# DATA MINING 2 Time Series - Similarities & Distances

Riccardo Guidotti

a.a. 2023/2024

Slides edited from Keogh Eamonn's tutorial



# **Distances and Similarities**

- Time series problems such as classification, forecasting, clustering, etc. require the usage of a notion of distance or similarity.
- What is similarity?
- It is the quality or state of being similar, likeness, resemblance, as a similarity of features.
- In TSA we recognize two types of similarity measures depending on the data representation considered:
  - shape-based similarity
  - structural-based similarity



# Shape vs Structural Similarities

#### **Shape-based Similarity**

- The original values of the time series are compared taking time into account.
- Better for short time series.

#### **Structural Similarity**

- Time series are transformed into an alternative representation where the novel features are time-independent.
- Better for long time series.



	min	max	mean	std
Q	1.8	2.9	2.0	1.3
С	0.0	1.0	0.2	1.2

#### **Euclidean Distance**

- Given two time series:
  - $Q = q_1 \dots q_n$
  - $C = c_1 \dots c_n$

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

• 
$$T1 = < 56$$
, 176, 110, 95 >  
•  $T2 = < 36$ , 126, 180, 80 >

 $D(T1,T2) = sqrt [ (56-36)^2 + (176-126)^2 + (110-180)^2 + (95-80)^2 ]$ 



# Problems with Euclidean Distance

- Euclidean distance is very sensitive to "distortions" in the data.
- These distortions are dangerous and should be removed.
- Most common distortions:
  - Offset Translation
  - Amplitude Scaling
  - Linear Trend
  - Noise
- They can be removed by using the appropriate normalization.



## Further Problems with Euclidean Distance

• Even after normalization, the Euclidean distance may still be unsuitable for some time series domains since it does not allow for acceleration and deceleration along the time axis.



# **Dynamic Time Warping**

 Sometimes two time series that are conceptually equivalent evolve at different speeds, at least in some moments.



• Euclidean distance - Fixed Time Axis: Sequences are aligned "one to one". Greatly suffers from the misalignment in data.



• Dynamic Time Warping - Warped Time Axis: Nonlinear alignments are possible. Can correct misalignments in data.





https://izbicki.me/blog/converting-images-into-time-series-for-data-mining.html8

# How is DTW Calculated?

- Every possible warping between two time series, is a path through the matrix.
- The constrained sequence of comparisons performed:
  - Start from pair of points (0,0)
  - After point (*i*,*j*), either *i* or *j* increase by one, or both of them
  - End the process on (n,m)





# **Dynamic Programming Approach**



**Step 3**: find the path with the lowest values, i.e., the best alignment between Q and C



# **Dynamic Programming Approach**

 $\gamma(i,j) = d(q_i,c_j) + \min\{\gamma(i-1,j-1), \gamma(i-1,j), \gamma(i,j-1)\}$ Step 2: compute the matrix of all path costs  $\gamma(i,j)$ • Start from cell (1,1)  $\gamma(1,1) = d(q_1,c_1) + \min\{\gamma(0,0), \gamma(0,1), \gamma(1,0)\}$ D(1,1)  $= d(q_1, c_1)$ = D(1,1)• Compute (2,1), (3,1), ..., (n,1)  $\gamma(i,1) = d(q_i,c_1) + \min\{\gamma(i-1,0), \gamma(i-1,1), \gamma(i,0)\}$ D(i,1)  $= d(q_i, c_1) + \gamma(i-1, 1)$  $= D(i,1) + \gamma(i-1,1)$ • Repeat for columns 2, 3, ..., n min +  $\overrightarrow{D(i,1)}$ The general formula applies















#### Point-to-point costs Ч t2 t1 4 Result: 4

























## DTW – A Real Example

- This example shows 2 oneweek periods from the power demand time series.
- Note that although they both describe 4-day work weeks, the blue sequence had Monday as a holiday, and the red sequence had Wednesday as a holiday.



