

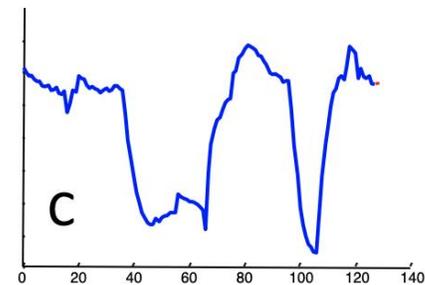
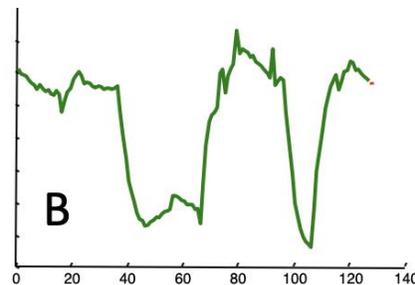
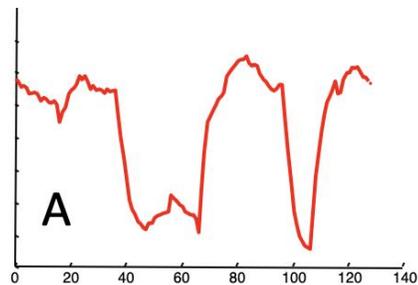
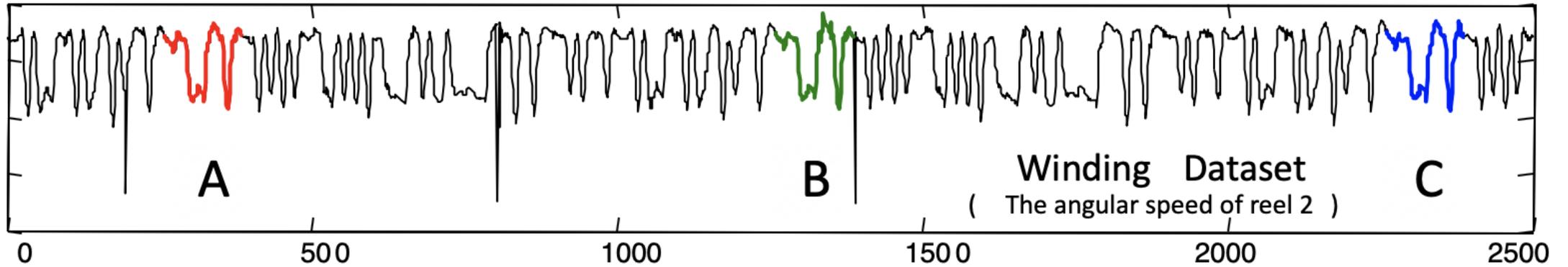
Time Series - Shapelet/Motif Discovery



Motif

Time Series Motif Discovery

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

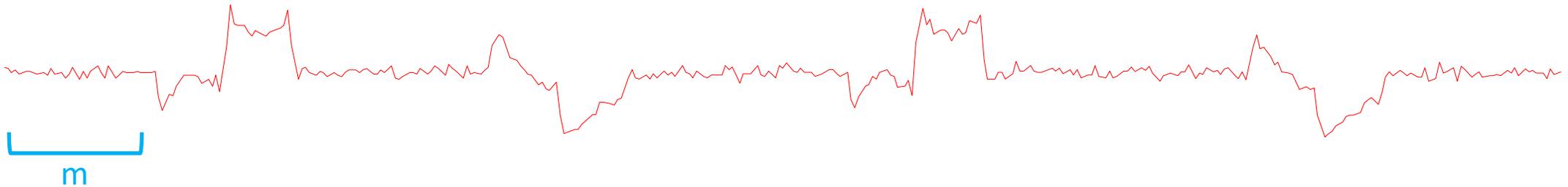


Why Find Motifs?

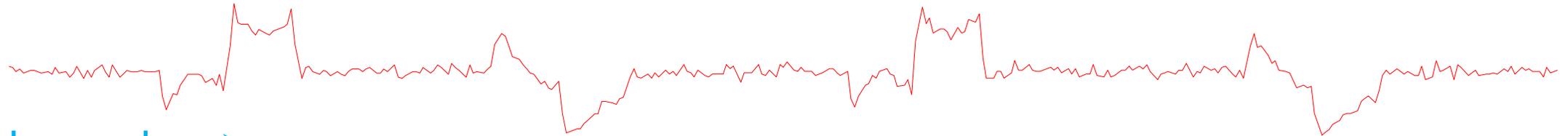
- Mining **association rules** in TS requires the discovery of motifs. These are referred to as primitive shapes and frequent patterns.
- Several **TS classifiers** work by constructing typical prototypes of each class. These prototypes may be considered motifs.
- Many **TS anomaly detection** algorithms 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.

Matrix Profile

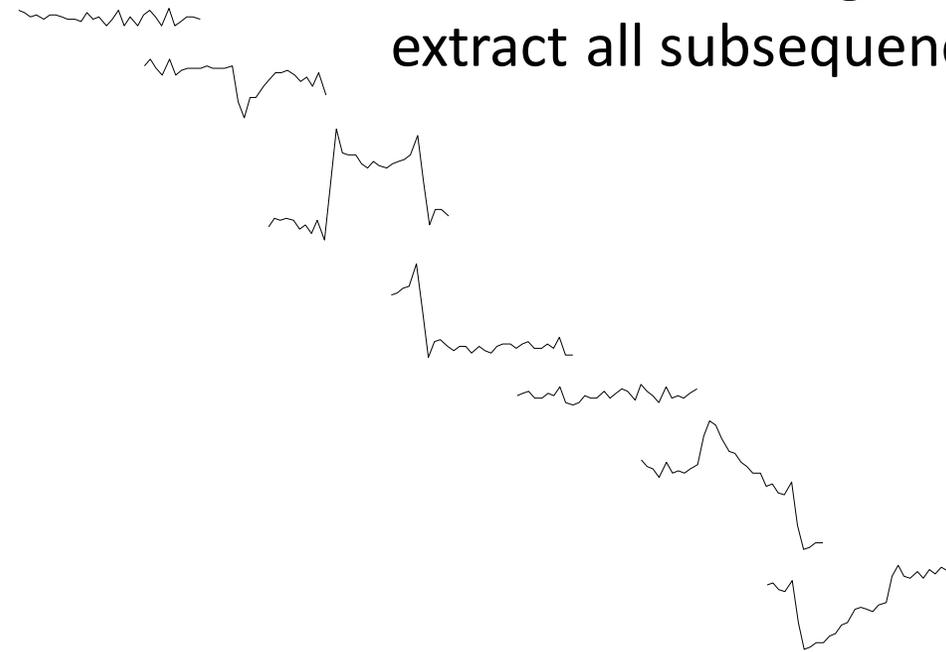
- The Matrix Profile (MP) is a data structure that annotates a TS and can be exploited for many purposes: e.g. efficient Motif Discovery.
- Given a time series, T and a desired subsequence length, m .



Matrix Profile



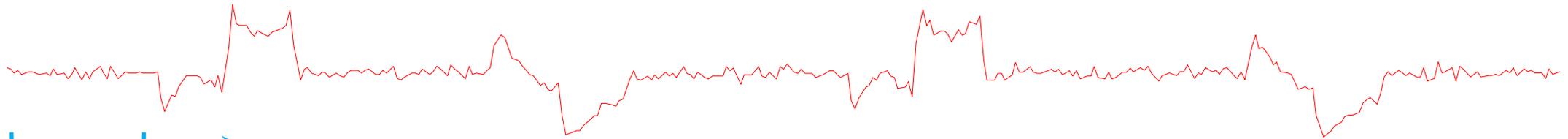
We can use sliding window of length m to extract all subsequences of length m .



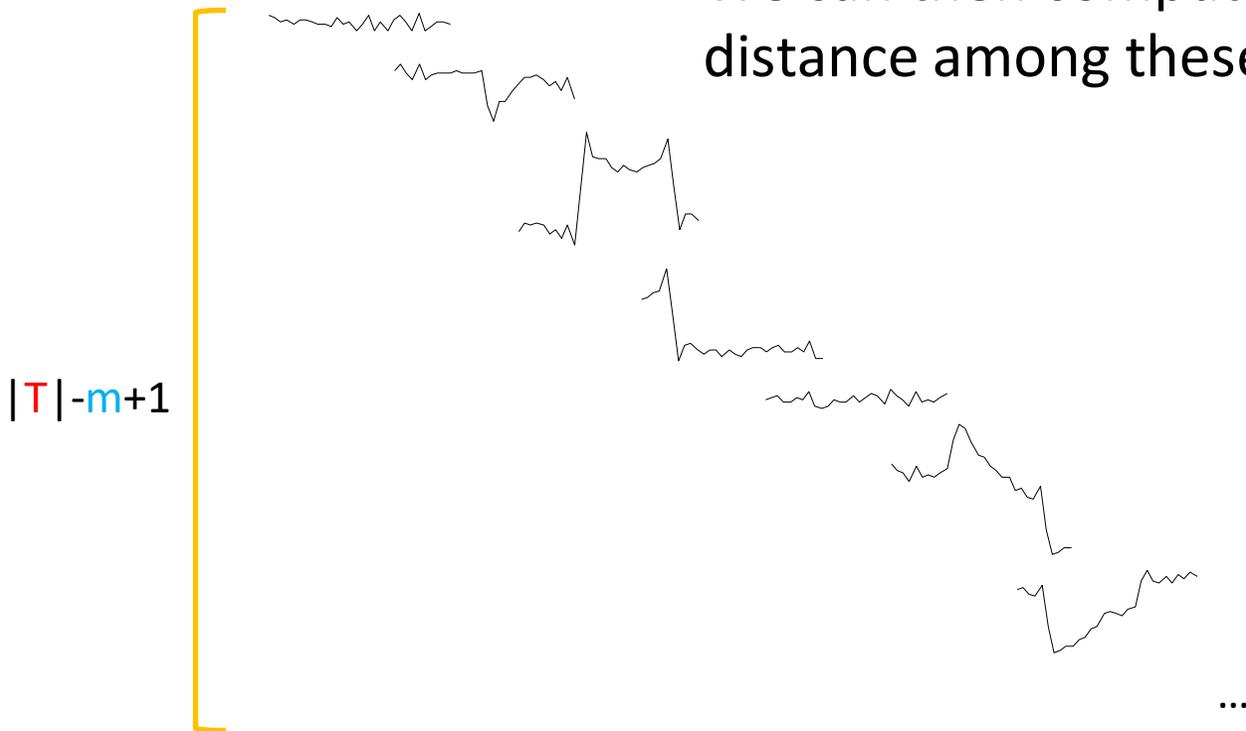
...

