

Structuring Typical Evolutions using Temporal-Driven Constrained Clustering

8 November 2012

Marian-Andrei Rizoiu Julien Velcin Stéphane Lallich ERIC Laboratory Université Lumière Lyon 2 France

> UNIVERSITÉ LUMIÈRE LYON 2 UNIVERSITÉ DE LYON



ProblemProposed SolutionsExperimentsConclusionDataset:the values for a certain number of numerical
features (x^d) for multiple entities (φ) at different
moments of time (t)

Dataset:

Dataset:



Dataset:





Dataset:





Dataset:





Dataset:

the values for a certain number of numerical features (x^d) for multiple entities (φ) at different moments of time (t)





Dataset:

the values for a certain number of numerical features (x^d) for multiple entities (φ) at different moments of time (t)



Dataset:

the values for a certain number of numerical features (x^d) for multiple entities (φ) at different moments of time (t)



Goal: Detect typical evolution patterns of individuals in the dataset

Goal: Detect typical evolution patterns of individuals in the dataset



a) the phases through which the entity collection went over time

Goal: Detect typical evolution patterns of individuals in the dataset

φ₁,φ₃

a) the phases through which the entity collection went over time

h

$$\mu_2$$

 μ_2
 μ_4
 μ_5
 μ_4
 μ_4
 μ_5
 μ_4
 μ_4
 μ_5
 μ_6
 μ_5
 μ_6
 μ_7
 μ_8
 μ_8

φ,

t

b) the trajectory of entities through the different phases $x_{1} = (\phi_{1}, t_{1}, x_{1}^{d}) + x_{2} = (\phi_{3}, t_{1}, x_{2}^{d})$ $\phi_{1} = \phi_{1} = \mu_{3} + \phi_{1} = \mu_{5}$ $\phi_{2} = \phi_{2} = \mu_{2} + \phi_{3} = \mu_{4} + \phi_{3} = \mu_{4} + \phi_{3} = \mu_{4}$

Summary:

1. Problem 1.1 Data 1.2 Goal

2. Proposed solutions:

- 2.1 A clustering solution
- 2.2 Temporal-Aware Dissimilarity Measure
- 2.3 Contiguity Penalty Measure
- 2.4 TDCK-Means algorithm
- 2.5 Evaluation measures
- 3. Experiments
 - 3.1 Qualitative evaluation
 - 3.2 Quantitative evaluation

4. Conclusion and perspectives

The resulted partition must ensure:

- → the descriptive coherence of clusters;
 → the temporal coherence of clusters;
- → continuous segmentation of observations belonging to an entity.

The resulted partition must ensure:

→ the descriptive coherence of clusters;
→ the temporal coherence of clusters;

→ continuous segmentation of observations belonging to an entity.

Temporal-aware dissimilarity measure

Contiguity penalty measure

The resulted partition must ensure:

→ the descriptive coherence of clusters;
→ the temporal coherence of clusters;

 Contiguity penalty measure

Temporal-aware

dissimilarity measure

K-Means like algorithm. Objective function to minimize:

$$J = \sum_{\mu_j \in M} \sum_{x_i \in C_j} \left(\|x_i - \mu_j\|_{TE} + \sum_{(x_k \notin C_j) \land (x_k^{\varphi} = x_i^{\varphi})} w(x_i, x_k) \right)$$

Contiguity penalty

measure

The resulted partition must ensure:

→ the descriptive coherence of clusters;
→ the temporal coherence of clusters;
Temporal-aware 1 dissimilarity measure

→ continuous segmentation of observations belonging to an entity.

K-Means like algorithm. Objective function to minimize:

$$J = \sum_{\mu_{j} \in M} \sum_{x_{i} \in C_{j}} \left(\|x_{i} - \mu_{j}\|_{TE} + \sum_{(x_{k} \notin C_{j}) \land (x_{k}^{\varphi} = x_{i}^{\varphi})} w(x_{i}, x_{k}) \right)$$

►

Euclidean distance

distance in the description space

Euclidean distance

distance in the description space

Temporal-aware dissimilarity measure distance in both description space and temporal space

Euclidean distance

distance in the description space

Temporal-aware dissimilarity measure distance in both description space and temporal space

$$||x_{i} - x_{j}||_{TE} = 1 - \left(1 - \frac{||x_{i}^{d} - x_{j}^{d}||^{2}}{\Delta x_{max}}\right) \left(1 - \frac{|x_{i}^{t} - x_{j}^{t}|^{2}}{\Delta t_{max}}\right)$$

Euclidean distance

distance in the description space

Temporal-aware dissimilarity measure distance in both description space and temporal space

$$||x_{i} - x_{j}||_{TE} = 1 - \left(1 - \frac{||x_{i}^{d} - x_{j}^{d}||^{2}}{\Delta x_{max}}\right) \left(1 - \frac{|x_{i}^{t} - x_{j}^{t}|^{2}}{\Delta t_{max}}\right)$$

Properties:

$$\Rightarrow ||x_i - x_j||_{TE} \in [0,1], \forall x_i, x_j \in X$$

$$\Rightarrow ||x_i - x_j||_{TE} = 0 \Leftrightarrow x_i^d = x_j^d \wedge x_i^t = x_j^t$$

$$\Rightarrow ||x_i - x_j||_{TE} = 1 \Leftrightarrow ||x_i^d - x_j^d|| = \Delta x_{max} \lor |x_i^t - x_j^t| = \Delta t_{max}$$

