# CS-541
# Wireless Sensor Networks

**Lecture 10: Time-series analysis**

Spring Semester 2017-2018

Prof Panagiotis Tsakalides, Dr Athanasia Panousopoulou, Dr Gregory Tsagkatakis

CS-541 Wireless Sensor Networks
University of Crete, Computer Science Department

FORTH
Institute of Computer Science

# Overview

- Time series analysis

- Intro to Machine learning

CS-541 Wireless Sensor Networks
University of Crete, Computer Science Department

FORTH
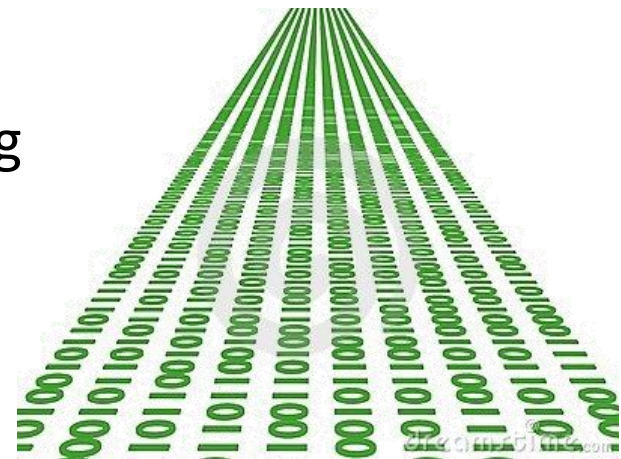Institute of Computer Science

# Stream Data Processing

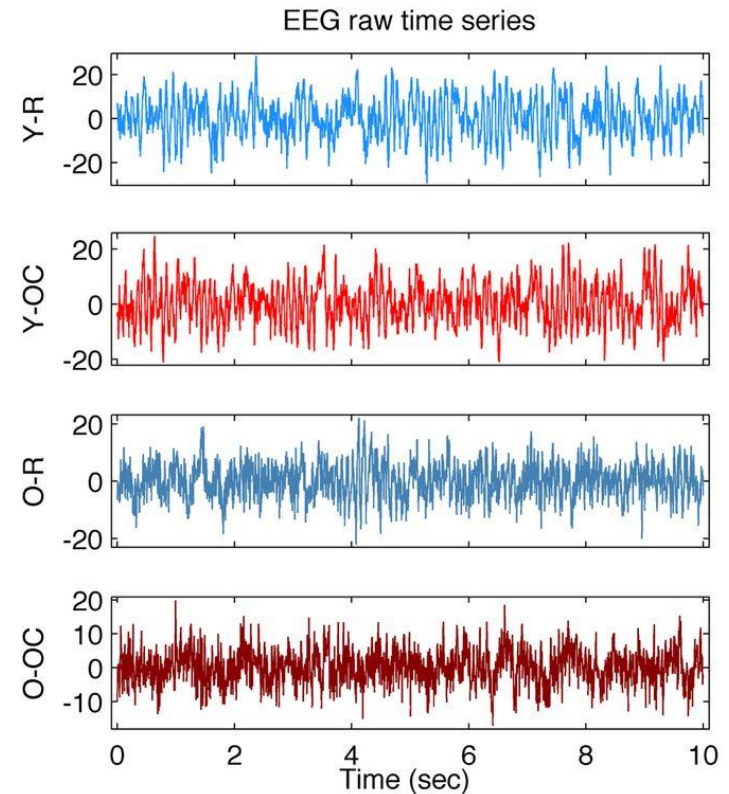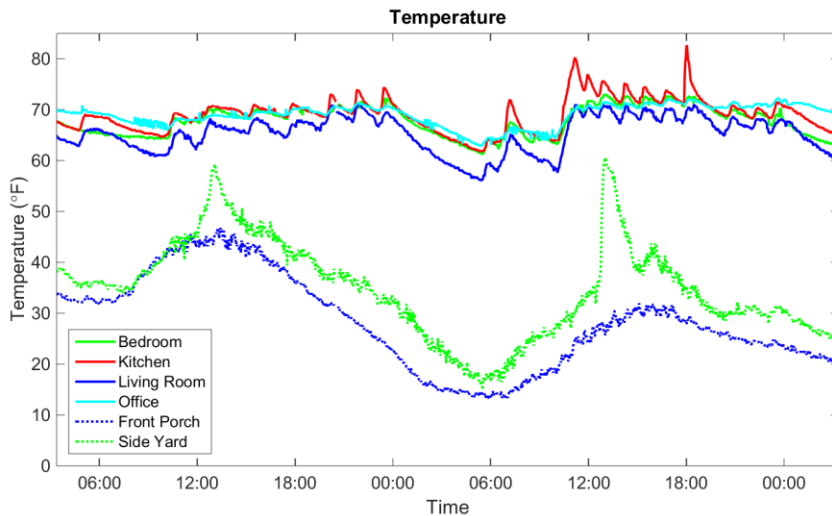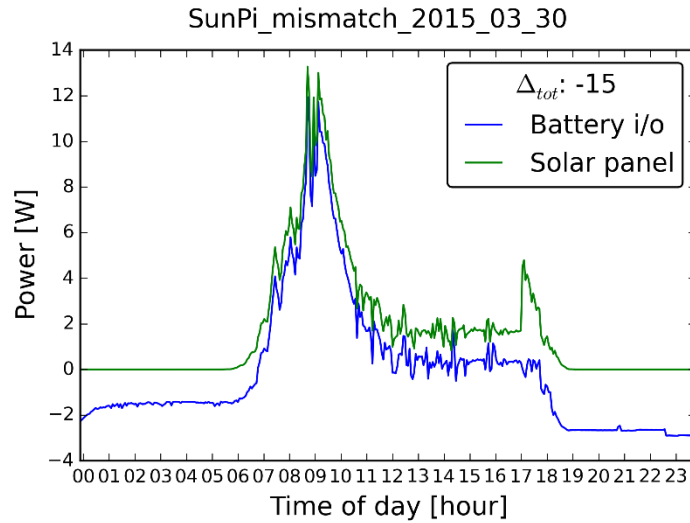Data streams—continuous, ordered, changing, fast, huge amount

- Huge *volumes* of continuous data, possibly infinite
- Fast *changing* and requires fast, real-time response

Applications

- Telecommunication records
- Network monitoring and traffic engineering
- Industrial processes: power & manufacturing
- Sensor, monitoring & surveillance

# Time-series in WSN

SunPi_mismatch_2015_03_30



EEG raw time series



Temperature

CS-541 Wireless Sensor Networks
University of Crete, Computer Science Department

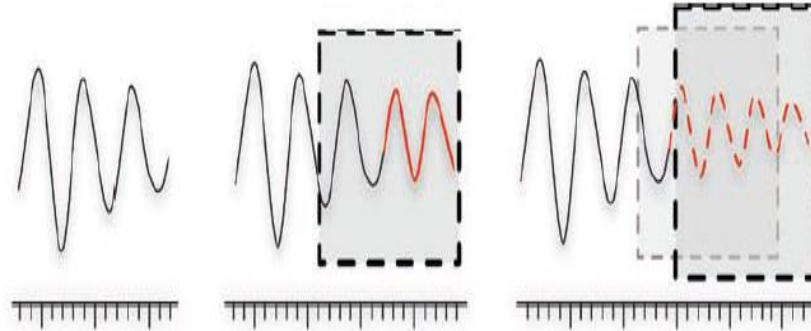FORTH
Institute of Computer Science

# Problems

- *Type 1*:  patterns, periodicities, and/or compress
  - Wearable, Smart city
- *Type 2*:  forecast, find motifs, quantify similarity
  - Activity recognition
- *Type 3*: Multiple time series analysis
  - Sensor networks

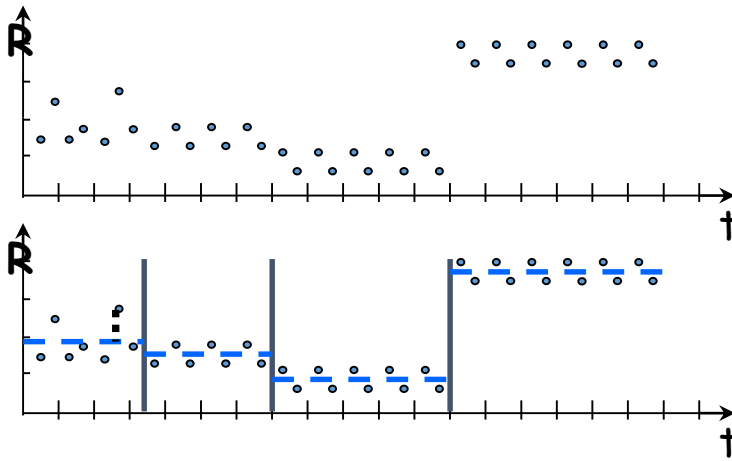"Predictions are very difficult… especially about the future"
Niels Bohr

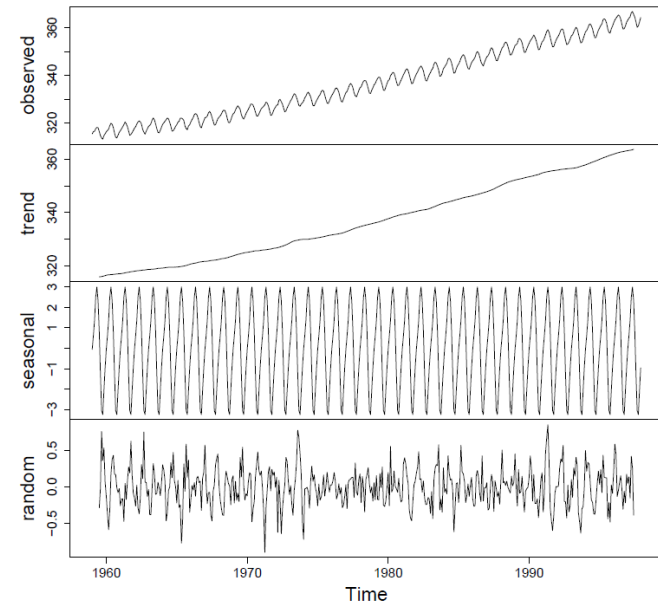CS-541 Wireless Sensor Networks
University of Crete, Computer Science Department

FORTH
Institute of Computer Science

# Applications

## Prediction - Forecasting



## Segmentation - Clustering



## Analysis

FORTH
Institute of Computer Science

# Time-Series data

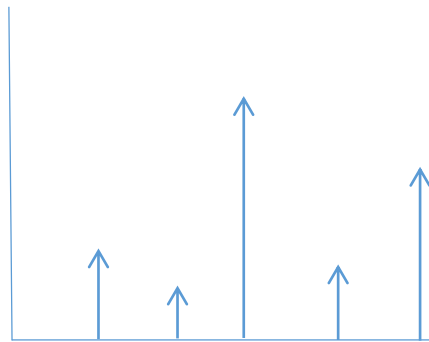Time series: sequence of observations $s_t \in R$ ordered in time t=1...N

Applications

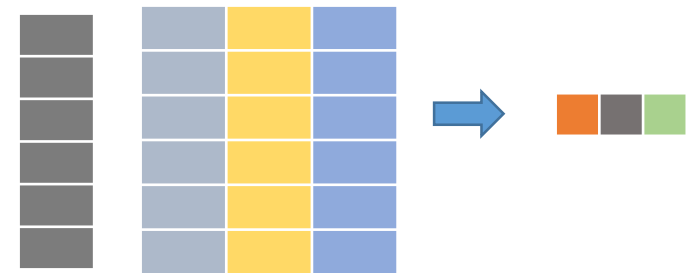• Weather, economic, marketing, web, envirometrics, sensor networks
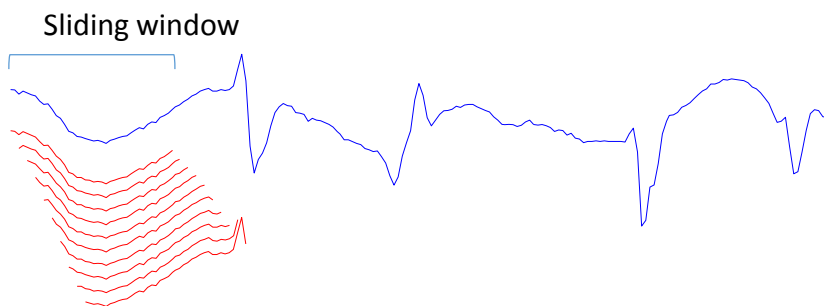
Representations