## Comparison of Euclidean Distance and DTW



**Word Spotting** 

# Comparison of Euclidean Distance and DTW

- Classification using 1-NN
- Class(x) = class of most similar training object
- Leaving-one-out evaluation
- For each object: use it as test set, return overall average

#### Accuracy

Dataset	Euclidean	DTW
Word Spotting	0.95	0.99
Sign language	0.71	0.74
GUN	0.95	0.99
Nuclear Trace	0.89	1.00
Leaves <sup>#</sup>	0.67	0.96
(4) Faces	0.94	0.97
Control Chart*	0.93	1.00
2-Patterns	0.99	1.00

# Comparison of Euclidean Distance and DTW

- Classification using 1-NN
- Class(x) = class of most similar training object
- Leaving-one-out evaluation
- For each object: use it as test set, return overall average
- DTW is two to three orders of magnitude slower than Euclidean distance.

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

#### Milliseconds

# Problems with Dynamic Time Warping

- Dynamic Time Warping gives much better results than Euclidean distance on many problems.
- Dynamic Time Warping is very very slow to calculate!
- Is there anything we can do to speed up similarity search under DTW?
# **Global Constraints**

- Slightly speed up the calculations
- Prevent pathological warpings





# **Global Constraints**

- A global constraint constrains the indices of the warping path  $w_k = (i,j)_k$  such that  $j-r \le i \le j+r$ , where r is a term defining allowed range of warping for a given point in a sequence.
- r can be considered as a *window* that reduces the number of calculus.



Sakoe-Chiba Band

Itakura Parallelogram

### Accuracy vs. Width of Warping Window



# Fast Approximations to DTW

• Approximate the time series with some compressed or downsampled representation and do DTW on the new representation.





# Fast Approximations to DTW

- There is strong visual evidence to suggests it works well
- In the literature there is good experimental evidence for the utility of the approach on clustering, classification, etc.



# **Distances and Normalizations**

- If measuring a distance to account for a shape-based similarity it is important to consider the level then the level, i.e., the mean, should not be removed.
- This kind of reasoning applies also to other features of the TS.

# **Global Structural Features**

# Structure-based Similarity

- For long time series, shape-based similarity typically give poor results.
- Structure-based similarity measure similarly of TS based on high level structure.
- The basic idea is to:
  - 1. extract *global* features from the time series,
  - 2. create a feature vector, and
  - 3. use it to measure similarity with Euclidean distance
- Example of features:
  - mean, variance, skewness, kurtosis,
  - 1<sup>st</sup> derivative mean, 1<sup>st</sup> derivative variance, ...
  - parameters of regression, forecasting, Markov model



Feature\Time Series	Α	В	С
Max Value	11	12	19
Mean	5.3	6.4	4.8
Min Value	3	2	5
Autocorrelation	0.2	0.3	0.5
	•••	•••	

# **Simple Standard Features**

- Mean
- Standard Deviation
- Variance
- Median
- 10th Percentile
- 25th Percentile
- 75th Percentile
- 90th Percentile
- IQR
- Covariance
- Skewness
- Kurtosis
- Min
- Max





- abs\_energy Returns the absolute energy of the time series which is the sum over the squared values
- absolute\_maximum Calculates the highest absolute value of the time series x.
- absolute\_sum\_of\_changes Returns the sum over the absolute value of consecutive changes in the series x
- agg\_autocorrelation Descriptive statistics on the autocorrelation of the time series.
- agg\_linear\_trend Calculates a linear least-squares regression
  for values of the time series that were aggregated over chunks versus the sequence from 0 up to the number of chunks minus one.
- approximate\_entropy Implements a vectorized Approximate entropy algorithm.
- ar\_coefficient This feature calculator fits the unconditional maximum likelihood of an autoregressive AR(k) process.
- augmented\_dickey\_fuller Does the time series have a unit root?
- autocorrelation Calculates the autocorrelation of the specified lag
- benford\_correlation Useful for anomaly detection applications. Returns the correlation from first digit distribution when
- binned\_entropy First bins the values of x into max\_bins equidistant bins.

