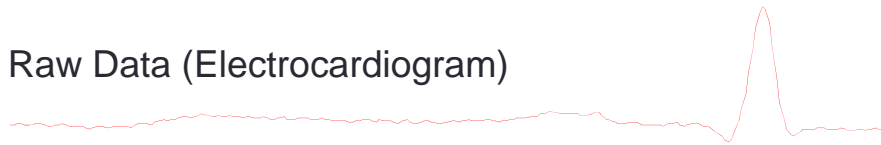


Adaptive Piecewise Constant Approximation

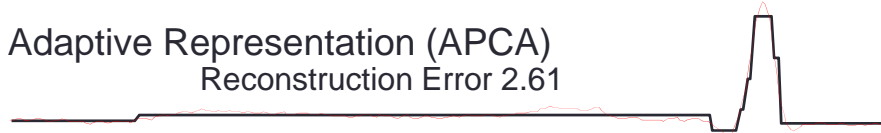


Generalize PAA to allow the piecewise constant segments to have arbitrary lengths. Note that we now need 2 coefficients to represent each segment, its value and its length.

Raw Data (Electrocardiogram)

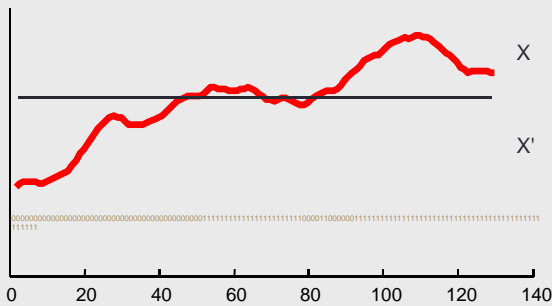


Adaptive Representation (APCA)
Reconstruction Error 2.61



- Many advantages, but challenging to index - implementation exists, but is very challenging.
- Very fast $O(n)$
- More efficient as other approaches
- Support series of arbitrary lengths.
- Supports non Euclidean measures.

Clipped data



Find the mean of a time series, convert values above the line to “1” and values below the line to “0”.

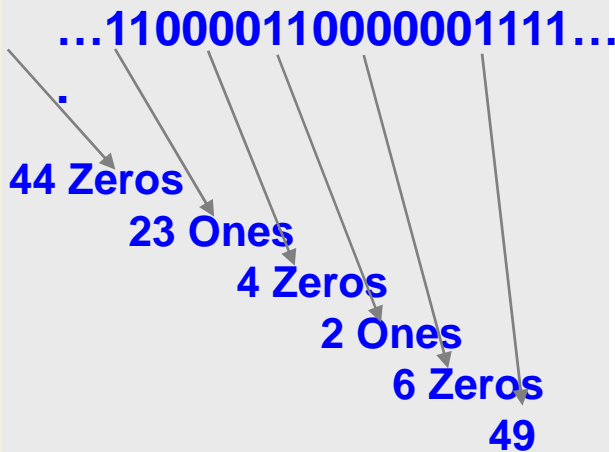


Tony Bagnall

Runs of “1”s and “0”s can be further compressed with run length encoding if desired.

This representation does allow lower bounding.

Ultra compact representation which may be particularly useful for specialized hardware.

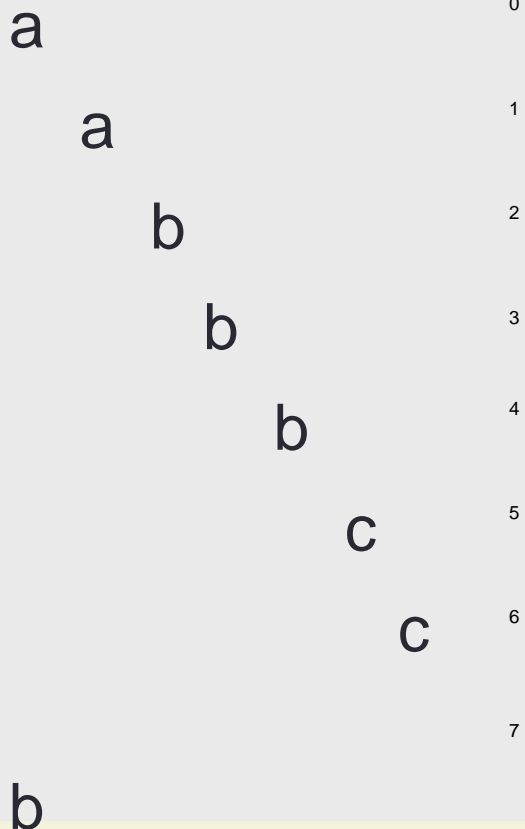
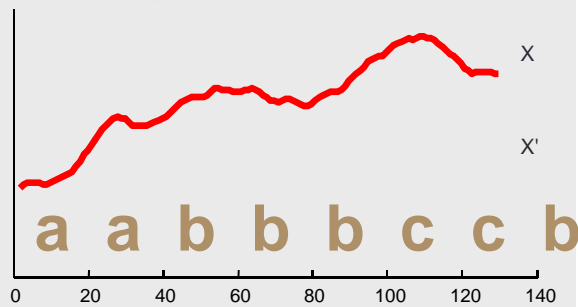


Ones

44

Zeros|23|4|2|6|49

Symbolic Approximation



Convert the time series into an alphabet of discrete symbols. Use string indexing techniques to manage the data.

- We could take advantage of a wealth of techniques from the very mature field of string processing and bioinformatics.
- How we should discretize the times series (discretize the values, the slope, shapes)?

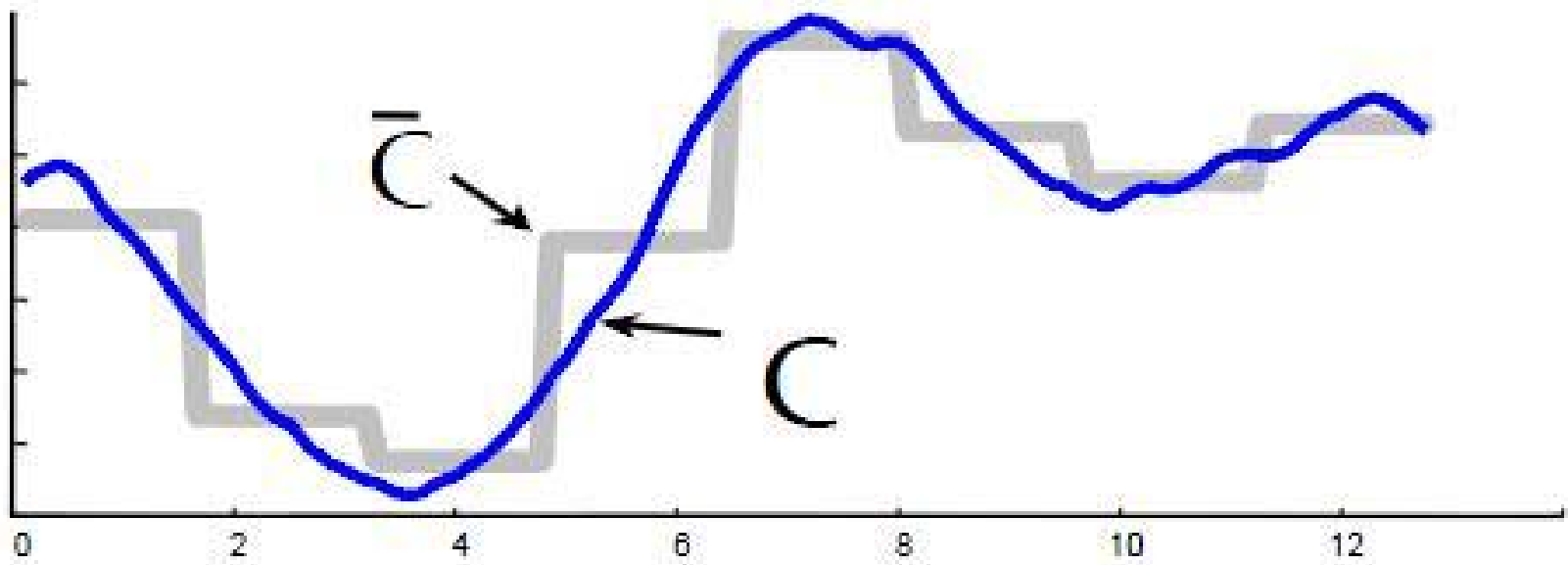
Symbolic Aggregate ApproXimation

- Algorithm for symbolic approximation
- 2002, Eamonn Keogh, University of California
- Significantly reduces the number of dimensions of the original time series
- Lower bounding of Euclidean distance

