# DATA MINING 2
# Time Series - ShapeIt/Motif Discovery

Riccardo Guidotti

a.a. 2019/2020
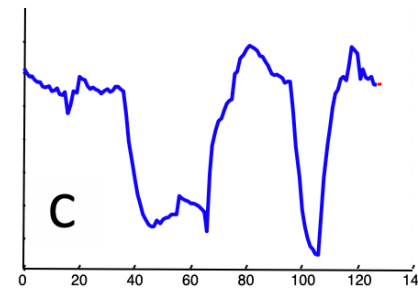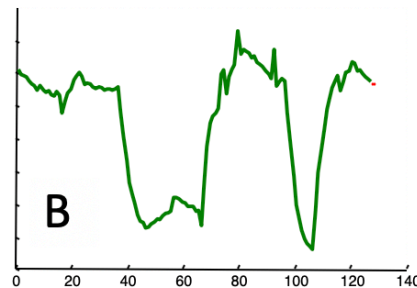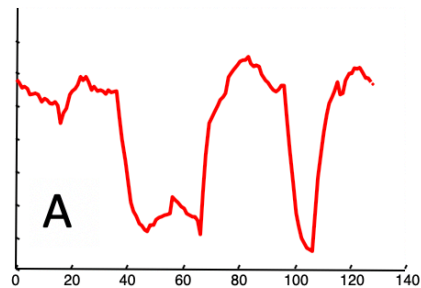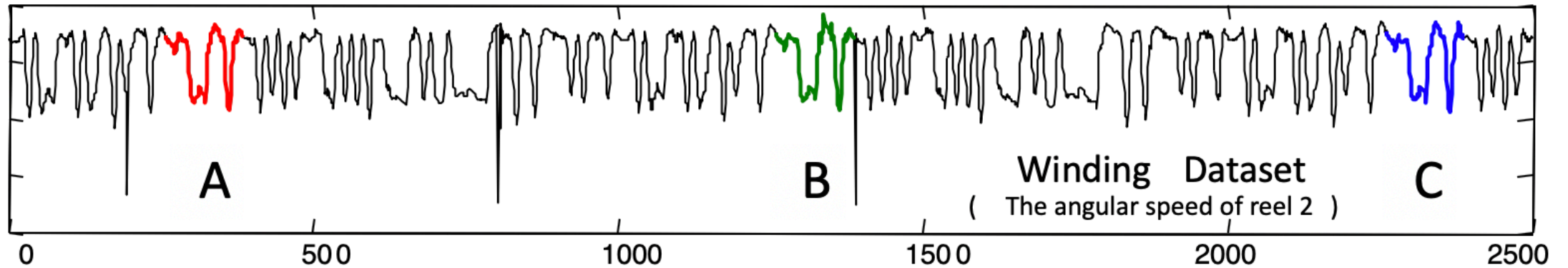
Slides edited from Keogh Eamonn's tutorial

UNIVERSITÀ DI PISA

# Motif

# Time Series Motif Discovery

- Finding repeated patterns, i.e., pattern mining.
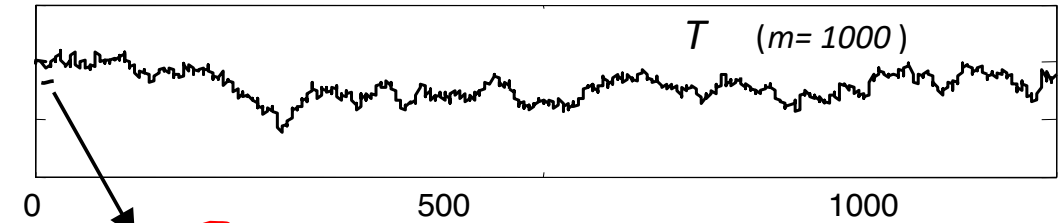- Are there any repeated patterns, of length $m$ in the TS?

# Why Find Motifs?

- Mining **association rules** in TS requires the discovery of motifs. These are referred to as primitive shapes and frequent patterns.

- Several **TS classifiers** work by constructing typical prototypes of each class. These prototypes may be considered motifs.

- Many **TS anomaly detection** algorithms consist of modeling normal behavior with a set of typical shapes (which we see as motifs), and detecting future patterns that are dissimilar to all typical shapes.

# How do we find Motifs?

- Given a predefined motif length $m$, a brute-force method searches for motifs from all possible comparisons of subsequences.

- It is obviously very slow and computationally expensive.

- The most reference algorithm is based on a hot idea from bioinformatics, random projection* and the fact that SAX allows use to lower bound discrete representations of TSs.

- J Buhler and M Tompa. Finding motifs using random projections. In RECOMB'01. 2001.