Semi-Supervised _____ clustering [Wagstaff & Cardie '00] pair-wise constraints apply penalty when constraints are broken

Semi-Supervised _____ clustering [Wagstaff & Cardie '00] pair-wise constraints apply penalty when constraints are broken

Segmentation contiguity soft MUST-LINK pair-wise constraints time-dependent Contiguity Penalty Function

Semi-Supervised _____ clustering [Wagstaff & Cardie '00] pair-wise constraints apply penalty when constraints are broken

Segmentation contiguity

soft MUST-LINK pair-wise constraints time-dependent
Contiguity
Penalty Function

Contiguity Penalty Function:

$$w(x_i, x_j) = \beta * e^{\frac{-1}{2} \left(\frac{|x_i^t - x_j^t|}{\delta}\right)^2}$$

for $x_i^{\varphi} = x_j^{\varphi}$

Semi-Supervised _____ clustering [Wagstaff & Cardie '00] pair-wise constraints apply penalty when constraints are broken

Segmentation contiguity soft MUST-LINK pair-wise constraints time-dependent
Contiguity
Penalty Function

Contiguity Penalty Function:

$$w(x_i, x_j) = \beta * e^{\frac{-1}{2} \left(\frac{|x_i^t - x_j^t|}{\delta}\right)^2}$$

for $x_i^{\varphi} = x_j^{\varphi}$



The TDCK-MeansInspired from K-Means. Iterativelyalgorithm:recomputes centroids and assignments of
observations to clusters.

Uses the Temporal-Aware Dissimilarity Function and the Contiguity Penalty Function.

Centroids: (μ_i^t, μ_i^d)

The TDCK-MeansInspired from K-Means. Iterativelyalgorithm:recomputes centroids and assignments of
observations to clusters.

Uses the Temporal-Aware Dissimilarity Function and the Contiguity Penalty Function.

Centroids: (μ_j^t, μ_j^d)



Partition evaluation measures

→ descriptive coherence of clusters;
→ temporal coherence of clusters;

→ continuous segmentation of observations belonging to an entity.

Partition evaluation measures

→ descriptive coherence of clusters; }
→ temporal coherence of clusters; }
✓ variance
✓ Tvar

→ continuous segmentation of observations belonging to an entity.

Shannon Entropy $A \rightarrow B \rightarrow A \rightarrow B$???

Partition evaluation measures

→ descriptive coherence of clusters; } variance
→ temporal coherence of clusters; } Variance
✓ Tvar

→ continuous segmentation of observations belonging to an entity.

Shannon Entropy $A \rightarrow B \rightarrow A \rightarrow B$???

Proposal: Correct the Shannon entropy to penalize changes

Summary:

I. Problem 1.1 Data 1.2 Goal

2. Proposed solutions:
2.1 A clustering solution
2.2 Temporal-Aware Dissimilarity Measure
2.3 Contiguity Penalty Measure
2.4 TDCK-Means algorithm
2.5 Evaluation measures

3. Experiments

- 3.1 Qualitative evaluation
- 3.2 Quantitative evaluation

4. Conclusion and perspectives

Compared Political23 countries, 60 years, 207 political,Dataset Idemographic, social and economic variables.

Compared Political23 countries, 60 years, 207 political,Dataset Idemographic, social and economic variables.

Execution TDCK-Means (8 clusters, $\beta = 0.003$ and $\delta = 3$)



M-A. Rizoiu, J. Velcin and S. Lallich Structuring Typical Evolutions using Temporal-Driven Constrained Clustering 11

Compared Political23 countries, 60 years, 207 political,Dataset Idemographic, social and economic variables.

Execution TDCK-Means (8 clusters, $\beta = 0.003$ and $\delta = 3$)



Compared Political23 countries, 60 years, 207 political,Dataset Idemographic, social and economic variables.

Execution TDCK-Means (8 clusters, $\beta = 0.003$ and $\delta = 3$)



Compared Political23 countries, 60 years, 207 political,Dataset Idemographic, social and economic variables.

Execution TDCK-Means (8 clusters, $\beta = 0.003$ and $\delta = 3$)



Quantitative evaluation

5 algorithms:

→ K-Means [MacQueen '67];
→ tcK-Means [Lin and Hauptmann '10]

 → Temporal-Driven K-Means; (uses Temporal-Aware Measure)
 → Constrained K-Means; (uses Contiguity Penalty Function)

→ TDCK-Means;

(combines the two above)

3 measures:

- → MDvar
- → Tvar
- → ShaP

Problem Proposed Solutions

Experiments

Conclusion



Summary:

Problem 1.1 Data 1.2 Goal

2. Proposed solutions:

2.1 A clustering solution

2.2 Temporal-Aware Dissimilarity Measure

2.3 Contiguity Penalty Measure

2.4 TDCK-Means algorithm

2.5 Evaluation measures

3. Experiments

3.1 Qualitative evaluation3.2 Quantitative evaluation

4. Conclusion and perspectives

Conclusion:

- → Studied the detection of typical evolutions starting from a collection of observations corresponding to entities;
- → Proposed a new **Temporal-Aware Measure**;
- → Proposed a new **Contiguity Penalty Function**;
- → Proposed a new algorithm for detecting evolutions: TDCK-Means;
- → Other applications: political careers, life trajectories *etc*.

Perspectives:

- → Generating the evolution graph;
- → Automatic description of generated evolution phases (clusters);
- → Flexible configuring the ration between the descriptive component and the temporal component in the dissimilarity measure.

Thank you!

Questions?

Impact of parameters β and δ