Symbolic Aggregate ApproXimation

- It consists of three steps
 - Normalisation of time series
 - The mean value is 0
 - PAA transformation
 - Divide time series into multiple segments, calculate the average for each segment, this value will replace all points in the segment
 - Discrete symbolic representation
 - Conversion time series to the alphabet symbols

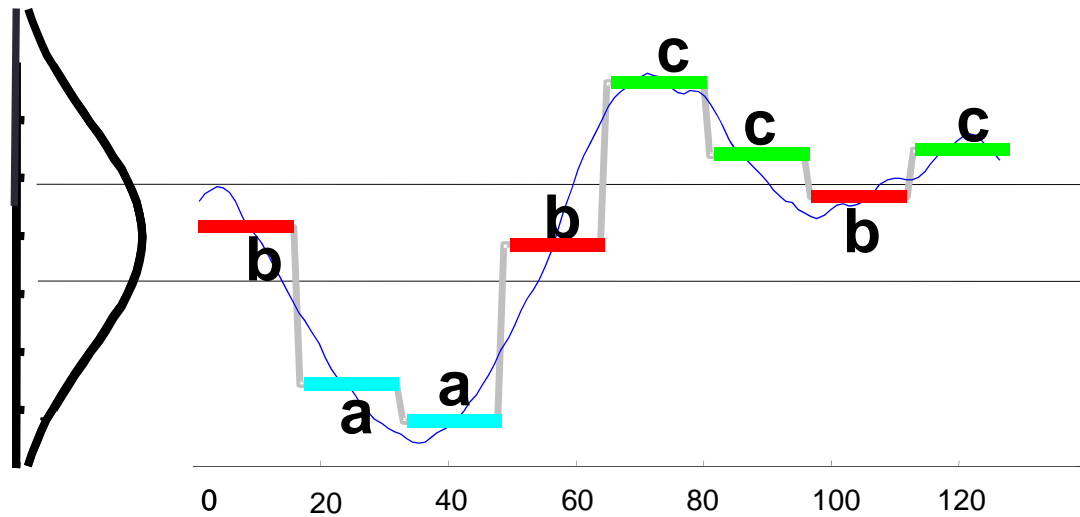
Symbolic Aggregate ApproXimation



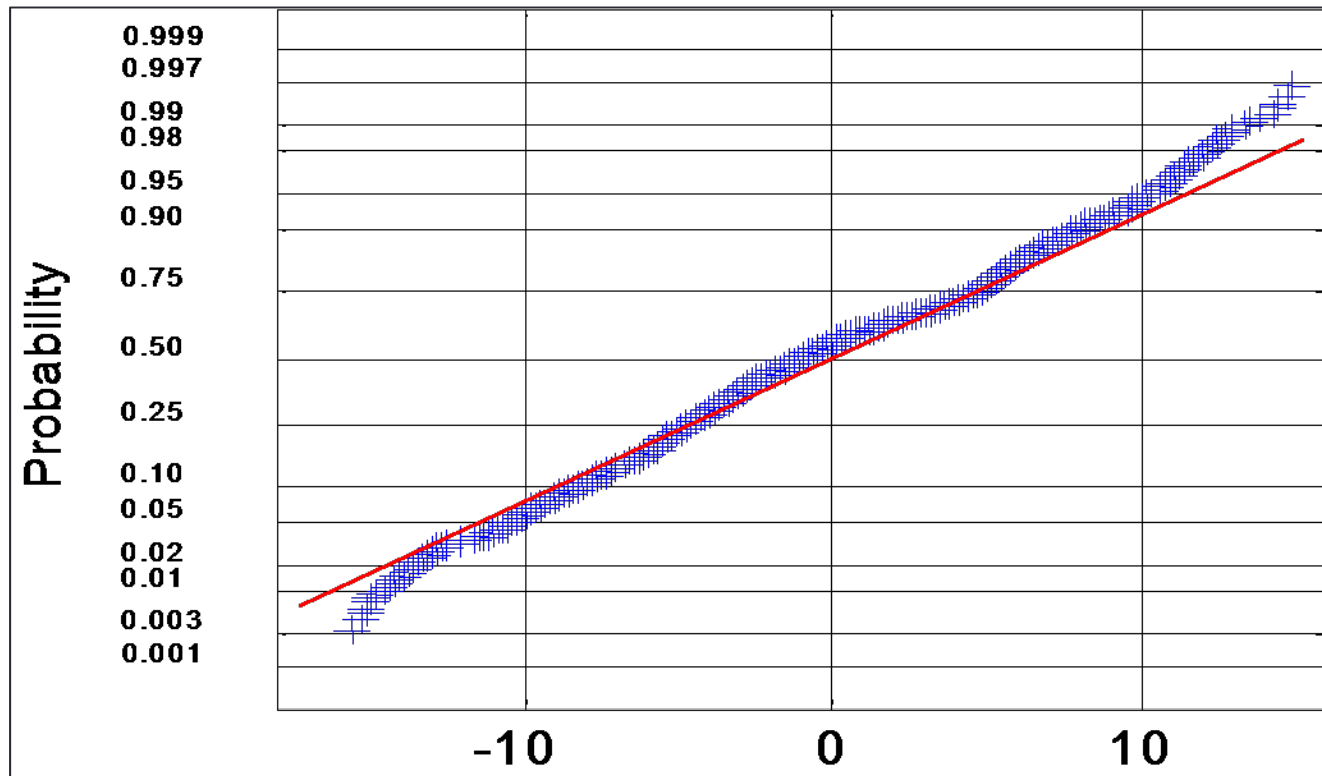
Symbolic Aggregate ApproXimation

- Normalized series has a normal distribution and shape of the Gaussian curve, we define the so-called breaking points, which divide the Gaussian curve into equal parts, discretization generates symbols with equal probability.