# Example of the Motif Discovery Algorithm

- Assume that we have a time series T of length 1,000, and a motif of length 16, which occurs twice, at time $T_1$ and time $T_{58}$.
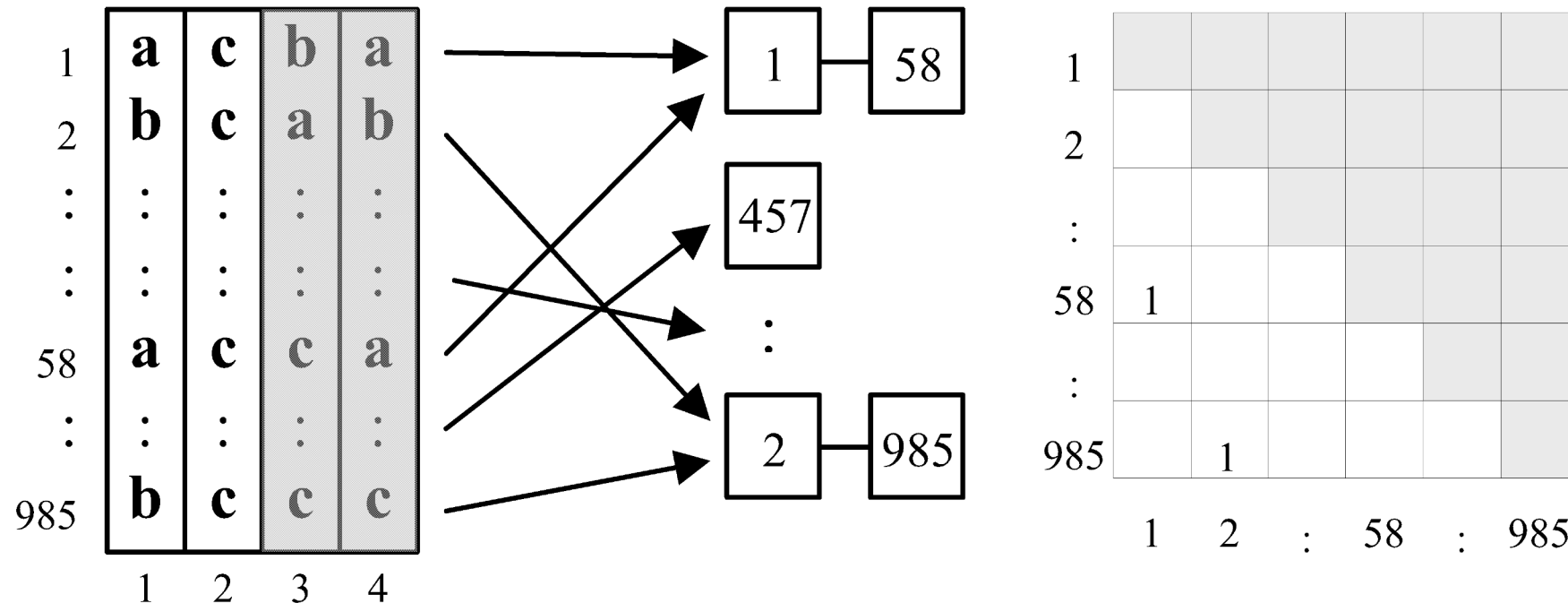


$T$   $(m = 1000)$

0          500          1000

$C_1$

$\hat{C}_1$   **a c b a**

$\hat{S}$

| | | | |
|---|---|---|---|
| **a** | **c** | **b** | **a** |
| b | c | a | b |
| : | : | : | : |
| : | : | : | : |
| a | c | c | a |
| : | : | : | : |
| b | c | c | c |

1
2
:
:
58
:
985

16

$a = 3$  { **a**,**b**,**c** }   *alphabet*
$n = 16$   *motif length*
$w = 4$   *sax window*

# Example of the Motif Discovery Algorithm

- 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.
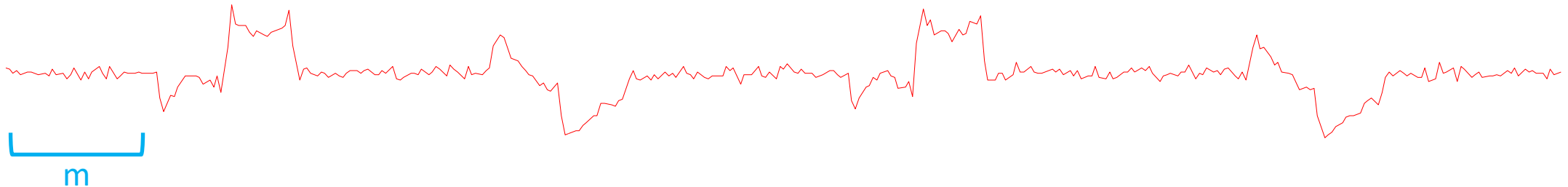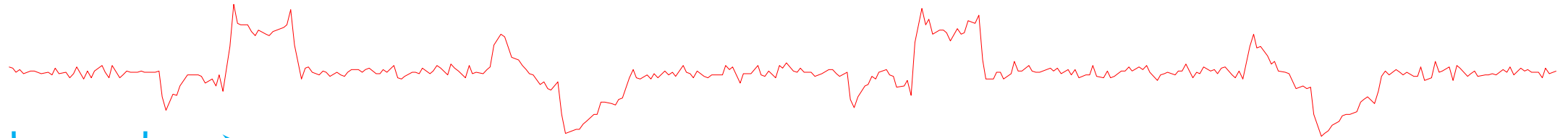
# Example of the Motif Discovery Algorithm
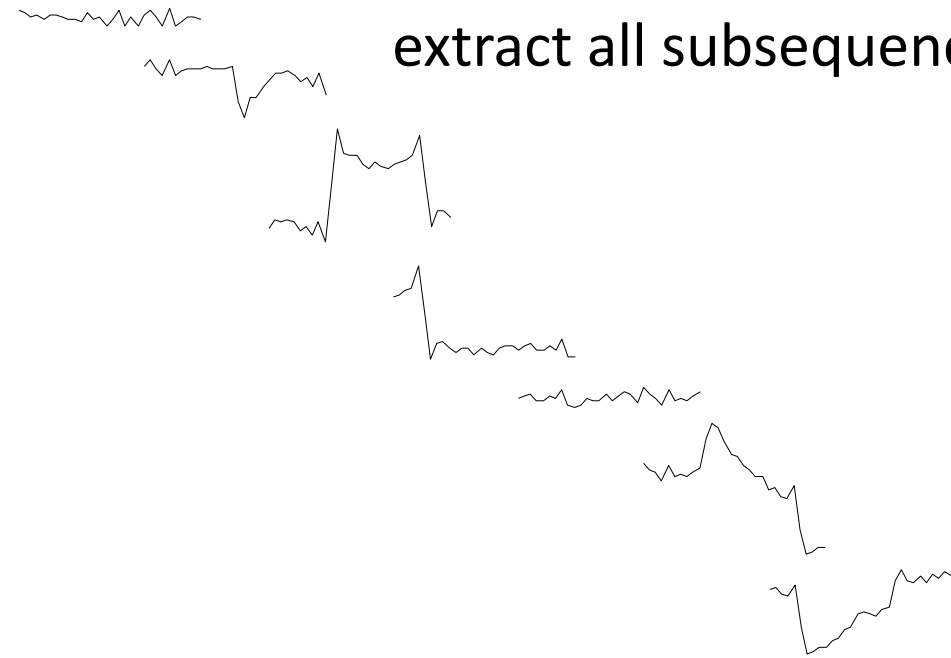
- A mask {2,4} was randomly chosen, so the values in columns {2,4} were used to project matrix into buckets.

