

Multivariate Time Series Multi-Cocustering. Application to Advanced Driving Assistance System Validation

Etienne Goffinet^{1,2}, Mustapha Lebbah¹, Hanane Azzag¹,
Loïc Giraldi³ and Anthony Coutant¹

1- Sorbonne Paris-Nord University - LIPN-UMR 7030
99 Avenue Jean Baptiste Clément, Villetaneuse - France
2- Groupe Renault SAS Avenue du Golf, Guyancourt - France
3- CEA, DES, IRESNE, DEC, 13108 Saint-Paul-Lez-Durance, France

Abstract. Driver assistance systems development remains a technical challenge for car manufacturers. Validating these systems requires to assess the assistance systems performances in a considerable number of driving contexts. *Groupe Renault* uses massive simulation for this task, which allows reproducing the complexity of physical driving conditions precisely and produces large volumes of multivariate time series. We present the operational constraints and scientific challenges related to these datasets and our proposal of an adapted model-based multiple cocustering approach, which creates several independent partitions by grouping redundant variables. This method natively performs model selection, missing values inference, noisy samples handling, confidence interval production, while keeping a sparse parameter numbers. The proposed model is evaluated on a synthetic dataset, and applied to a driver assistance system validation use-case.

1 Introduction

Advanced driver-assistance systems (ADAS) are automated embedded systems that bring support to driving tasks (e.g., emergency braking, adaptive cruise control, automated lighting, navigation, ...). Given the high number of different vehicles, driving conditions, traffic laws, and given the expected reliability, it is today impossible to validate ADAS exhaustively only with physical tests "on tracks". *Groupe Renault* has made the technical choice to invest in massive driving simulation technology to circumvent this issue. For a given use case, a simulation protocol produces a multivariate time series dataset with varying dimensions: from 100 to 100 000 observations, and from 10 to 100 variables. The behavior of the simulated car cannot always be predicted. In consequences, it is impossible to use supervised methods. Our objective is to provide a synthetic view of the simulated dataset to the ADAS developer, that simultaneously highlights the driving patterns and discriminates groups of dependent sensors.

The model-based clustering methods [1] are a family of unsupervised probabilistic approaches that natively provide probabilistic estimations of cluster memberships, confidence intervals for probabilistic outlier detection, can infer missing values and model sample noise. In addition, they can handle data sets

of varying sizes in opposition to existing Deep Learning clustering approaches [2, 3] that require massive amounts of data. An instance of such model-based approach has been used recently in the ADAS validation context [4], but restricted to univariate time series analysis. In a multivariate context, the possible relations

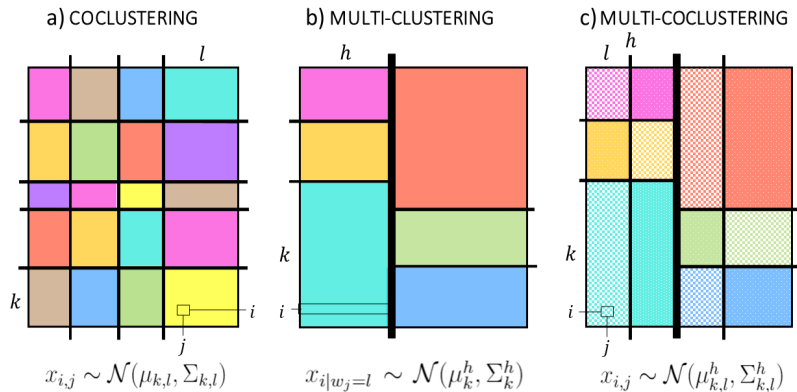


Fig. 1: Coclustering, Multi-Clustering and Multi-Coclustering, with k, l, h the row-cluster, correlated cluster and redundant cluster indices (respectively). Color and pattern designate block membership.