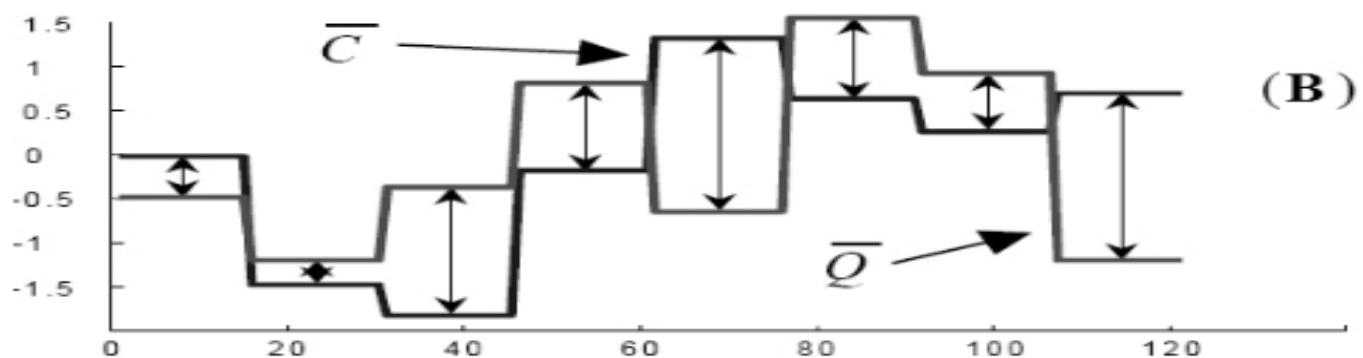
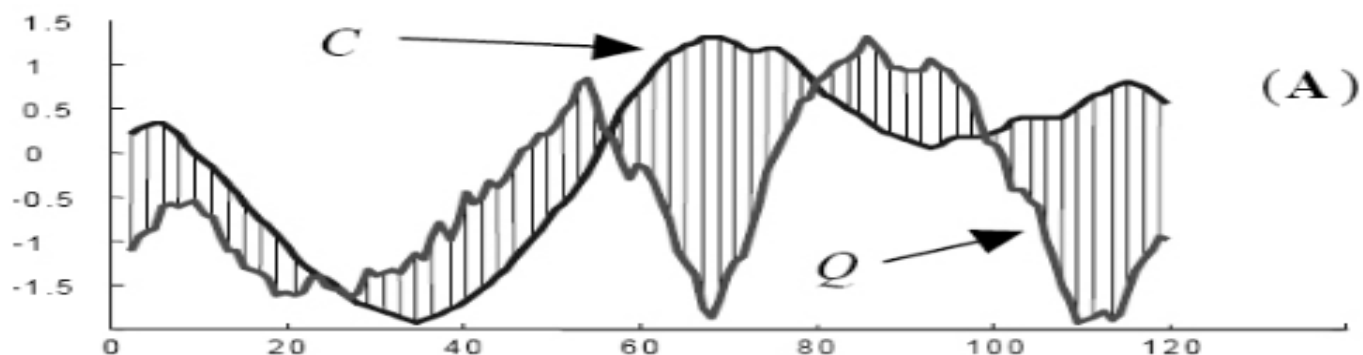
Symbolic Aggregate AppRoXimation



Symbolic Aggregate ApproXimation



Symbolic Aggregate ApproXimation

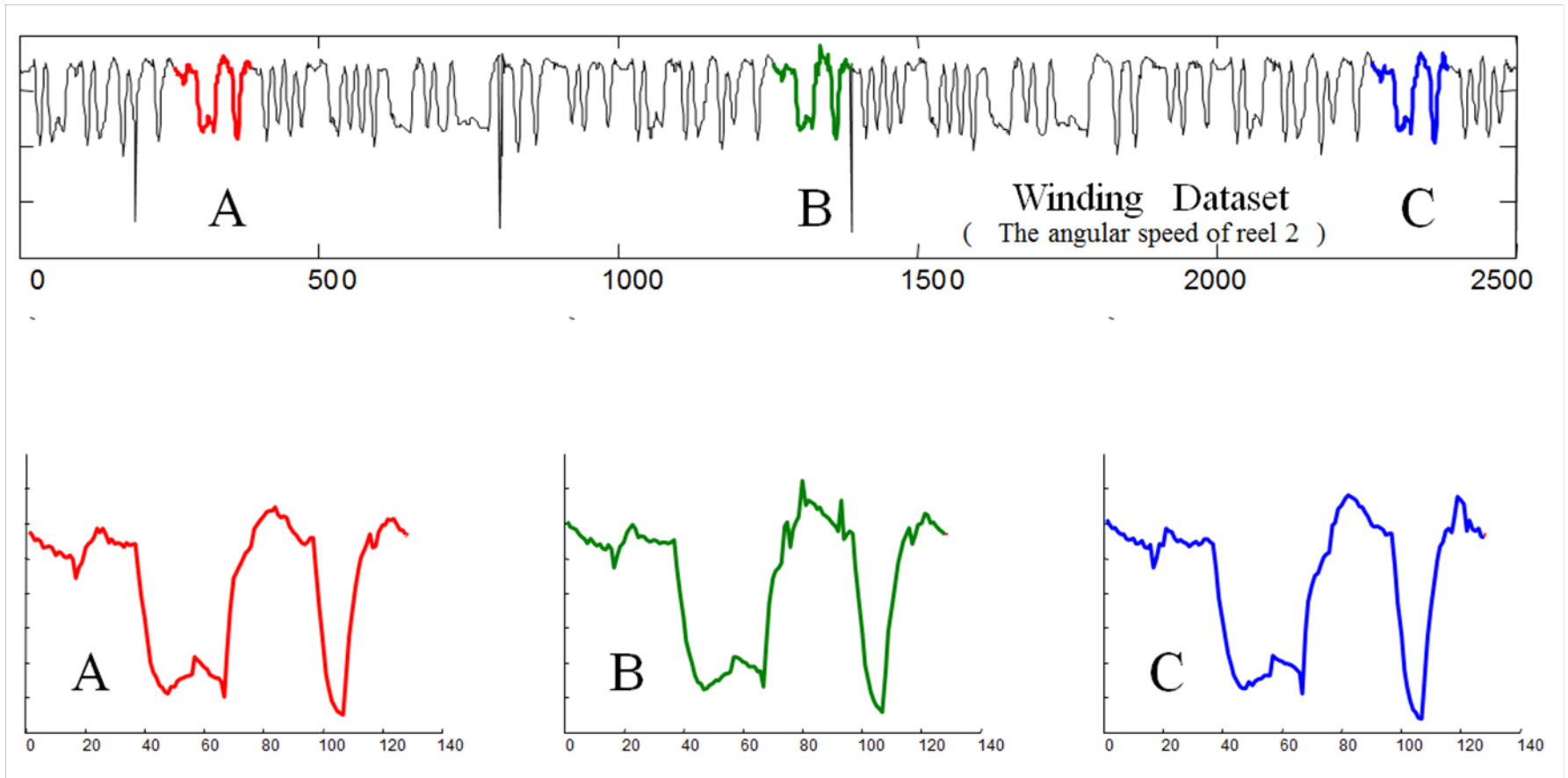


$\hat{C} = \mathbf{b a a b c c b c}$
 $\hat{Q} = \mathbf{b a b c a c c a}$

(C)

$$MINDIST(\hat{Q}, \hat{C}) \equiv \sqrt{\frac{n}{w}} \sqrt{\sum_{i=1}^w (dist(\hat{q}_i, \hat{c}_i))^2}$$