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

# Matrix Profile

- The Matrix Profile (MP) is a data structure that annotates a TS and can be exploited for many purposed: e.g. efficient Motif Discovery.
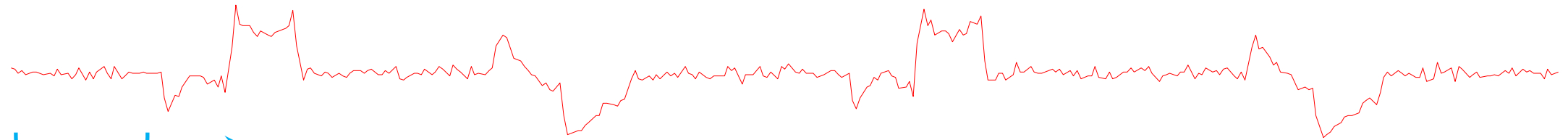
- Given a time series, T and a desired subsequence length, m.

# Matrix Profile

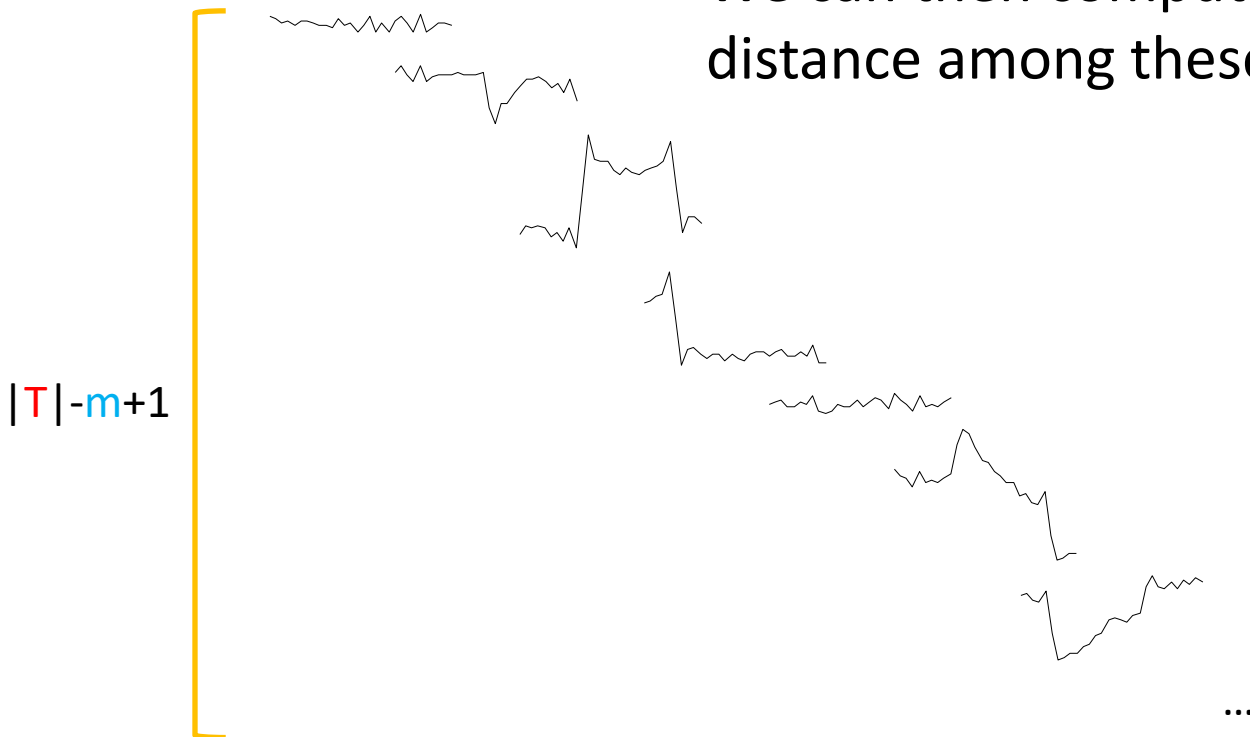We can use sliding window of length *m* to extract all subsequences of length *m*.