Sliding windows

Histograms
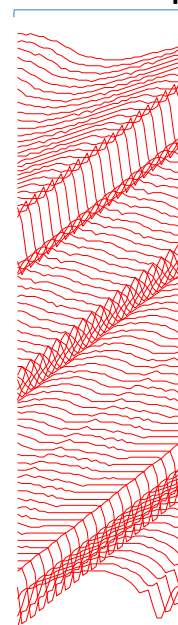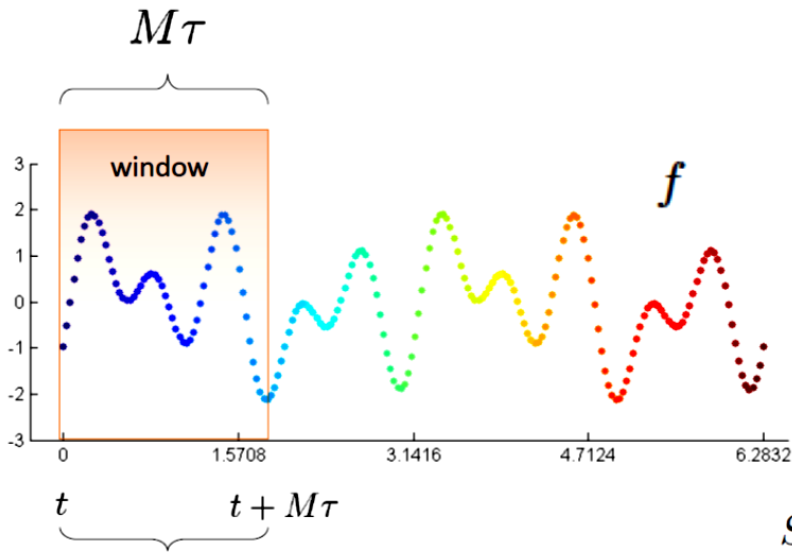
Transform coding

# Sliding window

- Given a time series, individual subsequences are extracted with a sliding window

Sliding window

All subsequences
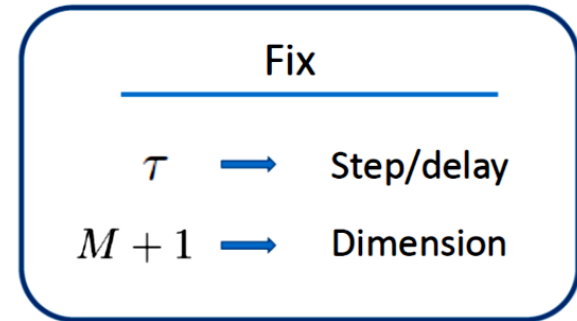
# Sliding windows embedding



$f(t), f(t+\tau), \ldots, f(t+M\tau)$

$$SW_{M,\tau}f(t) = \begin{bmatrix} f(t) \\ f(t+\tau) \\ \vdots \\ f(t+M\tau) \end{bmatrix} \in \mathbb{R}^{M+1}$$

**Fix**

| $\tau$ | $\Longrightarrow$ | Step/delay |
| $M+1$ | $\Longrightarrow$ | Dimension |

*Sliding Windows and Persistence: An application of topology to signal analysis,* J. Perea and J. Harer, 2015

**Data stream**



① Sensor stream

CS-541 Wireless Sensor Networks
University of Crete, Computer Science Department

FORTH
Institute of Computer Science

**Data stream**

① Sensor stream

② Temporal windowing

③ Hankelization process **H**

   ✔ $[n_1]$ lagged temporal windows of $[n_2]$ samples

**1st window**

$[n_2]$

$[n_1]$

| $h_0$ | $h_1$ | ... | | $h_{n_1}$ | ... | $h_{n_2-1}$ | $h_{n_2}$ |

**H**

CS-541 Wireless Sensor Networks
University of Crete, Computer Science Department

FORTH
Institute of Computer Science
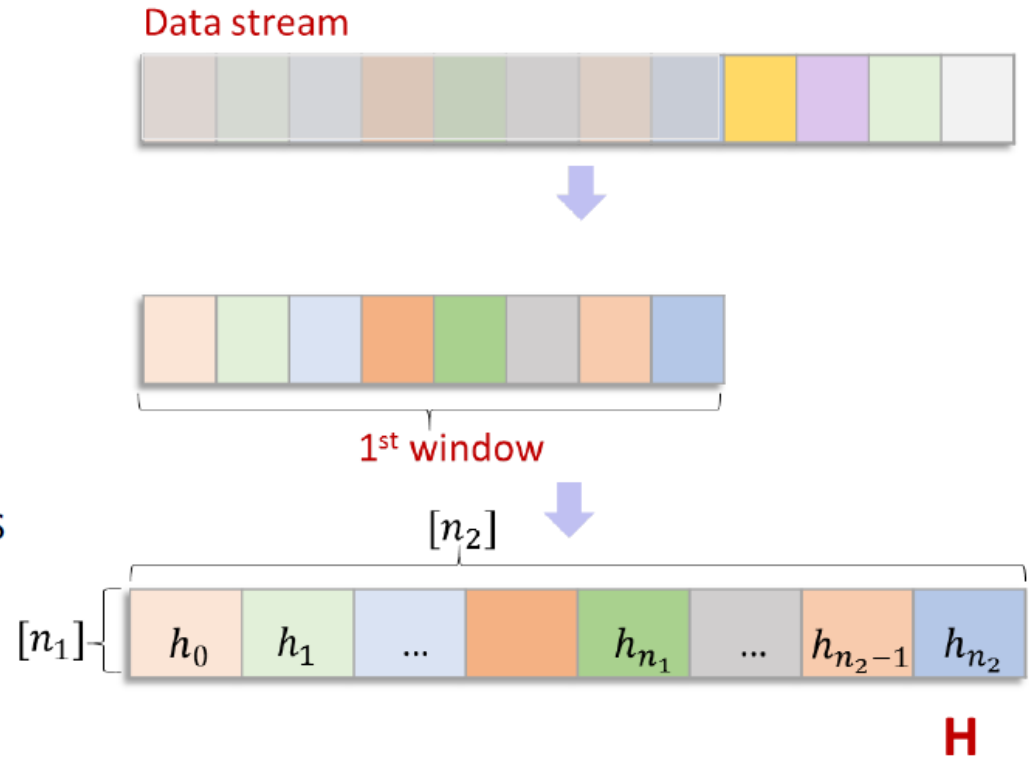
Data stream



① Sensor stream

② Temporal windowing

③ Hankelization process **H**
   ✔ $[n_1]$ lagged temporal windows of $[n_2]$ samples

lag

$1^{st}$ window

$[n_2]$

$[n_1]$ — | $h_0$ | $h_1$ | ... | | $h_{n_1}$ | ... | $h_{n_2-1}$ | $h_{n_2}$ |
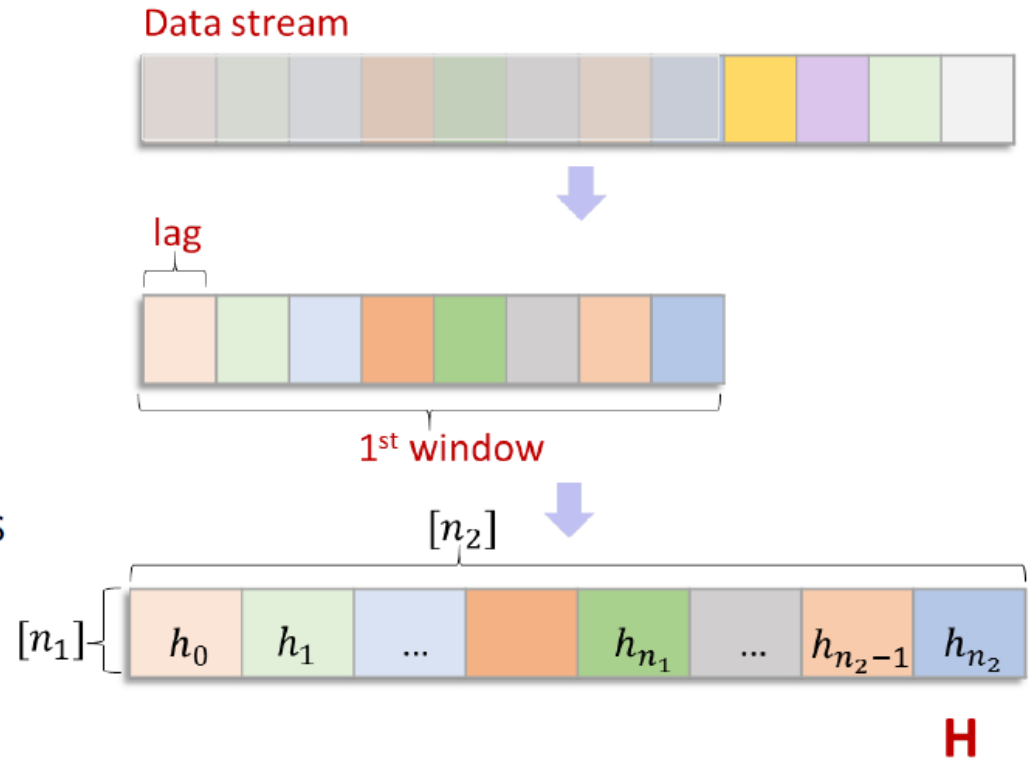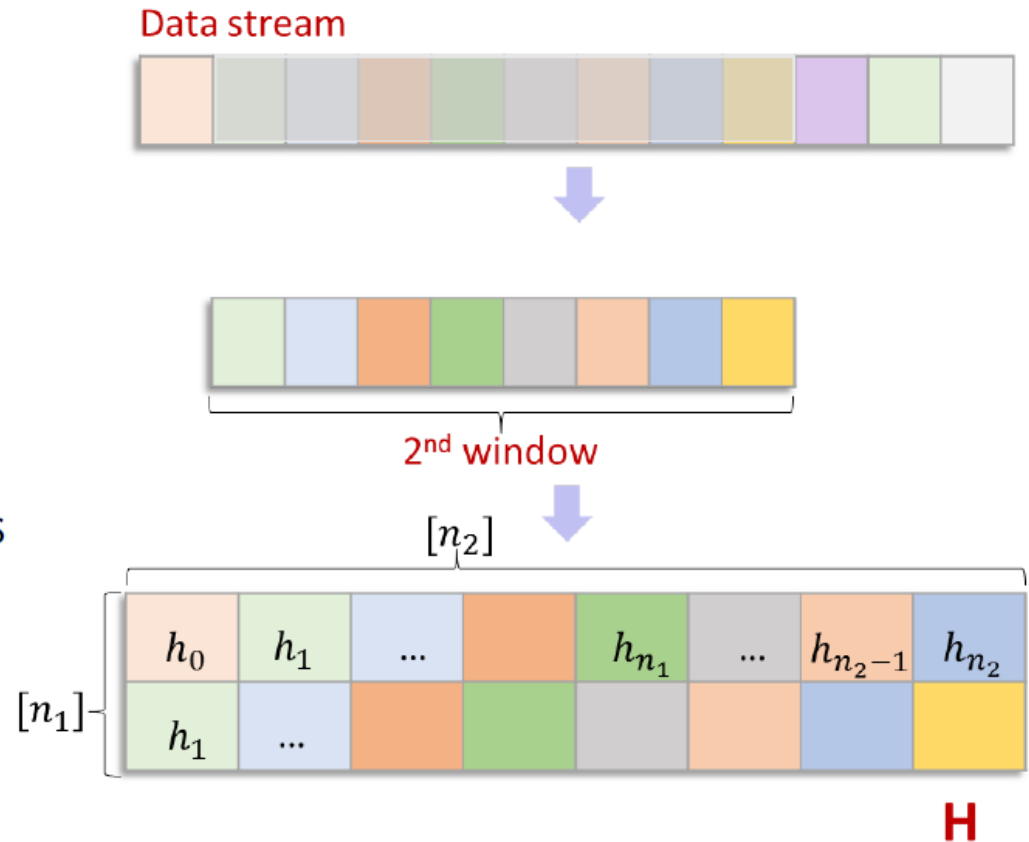
**H**

FORTH
Institute of Computer Science

Data stream



① Sensor stream

② Temporal windowing

③ Hankelization process **H**

✔ $[n_1]$ lagged temporal windows of $[n_2]$ samples

2nd window

$[n_2]$

$[n_1]$

| $h_0$ | $h_1$ | ... | | $h_{n_1}$ | ... | $h_{n_2-1}$ | $h_{n_2}$ |
| $h_1$ | ... | | | | | | |