between variables introduces additional difficulties. In most of real-life use cases, it is likely that each variable is associated to a different partition of the observation space (called *row-partition*). In the following, we assume the presence of a variable partition (called *redundant column-partition*), such that variables in a redundant column-cluster share the same row-partition. In our use cases, this assumption seems rational, as we know that many variables are related (e.g., position and acceleration, ADAS activation and braking). The *Multi-Clustering* [5] designates a set of methods that infers this multiple-partitions representation of a dataset, as illustrated in Fig. 1 - b). In a *parametric* model-based context [6, 7], using these approaches requires to know the partitions sizes a priori (rarely true in practice), or to perform a model selection step. However, the combinatorial nature of Multi-Clustering makes the model selection a complex task (there are, for example, 184755 different multi-clustering models with at most 10 clusters of variables and 10 clusters of observations). Heuristic approaches could be considered at the cost of assumptions on model structure (e.g. with greedy strategies). The *non-parametric* approach circumvents this issue by natively integrating a model selection. A non-parametric Bayesian Multi-Clustering model has been developed by [8] for continuous datasets. In this model, the rows belonging to a block follow an independent multivariate distribution (c.f. Fig. 1-b)). This Multi-Clustering method groups variables according to their row partition but cannot regroup variables with similar distributions. Moreover, estimating the block distributions parameters becomes problematic when the column cluster sizes are high (e.g., high dimensional covariance matrices). Each redundant col-

umn cluster is modeled with a coclustering structure to deal with this problem.

The model-based coclustering [9] (Figure 1 - a)) infers one row-partition and one column partition. This model assumes that all variables share the same row partition and that there exists groups of variables with common distribution. This assumption seems also natural, because our simulated variables are physically correlated (e.g., car speed and wheel rotation speed).

In this paper, we propose a new Non-Parametric Multi Coclustering (MCC) (c.f. Fig. 1-c) that combines the two previous approaches and produces the best of both worlds: the redundant partition layer extracts the true sparse number of row partition while the coclustering layers reduce the parameter set dimension and regroup correlated variables. To the best of our knowledge, there is only one comparable work by [10], which does not apply to multivariate time series.

2 Multi Coclustering of multivariate time series

2.1 MCC model definition

As a preprocessing step, the time series are first represented in a common basis, with a three steps transformations: 1) individual log-periodogram representation; 2) interpolation of the log-periodogram into a common frequency basis; 3) PCA on the frequency coefficients. This preprocessing outputs a dataset $X = (x_{i,j,s})_{n \times p \times d}$, with n observations described by p variables, and d the projection space dimension. We denote \mathbf{v} the redundant partition indicator, $(\mathbf{w}_h)_h$ the correlated partition indicator, and Z the row-partitions indicator matrix. The model is formally defined by:

$$\begin{aligned}
 x_{i,j} &| \{v_j = h, w_j^h = l, z_i^h = k, \theta_{k,l}^h\} \sim \mathcal{N}(\theta_{k,l}^h), \\
 \theta_{k,l}^h &\sim G_0, v_j \sim \text{Mult}(\eta), w_j^h \sim \text{Mult}(\rho_h), z_i^h \sim \text{Mult}(\pi_h), \\
 \eta_j(\mathbf{r}) &= r_j \prod_{j'=1}^{j-1} (1 - r_{j'}), \quad r_j \stackrel{\text{i.i.d.}}{\sim} \text{Beta}(1, \gamma), \gamma \sim \text{Gamma}(a_r, b_r) \\
 \rho_j^h(\mathbf{s}^h) &= s_j^h \prod_{j'=1}^{j-1} (1 - s_{j'}^h), \quad s_j^h \stackrel{\text{i.i.d.}}{\sim} \text{Beta}(1, \beta_h), \beta_h \sim \text{Gamma}(a_c, b_c) \\
 \pi_j^h(\mathbf{t}^h) &= t_j^h \prod_{j'=1}^{j-1} (1 - t_{j'}^h), \quad t_j^h \stackrel{\text{i.i.d.}}{\sim} \text{Beta}(1, \alpha_h), \alpha_h \sim \text{Gamma}(a_l, b_l),
 \end{aligned}$$

where each memberships proportions vector η , $(\rho^h)_h$ and $(\pi^h)_h$ follows a *stick-breaking* construction scheme [11] with Gamma distribution hyperprior. The block distributions are multivariate normal and G_0 is the associated conjugate Normal-inverse-Wishart prior. The complete hyperparameter set is noted χ .