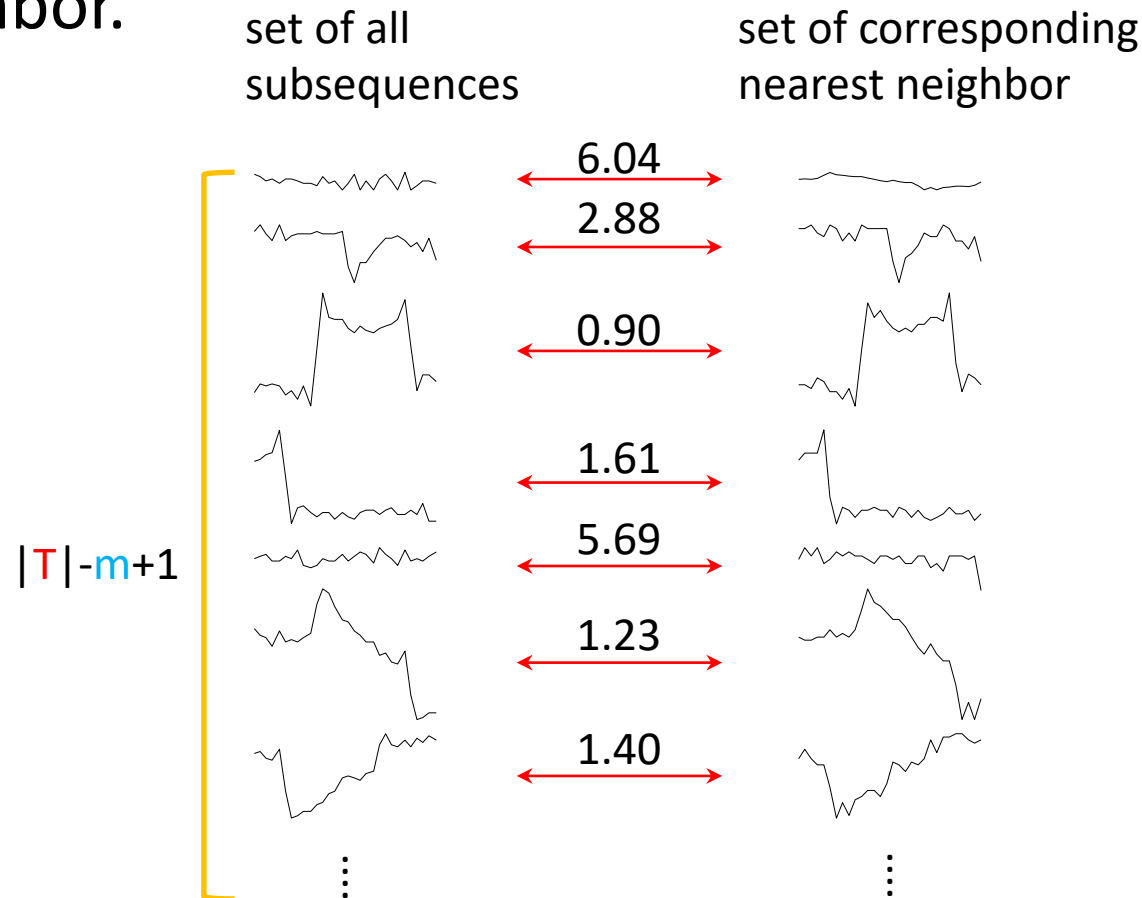
*m*

$|T|-m+1$

# Matrix Profile

We can then compute the pairwise distance among these subsequences.

$m$

$|T|-m+1$

| 0 | 7.6952 | 7.7399 | … |
|---|---|---|---|
| 7.6952 | 0 | 7.7106 | … |
| 7.7399 | 7.7106 | 0 | … |
| … | … | … | … |

…

# Matrix Profile

- For each subsequence we keep only the distance with the closest nearest neighbor.
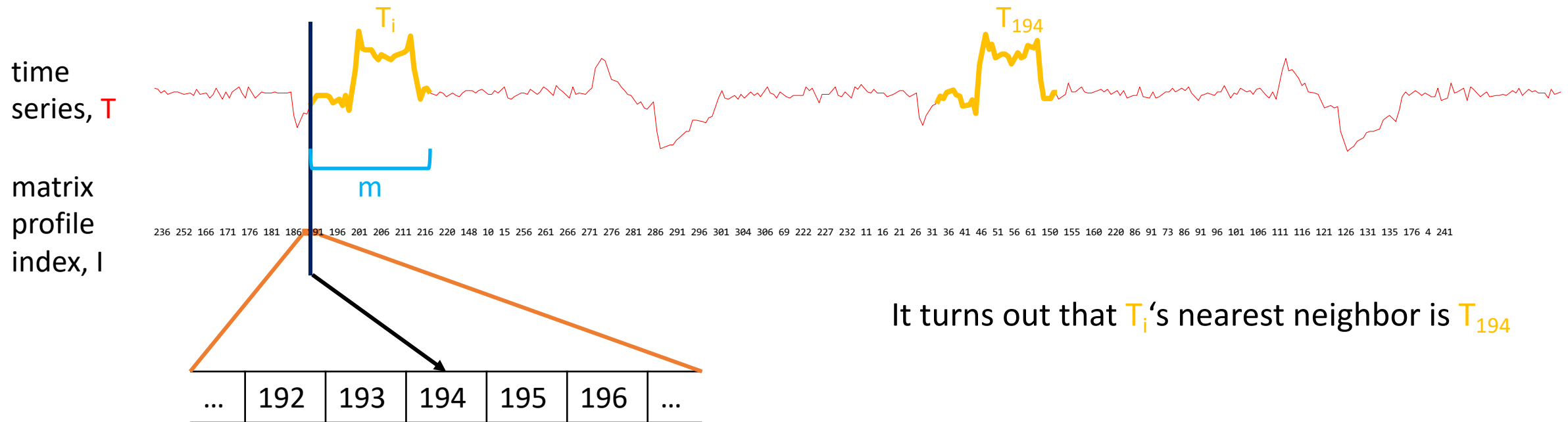
# Matrix Profile