$|T| - m + 1$

Matrix Profile



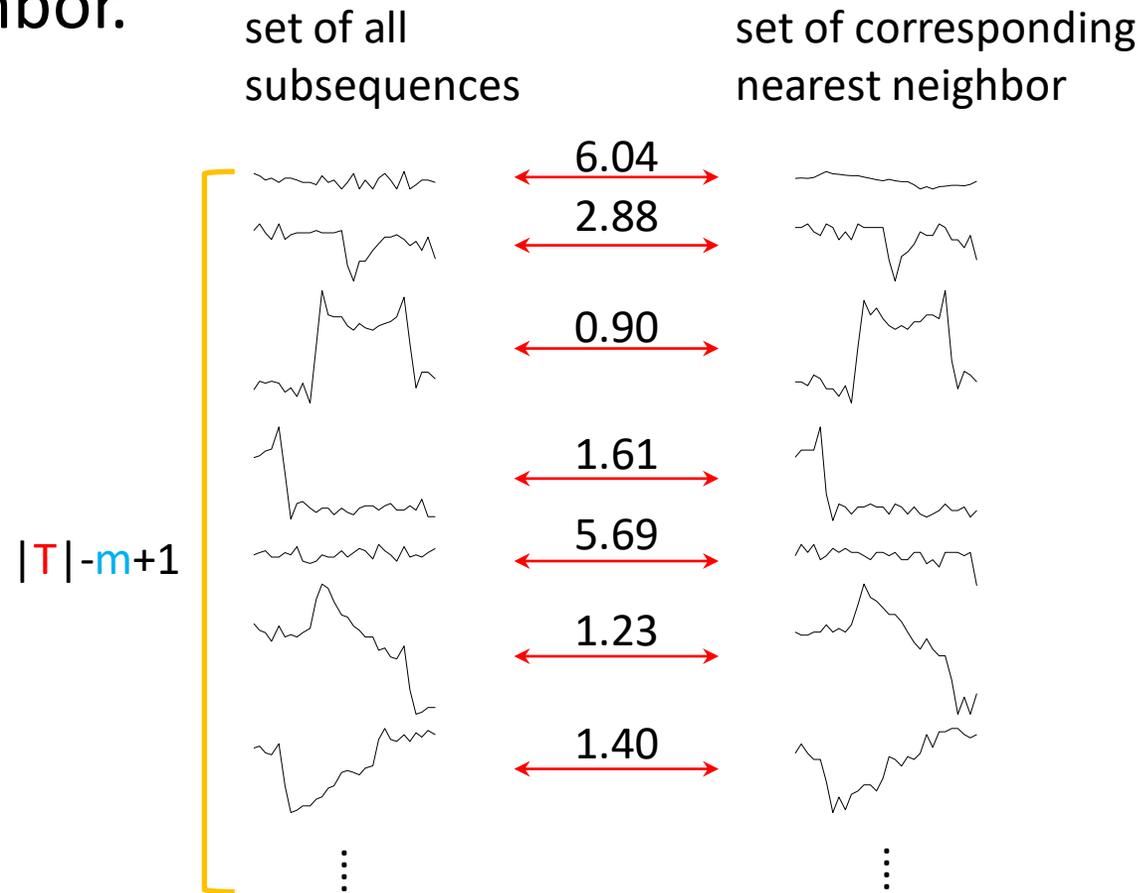
We can then compute the pairwise distance among these subsequences.



0	7.6952	7.7399	...
7.6952	0	7.7106	...
7.7399	7.7106	0	...
...

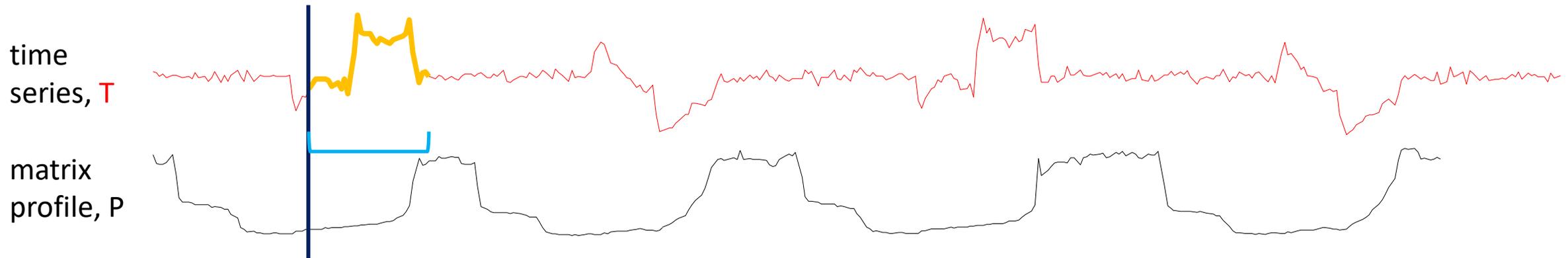
Matrix Profile

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



Matrix Profile

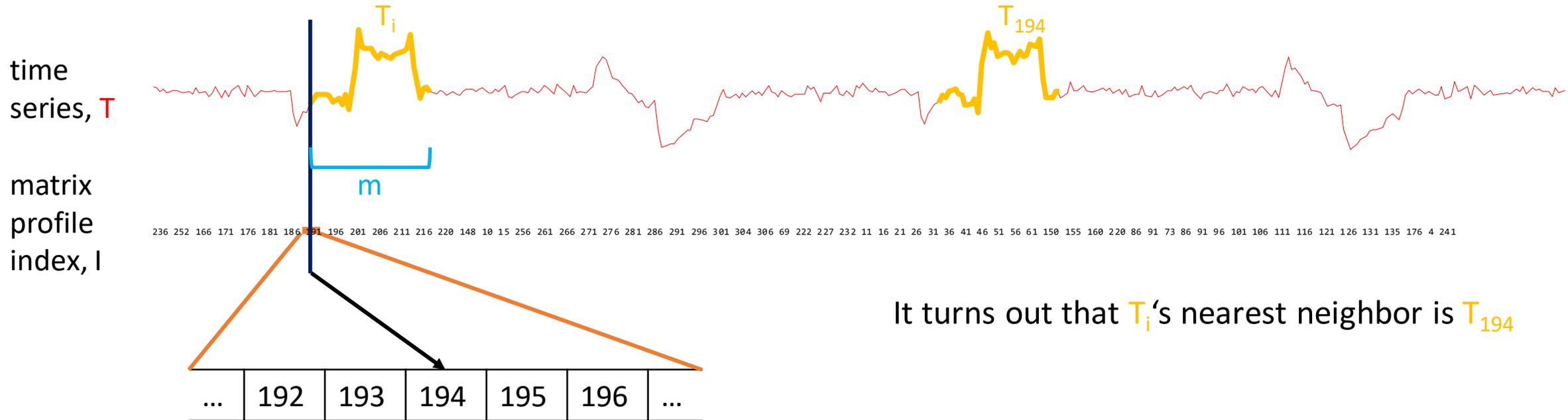
- The **distance** to the corresponding **nearest neighbor** of each subsequence can be stored in a vector called **matrix profile P**.



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**.

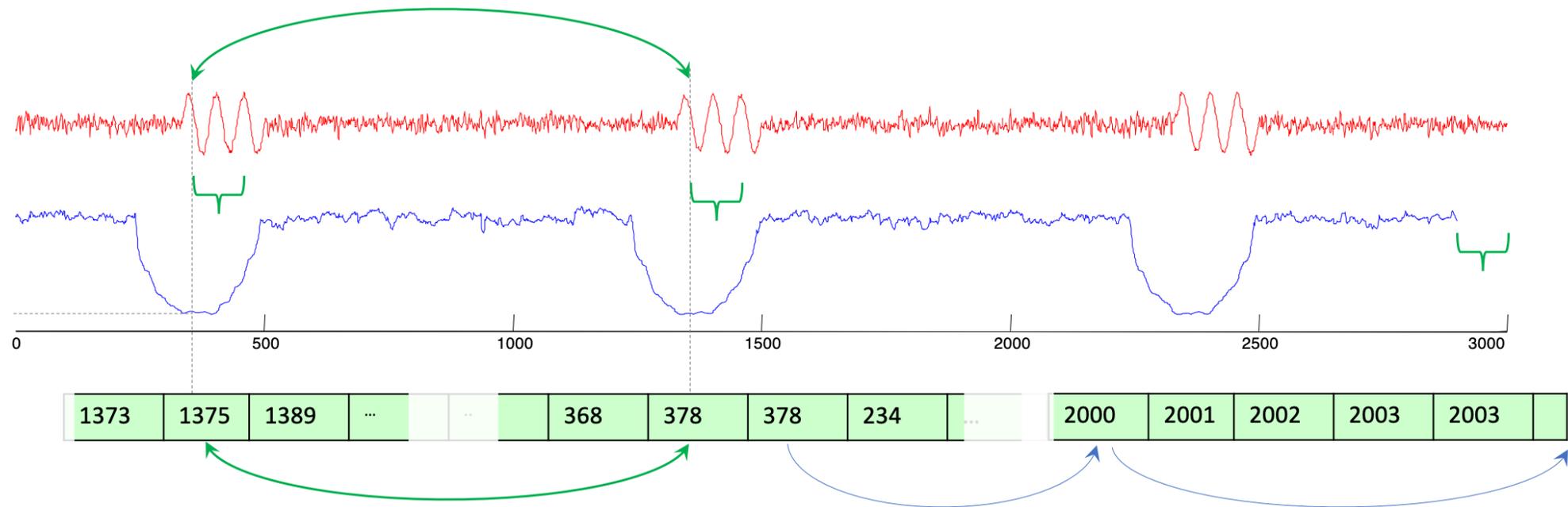


It turns out that T_i 's nearest neighbor is T_{194}

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

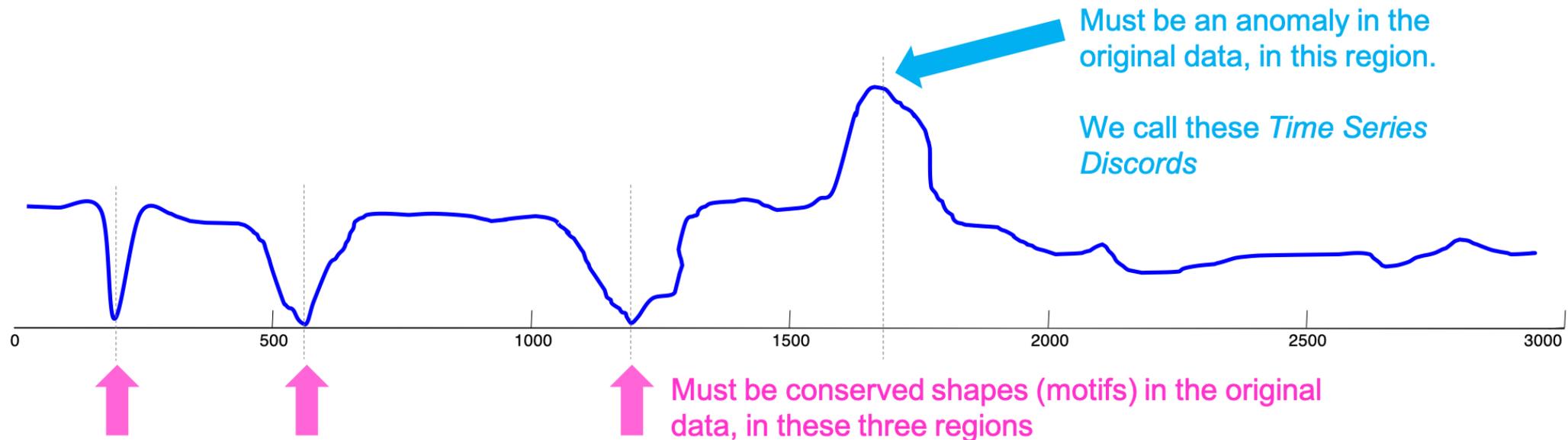
Matrix Profile

- 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.



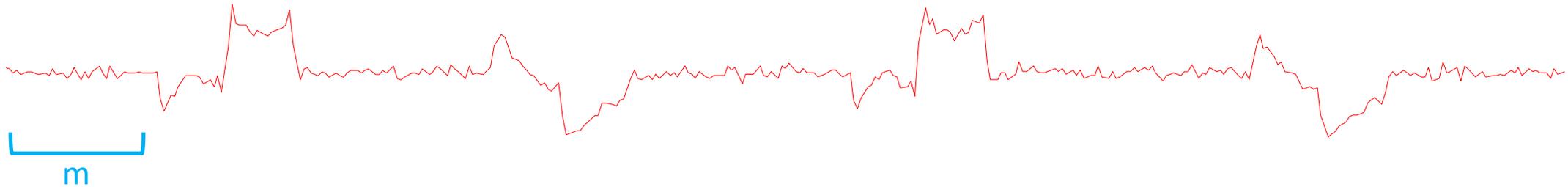
How to “read” a Matrix Profile

- 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).



How to Compute Matrix Profile?