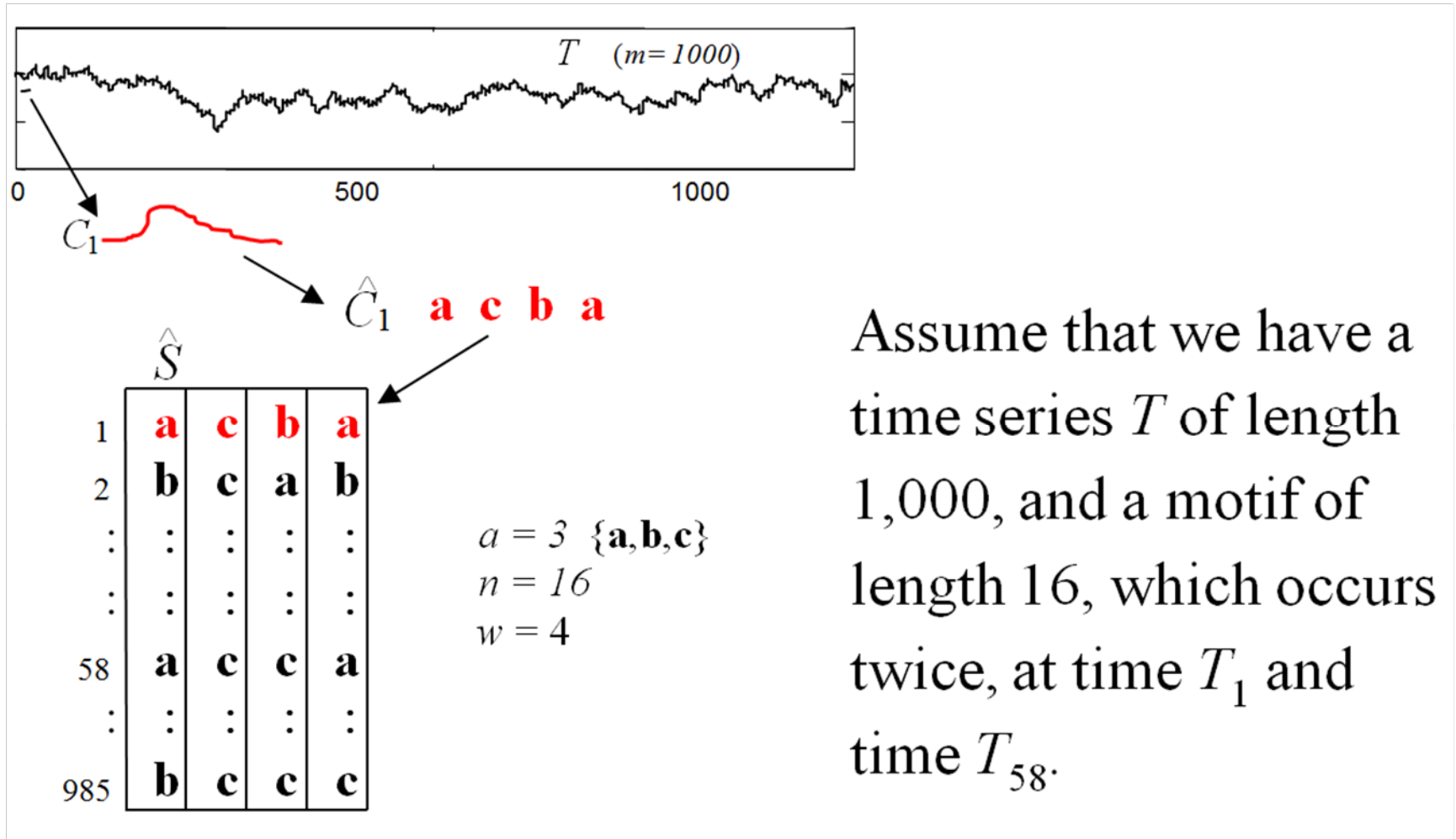
Motif Discovery



Motif Discovery

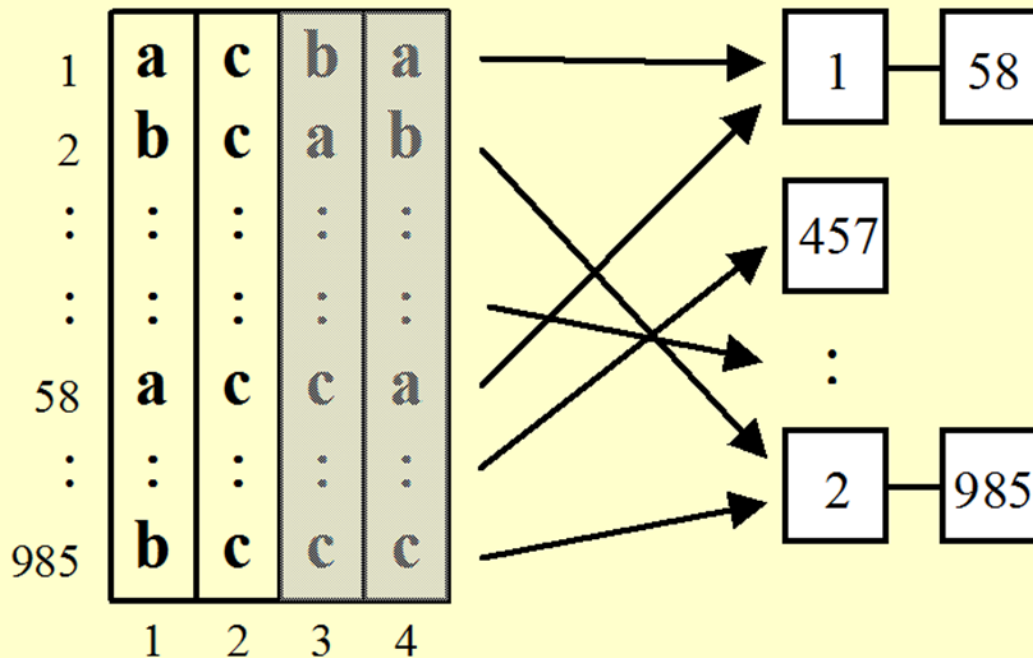
- We define motifs, but how do we find them?
- Brute force search algorithm is just too slow.
- The most reference algorithm is based on a idea from bioinformatics, random projection* and the fact that SAX allows use to lower bound discrete representations of time series.

Motif Discovery



Motif Discovery

A mask $\{1,2\}$ was randomly chosen, so the values in columns $\{1,2\}$ were used to project matrix into buckets.



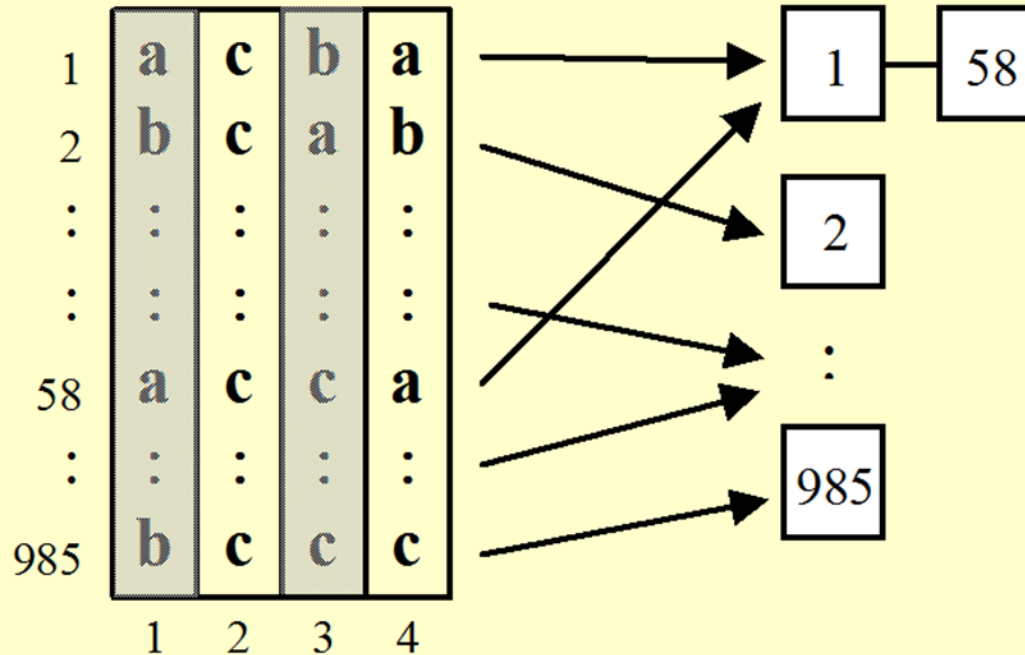
Collisions are recorded by incrementing the appropriate location in the collision matrix

The collision matrix is a grid with rows 1, 2, ..., 58, ..., 985 and columns 1, 2, ..., 58, ..., 985. Shaded cells indicate collisions. The collisions are recorded by incrementing the appropriate location in the collision matrix.