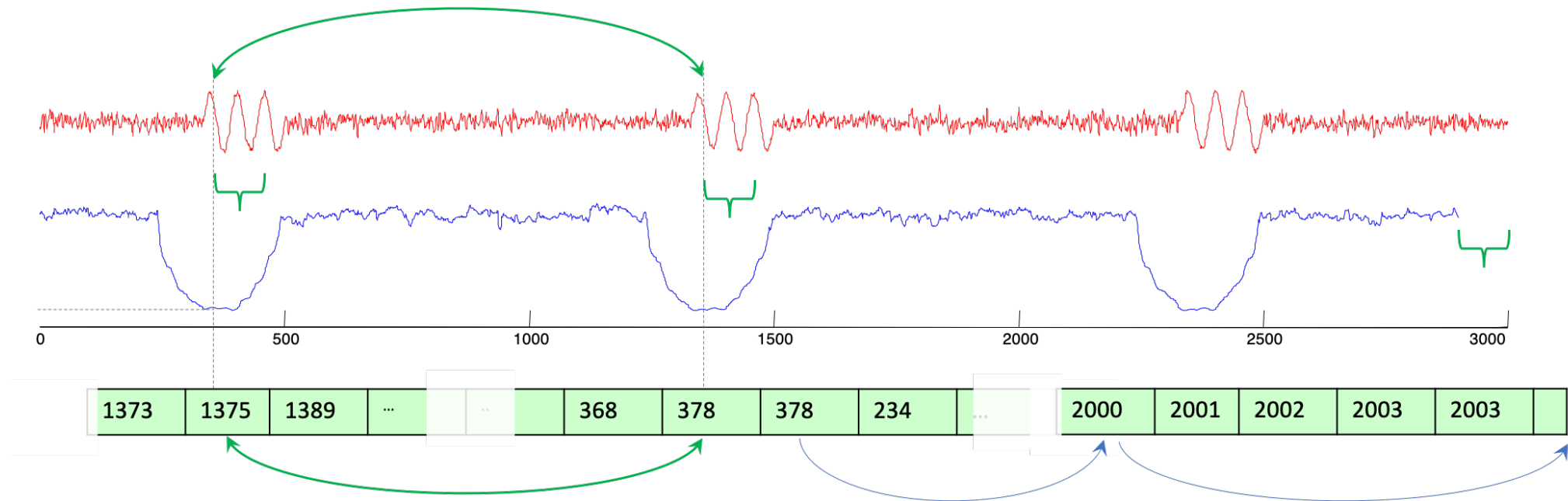- The distance to the corresponding nearest neighbor of each subsequence can be stored in a vector called **matrix profile P**.

time
series, T

matrix
profile, P

The matrix profile value at location *i* is the
distance between $T_i$ and its nearest neighbor

# Matrix Profile

- The index of corresponding nearest neighbor of each subsequence is also stored in a vector called matrix profile index.

$T_i$

$T_{194}$

time series, $T$

$m$

matrix profile index, $I$

236 252 166 171 176 181 186 191 196 201 206 211 216 220 148 10 15 256 261 266 271 276 281 286 291 296 301 304 306 69 222 227 232 11 16 21 26 31 36 41 46 51 56 61 150 155 160 220 86 91 73 86 91 96 101 106 111 116 121 126 131 135 176 4 241

It turns out that $T_i$'s nearest neighbor is $T_{194}$

| … | 192 | 193 | 194 | 195 | 196 | … |
|---|---|---|---|---|---|---|

The matrix profile value at location *i* is the distance between $T_i$ and its nearest neighbor

# Matrix Profile

- The MP index allows to find the nearest neighbor to any subsequence in constant time.

- Note that the pointers in the matrix profile index are not necessarily symmetric.

- If A points to B, then B may or may not point to A

- The classic TS motif: the two smallest values in the MP must have the same value, and their pointers must be mutual.

# How to "read" a Matrix Profile

- For relatively low values, you know that the subsequence in the original TS must have (at least one) relatively similar subsequence elsewhere in the data (such regions are "motifs")

- For relatively high values, you know that the subsequence in the original TS must be unique in its shape (such areas are anomalies).



Must be an anomaly in the original data, in this region.