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



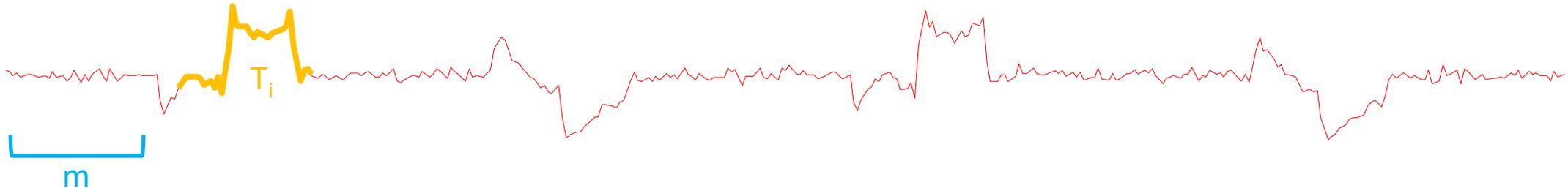
inf	inf	inf	inf	inf	inf	inf	inf	inf	inf	inf	inf	inf	inf	inf	inf	inf	inf	inf
-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----

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 .

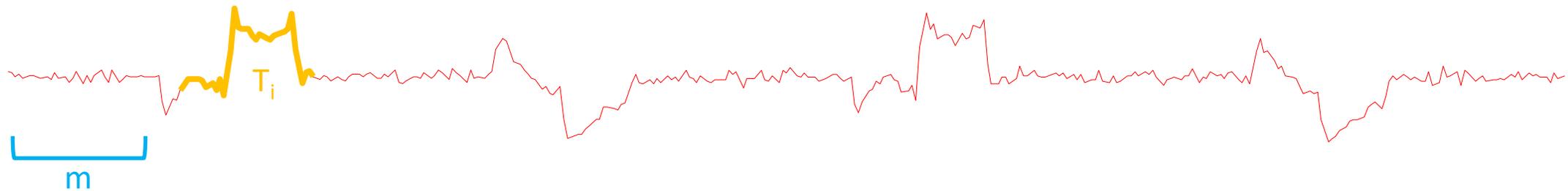


inf	inf	inf	inf	inf	inf	inf	inf	inf	inf	inf	inf	inf	inf	inf	inf	inf	inf	inf
-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----

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
-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----

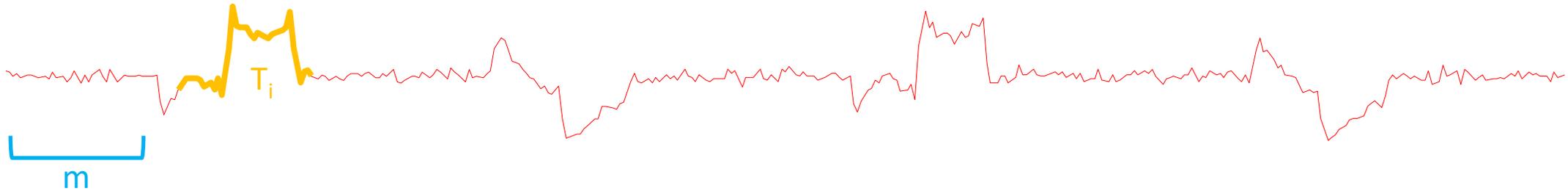
We compute the distances between T_i and every subsequences from T (time complexity = $O(|T| \log(|T|))$)
We then put the distances in a vector based on the position of the subsequences

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 .



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

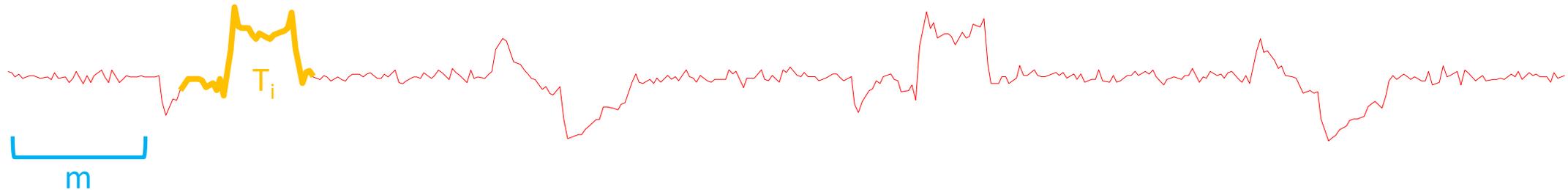
3	2	0	5	3	4	5	1	2	9	8	4	2	3	4	8	6	2	1
---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---



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, 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
-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----

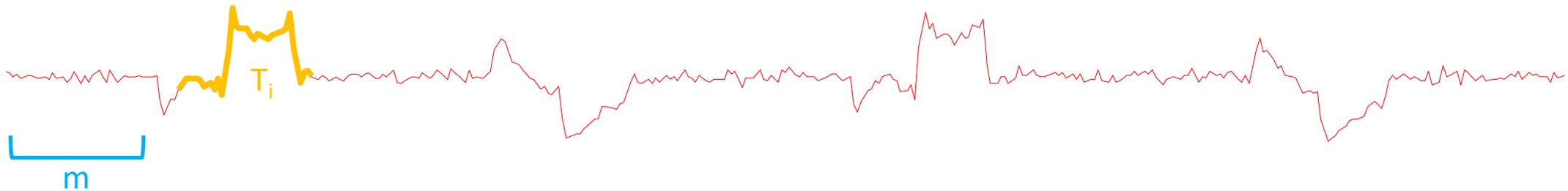
min
↕

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 .



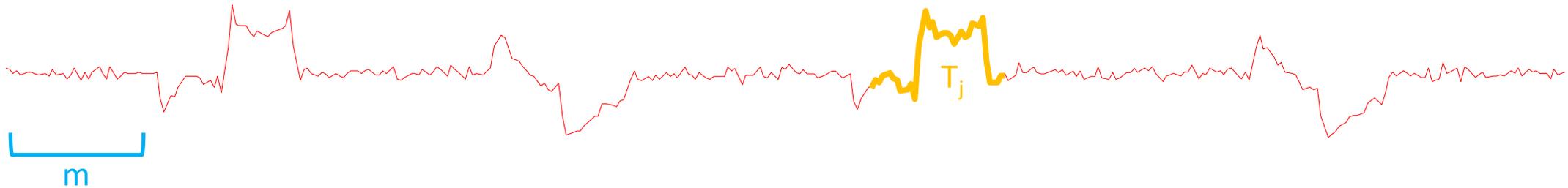
3	2	inf	5	3	4	5	1	2	9	8	4	2	3	4	8	6	2	1
---	---	-----	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---

After we finish to 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 .

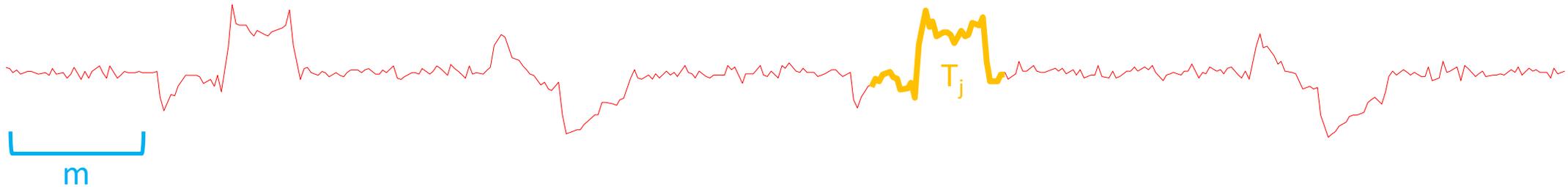


3	2	inf	5	3	4	5	1	2	9	8	4	2	3	4	8	6	2	1
---	---	-----	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---

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

How to Compute Matrix Profile?

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



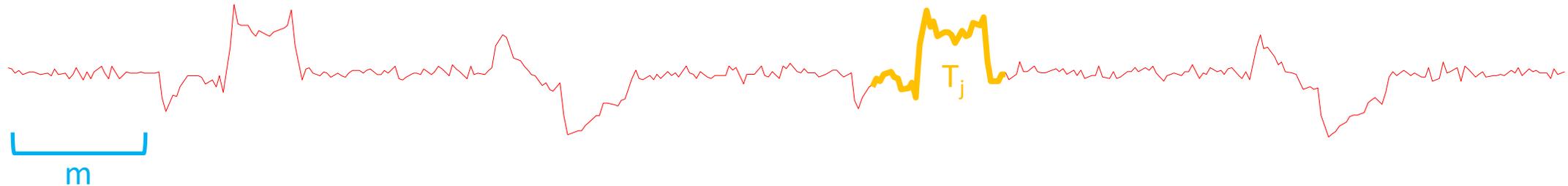
3	2	inf	5	3	4	5	1	2	9	8	4	2	3	4	8	6	2	1
---	---	-----	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---

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
---	---	-----	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---

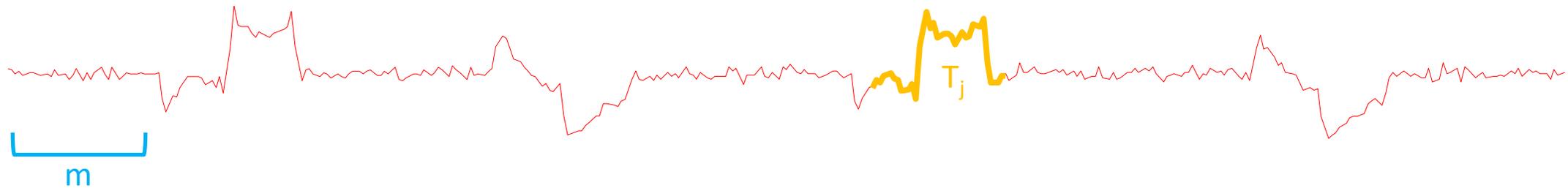
min

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
---	---	-----	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---

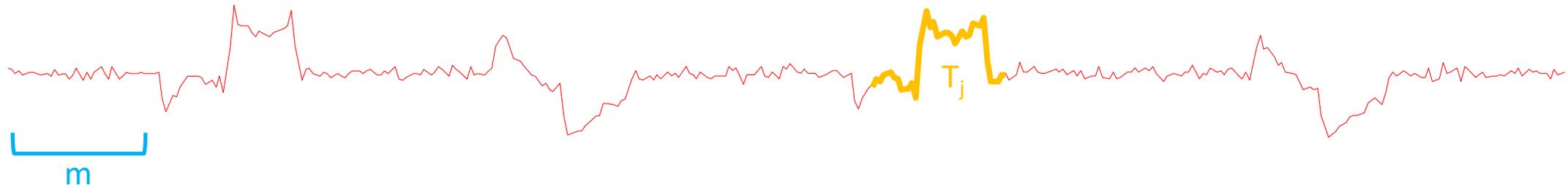
min \updownarrow

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
---	---	-----	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---

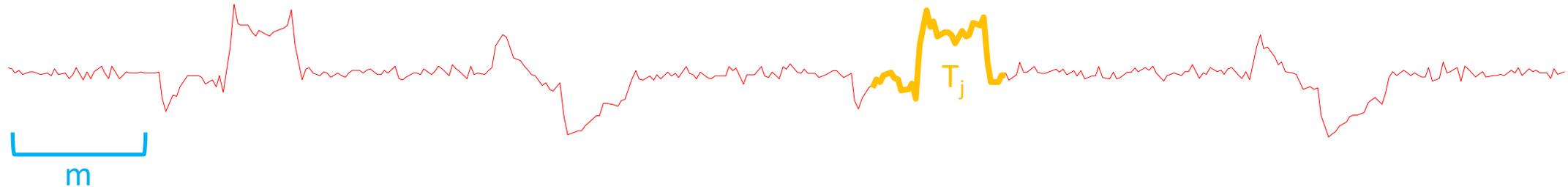
min \updownarrow

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

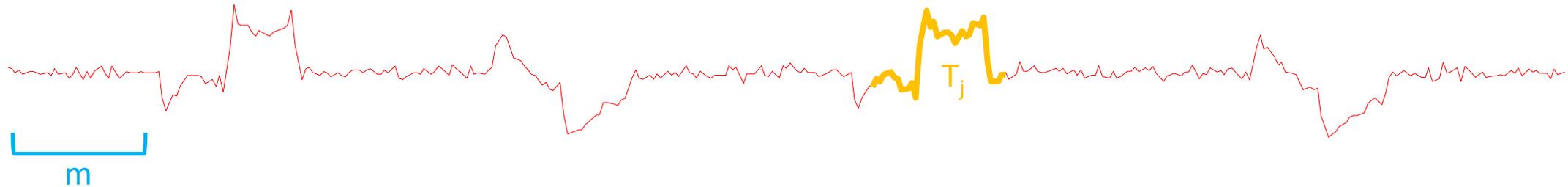
min

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
---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---

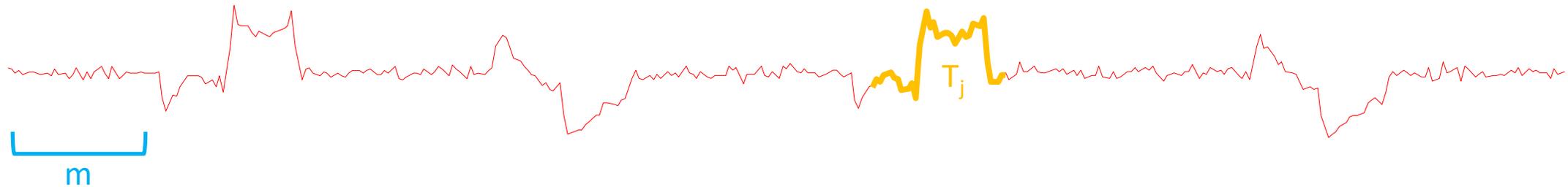
min \updownarrow

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 .