**H**

Data stream

① Sensor stream

② Temporal windowing

③ Hankelization process **H**
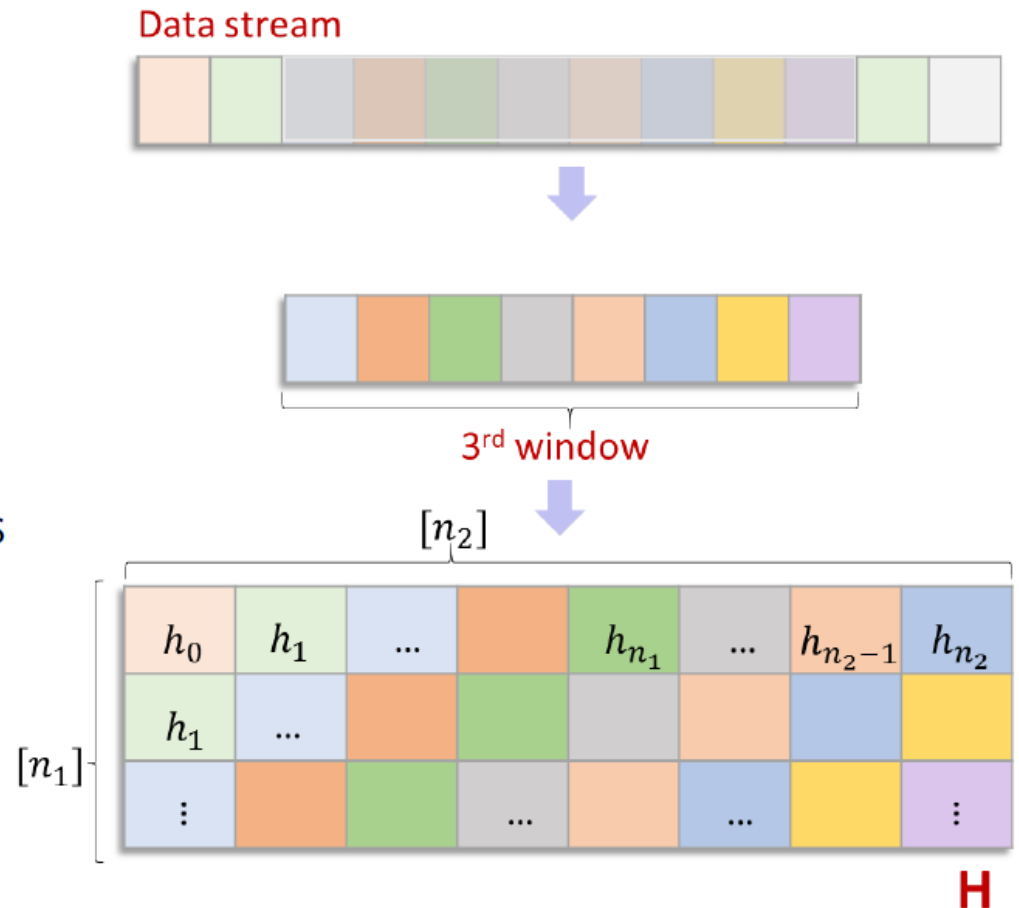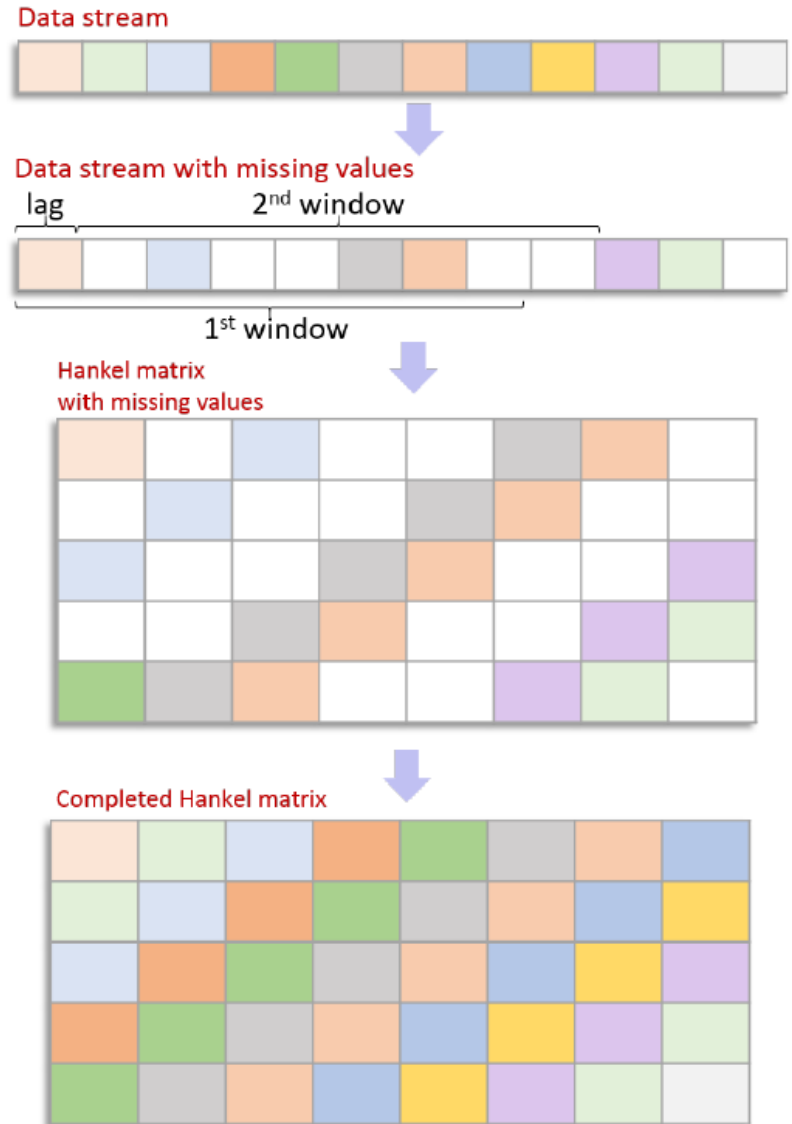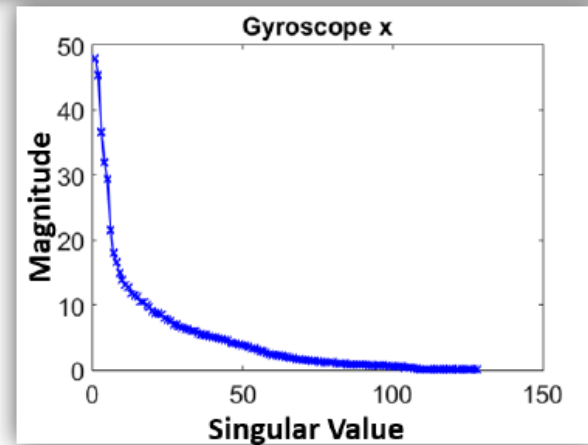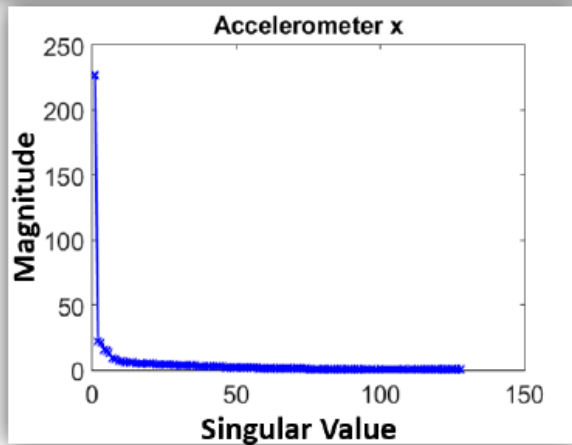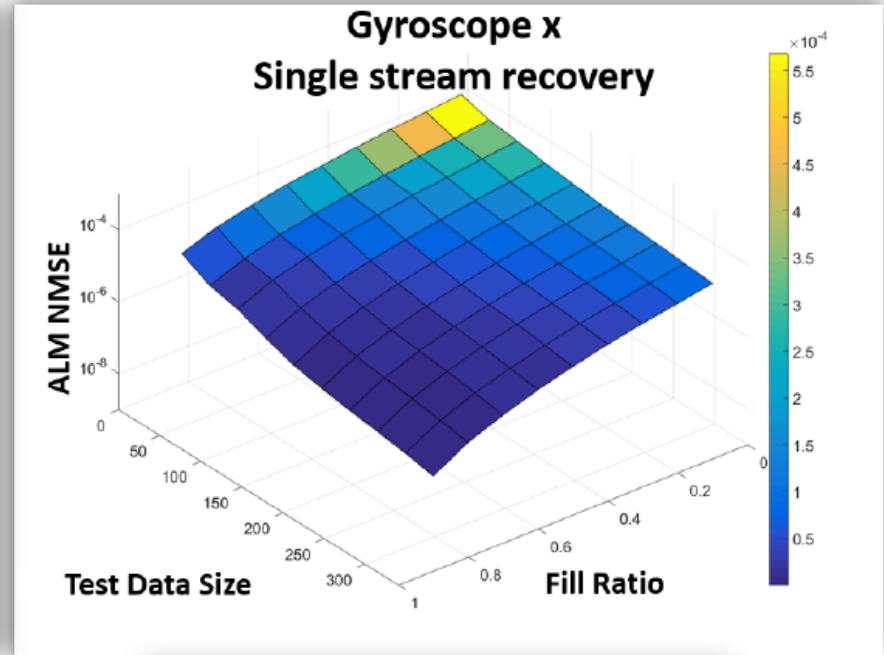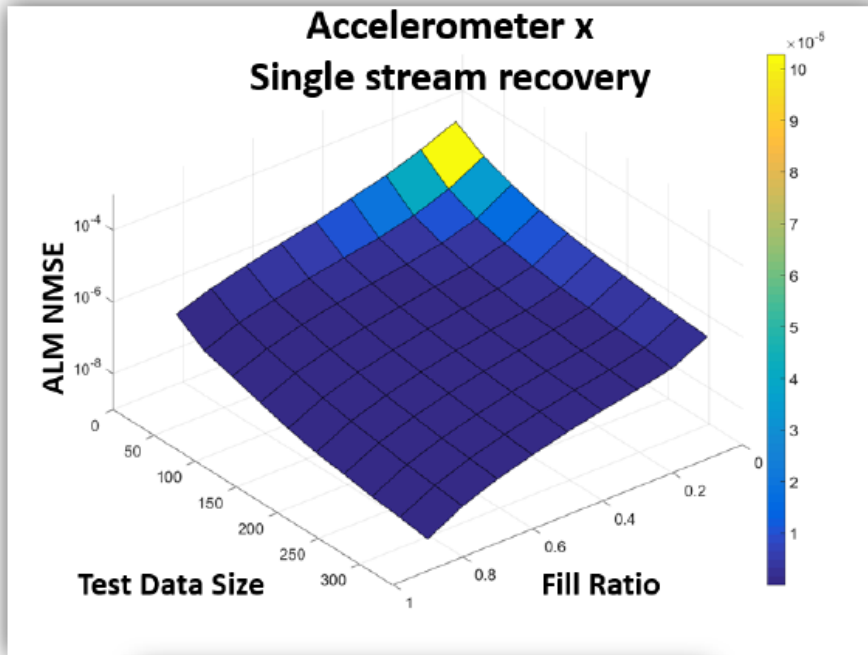
   ✔ $[n_1]$ lagged temporal windows of $[n_2]$ samples

3rd window

$[n_2]$

| $h_0$ | $h_1$ | ... | | $h_{n_1}$ | ... | $h_{n_2-1}$ | $h_{n_2}$ |
|---|---|---|---|---|---|---|---|
| $h_1$ | ... | | | | | | |
| ⋮ | | | ... | | ... | | ⋮ |

$[n_1]$

**H**

① **Test** sensor stream

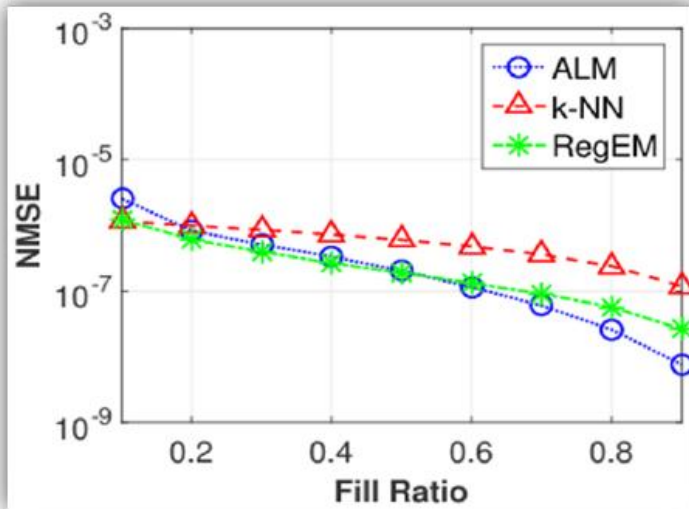② Introduction of missing values

③ Temporal windowing

④ Hankelization process **H**

⑤ Undersampled Hankel matrices that need to be reconstructed!

⇓

**Matrix Completion**

Data stream

Data stream with missing values

lag        2nd window

1st window

Hankel matrix with missing values

Completed Hankel matrix

CS-541 Wireless Sensor Networks
University of Crete, Computer Science Department

Accelerometer x
Single stream recovery

Gyroscope x
Single stream recovery

CS-541 Wireless Sensor Networks
University of Crete, Computer Science Department

# Autoregressive Models (AR)

Thus for *stationary* time series the mean value function is **constant** and the covariance function is only a **function of the distance in time** $(t - s)$

The "order" of the AR($p$) models is the number of prior values used in the model.