We call these *Time Series Discords*

Must be conserved shapes (motifs) in the original data, in these three regions

# How to Compute Matrix Profile?

- Given a time series, T and a desired subsequence length, m.



| inf | inf | inf | inf | inf | inf | inf | inf | inf | inf | inf | inf | inf | inf | inf | inf | inf | inf | inf |
|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|

Matrix profile is initialized as inf vector

This is just a toy example, so the values and the vector length does not fit the time series shown above

# How to Compute Matrix Profile?

- Given a time series, $T$ and a desired subsequence length, $m$.



| inf | inf | inf | inf | inf | inf | inf | inf | inf | inf | inf | inf | inf | inf | inf | inf | inf | inf | inf |
|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|

At the first iteration, a subsequence $T_i$ is randomly selected from $T$

# How to Compute Matrix Profile?

- Given a time series, $T$ and a desired subsequence length, $m$.



| inf | inf | inf | inf | inf | inf | inf | inf | inf | inf | inf | inf | inf | inf | inf | inf | inf | inf | inf |
|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|

We compute the distances between $T_i$ and every subsequences from $T$ (time complexity = $O(|T|\log(|T|))$)
We then put the distances in a vector based on the position of the subsequences

| 3 | 2 | 0 | 5 | 3 | 4 | 5 | 1 | 2 | 9 | 8 | 4 | 2 | 3 | 4 | 8 | 6 | 2 | 1 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|

The distance between $T_i$ and $T_1$ (first subsequence) is 3

# How to Compute Matrix Profile?

- Given a time series, $T$ and a desired subsequence length, $m$.



$T_i$

$m$

| inf | inf | inf | inf | inf | inf | inf | inf | inf | inf | inf | inf | inf | inf | inf | inf | inf | inf | inf |
|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|

We compute the distances between $T_i$ and every subsequences from $T$ (time complexity = $O(|T|\log(|T|))$)
We them put the distances in a vector based on the position of the subsequences

| 3 | 2 | 0 | 5 | 3 | 4 | 5 | 1 | 2 | 9 | 8 | 4 | 2 | 3 | 4 | 8 | 6 | 2 | 1 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|

Let say $T_i$ happen to be the third subsequences, therefore
the third value in the distance vector is 0

# How to Compute Matrix Profile?

- Given a time series, <span style="color:red">T</span> and a desired subsequence length, <span style="color:cyan">m</span>.



| inf | inf | inf | inf | inf | inf | inf | inf | inf | inf | inf | inf | inf | inf | inf | inf | inf | inf | inf | inf |
|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|

min ↕

Matrix profile is updated by apply elementwise minimum to these two vectors

| 3 | 2 | 0 | 5 | 3 | 4 | 5 | 1 | 2 | 9 | 8 | 4 | 2 | 3 | 4 | 8 | 6 | 2 | 1 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|

# How to Compute Matrix Profile?

- Given a time series, T and a desired subsequence length, m.



| 3 | inf | inf | inf | inf | inf | inf | inf | inf | inf | inf | inf | inf | inf | inf | inf | inf | inf | inf |
|---|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|

min

Matrix profile is updated by apply elementwise minimum to these two vectors

| 3 | 2 | 0 | 5 | 3 | 4 | 5 | 1 | 2 | 9 | 8 | 4 | 2 | 3 | 4 | 8 | 6 | 2 | 1 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|

# How to Compute Matrix Profile?

- Given a time series, T and a desired subsequence length, m.



| 3 | 2 | inf | 5 | 3 | 4 | 5 | 1 | 2 | 9 | 8 | 4 | 2 | 3 | 4 | 8 | 6 | 2 | 1 |

After we finish update matrix profile for the first iteration

| 3 | 2 | 0 | 5 | 3 | 4 | 5 | 1 | 2 | 9 | 8 | 4 | 2 | 3 | 4 | 8 | 6 | 2 | 1 |

# How to Compute Matrix Profile?

- Given a time series, $T$ and a desired subsequence length, $m$.



| 3 | 2 | inf | 5 | 3 | 4 | 5 | 1 | 2 | 9 | 8 | 4 | 2 | 3 | 4 | 8 | 6 | 2 | 1 |
|---|---|-----|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|

In the second iteration, we randomly select another subsequence $T_j$ and it happens to be the 12th subsequences

# How to Compute Matrix Profile?

- Given a time series, $T$ and a desired subsequence length, $m$.



| 3 | 2 | inf | 5 | 3 | 4 | 5 | 1 | 2 | 9 | 8 | 4 | 2 | 3 | 4 | 8 | 6 | 2 | 1 |
|---|---|-----|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|

Once again, we compute the distance between $T_j$ and every subsequences of $T$

| 2 | 3 | 1 | 4 | 4 | 3 | 6 | 2 | 1 | 5 | 8 | 0 | 2 | 3 | 5 | 9 | 4 | 2 | 2 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|

# How to Compute Matrix Profile?

- Given a time series, $T$ and a desired subsequence length, $m$.



| 3 | 2 | inf | 5 | 3 | 4 | 5 | 1 | 2 | 9 | 8 | 4 | 2 | 3 | 4 | 8 | 6 | 2 | 1 |
|---|---|-----|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|

min $\updownarrow$            The same elementwise minimum

| 2 | 3 | 1 | 4 | 4 | 3 | 6 | 2 | 1 | 5 | 8 | 0 | 2 | 3 | 5 | 9 | 4 | 2 | 2 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|

# How to Compute Matrix Profile?

- Given a time series, $T$ and a desired subsequence length, $m$.



| 2 | 2 | inf | 5 | 3 | 4 | 5 | 1 | 2 | 9 | 8 | 4 | 2 | 3 | 4 | 8 | 6 | 2 | 1 |
|---|---|-----|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|

min ⇅  The same elementwise minimum

| 2 | 3 | 1 | 4 | 4 | 3 | 6 | 2 | 1 | 5 | 8 | 0 | 2 | 3 | 5 | 9 | 4 | 2 | 2 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|

