Time Series Data Mining Challenges

Jose A. Lozano

Basque Center for Applied Mathematics (BCAM) University of the Basque Country UPV/EHU

EVIA, Sevilla 14 Junio, 2023



(日)

Outline of the presentation



Time Series Data Mining Activities

2 Clustering

- (Early) Supervised Classification
- Outlier/Anomaly Detection
- 5 Conclusions and Future Work



Outline of the presentation

Time Series Data Mining Activities

2 Clustering

- 3 (Early) Supervised Classification
- Outlier/Anomaly Detection
- 5 Conclusions and Future Work



Time series are all around













Definition and main characteristics

Definition

A time series is an ordered sequence of pairs of finite length L:

$$TS = \{(t_i, x_i) | i = 1, ..., L\},$$
(1)

where the timestamps $\{t_i\}_{i=1}^{L}$ take positive and ascending real values and the values of the time series (x_i) take univariate or multivariate real values.

Main characteristics

- Temporal correlation
- High dimensionality
- Noise

5.7

a center for soniied methemat

(日)

Time series forecasting



Examples

- Stock market prediction
- Temperature prediction



Time series data base: our object of study



- A set of time series (usually big)
- Different lengths
- Multidimensional



Time series clustering. Examples







Supervised classification of time series



Anomaly/outlier detection





Segmentation





・ロト ・ 理 ト ・ 理 ト ・ 理 ト



æ

Outline of the presentation

Time Series Data Mining Activities

2 Clustering

- 3 (Early) Supervised Classification
- Outlier/Anomaly Detection
- 5 Conclusions and Future Work



Time Series Data Mining Challenges Clustering

Time series clustering. Examples







▲ロト ▲舂 ト ▲ 恵 ト ▲ 恵 ト ● 恵 … 釣ぬ()

Differences with the classic clustering problem



Time Series Data Mining Challenges Clustering

Time series clustering: hierarchical, partitional



DISTANCE



◆ロ▶ ◆課 ▶ ◆理 ▶ ◆理 ▶ ─ 理 → ∽ !

Time Series Data Mining Challenges Clustering

Distance between time series

Rigid Distance

Flexible Distance





(日)



Euclidean Distance (ED)

$$D(X,Y) = \sqrt{\sum_{i=1}^{n} (x_i - y_i)^2}$$



- Easy to compute
- Only for series with the same distance \checkmark
- Does not consider the time
- Sensitivity to noise
- Requires the series to be sampled at the time stamps



Time Series Data Mining Challenges Clustering

Dynamic Time Warping (DTW)



- Takes into account the ordered sequence (time)
- It can deal with series of different sizes
- Does not require to be sampled in the same time stamps

Image: A math a math

 Computationally expensive O(min{m, n}²) ✓



Euclidean Distance vs Dynamic Time Warping





DTW

・ロト ・回ト ・回ト

taspe onle for appled mathematics

EUCLIDEAN

More on elastic distances

Cheap versions of dynamic time warping







▲ロト▲舂▶▲臣▶▲臣▶ 臣 のの

More on elastic distances

Edit distances for real sequences (LCSS, EDR, ERP)



Mori et al., 2016, R journal



ъ

・ ロ ト ・ 一 マ ト ・ 日 ト

More on elastic distances

On-line versions (Oregi et al 2019, PR)



(a) Incremental DTW computation.



(b) Weights over the DTW lattice.



Alternatives to calculating distances

- Represent each series by means of a set of features and calculate the distance between the features
- Learn a parametric model for each series and calculate the distance between the parameters



Last remarks on distances between series

Remarks

- There is no best distance (no free lunch)
- Each problem requires a different distance
- The distance to be used needs to be in agreement with our knowledge about what is far and what is close
- Hint: try several distances

Challenge:

Design a method to the (semi)automatic selection of a distance (e.g. Mori et al. 2016, TKDE)



...back to clustering: problems with k-means



basque center for applied mathematics

・ ロ ト ・ 一 マ ト ・ 日 ト

Remarks on clustering

- Recent papers on the computation of a mean series
- Alternate clustering methods: graph-based, spectral, model-based,...
- Multivariate time series clustering



Outline of the presentation

- 1 Time Series Data Mining Activities
- 2 Clustering
- (Early) Supervised Classification
 - 4 Outlier/Anomaly Detection
- 5 Conclusions and Future Work



Supervised classification of time series





(日)

General-purpose classifiers

- Each series is considered an instance
- Each time stamp is considered a feature

t ₁	t ₂	t ₃	 tn	С
<i>x</i> ₁₁	<i>x</i> ₁₂	<i>x</i> ₁₃	 x _{1n}	<i>C</i> ₁
<i>x</i> ₂₁	X ₂₂	<i>X</i> 23	 x _{2n}	<i>C</i> ₂
•••		•••	 	
<i>x</i> _{m1}	<i>x</i> _{m2}	x _{m3}	 x _{mn}	<i>C</i> ₂



General-purpose classifiers

- Each series is considered an instance
- Each time stamp is considered a feature

t ₂	t ₁	t ₃	 t _n	C
<i>x</i> ₁₂	<i>x</i> ₁₁	<i>x</i> ₁₃	 x _{1n}	<i>C</i> ₁
<i>x</i> ₂₂	<i>x</i> ₂₁	<i>x</i> ₂₃	 х _{2п}	<i>C</i> ₂
x _{m2}	<i>x</i> _{m1}	x _{m3}	 x _{mn}	<i>C</i> ₂



General-purpose classifiers

- Each series is considered an instance
- Each time stamp is considered a feature

<i>t</i> 2	t ₁	t ₃		tn	С
<i>x</i> ₁₂	<i>x</i> ₁₁	<i>x</i> ₁₃		x _{1n}	<i>C</i> ₁
<i>X</i> 22	<i>x</i> ₂₁	<i>x</i> ₂₃	•••	x _{2n}	<i>C</i> ₂
			• • •		
<i>x</i> _{m2}	<i>x</i> _{m1}	x _{m3}		x _{mn}	<i>C</i> ₂

CHALLENGE

When to use general-purpose and when time-series specific?