- c3 Uses c3 statistics to measure non linearity in the time series
- change\_quantiles First fixes a corridor given by the quantiles ql and qh of the distribution of x.
- cid\_ce This function calculator is an estimate for a time series complexity.
- count\_above Returns the percentage of values in x that are higher than t
- count\_above\_mean Returns the number of values in x that are higher than the mean of x
- count\_below Returns the percentage of values in x that are lower than t
- count\_below\_mean Returns the number of values in x that are lower than the mean of x
- cwt\_coefficients Calculates a Continuous wavelet transform for the Ricker wavelet, also known as the "Mexican hat wavelet" which is defined by
- energy\_ratio\_by\_chunks Calculates the sum of squares of chunk i out of N chunks expressed as a ratio with the sum of squares over the whole series.
- fft\_aggregated Returns the spectral centroid (mean), variance, skew, and kurtosis of the absolute fourier transform spectrum.

- fft\_coefficient Calculates the fourier coefficients of the onedimensional discrete Fourier Transform for real input by fast fourier transformation algorithm
- first\_location\_of\_maximum Returns the first location of the maximum value of x.
- first\_location\_of\_minimum Returns the first location of the minimal value of x.
- fourier\_entropy Calculate the binned entropy of the power spectral density of the time series (using the welch method).
- friedrich\_coefficients Coefficients of polynomial, which has been fitted to the deterministic dynamics of Langevin model
- has\_duplicate Checks if any value in x occurs more than once
- has\_duplicate\_max Checks if the maximum value of x is observed more than once
- has\_duplicate\_min Checks if the minimal value of x is observed more than once
- index\_mass\_quantile Calculates the relative index i of time series x where q% of the mass of x lies left of i.
- kurtosis Returns the kurtosis of x.
- large\_standard\_deviation Does time series have Targe standard deviation?

- last location\_of\_maximum Returns the relative last location of the maximum value of x.
- last\_location\_of\_minimum Returns the last location of the minimal value of x.
- lempel\_ziv\_complexity Calculate a complexity estimate based on the Lempel-Ziv compression algorithm.
- length Returns the length of x
- linear\_trend Calculate a linear least-squares regression for the values of the time series versus the sequence from 0 to length of the time series minus one.
- linear\_trend\_timewise Calculate a linear least-squares regression for the values of the time series versus the sequence from 0 to length of the time series minus one.
- longest\_strike\_above\_mean Returns the length of the longest consecutive subsequence in x that is bigger than the mean of x
- longest\_strike\_below\_mean Returns the length of the longest consecutive subsequence in x that is smaller than the mean of x

- fast fourier matrix\_profile Calculates the 1-D Matrix Profile[1]
  and returns Tukey's Five Number Set plus the mean of that Matrix Profile.
- max\_langevin\_fixed\_point Largest fixed point of dynamics :math:argmax\_x {h=0}` estimated from polynomial, which has been fitted to the deterministic dynamics of Langevin model
- maximum Calculates the highest value of the time series x.
- mean Returns the mean of x
- mean\_abs\_change Average over first differences.
- mean\_change Average over time series differences.
- mean n\_absolute max Calculates the arithmetic mean of the n absolute maximum values of the time series.
- mean\_second\_derivative\_central Returns the mean value of a central approximation of the second derivative
- median Returns the median of x
- minimum Calculates the lowest value of the time series x.
- number\_crossing\_m Calculates the number of crossings of x on m.
- number\_cwt\_peaks Number of different peaks in x.

- number\_peaks Calculates the number of peaks of at least support n in the time series x.
- partial\_autocorrelation Calculates the value of the partial autocorrelation function at the given lag.
- percentage\_of\_reoccurring\_datapoints\_to\_all\_datapoints Returns the percentage of non-unique data points.
- percentage\_of\_reoccurring\_values\_to\_all\_values Returns the percentage of values that are present in the time series more than once.
- permutation\_entropy Calculate the permutation entropy.
- quantile Calculates the q quantile of x.
- query\_similarity\_count This feature calculator accepts an input query subsequence parameter, compares the query (under z-normalized Euclidean distance) to all subsequences within the time series, and returns a count of the number of times the query was found in the time series (within some predefined maximum distance threshold).