# How to Compute Matrix Profile?

- Given a time series, $T$ and a desired subsequence length, $m$.



| 2 | 2 | inf | 5 | 3 | 4 | 5 | 1 | 2 | 9 | 8 | 4 | 2 | 3 | 4 | 8 | 6 | 2 | 1 |
|---|---|-----|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|

min ↕     The same elementwise minimum

| 2 | 3 | 1 | 4 | 4 | 3 | 6 | 2 | 1 | 5 | 8 | 0 | 2 | 3 | 5 | 9 | 4 | 2 | 2 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|

# How to Compute Matrix Profile?

- Given a time series, $T$ and a desired subsequence length, $m$.



| 2 | 2 | 1 | 5 | 3 | 4 | 5 | 1 | 2 | 9 | 8 | 4 | 2 | 3 | 4 | 8 | 6 | 2 | 1 |

min ↕   The same elementwise minimum

| 2 | 3 | 1 | 4 | 4 | 3 | 6 | 2 | 1 | 5 | 8 | 0 | 2 | 3 | 5 | 9 | 4 | 2 | 2 |

# How to Compute Matrix Profile?

- Given a time series, **T** and a desired subsequence length, **m**.



| 2 | 2 | 1 | 5 | 3 | 4 | 5 | 1 | 2 | 9 | 8 | 4 | 2 | 3 | 4 | 8 | 6 | 2 | 1 |

min ↕

| 2 | 3 | 1 | 4 | 4 | 3 | 6 | 2 | 1 | 5 | 8 | 0 | 2 | 3 | 5 | 9 | 4 | 2 | 2 |

We repeat the two steps (distance computation and update) until we have used every subsequences

# How to Compute Matrix Profile?

- Given a time series, T and a desired subsequence length, m.



$T_j$

m

| 2 | 2 | 1 | 5 | 3 | 4 | 5 | 1 | 2 | 9 | 8 | 4 | 2 | 3 | 4 | 8 | 6 | 2 | 1 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|

min ↕

| 2 | 3 | 1 | 4 | 4 | 3 | 6 | 2 | 1 | 5 | 8 | 0 | 2 | 3 | 5 | 9 | 4 | 2 | 2 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|

There are $|T|$ subsequences and the distance computation is $O(|T|\log(|T|))$

The overall time complexity is $O(|T|^2\log(|T|))$

# Motif Discovery From Matrix Profile



time series, T

matrix profile, P

Local minimums are corresponding to motifs

# Motif Discovery From Matrix Profile



- It is sometime useful to think of time series subsequences as points in m-dimensional space.

- In this view, dense regions in the m-dimensional space correspond to regions of the time series that have a low corresponding MP.

# Top-K Motifs

- We need a parameter R.
- 1 < R < (small number, say 3)
- Lets make R = 2 for now.
- We begin by finding the nearest pair of points, the *motif pair*....
- This the pair of subsequences corresponding to lowest pair of values in the MP

# Top-K Motifs

- We find the nearest pair of points are D1 apart.

- Lets draw a circle, D1 times R, around both points.

- Any points that are within either of these circles, are added to this motif, in this case just one.

- The Top-1 motif has three members, it is done.

# Top-K Motifs



- Now lets find the Top-2 motif. We find the **nearest pair of points**, excluding anything from the top motif.

- The nearest pair of points are D2 apart.

- Lets draw a circle D1 times R, around both points.

- Any points that are within either of these circles, is added to this motif, in this case there are two for a total of four items in the Top-2 Motif

# Top-K Motifs

- We are done with the Top-2 Motif
- Note that we will always have:
  - $D_1 < D_2 < D_3 \ldots D_K$
- **When to stop?** (what is K?)
- We could use MDL or a predefined K.

# Anomaly Discovery From Matrix Profile



- We need a parameter E of subseqeunces to exclude in the vicinity of the anomaly.

- Lets make E = 2 for now.

- We begin by finding the subsequence with the highest distance in the MP

- This corresponding to biggest anomaly

# Top-K Anomaly

- Then we look for the E closest subsequences to the anomaly.

- We remove all of them.

- We can use a predefined K or the MDL to stop.

# Shapelet

# Time Series Classification

- Given a set $X$ of $n$ time series, $X = \{x_1, x_2, ..., x_n\}$, each time series has m ordered values $x_i = \langle x_{t1}, x_{t2}, ..., x_{tm} \rangle$ and a class value $c_i$.

- The objective is to find a function $f$ that maps from the space of possible time series to the space of possible class values.

- Generally, it is assumed that all the TS have the same length $m$.

# Shapelet-based Classification

1. Represent a TS as a vector of distances with representative subsequences, namely shapelets.

2. Use it to as input for machine learning classifiers.



Urtica dioica

Verbena urticifolia

Verbena urticifolia

**Shapelet Dictionary**

5.1

Does Q have a subsequence within a distance 5.1 of shape | ?

**Leaf Decision Tree**

yes

no

0

1

Verbena urticifolia

Urtica dioica

| | |
|---|---|
| 3.2 | 8.7 |
| 1.4 | 7.9 |
| 6.7 | 4.2 |
| 9.2 | 3.4 |

Shapelet

Verbena urticifolia

Urtica dioica

# Time Series Shapelets

- Shapelets are TS subsequences which are maximally representative of a class.

- Shapelets can provide interpretable results, which may help domain practitioners better understand their data.

- Shapelets can be significantly more accurate/robust because they are *local features*, whereas most other state-of-the-art TS classifiers consider *global features*.