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

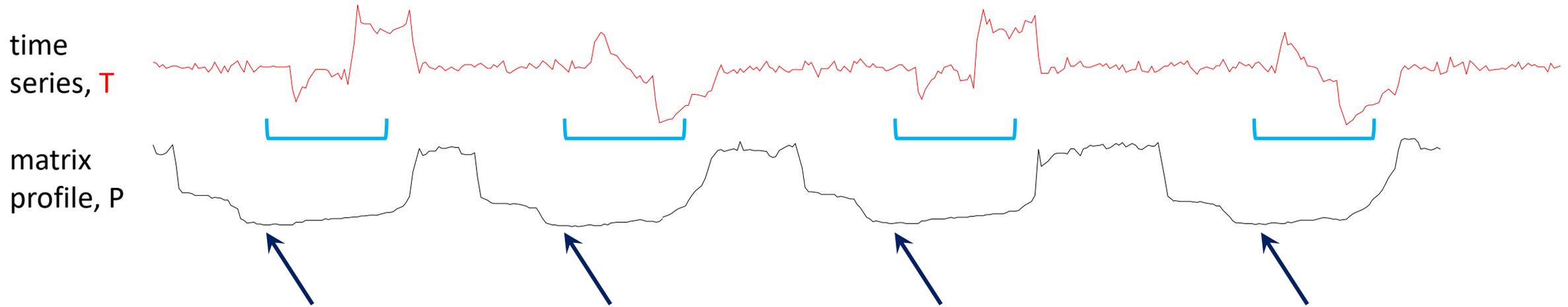
min

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 $O(|T|\log(|T|))$

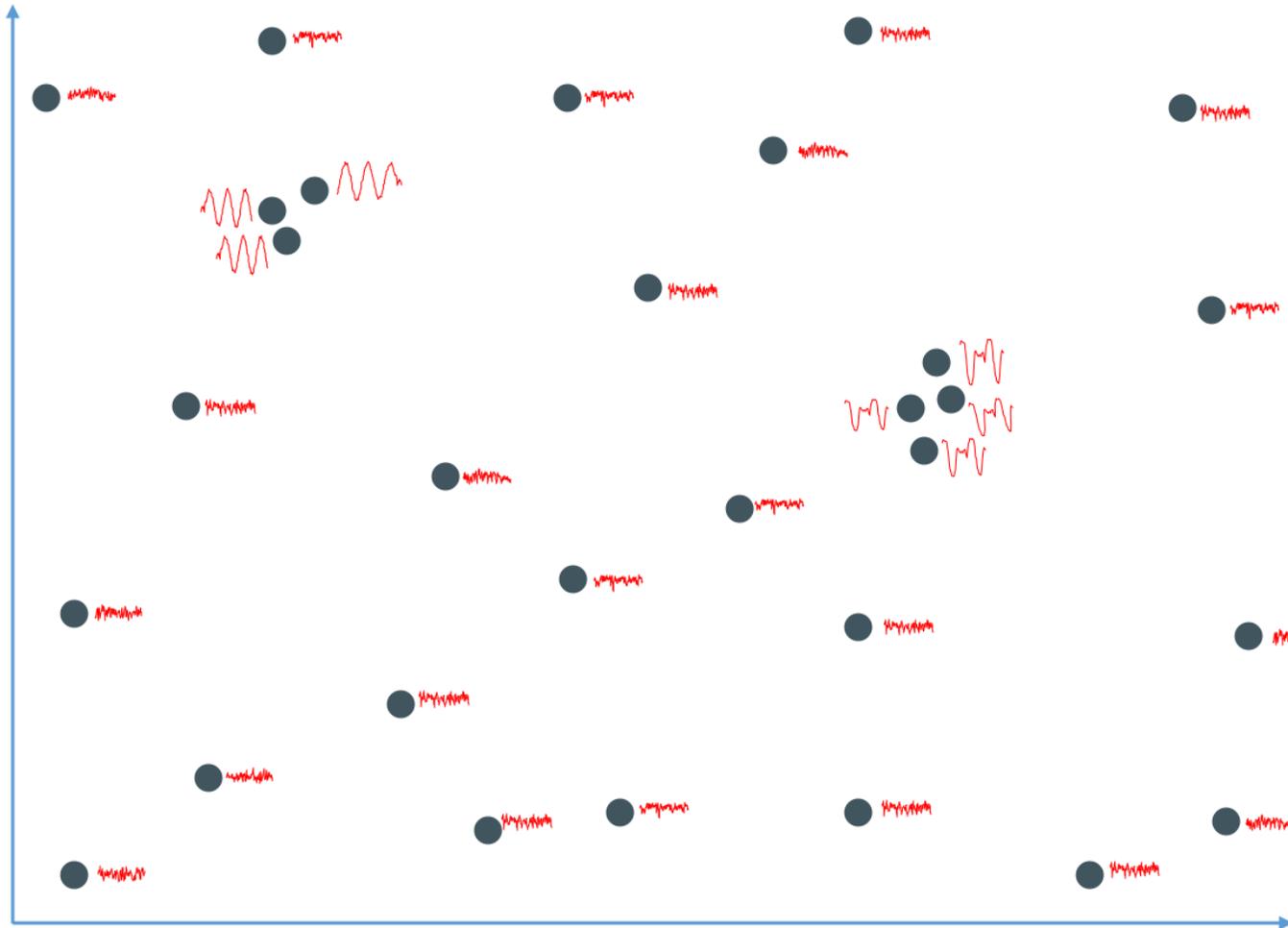
The overall time complexity is $O(|T|^2\log(|T|))$

Motif Discovery From Matrix Profile

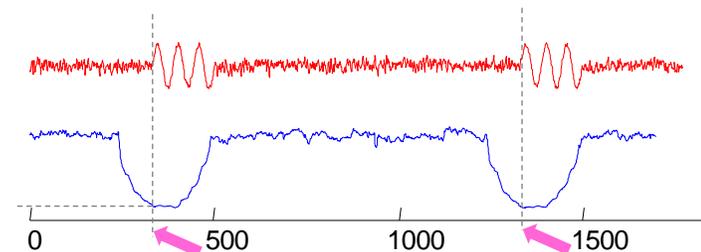


Local minimums are corresponding to motifs

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.



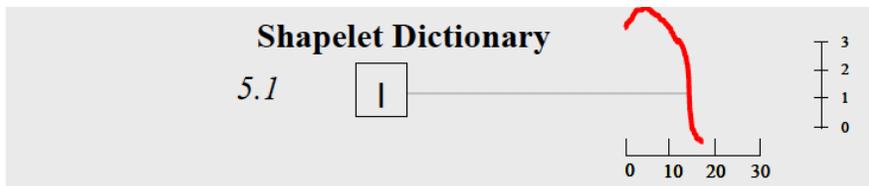
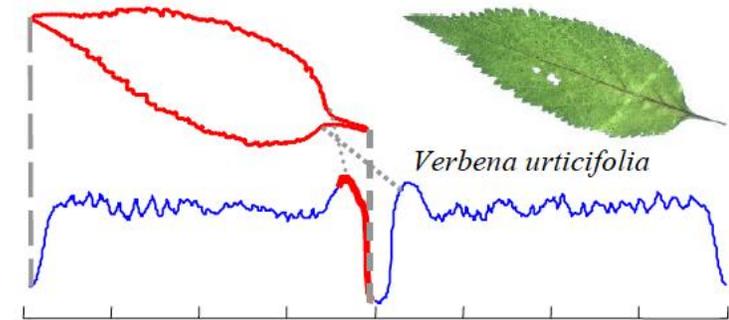
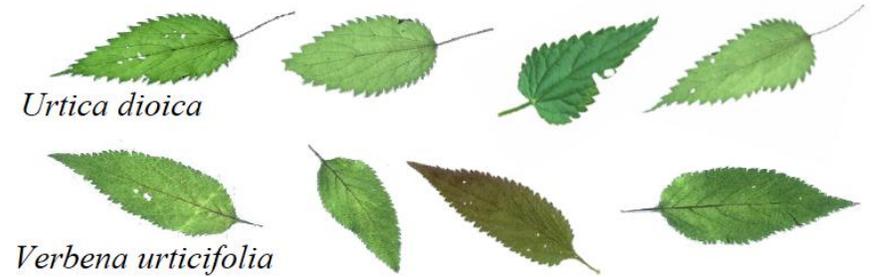
Shapelet

Time Series Classification

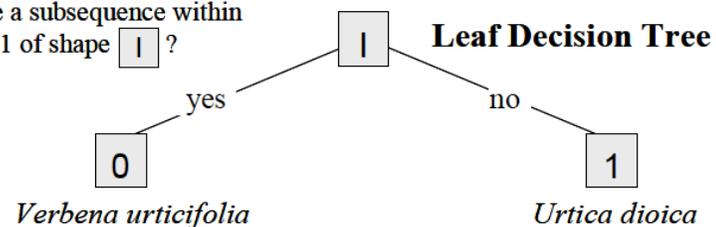
- Given a set X of n time series, $X = \{x_1, x_2, \dots, x_n\}$, each time series has m ordered values $x_i = \langle x_{t1}, x_{t2}, \dots, x_{tm} \rangle$ and a class value c_i .
- The objective is to find a function f that maps from the space of possible time series to the space of possible class values.
- Generally, it is assumed that all the TS have the same length m .

Shapelet-based Classification

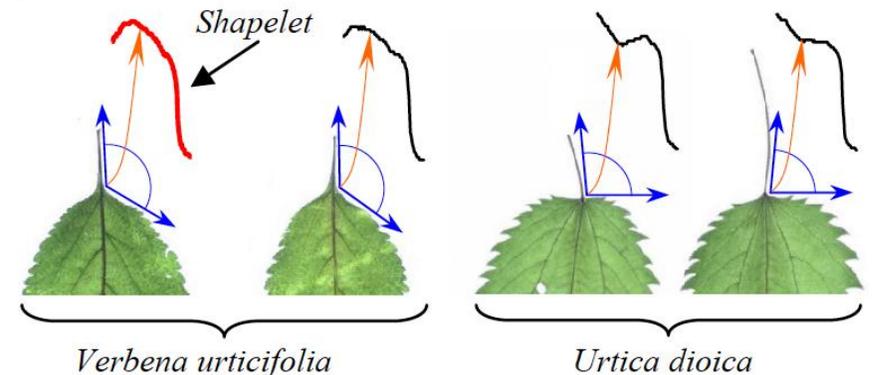
1. Represent a TS as a vector of distances with representative subsequences, namely shapelets.
2. Use it as input for machine learning classifiers.