- range\_count Count observed values within the interval [min, max).
- ratio\_beyond\_r\_sigma Ratio of values that are more than r \* std (so r times sigma) away from the mean of x.
- ratio\_value\_number\_to\_time\_series\_length Returns a factor which is 1 if all values in the time series occur only once, and below one if this is not the case.
- root\_mean\_square Returns the root mean square (rms) of the time series.
- sample\_entropy Calculate and return sample entropy of x.
- set\_property This method returns a decorator that sets the property key of the function to value
- skewness Returns the sample skewness of x (calculated with the adjusted Fisher-Pearson standardized moment coefficient G1).
- spkt\_welch\_density This feature calculator estimates the cross power spectral density of the time series x at different frequencies.
- standard\_deviation Returns the standard deviation of x
- sum\_of\_reoccurring\_data\_points Returns the sum of all data points, that are present in the time series more than once.
- sum\_of\_reoccurring\_values Returns the sum of all values, that

are present in the time series more than once.

- sum\_values Calculates the sum over the time series values
- symmetry\_looking Boolean variable denoting if the distribution of x looks symmetric.
- time\_reversal\_asymmetry\_statistic Returns the time reversal asymmetry statistic.

value\_count - Count occurrences of value in time series x.

- variance Returns the variance of x
- variance\_larger\_than\_standard\_deviation Is variance higher than the standard deviation?
- variation\_coefficient Returns the variation coefficient (standard error / mean, give relative value of variation around mean) of x

#### catch22: CAnonical Time-series CHaracteristics

- The catch22 feature set spans a diverse range of time-series characteristics representative of the diversity of interdisciplinary methods for TSA.
- Features in catch22 capture TS properties of the distribution of values in the TS, linear and nonlinear temporal autocorrelation properties, scaling of fluctuations, and others.
- Selected by applying the procedure describe in [Lubba 2019] to a set of 93 datasets containing over 147k TS and using a filtered version of the HCTSA feature library (4791 features).
- The reduction from 4791 to 22 features is associated with a 1000-fold reduction in computation time and near linear scaling with TS length, despite an average reduction in classification accuracy of just 7%.

Distribution DN\_HistogramMode\_5 DN\_HistogramMode\_10 Simple temporal statistics SB\_BinaryStats\_mean\_longstretch1 DN\_OutlierInclude\_p\_001\_mdrmd DN\_OutlierInclude\_n\_001\_mdrmd Linear autocorrelation CO\_f1ecac CO\_FirstMin\_ac SP\_Summaries\_welch\_rect\_area\_5\_1 SP\_Summaries\_welch\_rect\_centroid FC\_LocalSimple\_mean3\_stderr Nonlinear autocorrelation CO\_trev\_1\_num CO\_HistogramAMI\_even\_2\_5 IN\_AutoMutualInfoStats\_40\_gaussian\_fmmi Successive differences MD\_hrv\_classic\_pnn40 SB\_BinaryStats\_diff\_longstretch0 SB\_MotifThree\_quantile\_hh FC\_LocalSimple\_mean1\_tauresrat CO\_Embed2\_Dist\_tau\_d\_expfit\_meandiff Fluctuation Analysis SC\_FluctAnal\_2\_dfa\_50\_1\_2\_logi\_prop\_r1 SC\_FluctAnal\_2\_rsrangefit\_50\_1\_logi\_prop\_r1 OthersSB\_TransitionMatrix\_3ac\_sumdiagcov PD\_PeriodicityWang\_th0\_01

Mode of z-scored distribution (5-bin histogram) Mode of z-scored distribution (10-bin histogram)

Longest period of consecutive values above the mean Time intervals between successive extreme events above the mean Time intervals between successive extreme events below the mean

First 1/e crossing of autocorrelation function First minimum of autocorrelation function Total power in lowest fifth of frequencies in the Fourier power spectrum Centroid of the Fourier power spectrum Mean error from a rolling 3-sample mean forecasting

> Time-reversibility statistic,  $\langle (x_{t+1} - x_t)^3 \rangle_t$ Automutual information,  $m = 2, \tau = 5$ First minimum of the automutual information function

 $\begin{array}{c} \mbox{Proportion of successive differences exceeding 0.04 \sigma \ [20]} \\ \mbox{Longest period of successive incremental decreases} \\ \mbox{Shannon entropy of two successive letters in equiprobable 3-letter symbolization} \\ \mbox{Change in correlation length after iterative differencing} \\ \mbox{Exponential fit to successive distances in 2-d embedding space} \end{array}$ 

