

Dynamic Multi-Network Mining of Tensor Time Series

Kohei Obata, Koki Kawabata, Yasuko Matsubara, Yasushi Sakurai

SANKEN, Osaka University



Kohei Obata et al.

- Introduction
- Preliminaries
- Methodology
- Experiments
- Conclusion



Ubiquity of Tensor Time Series (TTS)

- IoT has facilitated the collection of TTS.
- TTS consists of multiple modes including Time.
 - e.g., Online activity data {Query, Country, Timestamp}, Air pollutant data {Pollutant, Site, Timestamp}, Automobile data {Sensor, Lap, Driver, Meter}.

Example of Online activity data:

The search amount of 6 Queries related to Covid19 taken from 10 Countries for 10 years (2013 ~ 2023).







Subsequence Clustering

- Important task for data mining.
 - uncover interesting patterns, useful for downstream tasks
- Interpretability of resulting clusters is also essencial.

Example of Online activity data:

The search amount of **6 Queries** related to Covid19 for **10 years** (2013 ~ 2023).





Subsequence Clustering

- Important task for data mining.
 - uncover interesting patterns, useful for downstream tasks
- Interpretability of resulting clusters is also essencial.
- **Dependency network** gives a clear explanation.
 - why a particular cluster distinguishes itself from another?
 - what happened during a period belonging to the cluster?





Existing Works of Subsequence Clustering

- Univariate time series clustering (e.g., K-means, DTW).
 - Focus on matching raw values.
 - Do not consider relationship between variables.
- Multivariate time series clustering (e.g., TICC¹¹, TAGM¹²).
 - <u>Characterize cluster with network</u> based on graphical lasso.
 - Discover clusters that other traditional methods cannot find.
- TTS clustering is more challenging.
 - It has an intricate dependency and huge data size due to the modes.
 - Applying MTS methods to TTS hinders interpretability and requires time.





[1] David Hallac, Sagar Vare, Stephen P. Boyd, and Jure Leskovec. 2017. Toeplitz Inverse Covariance-Based Clustering of Multivariate Time Series Data. In KDD. 215–223.

[2] Veronica Tozzo, Federico Ciech, Davide Garbarino, and Alessandro Verri. 2021. Statistical Models Coupling Allows for Complex Local Multivariate Time Series 6 Analysis.InKDD.1593–1603. Kohei Obata et al.

- Introduction
- Preliminaries
- Methodology
- Experiments
- Conclusion



Network & Graphical Lasso

Log-Likelihood

- Inverse covariance matrix (i.e. network) of Gaussian distribution encodes dependency between features.
 If there is an edge, then they are dependent given the rest of the variables.
- Graphical lasso infers a sparse network, which helps us understand the important relationship.





Problem Formulation

Given: (N + 1)th-order TTS $X \in \mathbb{R}^{D_1 \times \cdots \times D_N \times T}$ Find: \mathcal{M} that minimize the cost function $Cost_T(X; \mathcal{M})$

Cluster parameter $\mathcal{M} = \{\Theta, \mathcal{F}\}$ $- \left\{\begin{array}{c} \text{Cluster assignment set } \mathcal{F} = \{f_k\}_{k=1}^K \\ \text{e.g., } f_2 = \{cp_1, \dots, cp_2 - 1, cp_5, \dots, cp_6 - 1\} \\ \text{Model parameter set } \Theta = \{\theta_k\}_{k=1}^K \end{array}\right\}$





- Introduction
- Preliminaries
- Methodology
- Experiments
- Conclusion



DMM: Dynamic Multi-Network Mining

DMM achieves interpretable TTS clustering.

(1) Characterize a cluster with multiple networks.

- Extend graphical lasso to TTS.
- Each mode has a dependency network.

(2) Define the cost function (goodness of clustering).

- Based on the Minimum Description Length (MDL).
- Can determine any hyperparameters (e.g., λ , K).

(3) Propose the algorithm that minimize the cost.

- Based on the bottom-up algorithm.
- Scales linearly w.r.t. data size.
- It can be applied to long-range,

high-dimensional TTS.





 $=Cost_{A}(\mathcal{I})$



Coding length cost Model coding cost



 $Cost_T(X; M)$

Total description cost

iteration: 6, # of segments: 3



(1) Multimode Graphical Lasso

- We model TTS with *N* networks.
- The model is **interpretable** as each mode has a sparse network.





(2) Model Description Cost

- We define the criterion for the goodness of clustering based on Minimum Description Length (MDL).
- It can determine any hyperparameters (e.g., λ , K) by minimizing $Cost_T(X; M)$.









- Introduction
- Preliminaries
- Methodology
- Experiments
- Conclusion



DMM is Interpretable (Covid data)

• DMM helps us understand the real-world data.

The search amount of 6 Queries related to Covid19 taken from 10 Countries for 10 years (2013 ~ 2023).





DMM is Interpretable (Air data)

• DMM is <u>more interpretable</u> than existing MTS methods (TAGM and TICC).

Concentration of 6 Pollutants taken from 10 Locations in China for 4 years (2013 ~ 2017).





DMM is Accurate (Synthetic data)

• DMM accurately discovers cluster with different networks even at MTS.

Table 1. Macro-F1 score of Synthetic data.

Data	DMM	TAGM	TAGM [†]	TICC	TICC [†]
2 nd -order TTS A (MTS) (i) B D	0.955 0.926 0.956 0.960	0.915 <u>0.897</u> 0.770 0.907	$\begin{array}{c} 0.915 \\ 0.756 \\ \underline{0.811} \\ 0.912 \end{array}$	0.997 0.884 0.725 0.857	0.997 0.825 0.756 0.952
3 rd -order TTS A (ii) B C D	0.961 0.962 0.941 0.980	0.514 0.462 0.359 0.438	0.514 0.431 0.396 0.432	$ \begin{array}{r} 0.932 \\ \hline 0.844 \\ \hline 0.704 \\ \hline 0.838 \\ \end{array} $	0.923 0.770 0.594 0.741

All the hyperparameters, including # of clusters, are determined by the cost function.

• DMM is not affected by the number of variables.





DMM is Scalable

• DMM scales linearly w.r.t. the input data size.



• DMM is up to x300 faster than the existing methods.

Table 3: The data size and attributes for each dataset.					
ID	Dataset	Size	Description		
#1	E-commerce	(11, 10, 1796)	3 rd -order TTS		
#2	VoD	(8, 10, 1796)	(query, state, day)		
#3	Sweets	(9, 10, 1796)			
#4	Covid	(6, 10, 3652)	(
#5	GAFAM	(5, 10, 1796)	(query, country, day)		
#6	Air	(6, 12, 1461)	(pollutant, site, day)		
#7	Car-A	(6, 10, 4, 3241)	4th-order TTS (sensor, lap, driver, meter)		
#8	Car-H	(6, 10, 4, 4000)			





- Introduction
- Preliminaries
- Methodology
- Experiments
- Conclusion



Conclusion

- DMM is Interpretable, Accurate, and Scalable.
- DMM is a useful tool for **TTS subsequence clustering** that enables multifaceted analysis and understanding of TTS.
- Possible applications could be



Traffic Congestion



Factory



Healthcare



Thank you for listening!

Code & Data are available [1].





[1] https://github.com/KoheiObata/DMM

Kohei Obata et al.