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

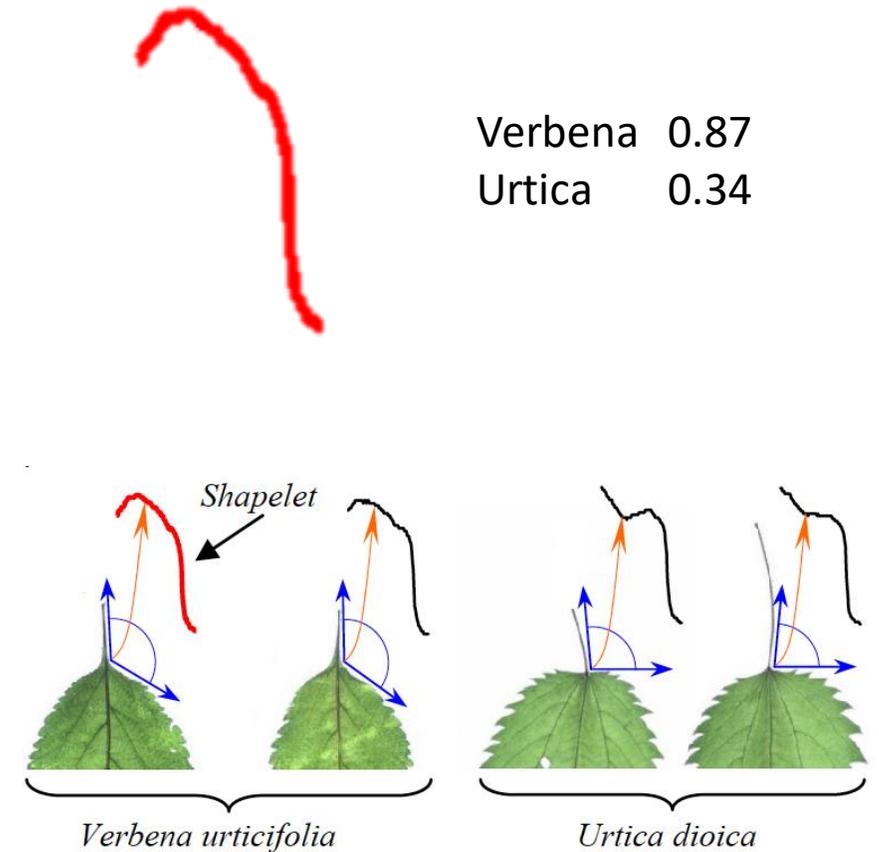


3.2	8.7
1.4	7.9
6.7	4.2
9.2	3.4



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*.



Finding Shapelets

FindingShapeletBF (dataset \mathbf{D} , $MAXLEN$, $MINLEN$)

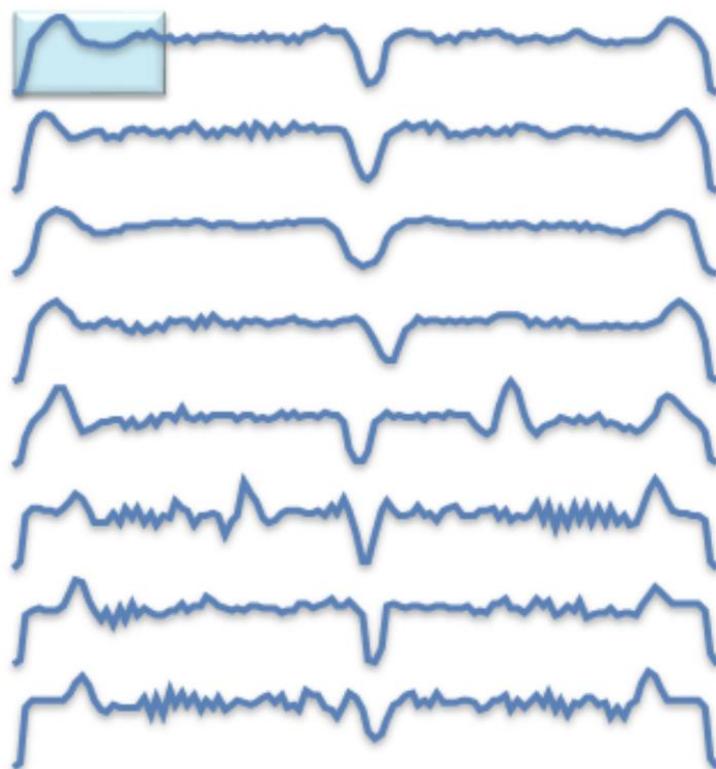
```
1  candidates ← GenerateCandidates( $\mathbf{D}$ ,  $MAXLEN$ ,  $MINLEN$ )
2  bsf_gain ← 0
3  For each  $S$  in candidates
4      gain ← CheckCandidate( $\mathbf{D}$ ,  $S$ )
5      If gain > bsf_gain
6          bsf_gain ← gain
7          bsf_shapelet ←  $S$ 
8      EndIf
9  EndFor
10 Return bsf_shapelet
```

Generate Candidate

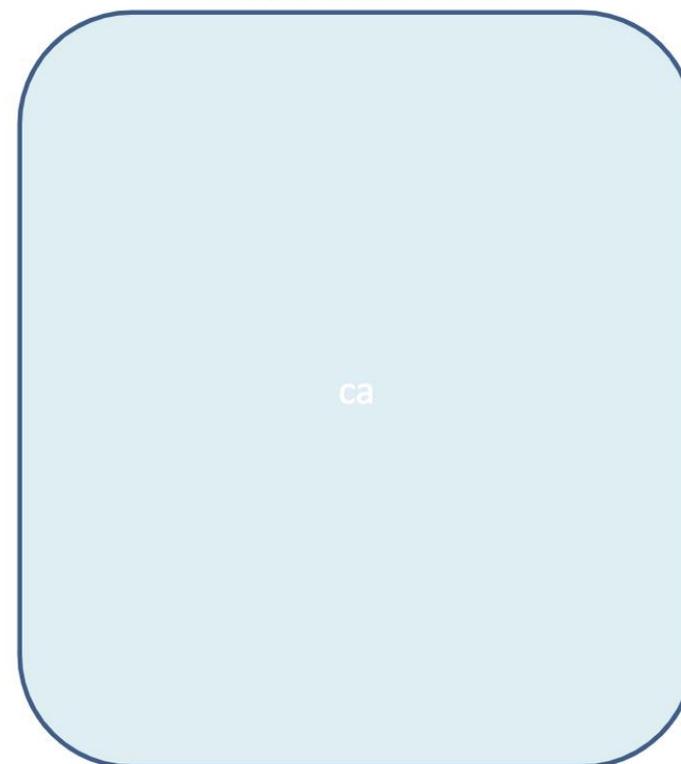
Sliding a **window of size l** across all of the time series objects in the dataset \mathbf{D} , extracts all of the possible candidates and adds them to the pool

GenerateCandidates (dataset \mathbf{D} , $MAXLEN$, $MINLEN$)	
1	$pool \leftarrow \emptyset$
2	$l \leftarrow MAXLEN$
3	While $l \geq MINLEN$
4	For T in \mathbf{D}
5	$pool \leftarrow pool \cup \mathbf{S}_T^l$
6	EndFor
7	$l \leftarrow l - 1$
8	EndWhile
9	Return $pool$

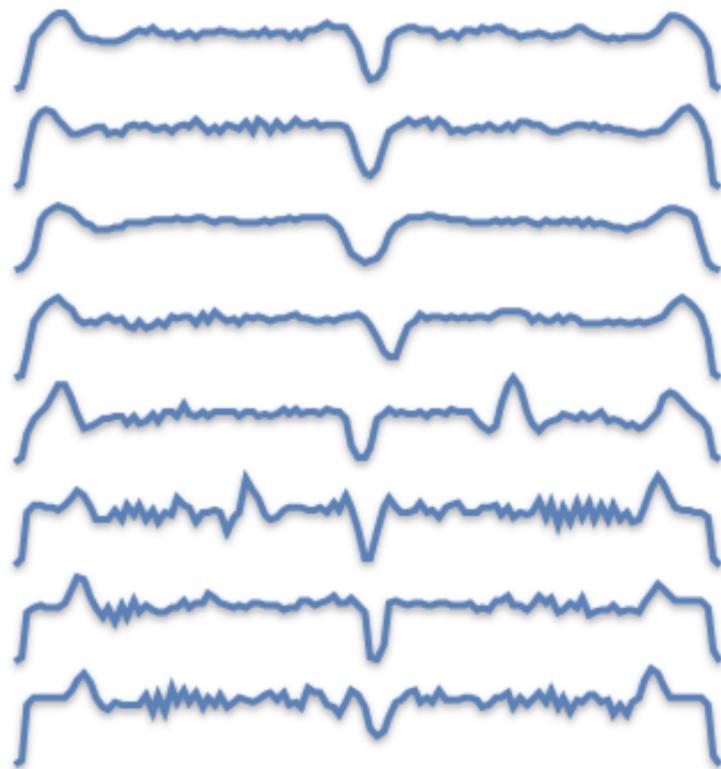
Extract Subsequences of all Possible Lengths



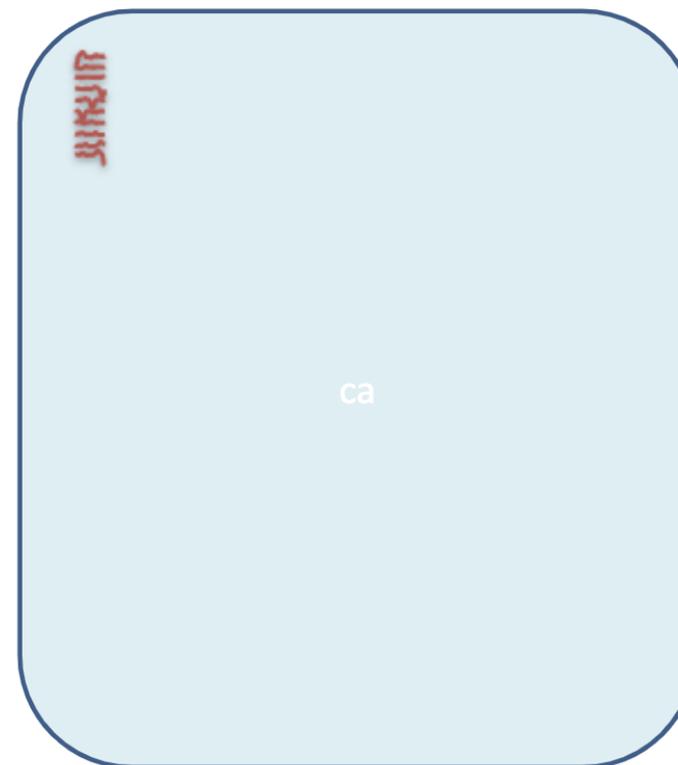
Candidates Pool



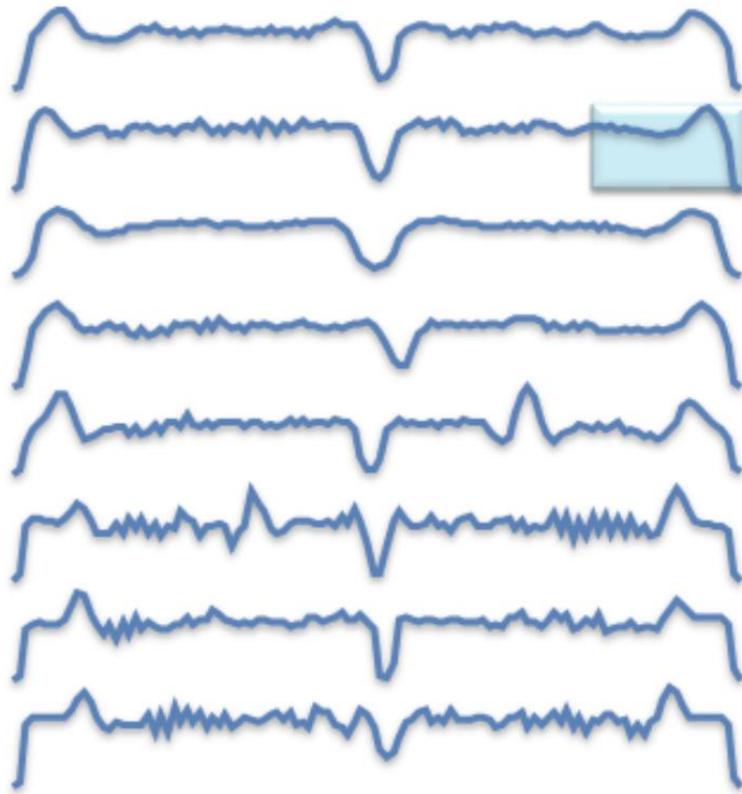
Extract Subsequences of all Possible Lengths