Proportion of slower timescale fluctuations that scale with DFA (50% sampling) Proportion of slower timescale fluctuations that scale with linearly rescaled range fits

Trace of covariance of transition matrix between symbols in 3-letter alphabet Periodicity measure of [31]

#### **Overview of Global Features and Relationships**



### **Global Feature-based Predictor**



#### Features, Approximations, Distances and Normalizations

- Normalizations can be applied before global features extraction depending on the objective of your TSA task.
- Time-Dependent approximations can be applied before global features extraction depending on the objective of your TSA task.
- It does not make any sense to use Time-Independent approximation after that global features have been extracted.
- It does not make any sense to use a distance function accounting for time like DTW after that global features have been extracted.

# Summary of Time Series Similarity

- If you have short time series
  - use DTW after searching over the warping window size
  - try also to approximate to speed up the calculus
- If you have long time series
  - if you do know something about your data => extract features
  - (and you know nothing about your data => try compression/approximation-based dissimilarity)

## References

- Forecasting: Principles and Practic. Rob J Hyndman and George Athanasaopoulus. (<u>https://otexts.com/fpp2/</u>)
- Time Series Analysis and Its Applications. Robert H. Shumway and David S. Stoffer. 4<sup>th</sup> edition.(<u>http://www.stat.ucla.edu/~frederic/415/S23/tsa4.pdf</u>)
- Mining Time Series Data. Chotirat Ann Ratanamahatana et al. 2010. (<u>https://www.researchgate.net/publication/227001229\_Mining\_Time\_Series\_Data</u>)
- Dynamic Programming Algorithm Optimization for Spoken Word Recognition. Hiroaki Sakode et al. 1978.
- Experiencing SAX: a Novel Symbolic Representation of Time Series. Jessica Line et al. 2009
- Compression-based data mining of sequential data. Eamonn Keogh et al. 2007.



#### Time Series Analysis and Its Applications

With R Examples

Fourth Edition

Description Springer

# **Exercises DTW**

# DTW – Exercise 1

• Given the following input time series:

- A) Compute the distance between "t1" and "t2", using the DTW with distance between points computed as d(x,y) = |x y|.
- B) If we repeat the computation of point (A) above, this time with a Sakoe-Chiba band of size r=1, does the result change? Why?
- C) If we compute DTW(T1,T2), where T1 is equal to t1 in reverse order (namely T1=<0,1,6,3,4>) and similarly for T2 (namely T2=<1,0,7,6,3>), is it true that DTW(T1,T2) = DTW(t1,t2)? Discuss the problem without providing any computation.








































- B) No. Because the DTW optimal path remains inside the band of size r=1
- C) Yes. The optimal path in one direction is the same in the opposite direction. Though, the cumulative costs matrix might look different.

# DTW – Exercise 2

Given the following time series:
t = < 2, 6, 9, 1, 6, 2 >
q = < 5, 1, 5, 5, 8, 4 >

compute

- (i) their Manhattan and Euclidean distance,
- (ii) their DTW, and (iii) their DTW with Sakoe-Chiba band of size r=1 (i.e. all cells at distance <= 1 from the diagonal are allowed).</li>
- For points (ii) and (iii) show the cost matrix and the optimal path found.

• Euclidean = sqrt(74) = 8.6, Manhattan = 20





## DTW – Exercise 3

• Given the following time series:

ID	Time series
W	< 6, 11, 13, 15 >
Х	< 10, 7, 7, 12, 14, 17 >
Y	< 9, 11, 14, 13, 20 >

 Compute the distances among all pairs of time series adopting a Dynamic Time Warping distance, and computing the distances between single points as d(x,y) = | x - y |. For each pair of time series compared also show the matrix used to compute the final result.

ID	Time series
W	< 6, 11, 13, 15 >
Х	< 10, 7, 7, 12, 14, 17 >
Y	< 9, 11, 14, 13, 20 >

#### W – X



#### **W** – Y



X – Y

[1,] [2,] [3,] [4,] [5,] [6,]