**Univariate AR** model

- **AR(1)**➜ $x_t = b_0 + b_1 x_{t-1} + \varepsilon_t$

- **AR(2)**➜ $x_t = b_0 + b_1 x_{t-1} + b_2 x_{t-2} + \varepsilon_t$

- **AR(p)**➜ $X_t = \sum_{i=1}^{p} \varphi_i X_{t-i} + \varepsilon_t.$

Solutions: Yule–Walker equations

Estimation of autocovariances, least squares regression

# Matrix formulation

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



$$\text{time} \quad \begin{bmatrix} \underline{X_{11}, X_{12}, \cdots, X_{1w}} \\ X_{21}, X_{22}, \ldots, X_{2w} \\ \vdots \\ \vdots \\ \vdots \\ X_{N1}, X_{N2}, \ldots, X_{Nw} \end{bmatrix} \times \begin{bmatrix} a_1 \\ a_2 \\ \vdots \\ a_w \end{bmatrix} = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ \vdots \\ y_N \end{bmatrix}$$
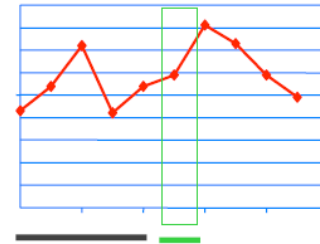
Ind-var1      Ind-var-w

# Matrix formulation

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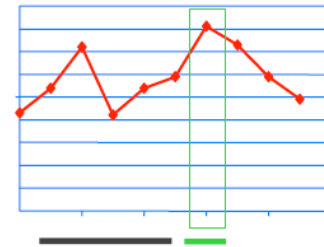


$$
\text{time} \quad \begin{bmatrix} X_{11}, X_{12}, \cdots, X_{1w} \\ X_{21}, X_{22}, \ldots, X_{2w} \\ \vdots \\ \vdots \\ \vdots \\ X_{N1}, X_{N2}, \ldots, X_{Nw} \end{bmatrix} \times \begin{bmatrix} a_1 \\ a_2 \\ \vdots \\ a_w \end{bmatrix} = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ \vdots \\ y_N \end{bmatrix}
$$

Ind-var1      Ind-var-w

CS-541 Wireless Sensor Networks
University of Crete, Computer Science Department

FORTH
Institute of Computer Science

# Vector Autoregressive Models (VAR)

**Vector AR (VAR)** extension to multiple time series

$$y_t = c + A_1 y_{t-1} + A_2 y_{t-2} + \cdots + A_p y_{t-p} + e_t,$$

- least squares: $B_1 = (Y_{t-1}^\mathsf{T} Y_{t-1})^{-1} Y_{t-1}^\mathsf{T} Y_t$ (under conditions)
- Determination of lag length is a trade-off

**Granger causality:** statistical hypothesis test for determining whether one time series X is useful in forecasting another time series Y, ('60)

$$Y_t = \alpha + \phi_1 Y_{t-1} + \beta_1 X_{t-1} + e_t$$

"**if** $\beta_1 = 0$ **then past values of X have no explanatory power for Y beyond that provided by past values of Y**".

CS-541 Wireless Sensor Networks
University of Crete, Computer Science Department

FORTH
Institute of Computer Science

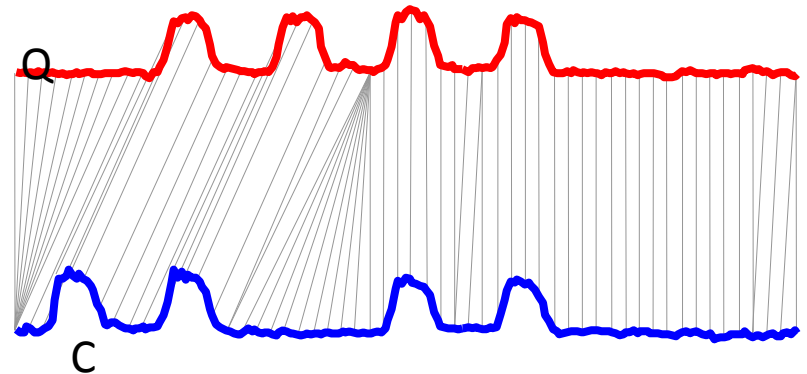# Similarity between time-Series
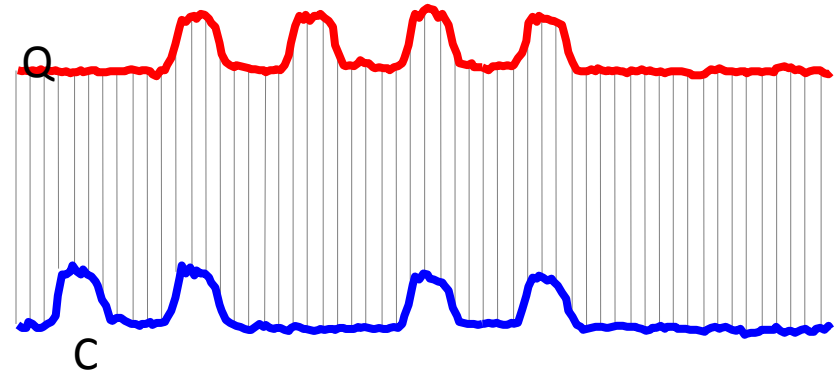
## Euclidean Distance

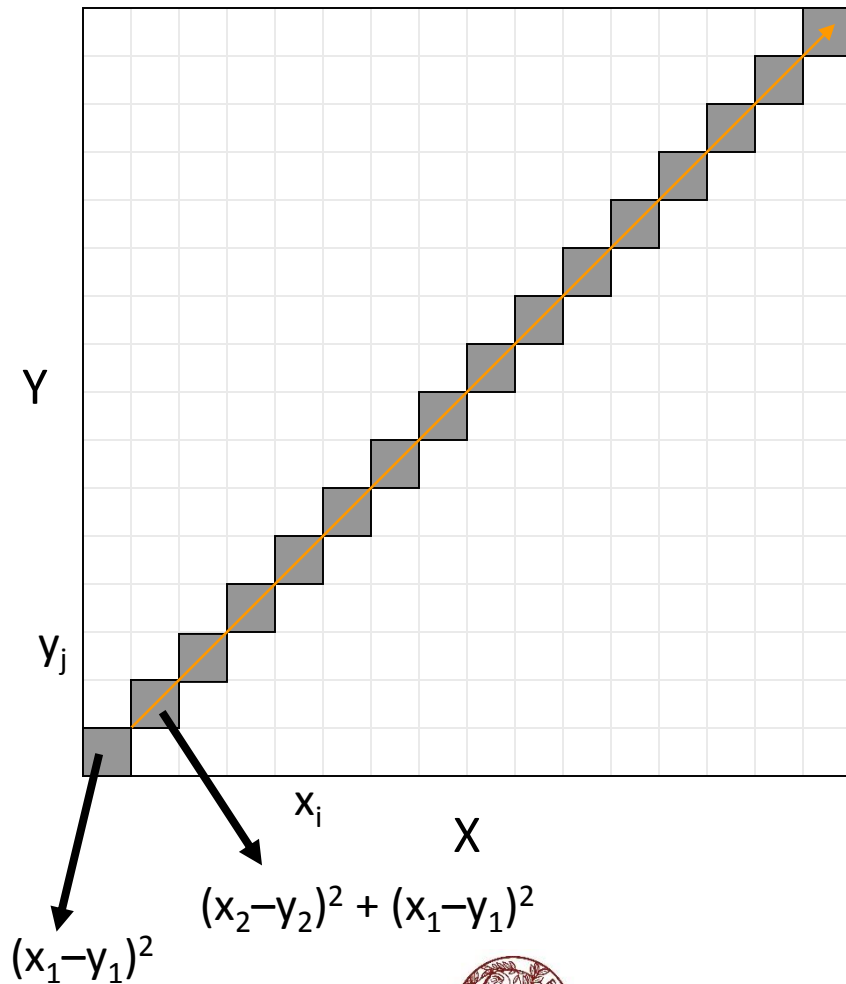$$D(\vec{x}, \vec{y}) = \sum_{i=1}^{n} (x_i - y_i)^2$$

(+) *Efficient computation*

(-) *Time shift, scaling*

## Dynamic Time Warping

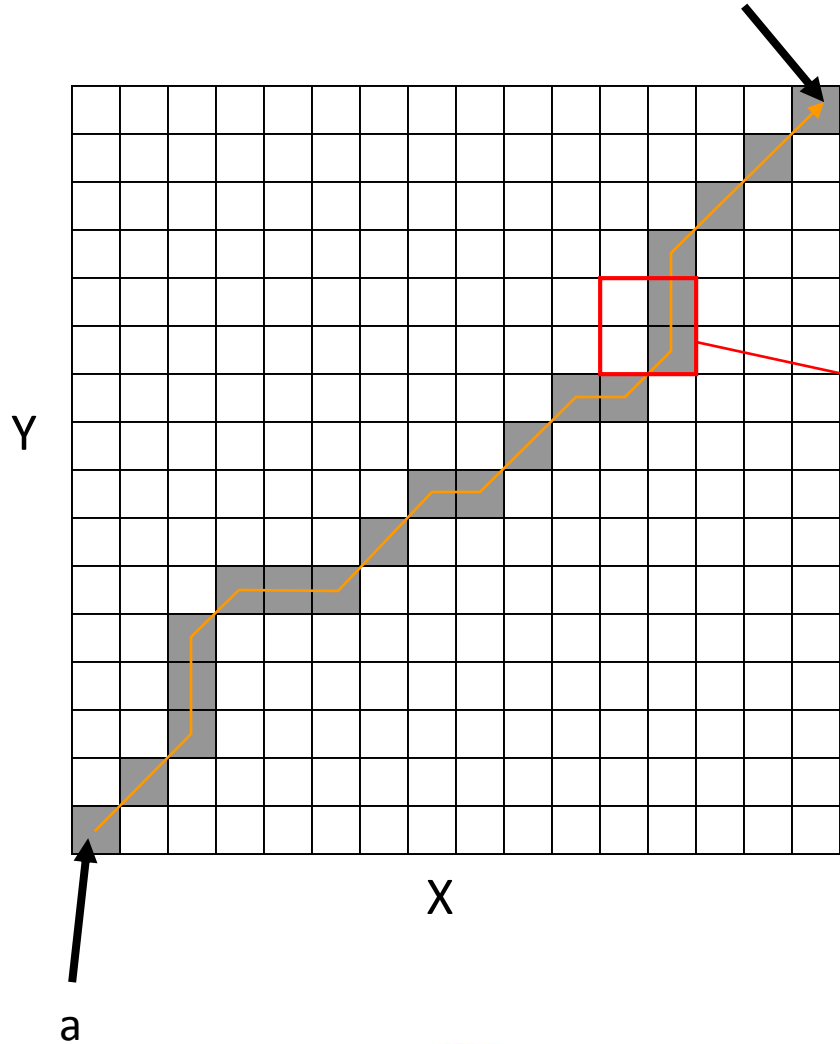• *Nonlinear alignments are possible.*

CS-541 Wireless Sensor Networks
University of Crete, Computer Science Department

FORTH
Institute of Computer Science

# DTW: Euclidean Distance



$Y$

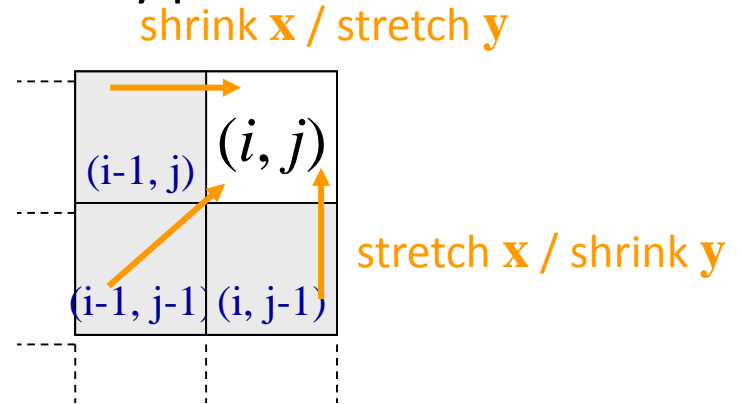$y_j$

$x_i$

$X$

$(x_2-y_2)^2 + (x_1-y_1)^2$

$(x_1-y_1)^2$

- Each cell $c = (i, j)$ is a pair of indices whose corresponding values will be computed, $(x_i-y_j)^2$, and included in the sum for the distance.

- Euclidean path:
  - $i = j$ always.
  - Ignores off-diagonal cells.

