Advanced Data Mining

Piotr Lipiński

- □ How to cluster time series?
 - see the jupyter python notebook with some examples



- □ How to cluster time series?
 - see the jupyter python notebook with some examples



• **APPROACH 1:** (for regular time series)

- □ daily (or weekly, monthly, annual) average profiles may define patterns
- each day (or week, month, year) time series can be matched to one of these patterns
- the time series matched to the same pattern define the cluster
- how to match time series to the patterns?

- □ How to cluster time series?
 - see the jupyter python notebook with some examples



APPROACH 2:

□ consider the time series as a vector of numbers, apply one of classic clustering algorithms, e.g. k-means, DBScan, etc.

Difficulties:

- time series may be of various length
- similar, but shifted or rescaled time series will lead to dissimilar vectors

APPROACH 3:

stay with original time series, but introduce a distance measure over them

□ How to measure similarities between time series?



- □ How to measure similarities between time series?
 - Euclidean distance



- □ How to measure similarities between time series?
 - Euclidean distance





7

- □ How to measure similarities between time series?
 - Euclidean distance / Manhattan distance / cosine distance / etc.
 - Euclidean distance with feature-based representation



8

- □ How to measure similarities between time series?
 - Euclidean distance / Manhattan distance / cosine distance / etc.
 - Euclidean distance with feature-based representation
 - Dynamic Time Warping (DTW)
 - DTW tries to synchronize time series before comparing



Dynamic Time Warping (DTW)

- Dynamic Time Warping (DTW)
 - \Box consider s[1..n] and t[1..m]
 - for each i = 1..n, s[i] must be matched with one or more elements of t
 - for each k = 1..m, t[k] must be matched with one or more elements of s
 - s[1] must be matched with t[1]
 (but may be also matched with other elements of t, t[1] similarly)
 - s[n] must be matched with t[m]
 (but may be also matched with other elements of t, t[m] similarly)
 - the mapping of elements of s to elements of t must be monotonically increasing,
 - for each i, j = 1..n, if i < j, then for each k, l = 1..m if s[i] matches t[k] and s[j] matches t[l], then k <= l</p>



10

```
Dynamic Time Warping (DTW)
```

```
DTW-Distance(s: array[1..n], t: array[1..m])
```

```
DTW := array[0..n, 0..m]

for i := 0 to n

for j := 0 to m

DTW[i, j] := infinity
```

```
DTW[0, 0] := 0

for i := 1 to n

for j := 1 to m

DTW[i, j] := d(s[i], t[j]) + min(

DTW[i-1, j],

DTW[i, j-1],

DTW[i-1, j-1])
```

```
return DTW[n, m]
```

□ How to cluster time series?

APPROACH 3:

stay with original time series, but introduce a distance measure over them

Difficulties:

- **DTW:** may be inefficient for larger time series
 - O(nm)
 - possible improvements, e.g. PrunedDTW, SparseDTW, FastDTW, MultiscaleDTW, etc.
- **K-means:** seemingly easy, but
 - how to define the center of the cluster for DTW distance?
 - minimizes the quantisation error, related to Euclidean distance, not to DTW distance

REMINDER: k-means

- □ Let $D = {x_1, x_2, ..., x_N}$ be the dataset of N data vectors $x_1, x_2, ..., x_N$. Let K be the number of cluster to define.
- Each cluster C_k is defined by a point \mathbf{r}_k , called the center of the cluster. Each data vector is assigned to the cluster of the closest center.
 - in the case of equal distances to a few centers, they may be considered in a predefined order or by random
- The goal is to discover a partition $C = \{C_1, C_2, ..., C_K\}$ of the set D of the size K (i.e. K parwise disjoint sets $C_1, C_2, ..., C_K$ such that $C_1 \cup C_2 \cup ... \cup C_K = D$) minimizing the criterium function

$$F(C) = \sum_{k=1}^{K} \sum_{\mathbf{x} \in C_k} \| \mathbf{x} - \mathbf{r}_k \|^2$$

□ k-means is one of the algorithms solving such a problem

REMINDER: k-means

□ Minimization of the criteria function

$$F(C) = \sum_{k=1}^{K} \sum_{\mathbf{x} \in C_k} \| \mathbf{x} - \mathbf{r}_k \|^2$$

may be performed iteratively in two steps:

- having vectors r_k, find the optimal assignment of data vector to the clusters obviously: each data vector should be assigned to the cluster of the closest center r_k
- having the assignment of data vector to the clusters, find vectors \mathbf{r}_k
 - it is less obvious mathematical approach based on derivatives
 - \square the solution is to set \mathbf{r}_k in the barycenter of the set of vectors forming the cluster

REMINDER: k-means

□ k-means

FOR k = 1, 2, ..., K r_k = a random point from D WHILE there are still changes in C_k FOR k = 1, 2, ..., K $C_k = \{x \in D : d(x, r_k) < d(x, r_l) \text{ for each } l = 1, 2, ..., K, l \neq k\}$ FOR k = 1, 2, ..., K r_k = barycenter of C_k

□ How to cluster time series?

APPROACH 3:

stay with original time series, but introduce a distance measure over them

Difficulties:

- **DTW:** may be inefficient for larger time series
 - O(nm)
 - possible improvements, e.g. PrunedDTW, SparseDTW, FastDTW, MultiscaleDTW, etc.
- **K-means:** seemingly easy, but
 - how to define the center of the cluster for DTW distance?
 - minimizes the quantisation error, related to Euclidean distance, not to DTW distance

APPROACH 4:

DTW Barycenter Averaging k-means (**DBA-k-means**)

DTW Barycenter Averaging k-means (DBA-k-means)



F. Petitjean, A. Ketterlin, P. Gancarski, "A global averaging method for dynamic time warping, with applications to clustering". Pattern Recognition, 3(44), 2011, pp.678-693. Piotr Lipiński, Advanced Data Mining