Verbena   0.87
Urtica       0.34

Shapelet

Verbena urticifolia

Urtica dioica

# Extract Subsequences of all Possible Lengths



Candidates Pool

# Extract Subsequences of all Possible Lengths



Candidates Pool

# Extract Subsequences of all Possible Lengths



Candidates Pool

# Extract Subsequences of all Possible Lengths

# Extract Subsequences of all Possible Lengths



Candidates Pool

# Distance with a Subsequence

- Distance from the TS to the subsequence *SubsequenceDist(T, S)* is a distance function that takes time series *T* and subsequence *S* as inputs and returns a nonnegative value *d*, which is the distance from *T* to *S*.

- *SubsequenceDist(T, S) = min(Dist(S, S')), for S' $\in S_T^{|S|}$*

- where $S_T^{|S|}$ is the set of all possible subsequences of *T*

- Intuitively, it is the distance between *S* and its best matching location in *T*.

# Testing The Utility of a Candidate Shapelet

- Arrange the TSs in the dataset *D* based on the distance from the candidate.

- Find the optimal split point that maximizes the information gain (same as for Decision Tree classifiers)

- Pick the candidate achieving best utility as the shapelet

# Entropy



- A TS dataset *D* consists of two classes, A and B.

- Given that the proportion of objects in class A is *p(A)* and the proportion of objects in class B is *p(B)*,

- The **Entropy** of D is: *I(D) = -p(A)log(p(A)) -p(B)log(p(B))*.

- Given a strategy that divides the *D* into two subsets $D_1$ and $D_2$, the information remaining in the dataset after splitting is defined by the weighted average entropy of each subset.

- If the fraction of objects in $D_1$ is $f(D_1)$ and in $D_2$ is $f(D_2)$,

- The total entropy of *D* after splitting is $\hat{I}(D) = f(D_1)I(D_1) + f(D_2)I(D_2)$.

# Information Gain



Split point distance from shapelet = 5.1

- Given a certain split strategy *sp* which divides *D* into two subsets $D_1$ and $D_2$, the entropy before and after splitting is *I(D)* and *Î(D)*.

- The **information gain** for this splitting rule is:

- *Gain(sp) = I(D) - Î(D) =*

-         $= I(D) - f(D_1)I(D_1) + f(D_2)I(D_2)$.



- We use the distance from *T* to a shapelet *S* as the splitting rule *sp*.

# Problem

- The total number of candidate is

$$\sum_{l=MINLEN}^{MAXLEN} \sum_{T_i \in D} (|T_i| - l + 1)$$

- For each candidate you have to compute the distance between this candidate and each training sample

- For instance
  - 200 instances with length 275
  - 7,480,200 shapelet candidates

# Speedup

- Distance calculations form TSs to shapelet candidates is expensive.
- Reduce the time in two ways
- Distance Early Abandon
  - reduce the distance computation time between two TS
- Admissible Entropy Pruning
  - reduce the number of distance calculations

# Distance Early Abandon

- We only need the minimum distance.

- Method
  - Keep the best-so-far distance
  - Abandon the calculation if the current distance is larger than best-so-far.

# Admissible Entropy Prunining

- We only need the best shapelet for each class

- For a candidate shapelet
  - We do not need to calculate the distance for each training sample
  - After calculating some training samples, the upper bound of information gain < best candidate shapelet
  - Stop calculation
  - Try next candidate

# An Alternative Way for Extracting Shapelets

- The minimum distances (M) between Ts and Shapelets can be used as predictors to approximate the TSs label (Y) using a linear model (W):

$$\hat{Y}_i \;=\; W_0 + \sum_{k=1}^{K} M_{i,k} W_k, \quad \forall i \in \{1, \dots, I\}$$

- A logistic regression loss can measure the quality of the prediction:

$$\mathcal{L}(Y, \hat{Y}) \;=\; -Y \ln \sigma(\hat{Y}) - (1 - Y) \ln \left(1 - \sigma(\hat{Y})\right)$$

- The objective is to minimize a regularized loss function accross all the instances (I) :

$$\underset{S,W}{\operatorname{argmin}} \; \mathcal{F}(S, W) = \underset{S,W}{\operatorname{argmin}} \sum_{i=1}^{I} \mathcal{L}(Y_i, \hat{Y}_i) + \lambda_W \|W\|^2$$

- We can find the optimal shapelet for the objective function in a NN fashion by updating the shapelets in the minimum direction of the objective, hence the first gradient. Similarly, the weights can be jointly updated towards minimizing the objective function.

# Motif/Shapelet Summary

- A **motif** is a repeated pattern/subsequence in a given TS.



- A **shapelet** is a pattern/subsequence which is maximally representative of a class with respect to a given dataset of TSs.



Shapelet

Verbena urticifolia

Urtica dioica

# References

- Matrix Profile I: All Pairs Similarity Joins for Time Series:  A Unifying View that Includes Motifs, Discords and Shapelets. Chin-Chia Michael Yeh et al. 1997

- Time Series Shapelets: A New Primitive for Data Mining. Lexiang Ye and Eamonn Keogh. 2016.

- Josif Grabocka, Nicolas Schilling, Martin Wistuba, Lars Schmidt-Thieme (2014): Learning Time-Series Shapelets, in Proceedings of the 20th ACM SIGKDD Conference on Knowledge Discovery and Data Mining, KDD 2014

- Deep learning for time series classication: a review. Hassan Ismail Fawaz et al. 2019.