CS-541 Wireless Sensor Networks
University of Crete, Computer Science Department

# DTW: Dynamic time warping

b

DTW allows any path.

shrink **x** / stretch **y**

$(i, j)$

$(i-1, j)$

stretch **x** / shrink **y**

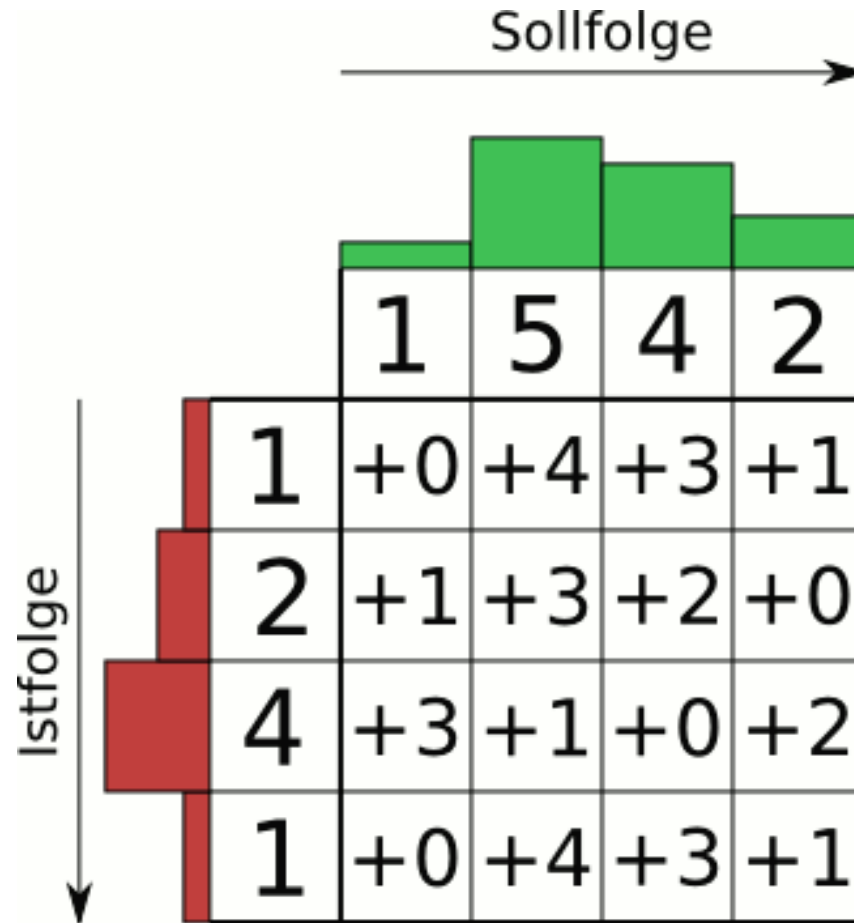$(i-1, j-1)$  $(i, j-1)$

Y

X

a

Dynamic Programming approach

$$D(i, j) = |\, x_i - y_j \,| + \min \{\quad D(i-1, j),$$
$$D(i-1, j-1),$$
$$D(i, j-1) \}$$

- Extend sequences by repeating elements
- Euclidean distance between extended sequences
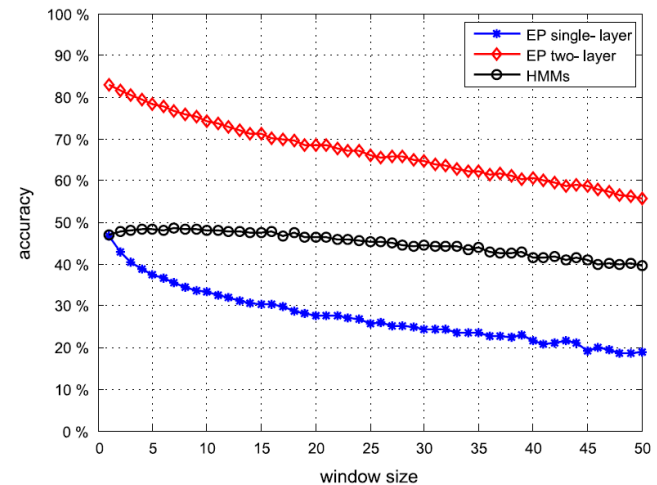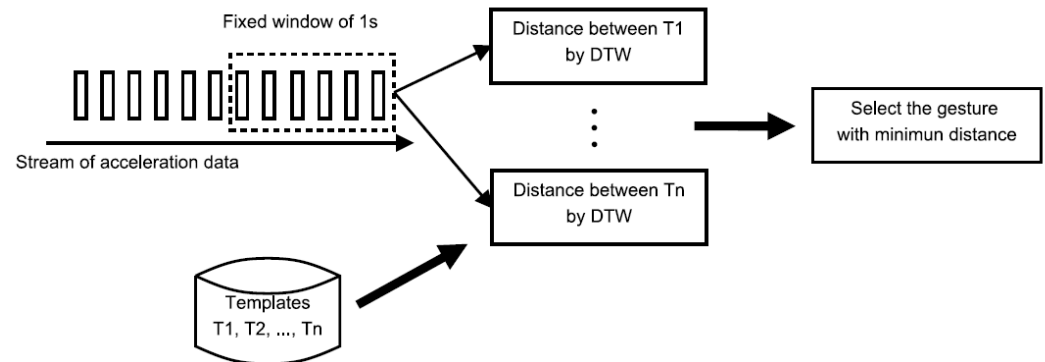
FORTH
Institute of Computer Science
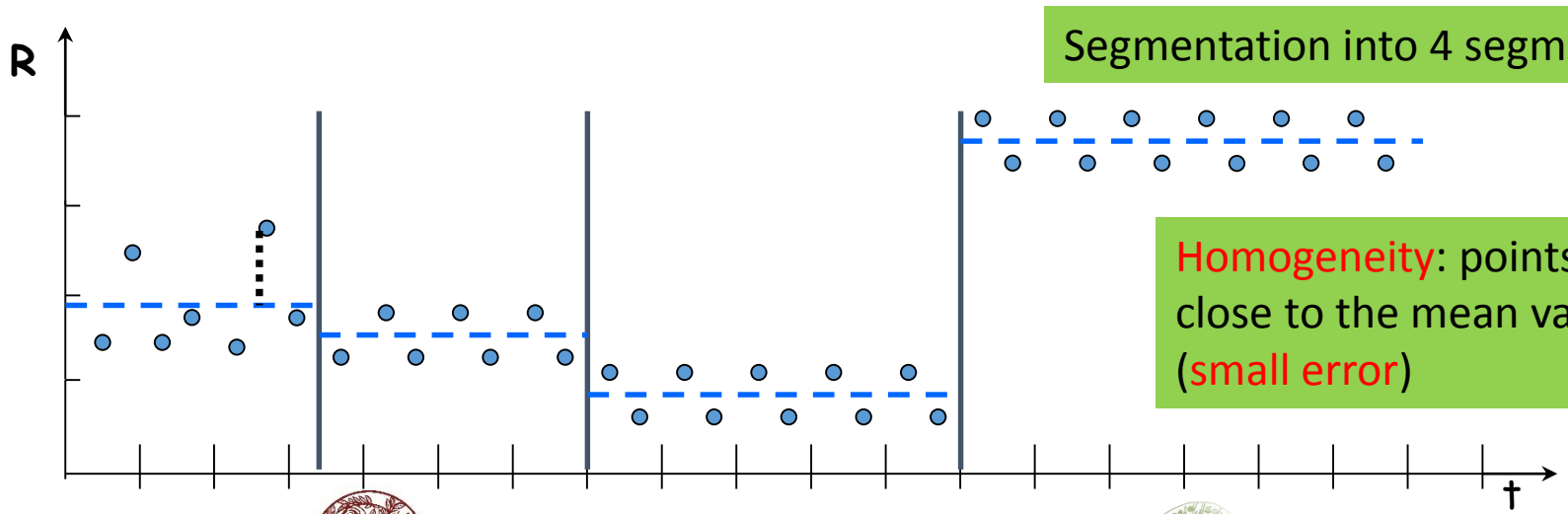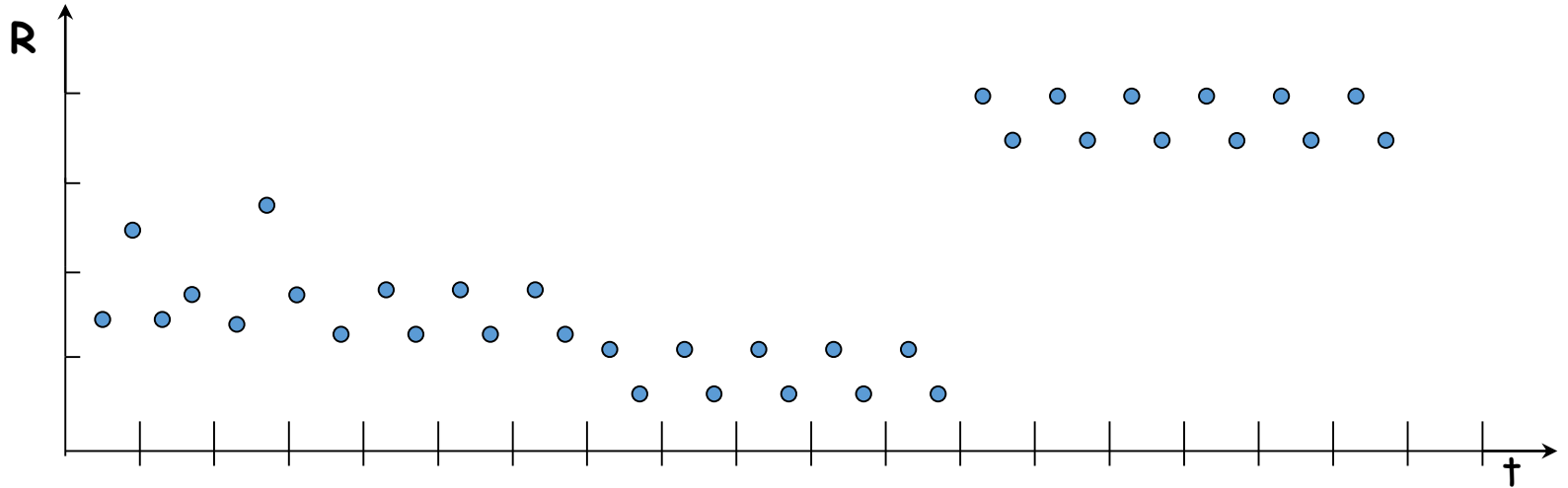
# DTW example

# DTW based activity recognition



Wang, Liang, et al. "A hierarchical approach to real-time activity recognition in body sensor networks." *Pervasive and Mobile Computing* 8.1 (2012): 115-130.

# Stream Data Processing



Segmentation into 4 segments

Homogeneity: points are close to the mean value (small error)

# The K-segmentation problem

- A K-segmentation S: a partition of T into K contiguous segments $\{s_1, s_2, ..., s_K\}$.

- Similar to K-means clustering, but now we need the points in the clusters to respect the order of the sequence

> Given a sequence **T** of length **N** and a value **K**, find a **K**-segmentation **S** = $\{s_1, s_2, ...,s_K\}$ of **T** such that the **SSE** error **E** is minimized.

Solve via Dynamic Programming:

- Construct the solution of the problem by using solutions to problems of smaller size

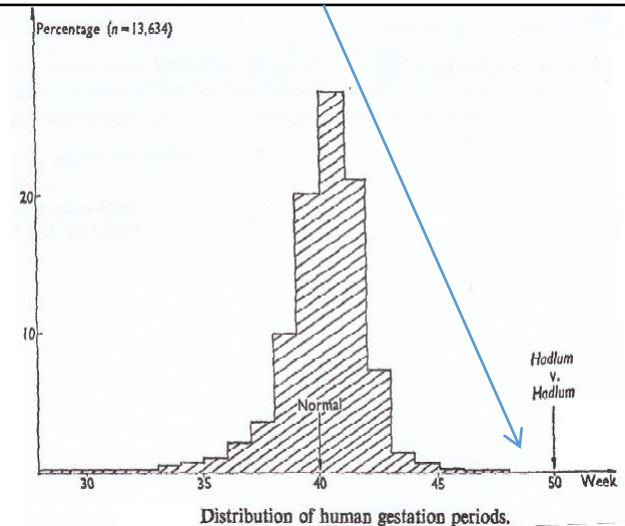- Build the solution bottom up from smaller to larger instances

# Outlier detection

Definition (anomaly/novelty detection)

"those measurements that significantly deviate from the normal pattern of the sensed data"

Types: Noise, Errors, Events & Attacks

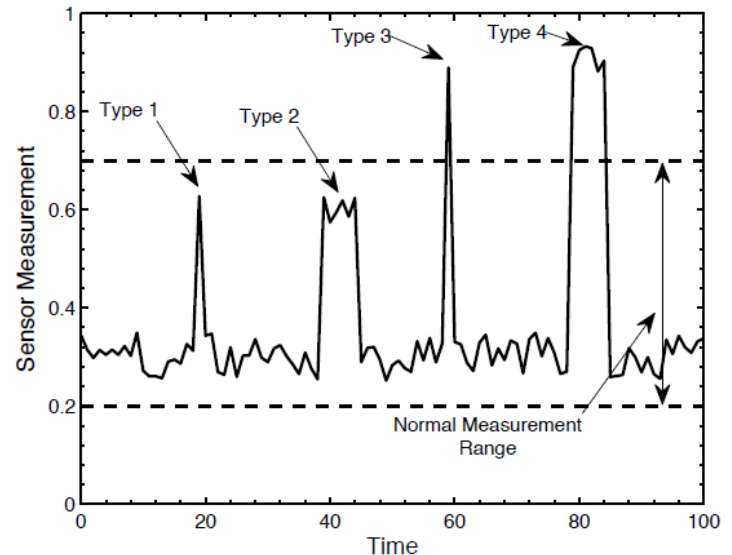| Outlier Detection | Event Detection |
|---|---|
| No prior knowledge | semantics |
| comparative | Threshold based |
| False alarms | Detection |

The birth of a child to Mrs. Hadlum happened 349 days (11,5 months) after Mr. Hadlum left for military service.