1						
2						
:						
:						
58	1					
:						
:						
985		1				
	1	2	:	58	:	985

Motif Discovery

A mask $\{2,4\}$ was randomly chosen, so the values in columns $\{2,4\}$ were used to project matrix into buckets.



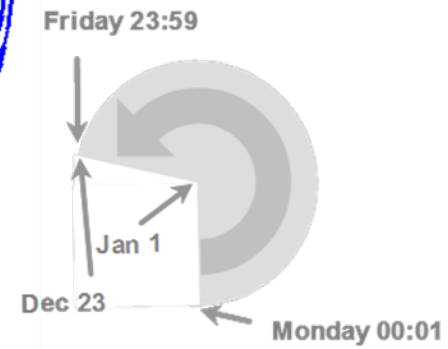
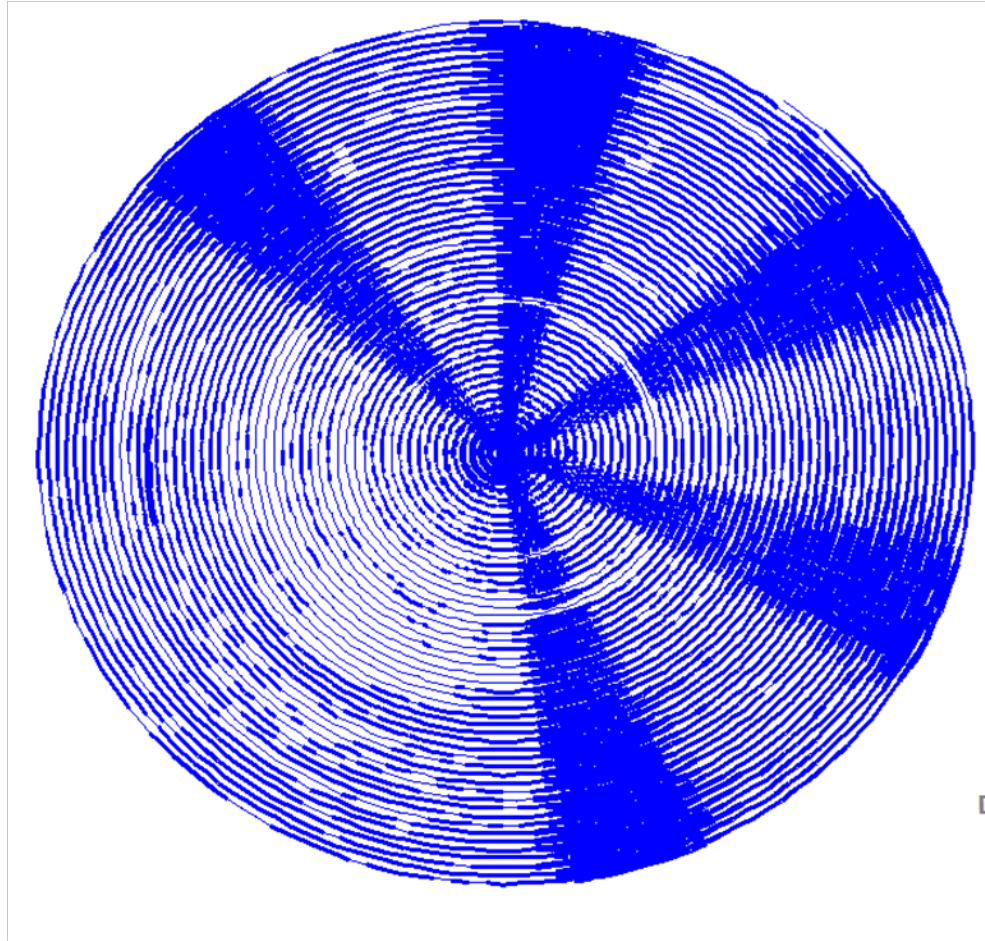
Once again, collisions are recorded by incrementing the appropriate location in the collision matrix

1						
2						
:						
58	2					
:						
985		1				
	1	2	:	58	:	985

Motif Discovery

1						
2	2					
:	1	3				
58	27	2	1			
:	3	2	2	1		
985	0	1	2	1	3	

Visualization

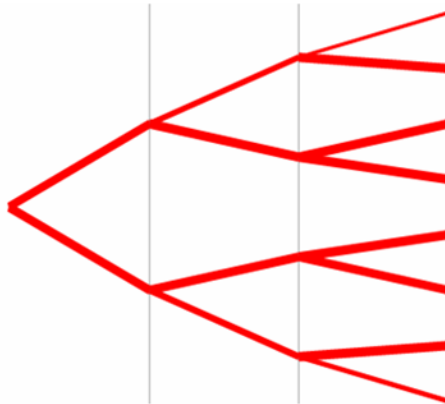


Visualization

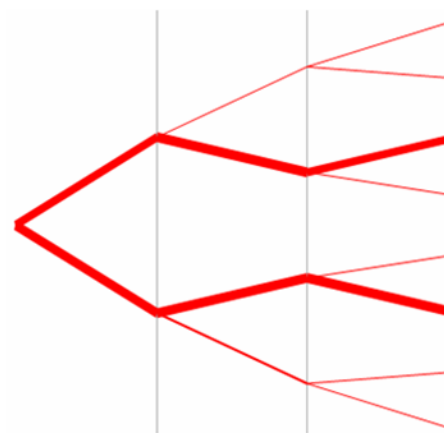
- Time Series Spiral
 - Simple and intuitive
 - Many extensions possible
 - Only useful on periodic data, and only then if you know the period

Visualization

```
010110010111100110100100001000101
00110110101110000101010111011110
001101101101111110100110010010001
101000111100110110100010111100010
110100110110011010000001001100010
011100000111010011001011000010100
10
```