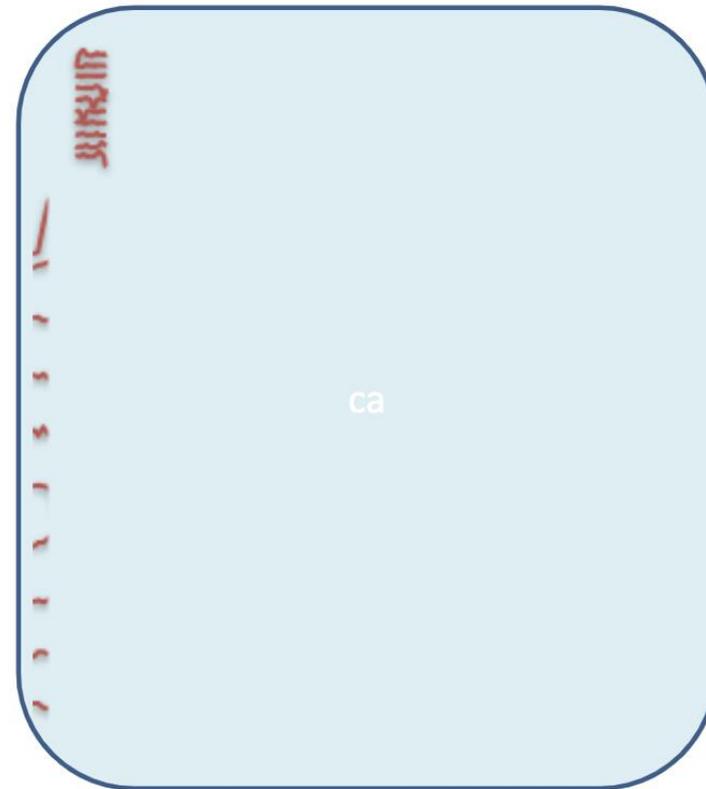
Candidates Pool



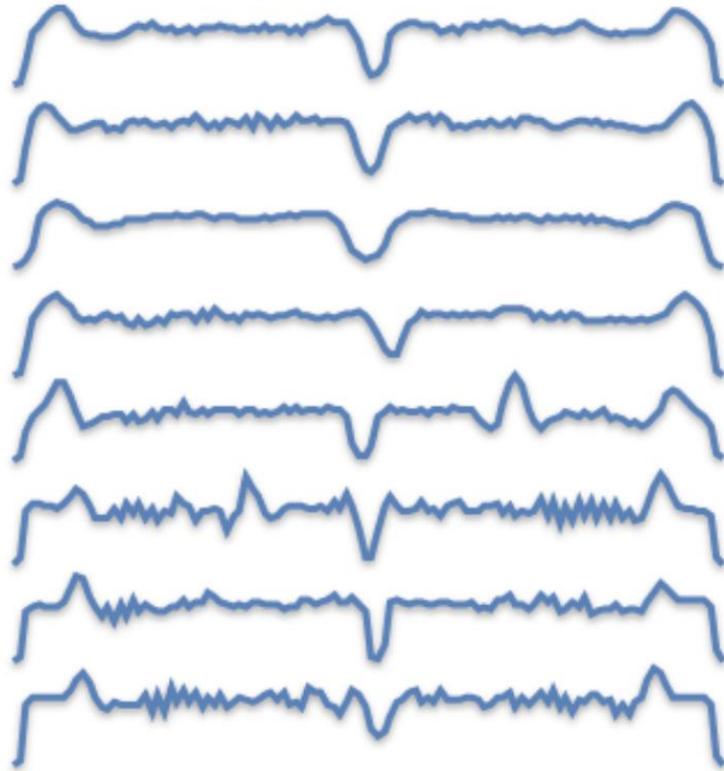
Extract Subsequences of all Possible Lengths



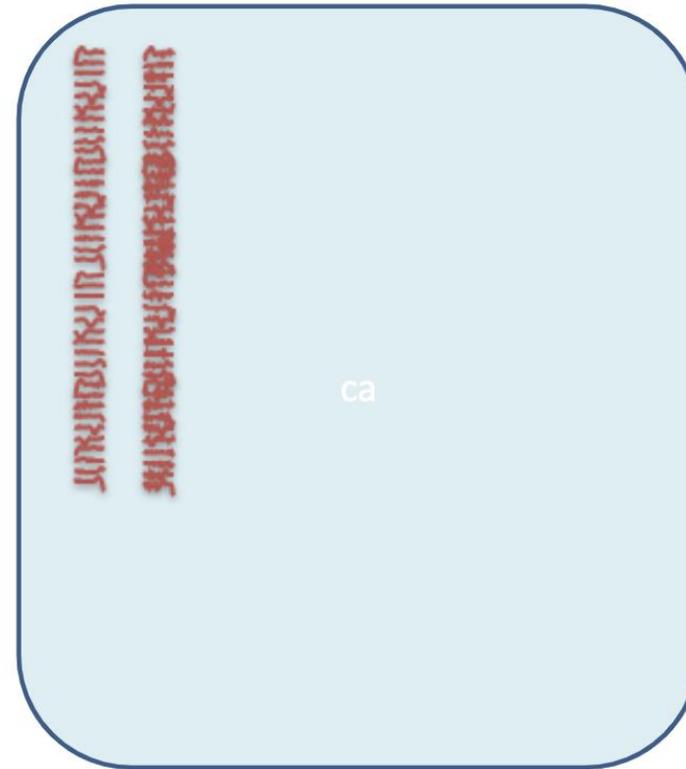
Candidates Pool



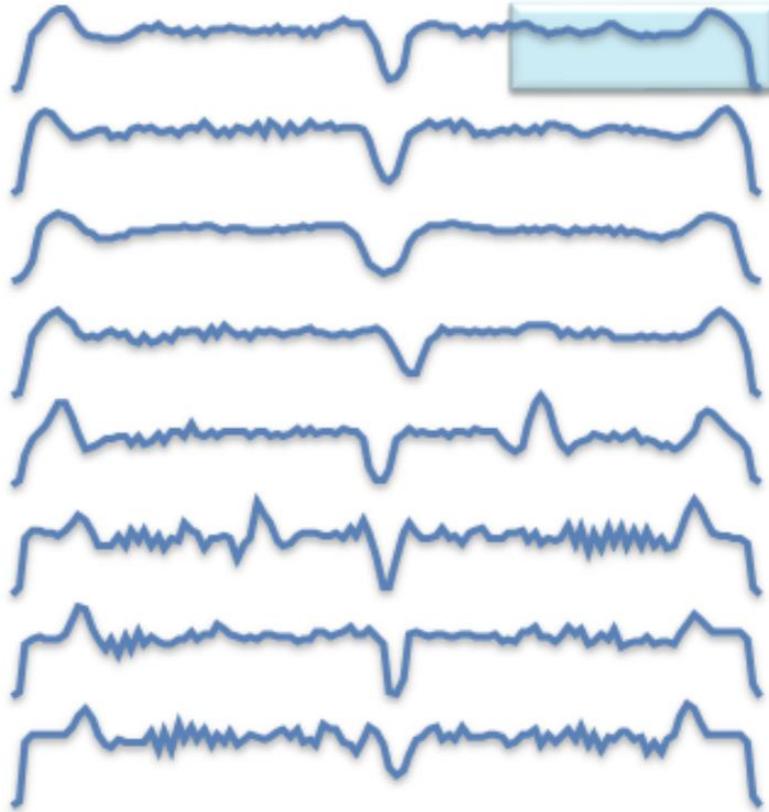
Extract Subsequences of all Possible Lengths



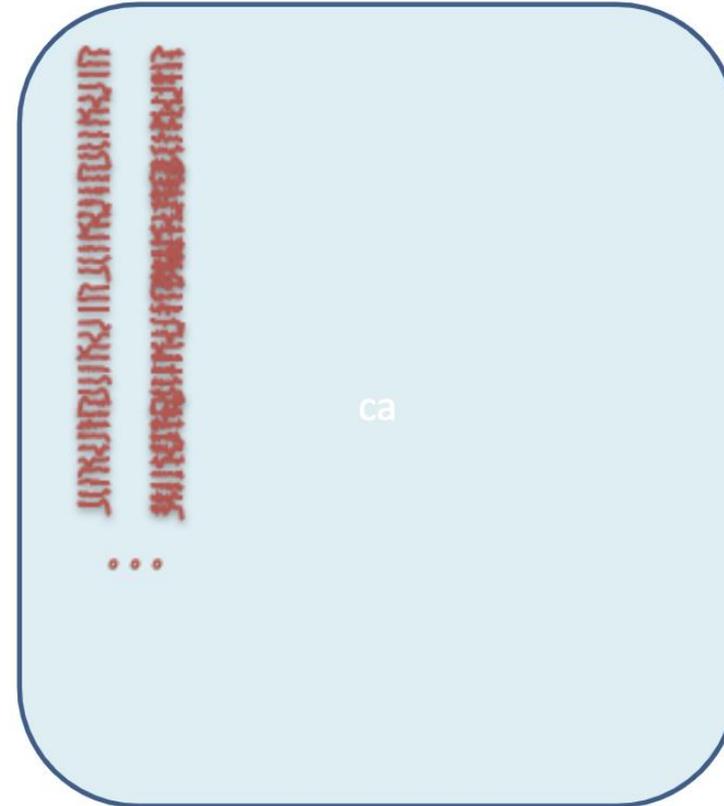
Candidates Pool



Extract Subsequences of all Possible Lengths



Candidates Pool

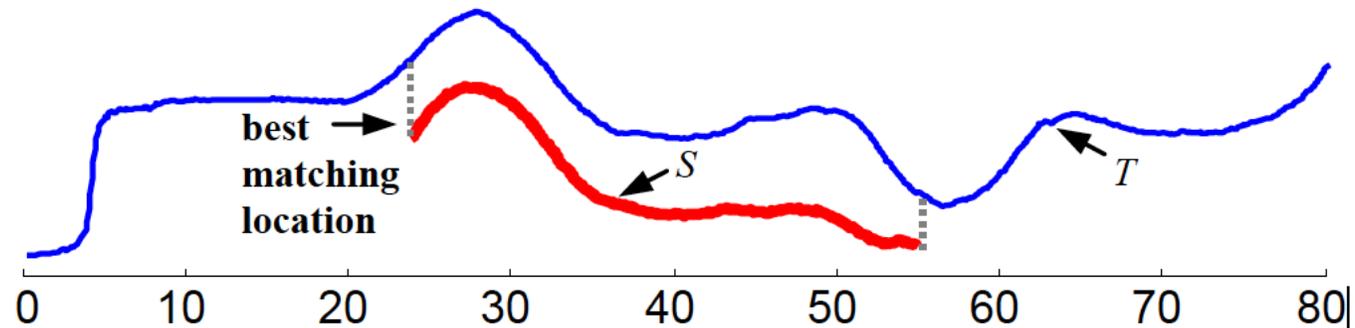


Check Candidates

CheckCandidate (dataset \mathbf{D} , shapelet candidate S)	
1	$objects_histogram \leftarrow \emptyset$
2	For each T in \mathbf{D}
3	$dist \leftarrow \text{SubsequenceDist}(T, S)$
4	insert T into $objects_histogram$ by the key $dist$
5	EndFor
6	Return CalculateInformationGain($objects_histogram$)