Distribution of human gestation periods.

# Types of outliers

- First Order Anomalies:
  - Partial data measurements are anomalous at a sensor node

- Second Order Anomalies:
  - All data measurements at a sensor node are anomalous

- Third Order Anomalies:
  - Data from a set of sensor nodes are anomalous

Type 1: Incidental absolute errors:
- A short-term extremely high anomalous

Type 2: Clustered absolute errors:
- A continuous sequence of *type 1* errors

Type 3: Random errors:
- Short-term observations outside normal range

Type 4: Long term errors:
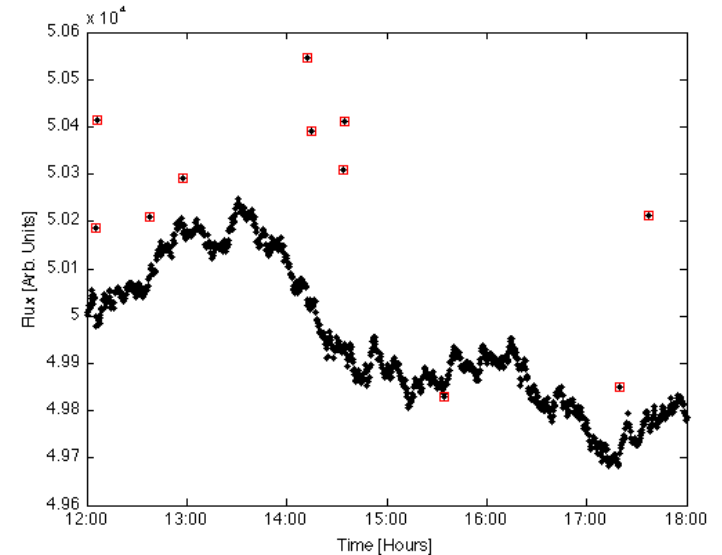- A continuous sequence of *type 3* errors

# Outlier detection in WSNs

Objectives

➢ Data reliability

➢ Quality of Service

➢ Communications overhead

➢ Adaptive sampling rates
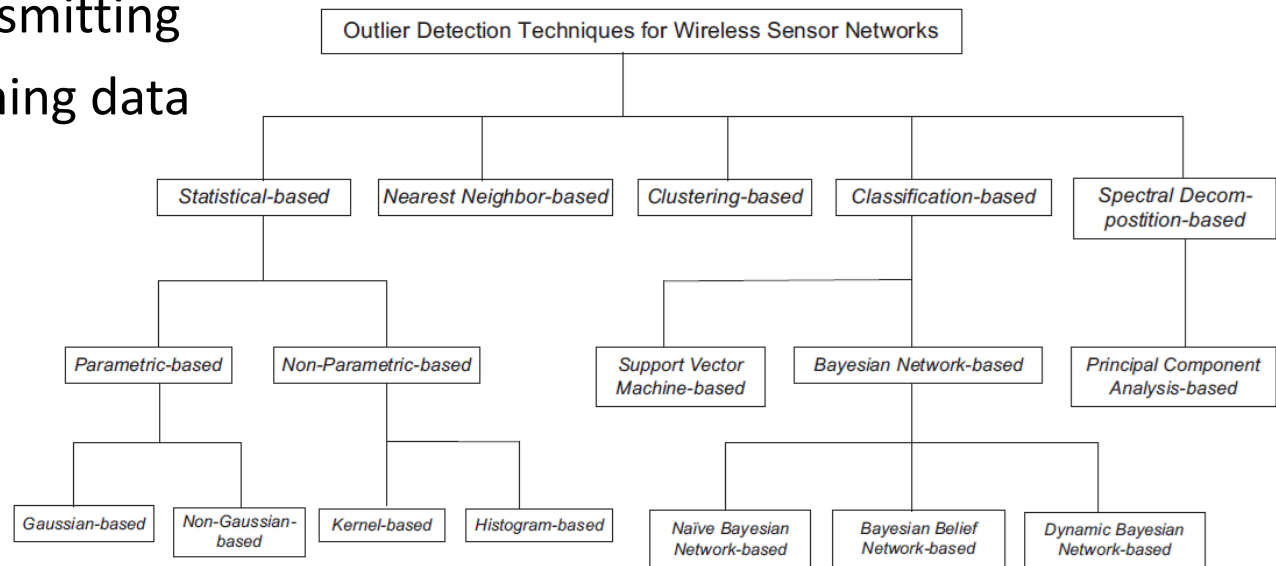
➢ Security alert

Applications

• Environmental monitoring (e.g. fire)

• Health monitoring (e.g. heart attack)

• Industrial monitoring (e.g. malfunctions)

# Outlier detection in WSNs

Challenges

- Low cost & quality
- Processing vs Transmitting
- Distributed streaming data
- Network topology
  - Failures,
  - Disconnections,
  - Mobility
- Deployment scale
- Type detection

# Statistical

Gaussian-based models

• Send measurements -> model
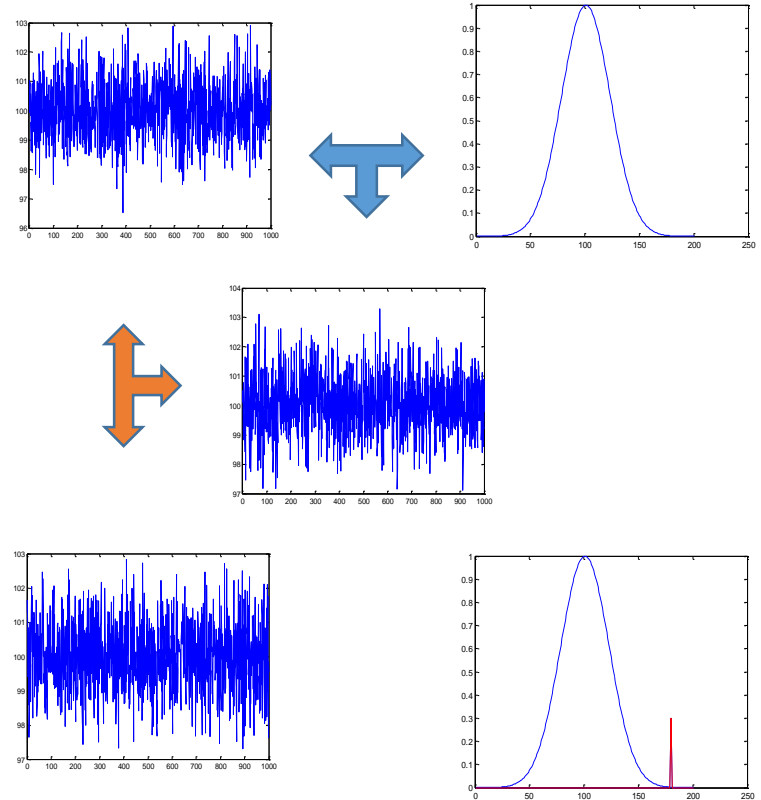
• Build model -> send parameters

Non-Gaussian

• Symmetric α-stable distributions

Mixtures

Clusters

**Detection Thresholds**

CS-541 Wireless Sensor Networks
University of Crete, Computer Science Department

# Non-parametric modeling

## Histogram based

1. Obtain $v_{min}$ and $v_{max}$ information

2. Collect histogram

3. Collect outliers and potential outliers

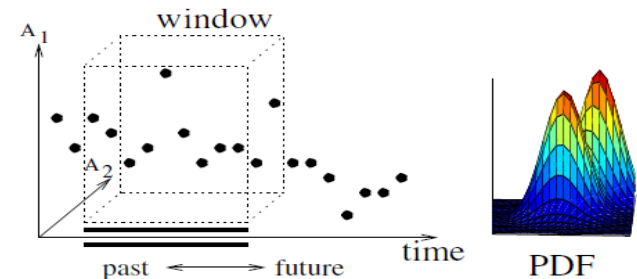4. Diffuse potential outliers and count the number of neighbors within d
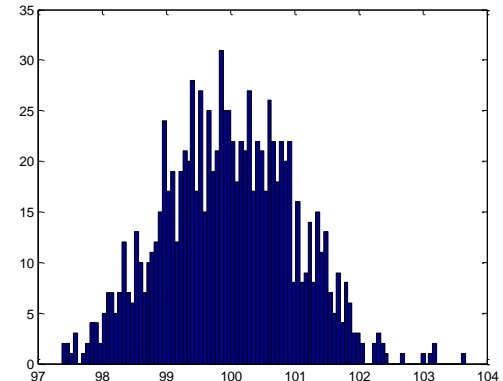
- Number of bins

- Thresholds

## Kernel Density Estimation

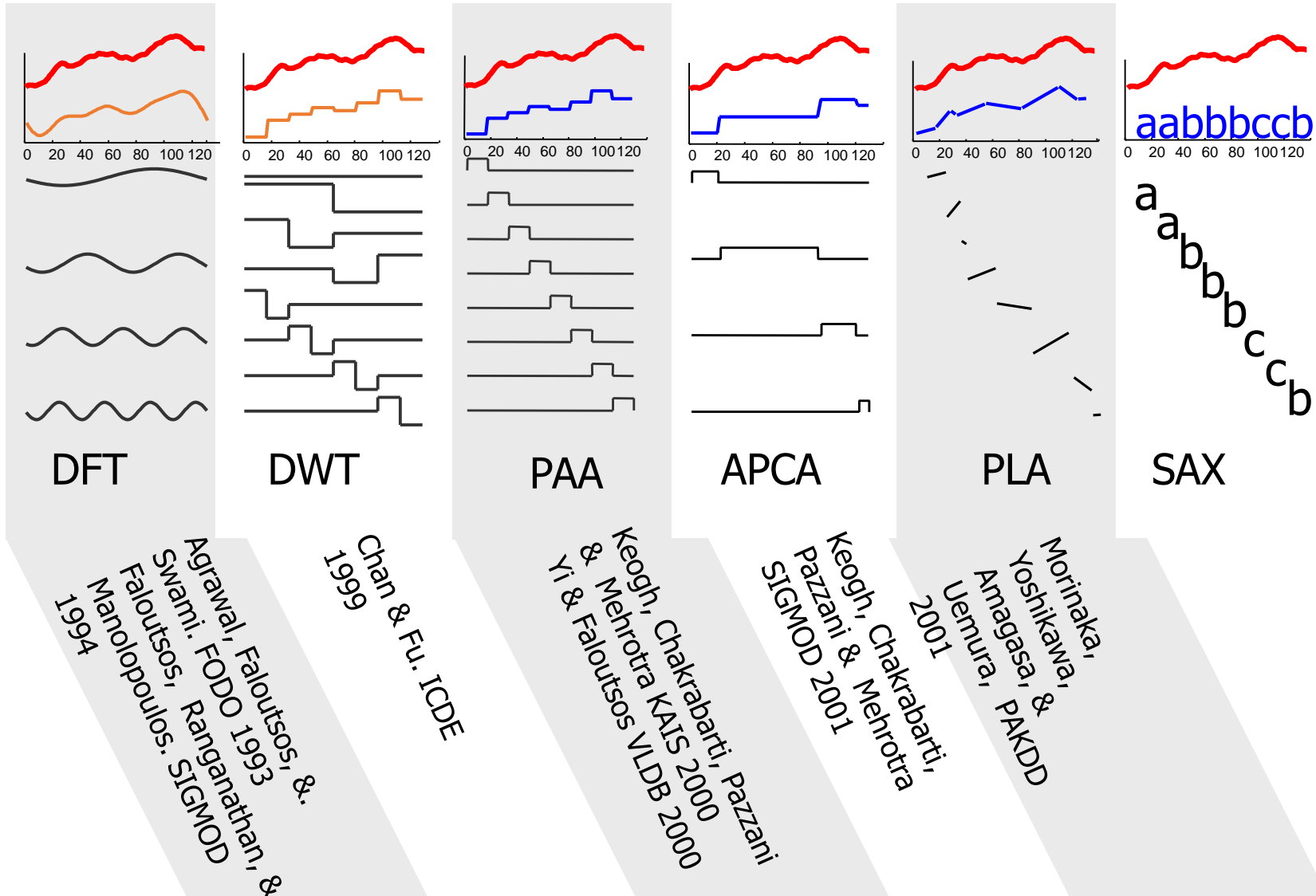$$f(x) = \frac{1}{n} \sum_{i=1}^{n} \frac{1}{h_i} K(\frac{x - x_i}{h_i})$$