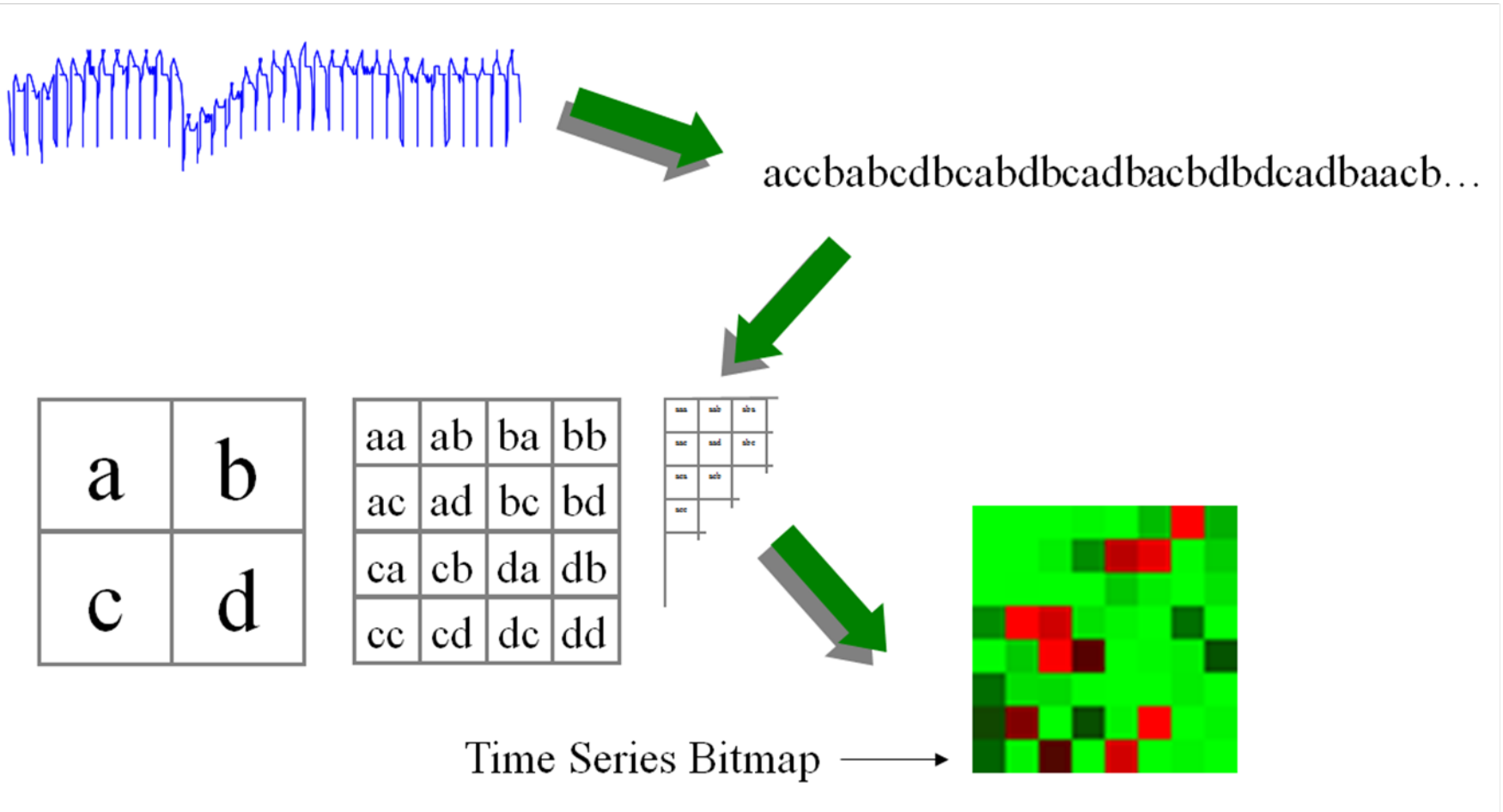


```
10001000101001000101010100001010
100010101110111101011010010111010
010101001110101010100101001010101
110101010010101010110101010010110
010111011110100011100001010000100
111010100011100001010101100101110
101
```



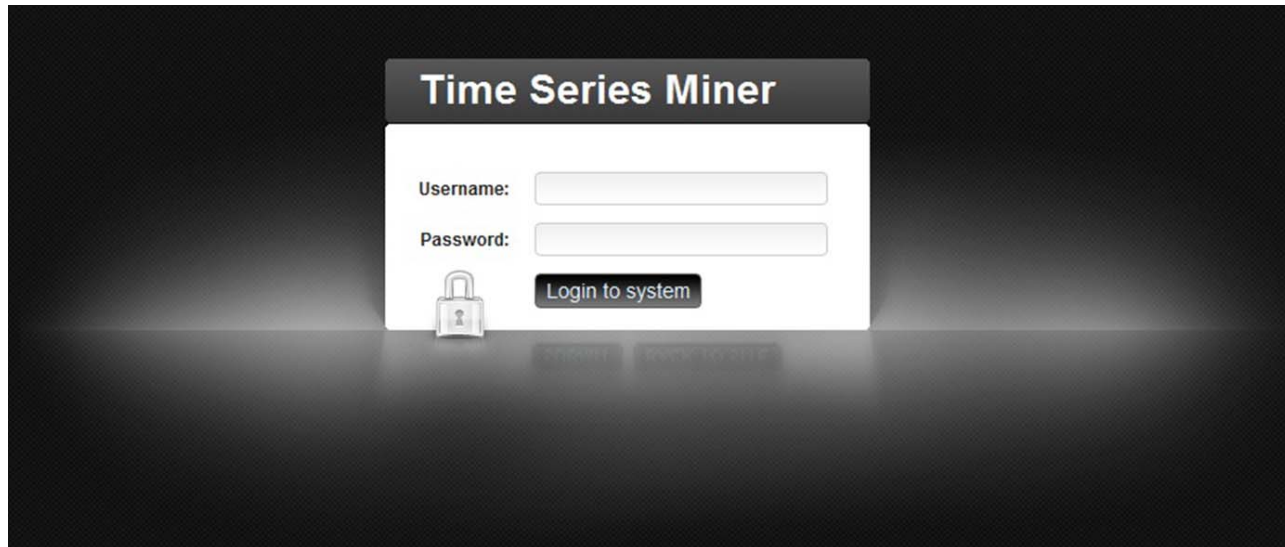
Lets put the sequences into a depth limited suffix tree, such that the frequencies of all triplets are encoded in the thickness of branches...

Visualization



TSMiner

- www.tsminer.cz, www.tsminer.com, www.tsminer.eu
- Multi-tier application for time series datamining
 - MS SQL Server 2008 R2
 - ASP.NET, .NET Framework 4.0



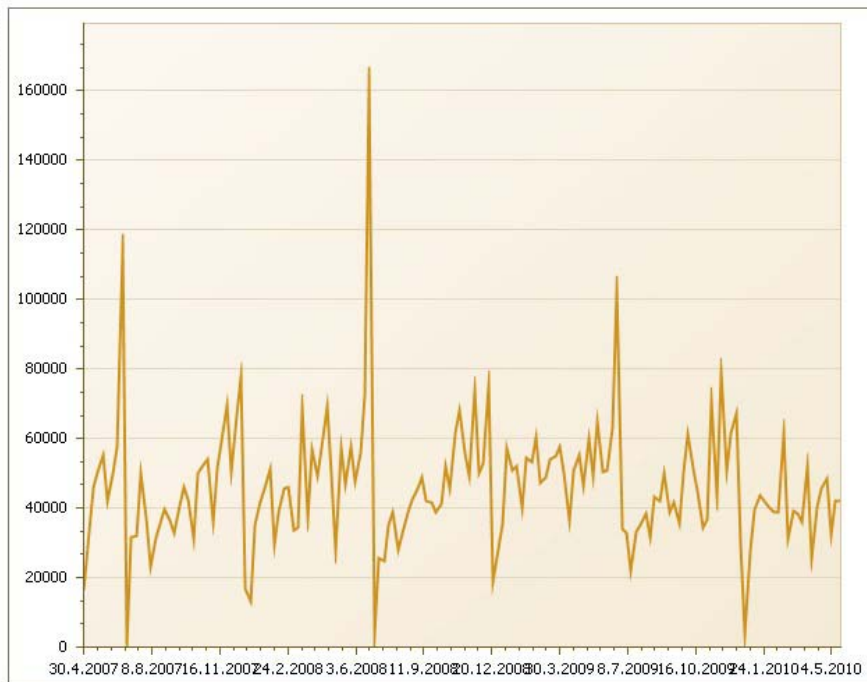
TSMiner

Časové řady

Vyberte projekt:

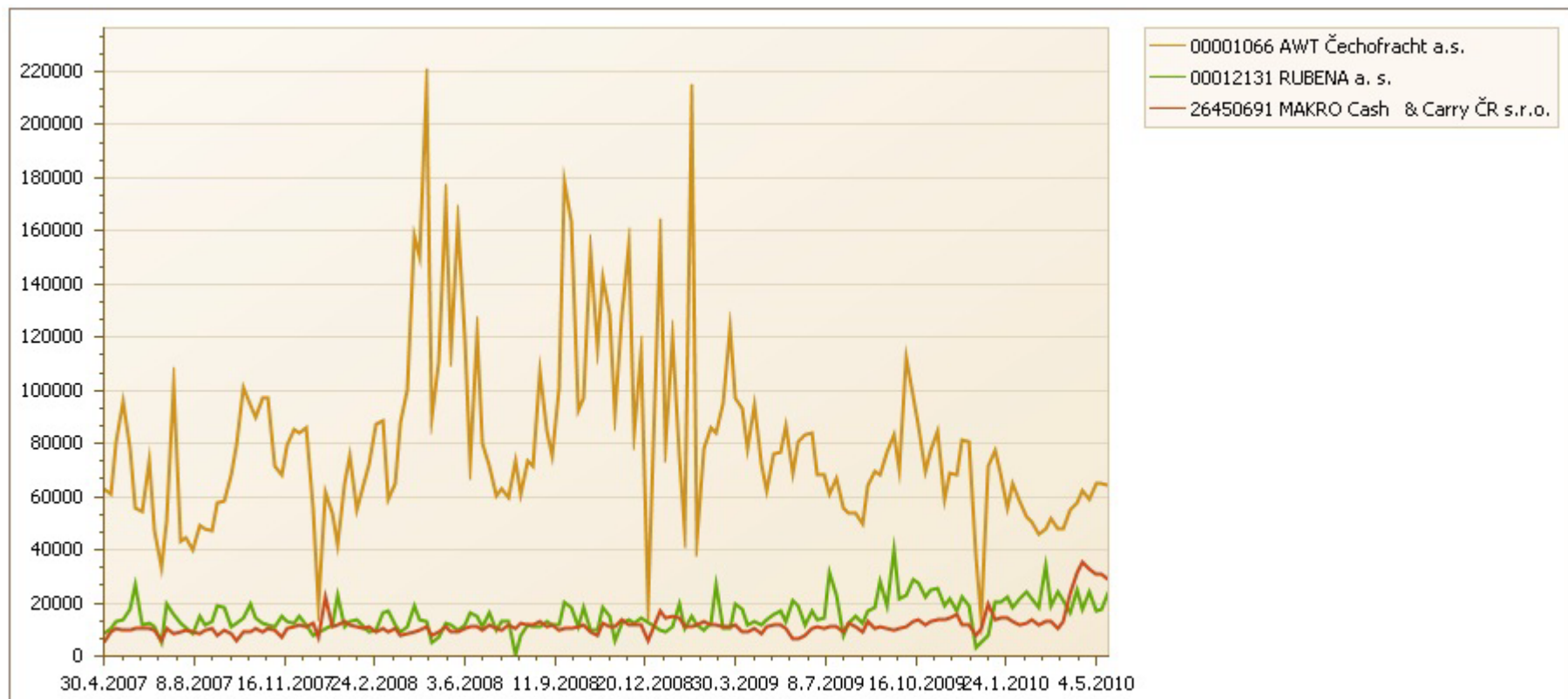
Prodejní databáze

- 00000001 Daniel Holásek
- 00001066 AWT Čechofracht a.s.
- 00012131 RUBENA a.s.
- 00559709 GEODIS BRNO, spol. s r.o.
- 00570664 Servant, a.s.
- 00676853 HOSPIMED, spol. s r.o.
- 10356592 Ing. Vladimír Výborný - LAN - projekt
- 10625658 Ing. Leo Stoklasa
- 11278528 Antonín Kolinger ASON
- 12164739 COMPEX, spol. s r.o.
- 13389271 Ing. Petr Kuba
- 15059278 MATEZA spol. s r.o.
- 15609391 Ing. Jaroslav Chutný
- 15887405 Motorsport, spol. s r.o.
- 15887791 AUTO KELLY, A.S.
- 16192648 DERMATEX, spol. s r.o.
- 16193407 WELLA CZ s.r.o.**
- 16556402 Straumann s.r.o.
- 18381201 AUTO - COLOR spol. s r.o.
- 18568092 Ing. Petr Prokúpek
- 18627722 ORIFLAME CZECH REPUBLIC spol. s r.o.
- 18630774 KODYS spol. s r.o.
- 18828507 REDA a.s.
- 19015909 Schindler, spol. s r.o.
- 25036661 Karton s.r.o.
- 25044516 FCC průmyslové systémy s.r.o.
- 25046225 ELTEQ, spol. s r. o.
- 25083163 Dentamed (ČR), spol. s r.o.
- 25099418 PRIMAVERA ANDORRANA s.r.o.
- 25101625 HUSKY CZ s.r.o.
- 25123998 TOMKET, s.r.o.
- 25130340 BRITEX CZ, S.R.O.
- 25140388 FISCHER INTERNATIONAL S.R.O.
- 25158694 PROFIMED s.r.o.
- 25262785 Smart Print s.r.o.
- 25266276 CLIPET s.r.o.



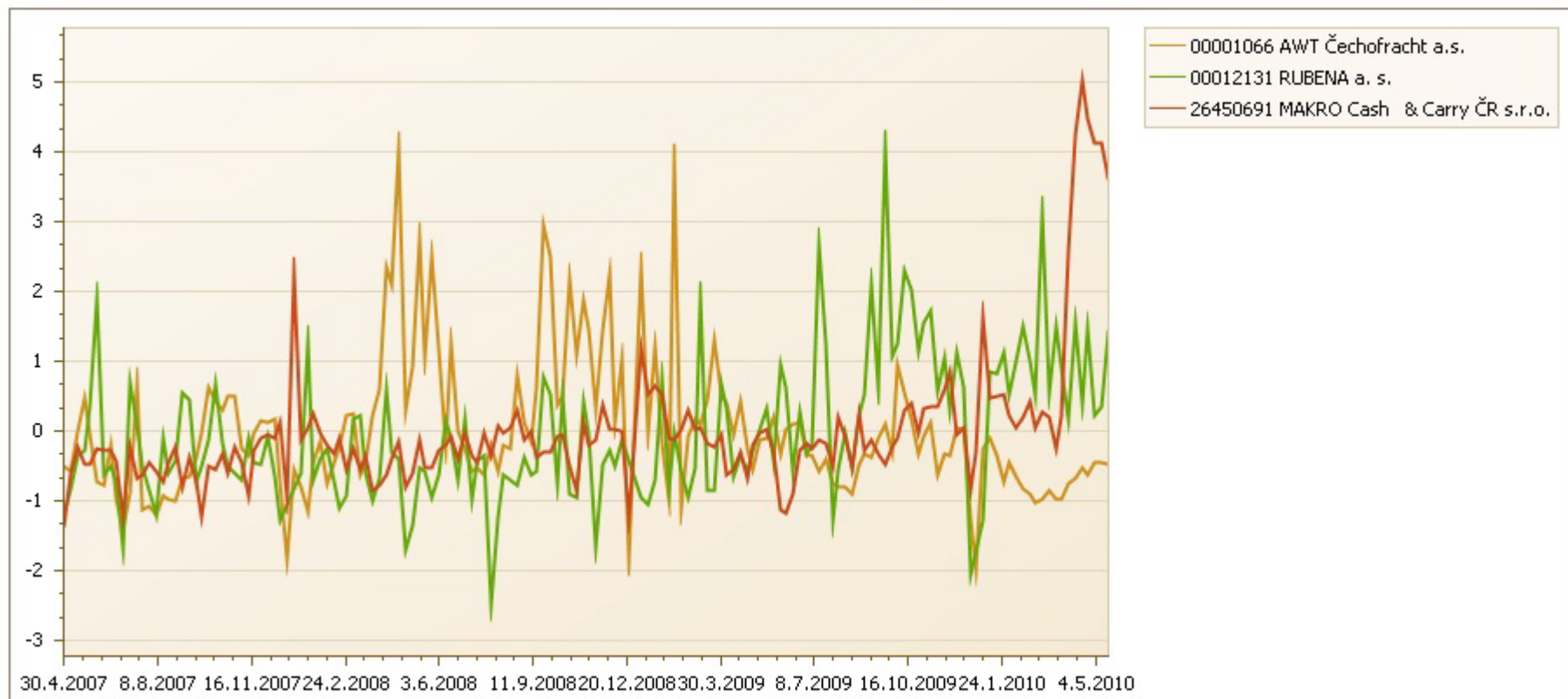
Počet prvků:	159	Součet řady:	7349670,81	Průměrná hodnota:	46224,34
Minimum:	360,57	Maximum:	163098,38	Směrodatná odchylka:	18201,85