2.2 Inference

We propose an iterative two-steps inference: 1) update of \mathbf{v} given Z ; 2) update of Z and \mathbf{w} given \mathbf{v} . In the first step, each membership v_j is sampled sequentially

according to:

$$p(v_j | \mathbf{v}_{-j}, Z, \mathbf{x}_{\cdot,j}, \chi) \propto \begin{cases} \frac{p_h}{p-1+\gamma} p(\mathbf{x}_{\cdot,j} | \mathbf{z}^h, \chi), & \text{existing cluster } h, \\ \frac{\gamma}{n-1+\gamma} p(\mathbf{x}_{\cdot,j} | \chi), & \text{new cluster,} \end{cases} \quad (1)$$

$$(2)$$

with $\mathbf{v}_{-j} = \{v_i : i \neq j\}$ and $p(\mathbf{x}_{\cdot,j} | \mathbf{z}^h)$ the distribution of $x_{\cdot,j}$ in \mathbf{z}^h . In Eq. (1), the predictive distributions $p(\mathbf{x}_{\cdot,j} | \mathbf{z}^h, \chi)$ reduces to a product of multivariate t-student densities [12] thanks to an appropriate choice of conjugate prior G_0 . In Eq. (2), $p(\mathbf{x}_{\cdot,j} | \chi)$ is estimated with an univariate Dirichlet Process Mixture (once per variable, prior to the inference of the MCC), which also produces the optimal $\hat{\mathbf{z}}^j$ row partition associated with the new cluster. After \mathbf{v} 's update, the concentration parameter γ is updated based on [13]. In the second step, an NPLBM [14] is inferred on each sub-dataset $(X^h)_h$ defined by the partition \mathbf{v} with a Gibbs algorithm that simulates $p(\mathbf{z}^h, \mathbf{w}^h | \mathbf{v}, X^h, \chi)$ (c.f. [14]).

3 Applications

3.1 Experiments on synthetic data

We compare the MCC's performances to other methods adapted to the use case constraints. The experiments dataset is composed of observations generated from a ground truth row-partition, that we seek to estimate. This ground truth information is "corrupted" by adding a proportion of other variables, uninformatives (generated from a one-cluster partition) or misleading (generated from another row-partition than the ground truth), noted p_u and p_m . The base-lines methods are non-parametric model-based approaches to Multi-Clustering: Dirichlet Process Mixture model (DPM) that infers a row-partition with only one column-cluster, non-parametric coclustering method (CC) [14] displayed in Fig. 1, decoupled Multi-Clustering based on a two-layers nested DPM (NDPM) and non-parametric Multi-clustering (MCD) equivalent to MCC without the co-clustering layer. The results, displayed in Table 1, shows the interest of using MCC in the presence of uninformatives or misleading variables.

3.2 Emergency lane Keeping assist use case

After validating MCC on a synthetic dataset, we apply it on a real ADAS validation dataset. In a straight lane scenario, the vehicle under test is drifting towards an oncoming car on the other lane. The ELK system is expected to detect the drifting, the oncoming car and put it back to its lane center with an emergency maneuver. The simulation generates a dataset of $n = 500$ described by $p = 150$ temporal variables, totalling 75000 time series. The objective is to find driving patterns, to discriminate relevant groups of sensors and to isolate the ADAS activation context. The Fig. 2-a) represents the row partitions overlaid on the dataset, with color representing redundant column clusters, and transparency indicating block partitions. Fig. 2-b) shows the car trajectories (surrounded

Table 1: Row-partition quality described by ARI, RI and row-cluster number \hat{K}

p_u	p_m	Scores	DPM	CC	NDPM	MCD	MCC
0	0	RI	0.84	1.0	0.57	1.0	1.0
0	100%		0.73	0.85	0.64	0.99	1.0
0	200%		0.79	0.85	0.76	0.88	1.0
100%	0		0.80	1.0	0.71	0.99	1.0
200%	0		0.51	1.0	0.60	0.98	1.0
0	0		ARI	0.67	1.0	0.27	1.0
0	100%	0.34		0.61	0.27	0.97	1.0
0	200%	0.47		0.61	0.39	0.70	1.0
100%	0	0.56		1.0	0.44	0.99	1.0
200%	0	0.13		1.0	0.27	0.96	1.0
0	0	\hat{K}		3	3	2	3
0	100%		7	6	3	3	3
0	200%		6	6	5	6	3
100%	0		3	3	2	3	3
200%	0		3	3	2	3	3