Kernel

Bandwidth

Gaussian    $K(u) = \frac{1}{\sqrt{2\pi}} e^{-\frac{1}{2}u^2}$

# Time series



DFT      DWT      PAA      APCA      PLA      SAX

aabbbccb

Agrawal, Faloutsos, & Swami. FODO 1993
Faloutsos, Ranganathan, & Manolopoulos. SIGMOD 1994
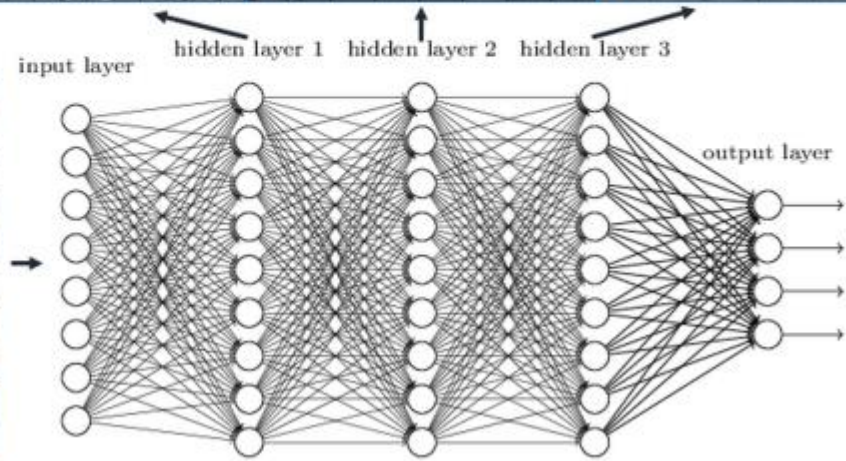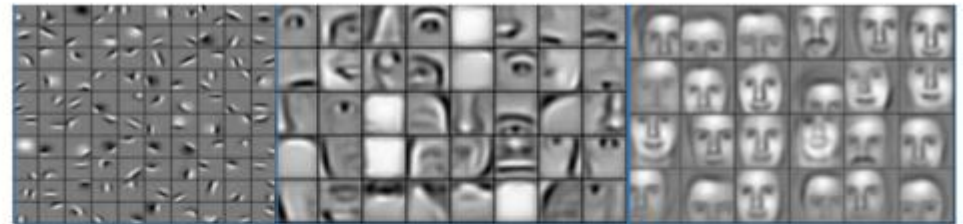
Chan & Fu. ICDE 1999

Keogh, Chakrabarti, Pazzani & Mehrotra KAIS 2000
Yi & Faloutsos VLDB 2000

Keogh, Chakrabarti, Pazzani & Mehrotra SIGMOD 2001

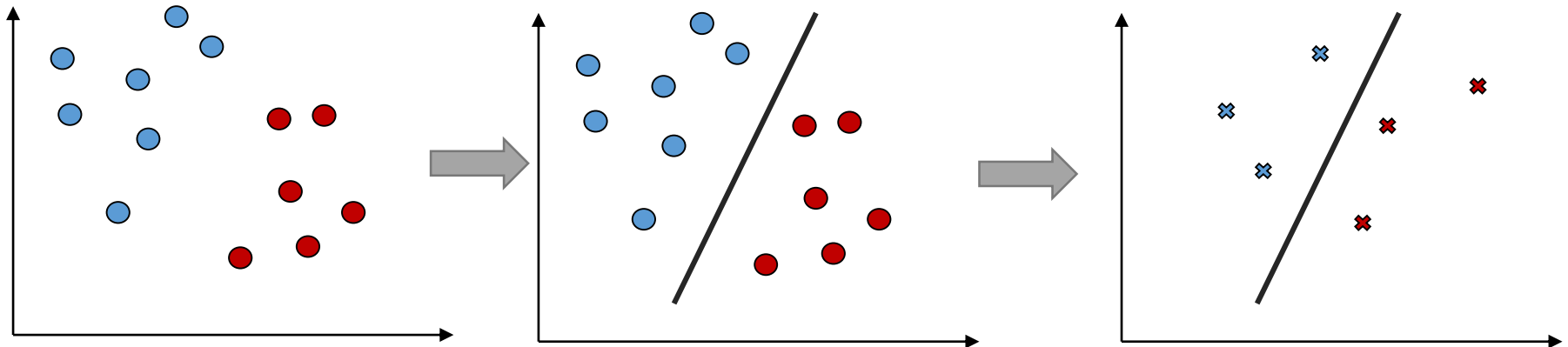Morinaka, Yoshikawa, Amagasa, & Uemura, PAKDD 2001

# Machine Learning



input layer — hidden layer 1 — hidden layer 2 — hidden layer 3 — output layer

# Machine Learning

Machine learning: construction and study of [algorithms](algorithms) that can [learn](learn) from data

- Models of example inputs (training data) → make predictions or decisions on new inputs (testing data)
- Data: characteristics
- Prior assumptions: a priori knowledge
- Representation: How do we represent the data
- Model / Hypothesis space: Hypotheses to explain the data
- Feedback / learning signal: Learning signal (delayed, labels)
- Learning algorithm: Model update
- Evaluation: Check quality

CS-541 Wireless Sensor Networks
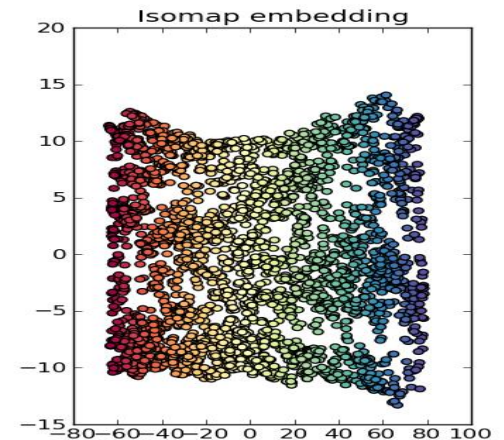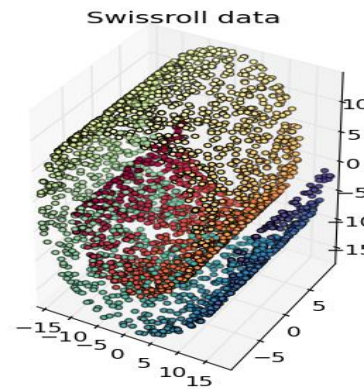University of Crete, Computer Science Department

# Types of ML

**Supervised learning:** present example inputs and their desired outputs (**labels**) →  learn a general rule that maps inputs to outputs.

CS-541 Wireless Sensor Networks
University of Crete, Computer Science Department

FORTH
Institute of Computer Science

# Types of ML

**Unsupervised learning:** no labels are given → find structure in input.



Swissroll data

Isomap embedding

CS-541 Wireless Sensor Networks
University of Crete, Computer Science Department

FORTH
Institute of Computer Science

# Types of ML

**Reinforcement learning:** system interacts with environment and must perform a certain goal without explicitly telling it whether it has come close to its goal or not.

CS-541 Wireless Sensor Networks
University of Crete, Computer Science Department

# Applications in WSNs

## Network performance optimization

- Routing

- Distributed regression framework

- Data Aggregation

- Localization and Objects Targeting

- Medium Access Control

## Data Mining

- Activity recognition

- Event Detection and Query Processing

# Unsupervised learning - Clustering

*What is a cluster*?

groups of data instances that are similar to each other in one cluster and data instances that are very different from each other into different clusters

Hard vs. Soft

- *Hard*: belong to single cluster
- *Soft*: belong to multiple clusters

Flat vs. Hierarchical

- *Flat*: clusters are flat
- *Hierarchical*: clusters form a tree

CS-541 Wireless Sensor Networks
University of Crete, Computer Science Department

FORTH
Institute of Computer Science

# K-means clustering

- K-means is a <span style="color:red">partitional clustering</span> algorithm
- Let the set of data points (or instances) *D* be

  $\{\mathbf{x}_1, \mathbf{x}_2, ..., \mathbf{x}_n\}$,

  where $\mathbf{x}_i = (x_{i1}, x_{i2}, ..., x_{ir})$ is a <span style="color:blue">vector</span> in a real-valued space $X \subseteq R^r$, and *r* is the number of attributes (dimensions) in the data.

- The *k*-means algorithm partitions the given data into *k* clusters.
  - Each cluster has a cluster **center**, called <span style="color:red">**centroid**</span>.
  - *k* is specified by the user

CS-541 Wireless Sensor Networks University of Crete, Computer Science Department

# K-means algorithm

Given *k*, the *k-means* algorithm works as follows:

1) Randomly choose *k* data points (seeds) to be the initial centroids, cluster centers
2) Assign each data point to the closest centroid
3) Re-compute the centroids using the current cluster memberships.
4) If a convergence criterion is not met, go to 2).

Stopping criteria

- no re-assignments of data points to different clusters
- no change of centroids
- minimum decrease in the $SSE = \sum_{j=1}^{k} \sum_{\mathbf{x} \in C_j} dist(\mathbf{x}, \mathbf{m}_j)^2$
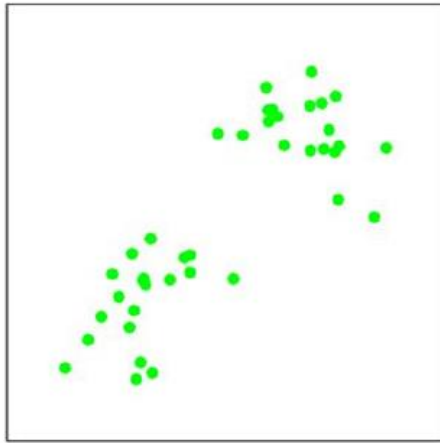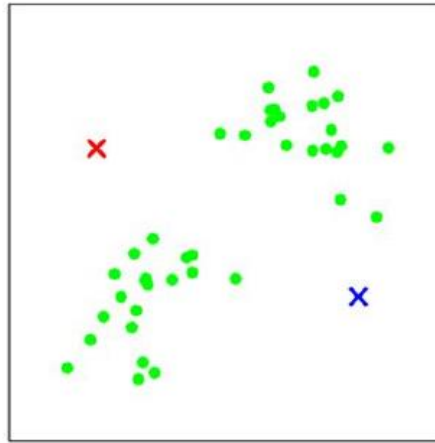
# K-means example

Complexity is O( n * K * I * d )
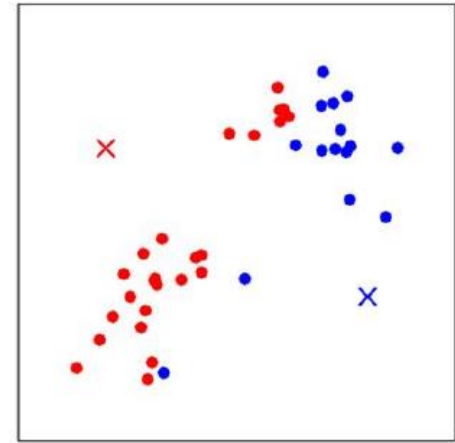n = number of points, K = number of clusters,
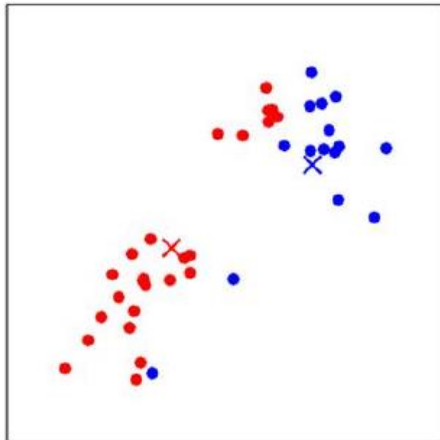I = number of iterations, d = dimensionality



(a)

(b)

(c)

(d)

(e)

(f)