TSMiner



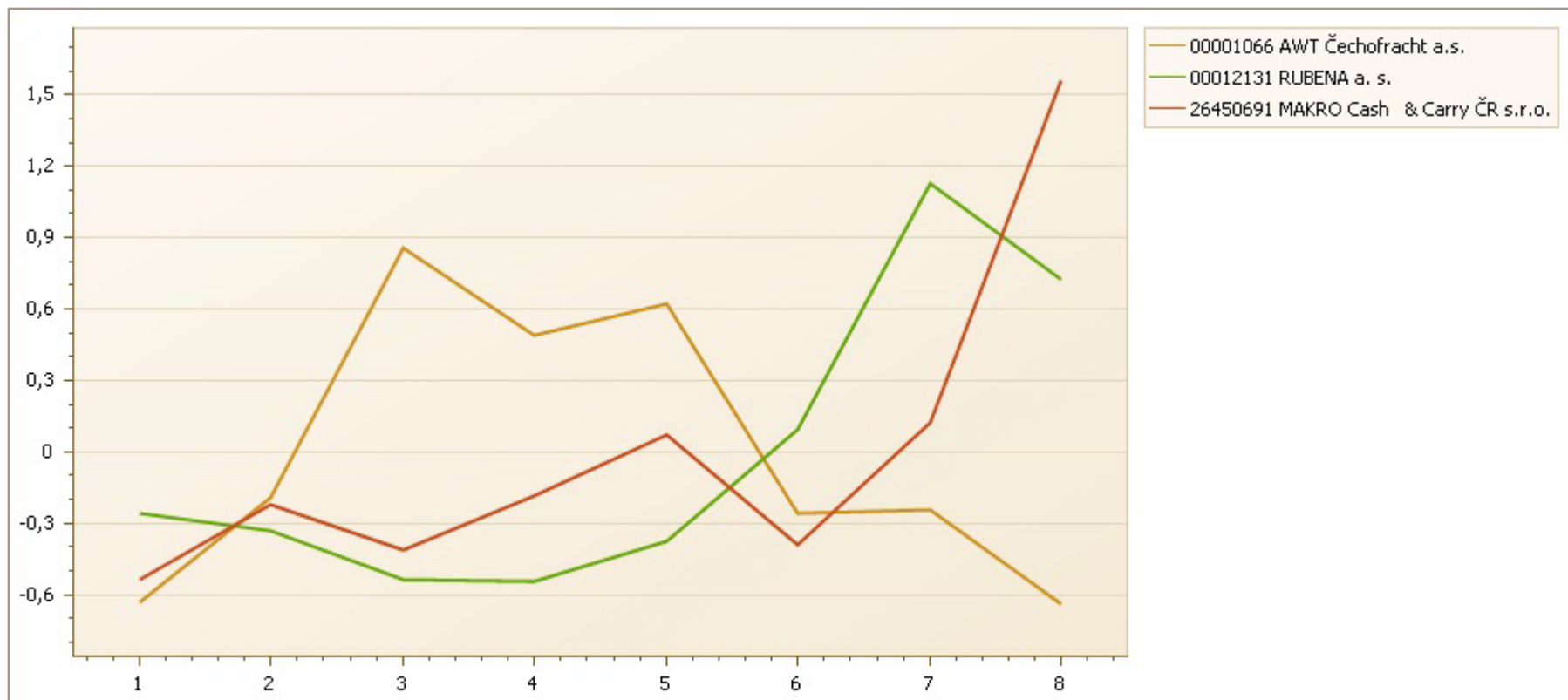
Vypnout transformace Normalizovat řadu Aproximace PAA Aproximace SAX

TSMiner



Vypnout transformace Normalizovat řadu Aproximace PAA Aproximace SAX

TSMiner



Vypnout transformace Normalizovat řadu Aproximace PAA Aproximace SAX

TSMiner

Vypnout transformace Normalizovat řadu Aproximace PAA Aproximace SAX

Počet segmentů:

Velikost abecedy:

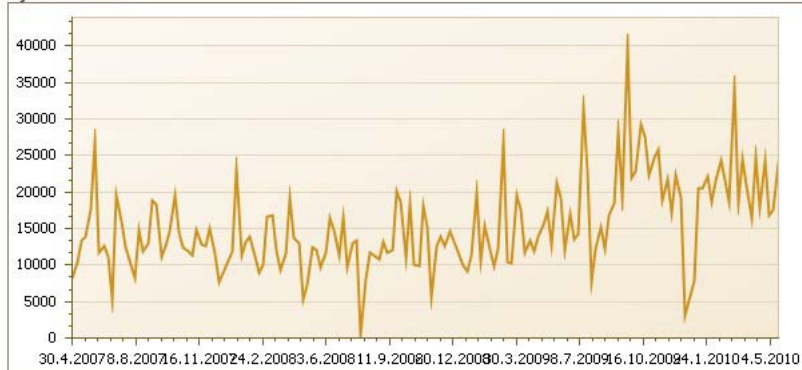
00001066 AWT Čechofracht a.s.	BBDCBBB
00012131 RUBENA a. s.	BBBBBCDD
26450691 MAKRO Cash & Carry ČR s.r.o.	BBBBCBCD

TSMiner

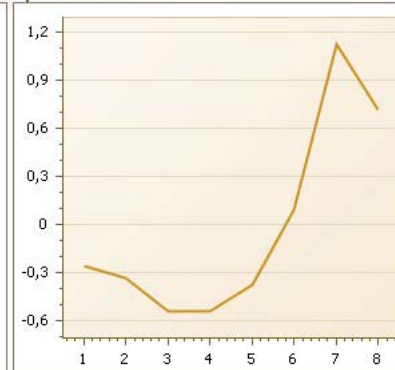
Klasifikace třídy

Vyberte řadu pro klasifikaci: Počet segmentů: Velikost abecedy: Počet sousedů:

Vybraná řada:



Aproximace:



Klasifikace:



Klasifikace: Třída A

Main problems of time series analysis

- Pattern search w/o prior parameter setting
- Clustering of streamed data
- Time series merging – finding all shared subsequences
- „Why“ analysis in classification and clustering, automatic generation of explanation
- Weighed representation of time series
- Visualization of large time series

(This slide was translated to English by V. Svátek)

Thanks for your attention