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

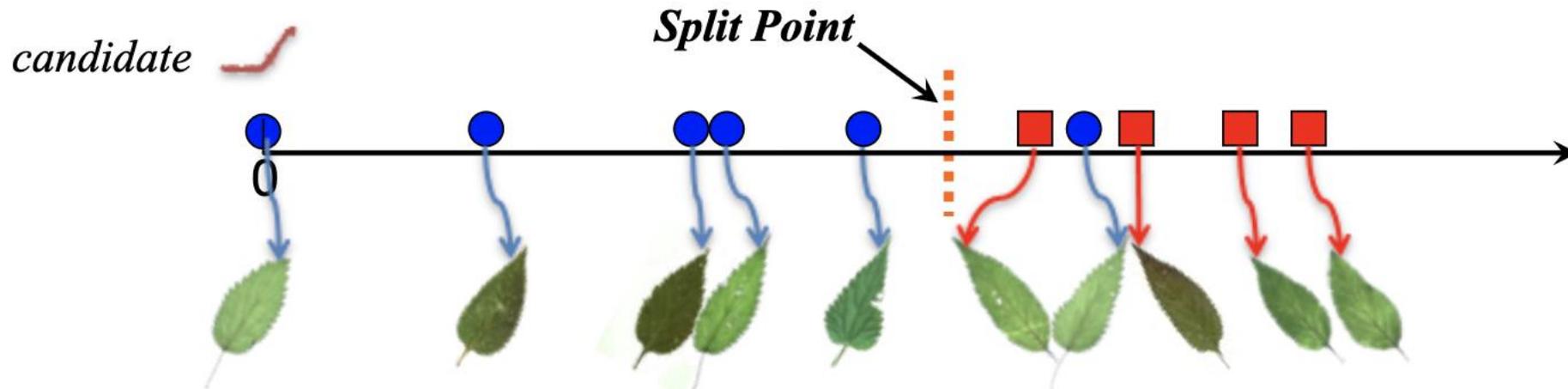


Check Candidates

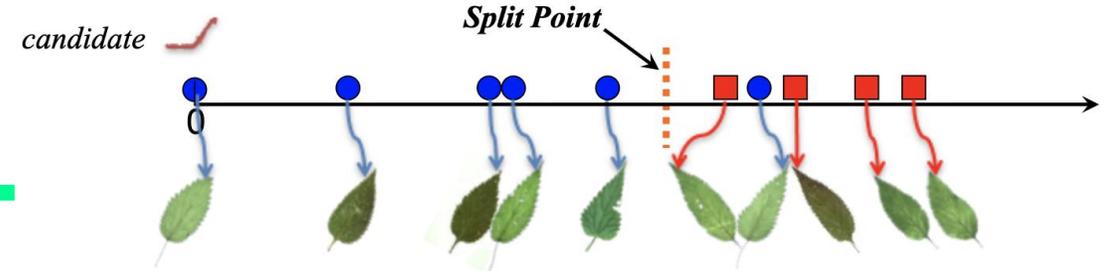
CheckCandidate (dataset \mathbf{D} , shapelet candidate S)	
1	$objects_histogram \leftarrow \emptyset$
2	For each T in \mathbf{D}
3	$dist \leftarrow \text{SubsequenceDist}(T, S)$
4	insert T into $objects_histogram$ by the key $dist$
5	EndFor
6	Return CalculateInformationGain($objects_histogram$)

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

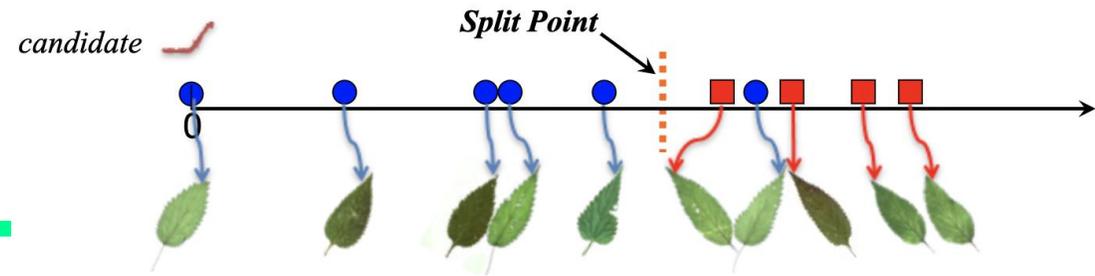


Entropy



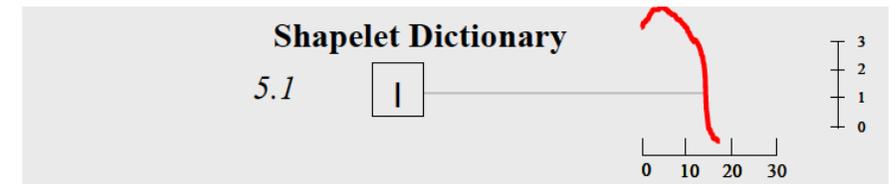
- A TS dataset D consists of two classes, A and B.
- Given that the proportion of objects in class A is $p(A)$ and the proportion of objects in class B is $p(B)$,
- The **Entropy** of D is: $I(D) = -p(A)\log(p(A)) - p(B)\log(p(B))$.
- Given a strategy that divides the D into two subsets D_1 and D_2 , 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

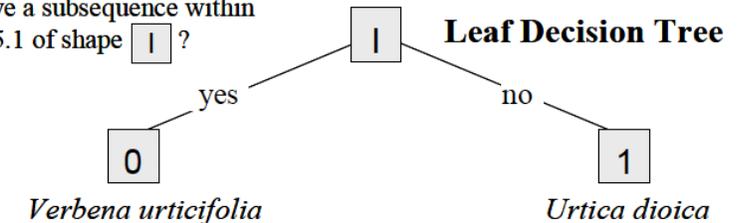


Split point
distance from
shapelet = 5.1

- 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 $\hat{I}(D)$.
- The **information gain** for this splitting rule is:
- $Gain(sp) = I(D) - \hat{I}(D) =$
- $= I(D) - f(D_1)I(D_1) + f(D_2)I(D_2).$
- We use the distance from T to a shapelet S as the splitting rule sp .



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

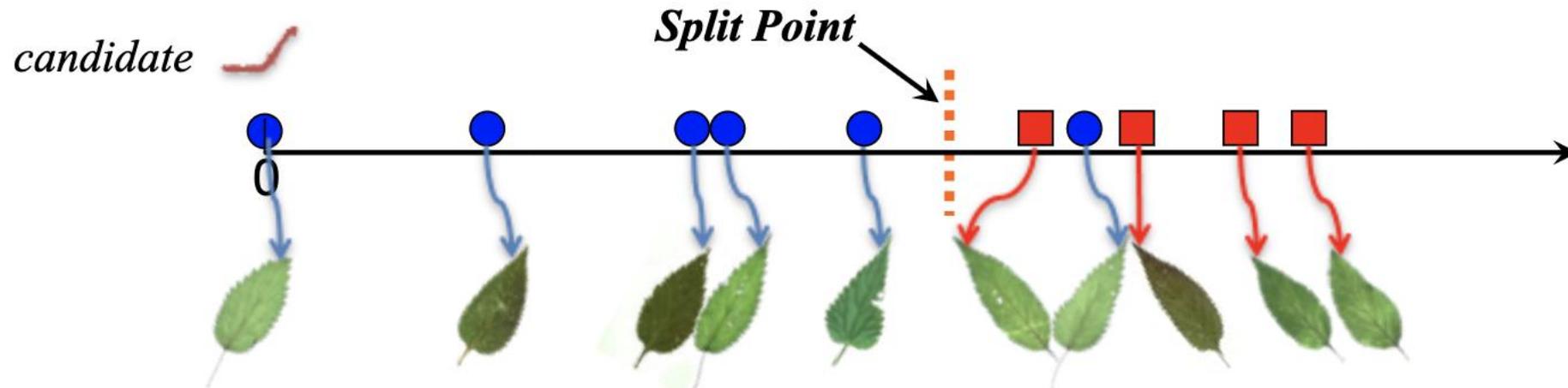


Problem

- 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 (space inefficiency)
- For instance
 - 200 instances with length 275
 - 7,480,200 shapelet candidates

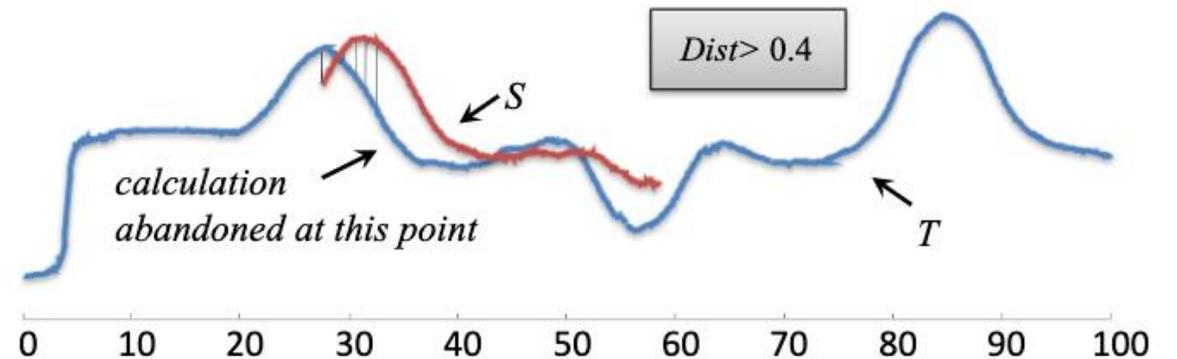
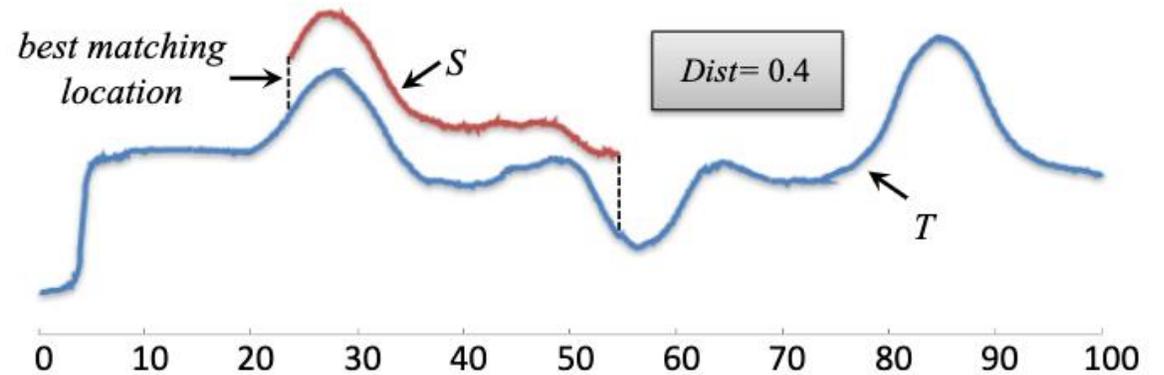
Speedup

- Distance calculations from 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