CS-541 Wireless Sensor Networks
University of Crete, Computer Science Department
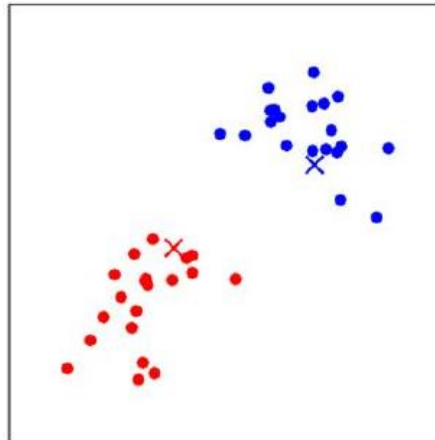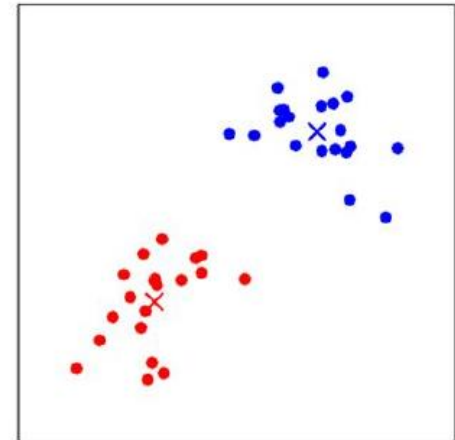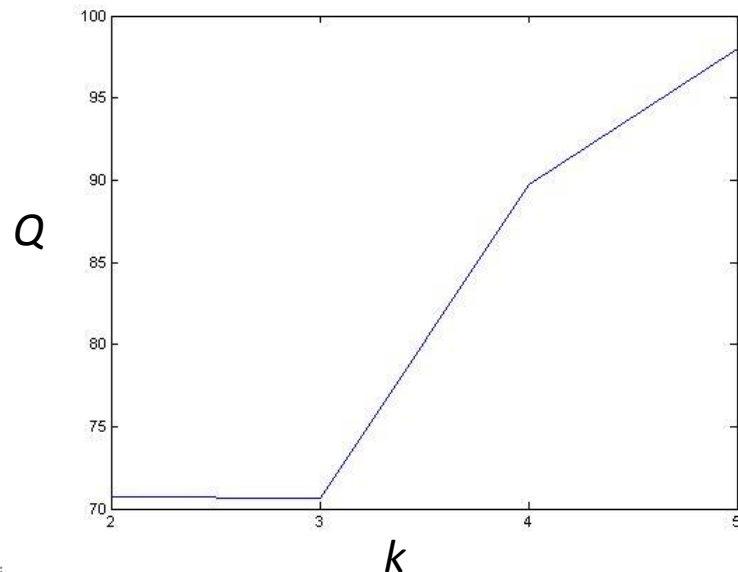
# Issues with K-means

- Random initialization -> different clusters each time
- Data points are assigned to only one cluster
- Implicit assumptions about the "shapes" of clusters
- You have to pick the number of clusters…

Cluster tightness

$$Q = \sum_{i=1}^{k} \frac{1}{|C_i|} \sum_{\boldsymbol{x} \in C_i} d(\boldsymbol{x}, \mu_i)$$

CS-541 Wireless Sensor Networks
University of Crete, Computer Science Department
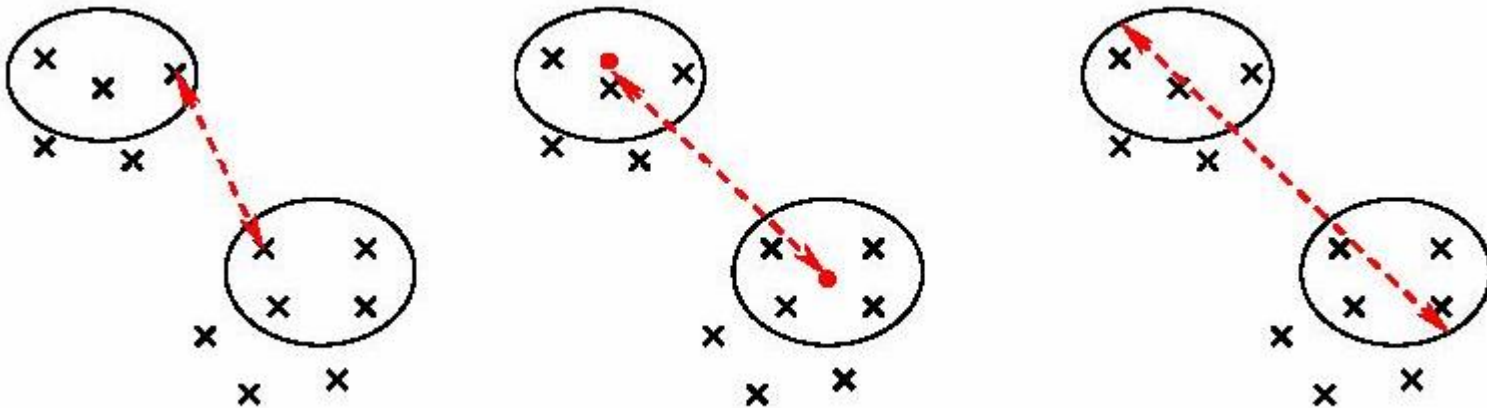
FORTH
Institute of Computer Science

# Distance Between Two Clusters

single-link clustering: distance between clusters -> shortest distance between any two members.
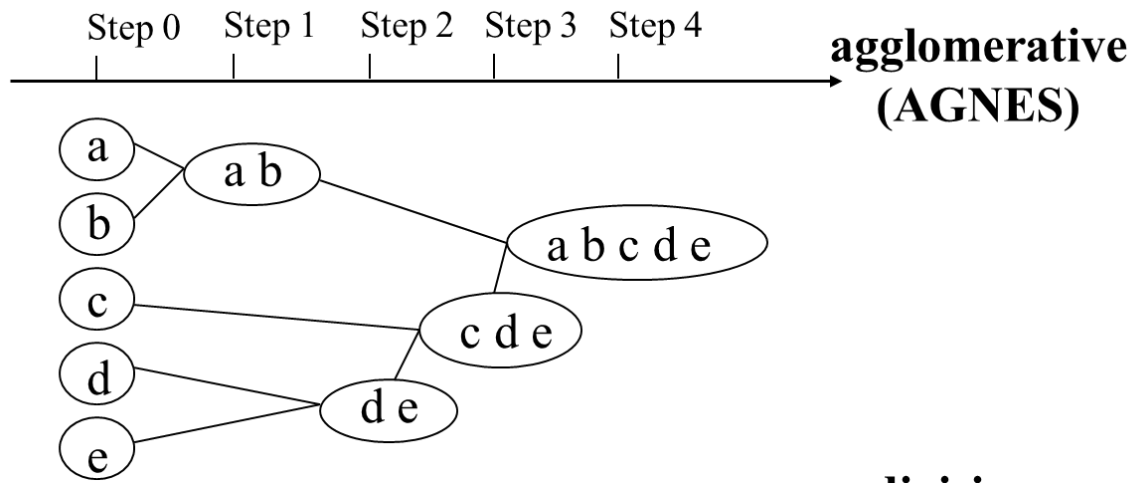
 complete-link clustering: distance between clusters -> longest distance between any two members.

average-link clustering: distance between clusters  -> average distance between any two members
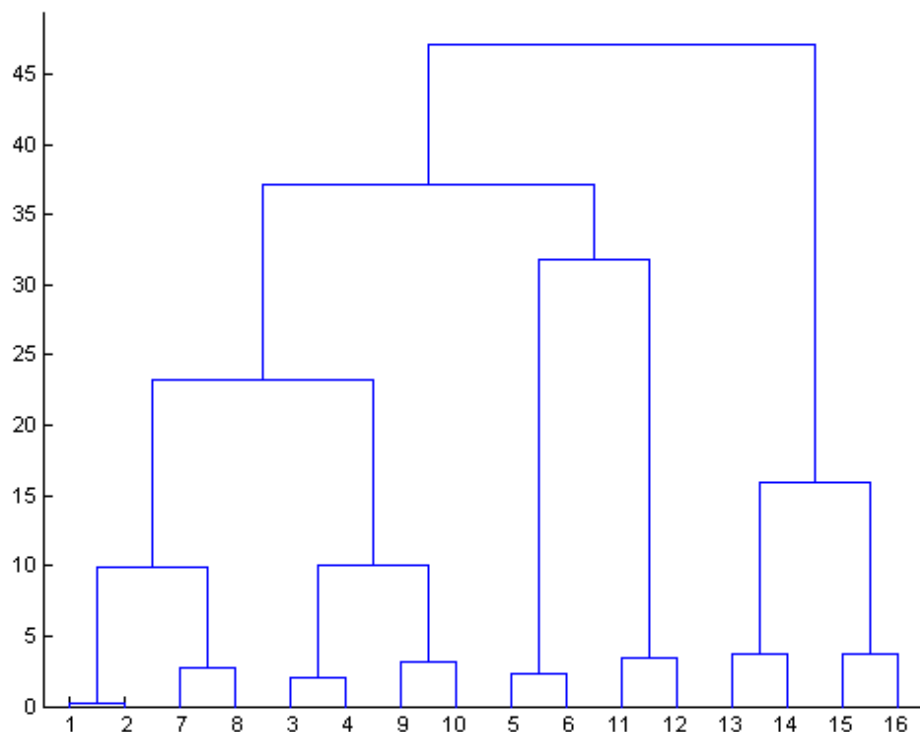
# Hierarchical Agglomerative Clustering

- We start with every data point in a separate cluster
- We keep merging the most similar pairs of data points/clusters until we have one big cluster left
- This is called a bottom-up or agglomerative method



CS-541 Wireless Sensor Networks
University of Crete, Computer Science Department
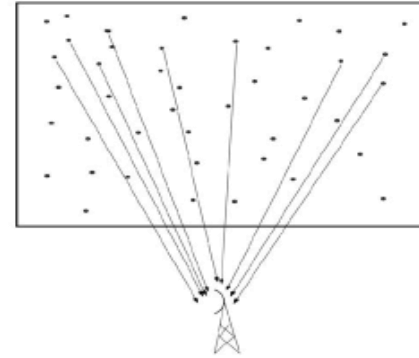
# Hierarchical Clustering (cont.)



- This produces a binary tree or **dendrogram**
- The final cluster is the root and each data item is a leaf
- The height of the bars indicate how close the items are

CS-541 Wireless Sensor Networks
University of Crete, Computer Science Department
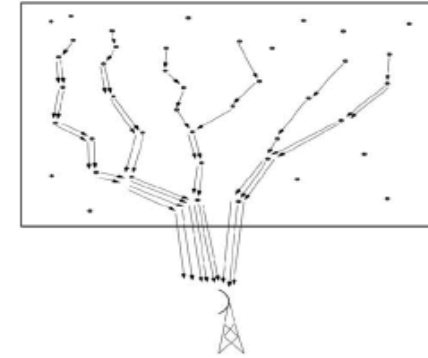
FORTH
Institute of Computer Science
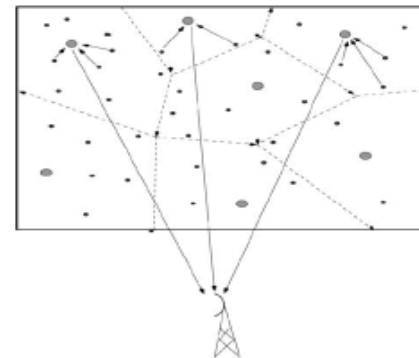
# Clustering in WSN

- Scalability:
  - Reduce routing tables to within cluster
- Data Aggregation
  - Energy reduction vs. full data transmission
  - CH based data fusion
  - multi-hop tree structure aggregation
- Load Balancing
  - Eliminate redundant data transmissions
  - Communications between CHs
- Energy reduction
  - Selective sampling within cluster
  - Short-range communications with CH
- Robustness & Fault tolerance
  - Support node failure/recovery
  - mobility of sensors
  - noisy measurements etc.
- Efficiency
  - Collision avoidance (intra vs. inter cluster communications)
  - Latency reduction by reducing hops
  - Network life-time maximization
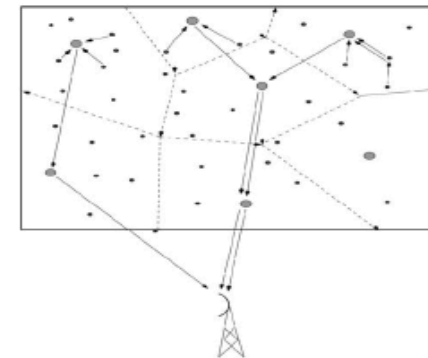  - Quality-of-service



(a)          (b)

(c)          (d)

CS-541 Wireless Sensor Networks
University of Crete, Computer Science Department

FORTH
Institute of Computer Science

# Reading List

- Esling, Philippe, and Carlos Agon. "Time-series data mining." *ACM Computing Surveys (CSUR)* 45.1 (2012): 12.

CS-541 Wireless Sensor Networks
University of Crete, Computer Science Department

FORTH
Institute of Computer Science