• • • • • • • • • •

What is relevant in TSC?

PROBLEM I



PROBLEM II





What is relevant in TSC?

PROBLEM I



SHAPE

PROBLEM II





◆□▶ ◆□▶ ◆□▶ ◆□▶ ●□ ●

What is relevant in TSC?

PROBLEM I



SHAPE

PROBLEM II





A taxonomy of time series classification methods

Taxonomy

- Distance-based classifiers
- Model-based classifiers
- Feature-based classifiers



Taxonomy of distance-based TSC (Abanda et al. 2019, DAMI)




1-Nearest Neighbour (1-NN)



1-Nearest Neighbour (1-NN)



Distance-based. Distance features. Global



< ロ > < 回 > < 回 > < 回 > < 回 >

э

Distance-based. Distance features. Global





<ロト <回 > < 注 > < 注 >

Distance-based. Distance features. Global



- Any general-purpose algorithm could be applied
- It depends on the number of series in training set

 Computationally expensive

Distance-based. Distance features. Local





Distance-based. Embedding



Training set





C₁

 C_3

 C_n

Distance-based. Embedding



 Many classifiers defined in Euclidean spaces

A B > A B >

- Computational complexity
- Prediction



Distance Kernels





Kernel matrix



$$f(x) = \sum \alpha_i \kappa(x, x_i)$$

Pattern function

asoue center for applied mathematics

UPV - EHU

Definite (PSD) Kernel

- All the SVM machinery works 🗸
- Difficult to define/check

Indefinite

- Theoretical properties are lost 🗸
- Easy to define
- Some methods can not be applied \checkmark

Distance Kernels. Indefinite

Gaussian Distance Substitution Kernels

$$GDS_d(x, x') = exp\left(-rac{d(x, x')^2}{\sigma^2}
ight)$$
 where $d = DTW, ...$



ъ

Feature-based time series classification





3

・ ロ ト ・ 一 マ ト ・ 日 ト

Feature-based time series classification





(日)

Feature-based time series classification



Features

- Statistics: mean, variance
- Autorregresive coefficients
- Fourier coefficients
- Shift, trend, ...
- tsfeatures (Yang et al. 2015)

A B > A B >



Feature-based time series classification



- Representation independent on the number of series
- Interpretable representation
- Challenge: what features to use?



Model-based time series classification



Model-based time series classification



- Good results with an appropriate model
- Choice of model
- Existence of model

(日)



Early time series classification

Examples

- Early activity recognition
- Early disease recognition in electrocardiograms
- Early detection of sepsis in newborn
- Early detection of failures in machines (predictive maintenance)



Early time series classification

TRAINING SET



Early time series classification (Mori et al 2017, DAMI, TNNLS)





・ロト ・回ト ・ヨト

Early time series classification





・ ロ ト ・ 一 マ ト ・ 日 ト

Early time series classification





Multivariate time series classification

CHALLENGE



・ロト ・ 四ト ・ ヨト ・ ヨト

Outline of the presentation

1 Time Series Data Mining Activities

2 Clustering

- 3 (Early) Supervised Classification
- Outlier/Anomaly Detection
- 5 Conclusions and Future Work



Outlier vs Anomaly





Type of outlier: point outlier





Type of outlier: subsequence outlier



tip Camp UPV - EHU

미 🛛 🖉 🕨 🖉 🖻 🖉 토 🖉 🗩 📃 🔊

Type of outlier: series outlier



$$|\mathbf{x}_t - \hat{\mathbf{x}}_t| > \tau$$



<ロ> < 団> < 団> < 豆> < 豆> < 豆> < 豆</p>

$$|\mathbf{x}_t - \hat{\mathbf{x}}_t| > \tau$$





tasque center for applied mathematics UPV - EHL

$$|\mathbf{x}_t - \hat{\mathbf{x}}_t| > au$$
 MAD



basque center for applied mathematics

$$|\mathbf{X}_t - \hat{\mathbf{X}}_t| > \tau$$



・ロト ・四ト ・ヨト ・ヨト





An overview of outlier/anomaly detection





Outline of the presentation

- 1 Time Series Data Mining Activities
- 2 Clustering
- 3 (Early) Supervised Classification
- 4 Outlier/Anomaly Detection
- 5 Conclusions and Future Work



Not too explored lands

Challenges

- Time series subset selection
- Learning in weakly environments: semi-supervised, multi-label, crowd learning
- Theoretical bounds on learning: assumptions on the generating model



Collaboration

- Usue Mori (UPV/EHU), Amaia Abanda (BCAM)
- Ane Blazquez (Ikerlan), Angel Conde (Ikerlan)
- Aritz Perez (BCAM), Izaskun Oregui (Tecnalia), Javier del Ser (Tecnalia)
- Josu Ircio (Ikerlan), Aizea Lojo (Ikerlan)



Time Series Data Mining Challenges

Jose A. Lozano

Basque Center for Applied Mathematics (BCAM) University of the Basque Country UPV/EHU

EVIA, Sevilla 14 Junio, 2023



(日)