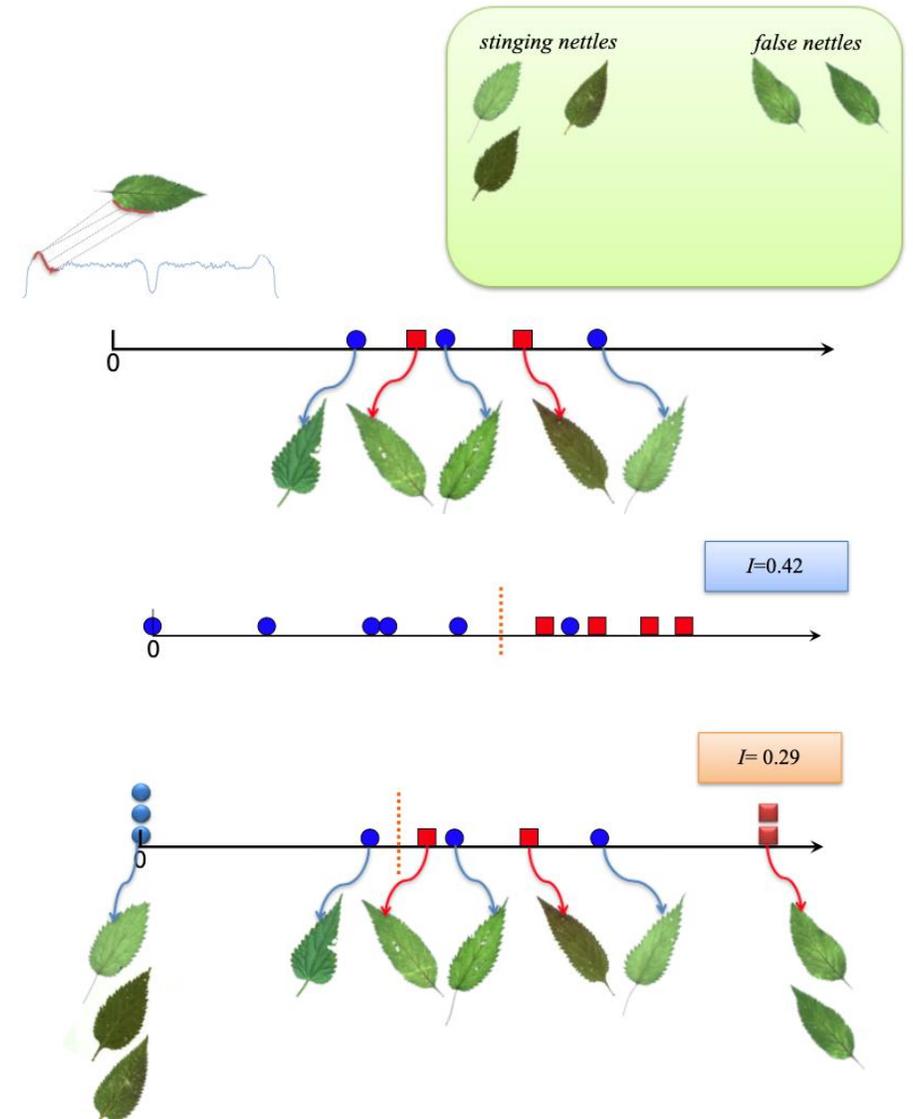
Distance Early Abandon

- We only need the minimum distance.
- Method
 - Keep the best-so-far distance
 - Abandon the calculation if the partial current distance is larger than best-so-far.
 - We can avoid to compute the full distance for S if the partial one is greater than the best so far



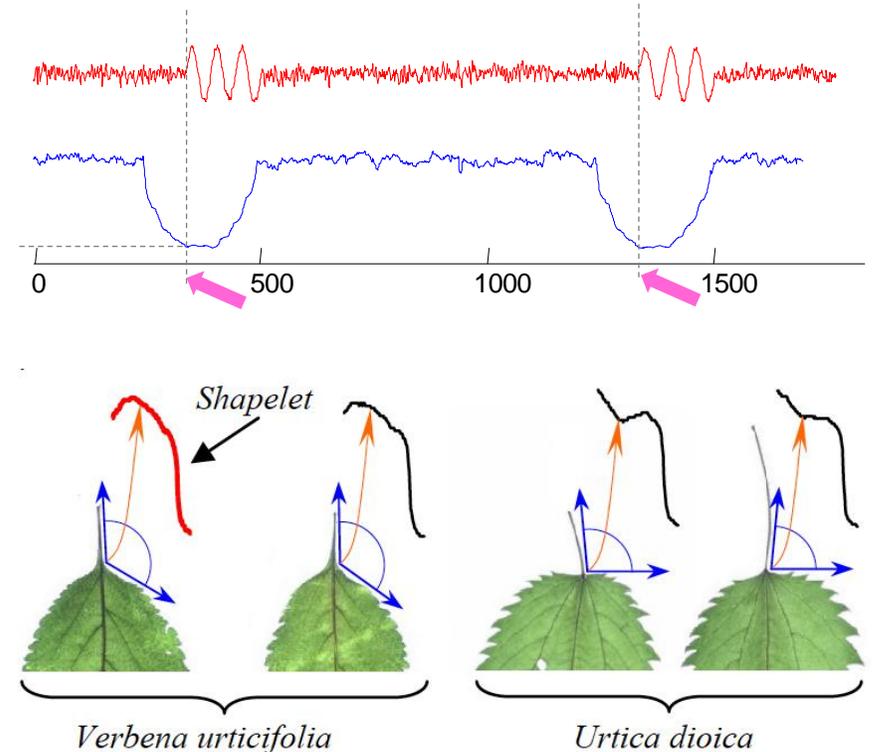
Admissible Entropy Pruning

- 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 (corresponding to the optimistic scenario) $<$ best candidate shapelet
 - Stop calculation
 - Try next candidate



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.



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

Matrix Profile I: All Pairs Similarity Joins for Time Series: A Unifying View that Includes Motifs, Discords and Shapelets

Chin-Chia Michael Yeh, Yan Zhu, Lindanila Ulanova, Nurjahan Begum, Yifei Ding, Hong Anli Diao, Diego Furtado Silva, Abhishall Mueen, and Eamonn Keogh

University of California, Riverside, University of São Paulo, University of New Mexico
(yeh003, yzhu01, hdu001, abg001, yding07, lulu001)@ucr.edu, d.furtado@usp.br, mueen@unm.edu, eamonn@cs.ucr.edu

Abstract—The all-pairs-similarity-search (or similarity join) problem has been extensively studied for text and a handful of other domains. However, surprisingly little progress has been made on similarity joins for time series subsequences. The lack of progress probably stems from the daunting nature of the problem. For even modest sized datasets the obvious nested-loop algorithm can take months, and the typical speed-up techniques in this domain (i.e., indexing, lower bounding, triangularization, pruning and early abandonment of leaf probes) are too costly, or require significant domain knowledge. In this work we introduce a novel scalable algorithm for time series subsequences all-pairs-similarity search. For exceptionally large datasets, the algorithm can be trivially cast as an anytime algorithm and produce high-quality approximate solutions in reasonable time. The exact similarity join algorithm computes the answer to the time series motif and time series discord problems in a side-effect, and our algorithm incidentally provides the fastest known algorithm for both these extensively-studied problems. We demonstrate the utility of our ideas for many time series data mining problems, including motif discovery, anomaly discovery, shapelet discovery, semantic segmentation, density estimation, and contrast set mining.

Keywords—Time Series, Similarity Joins, Motif Discovery

1. INTRODUCTION

The all-pairs-similarity-search (also known as similarity join) problem comes in several variants. The basic task is this: Given a collection of data objects, retrieve the nearest neighbor for each object. In the text domain, the algorithm has applications in a host of problems, including community discovery, duplicate detection, collaborative filtering, clustering, and query refinement [1]. While virtually all text processing algorithms have analogues in time series data mining, there has been surprisingly little progress on Time Series subsequences All-Pairs-Similarity-Search (TSAPSS).

We believe that this lack of progress stems not from a lack of interest in this useful primitive, but from the daunting nature of the problem. Consider the following example that reflects the needs of an industrial collaborator. A boiler at a chemical refinery reports pressure once a minute. After a year, we have a time series of length 525,600. A plant manager may wish to do a similarity self-join on this data with week-long subsequences (10,080) to discover operating regimes (constant vs. winter or light decline vs. heavy decline, etc). The obvious nested loop algorithm requires 132,899,092,500 Euclidean distance computations. If we assume each one takes 6,000 seconds, then the join will take 15.8 days. The core contribution of this work is to show that we can reduce this time to 6.3 hours, using an off-the-shelf desktop computer. Moreover, we show that this join can be computed and/or updated incrementally. This we would maintain this join essentially forever on a standard

desktop, even if the data arrival frequency was much faster than once a minute.

Our algorithm uses an ultra-fast similarity search algorithm under conventional Euclidean distance as a subroutine, exploiting the overlap between subsequences using the classic Fast Fourier Transform (FFT) algorithm.

Our method has the following advantages/features:

- It is exact, providing no false positives or false dismissals.
- It is simple and parameter-free. In contrast, the more general metric space APSS algorithms require building and tuning spatial access methods and/or hash functions.
- Our algorithm requires an inconsequential space overhead, just $O(N)$ with a small constant factor.
- While our exact algorithm is extremely scalable, for extremely large datasets we can compute the results in an anytime fashion, allowing ultra-fast approximate solutions.
- Having computed the similarity join for a dataset, we can incrementally update it very efficiently, in many domains this means we can effectively maintain exact joins on streaming data forever.
- Our method provides full joins, eliminating the need to specify a similarity threshold, which as we will show, is a near impossible task in this domain.
- Our algorithm is embarrassingly parallelizable, both on multicore processors and in distributed systems.

Time Series Shapelets: A New Primitive for Data Mining

Lexiang Ye
Dept of Computer Science & Engineering
University of California, Riverside, CA 92521
ylexiang@cs.ucr.edu

Eamonn Keogh
Dept of Computer Science & Engineering
University of California, Riverside, CA 92521
eamonn@cs.ucr.edu

ABSTRACT

Classification of time series has been attracting great interest over the past decade. Recent empirical evidence has strongly suggested that the simple nearest neighbor algorithm is very difficult to beat for some time series problems. While this may be considered good news, given the simplicity of implementing the nearest neighbor algorithm, there are some negative consequences of this. First, the nearest neighbor algorithm requires storing and searching the entire dataset, resulting in a time and space complexity that limits its applicability, especially in resource-limited settings. Instead, beyond more classification accuracy, we often wish to gain some insight into the data.

In this work we introduce a new time series primitive, time series shapelets, which addresses these limitations. Informally, shapelets are time series subsequences which are in some sense maximally representative of a class. As we shall show with extensive empirical evaluation on diverse domains, algorithms based on the time series shapelet primitive can be interpreted, more accurate and significantly faster than state-of-the-art classifiers.

Categories and Subject Descriptors

H.2.8 [Database Management] Database Applications - Data Mining

General Terms

Algorithms, Experimentation

1. INTRODUCTION

While the last decade has seen a huge interest in time series classification, to date the most accurate and robust method is the simple nearest neighbor algorithm [9][12][1]. While the nearest neighbor algorithm has the advantages of simplicity and not requiring extensive parameter tuning, it does have several important disadvantages. Chief among these are its space and time requirements, and the fact that it does not tell us anything about why a particular object was assigned to a particular class.

In this work we present a novel time series data mining primitive called time series shapelets. Informally, shapelets are time series subsequences which are in some sense maximally representative of a class. While it is believed that shapelets have many uses in data mining, one obvious application of them is to navigate the two weaknesses of the nearest neighbor algorithm noted above.

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. To copy otherwise, to republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee.
KDD '09, June 20-24, 2009, Paris, France.
Copyright 2009 ACM 978-1-60558-495-9/09/08...\$5.00.

Because we are defining and solving a new problem, we will take some time to consider a detailed motivating example. Figure 1 shows some examples of leaves from two classes, *Urtica dioica* (stinging nettle) and *Ferula arvensis*. These two plants are commonly confused, hence the colloquial name "false nettle" for *Ferula arvensis*.



Figure 1. Sample of leaves from two species. Note that some of leaves have the same leaf damage.

Suppose we wish to build a classifier to distinguish these two plants, what feature should we use? Note that the main variability of color and size within each class completely obscures the inter-variability between classes, our best hope is based on the shapes of the leaves. However, as we can see in Figure 1, the differences in the global shape are very subtle. Furthermore, it is very common for leaves to have distortions or "occlusions" due to insect damage, and these are likely to confound any global measures of shape. Instead, we attempt the following. We first convert each leaf into a one-dimensional representation as shown in Figure 2.

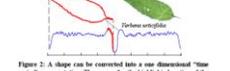


Figure 2. A shape can be converted into a one-dimensional "time series" representation. The reason for the highlighted section of the time series will be made apparent shortly.

Such representations have been successfully used for the classification, clustering and outlier detection of shapes in recent years [8]. However, here we find that using a nearest neighbor classifier with either the (optimal nearest neighbor) Euclidean distance or Dynamic Time Warping (DTW) distance does not significantly outperform random guessing. The reason for the poor performance of these otherwise very competitive classifiers seems to be due to the fact that the data is somewhat noisy (i.e. insect holes, and different stem lengths), and this noise is enough to swamp the subtle differences in the shapes.