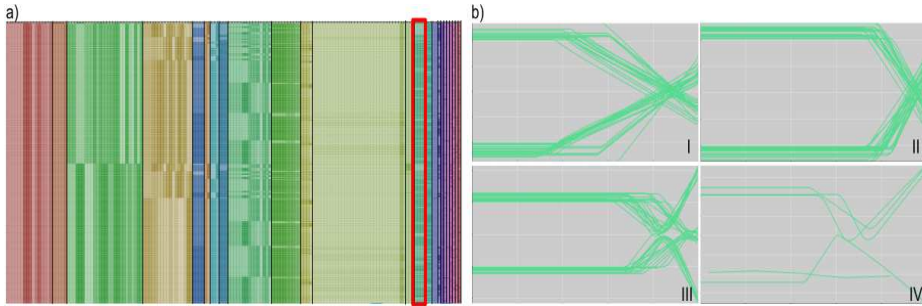


Fig. 2: a) MCC results. b) Car trajectories in the highlighted clusters

by a red rectangle in Fig. 2-a)) distributed in the row-clusters, where several scenarios are represented: the sub-graph I shows a late ADAS activation scenario, sub-graph II regroups cases when ELK did not activate. In sub-graph III, the systems correctly change the car trajectory and prevent the collision. The last graph IV shows several outlier simulations detected using the probability density of the MCC model. In addition to discriminating the non-informative variables (leftmost columns in red color), the Fig. 2-A) shows a hierarchical structure between several row partitions (leftmost "green", "ochre" and "blue" columns). These partitions are linked to the vehicle orientation (e.g., road angle or steering wheel angle), lane position and ADAS activation indicators. Based on these results, the next steps for the ADAS system developer consists in relating the partitions to the input simulation parameters and assessing the ADAS compliance with its specifications.

4 Conclusion

This paper describes a new Bayesian non-parametric based method designed for the exploration of multivariate time series datasets produced by driving simulations. This solution infers a multi-level dependency structure that highlights the relationships between sensors and discriminates driving simulations patterns. This method can be used in other domains, as long as the corresponding structure dependency assumptions hold. In the future we consider combining the multivariate time-series multi-coclustering with the classification of the simulation parameters in order to identify ADAS failure and success contexts.

References

- [1] Arthur P Dempster, Nan M Laird, and Donald B Rubin. Maximum likelihood from incomplete data via the em algorithm. *Journal of the Royal Statistical Society: Series B (Methodological)*, 39(1):1–22, 1977.
- [2] Naveen Sai Madiraju, Seid M Sadat, Dimitry Fisher, and Homa Karimabadi. Deep temporal clustering: Fully unsupervised learning of time-domain features. *arXiv preprint arXiv:1802.01059*, 2018.
- [3] Florent Forest, Alex Mourer, Mustapha Lebbah, Hanane Azzag, and Jérôme Lacaille. An invariance-guided stability criterion for time series clustering validation. In *International Conference on Pattern Recognition (ICPR)*, 2020.
- [4] Etienne Goffinet, Mustapha Lebbah, Hanane Azzag, and Loic Giraldi. Autonomous driving validation with model-based dictionary clustering. In *European Conference on Machine Learning, (ECML PKDD)*, pages 323–338. Springer International Publishing, 2021.
- [5] Juhua Hu and Jian Pei. Subspace multi-clustering: a review. *Knowledge and information systems*, 56(2):257–284, 2018.
- [6] Matthieu Marbac and Vincent Vandewalle. A tractable multi-partitions clustering. *Computational Statistics & Data Analysis*, 132:167–179, 2019.
- [7] Vincent Vandewalle. Multi-partitions subspace clustering. *Mathematics*, 8(4):597, 2020.
- [8] Yue Guan, Jennifer G Dy, Donglin Niu, and Zoubin Ghahramani. Variational inference for nonparametric multiple clustering. In *MultiClust Workshop, KDD-2010*, 2010.
- [9] Gérard Govaert and Mohamed Nadif. *Co-clustering: models, algorithms and applications*. John Wiley & Sons, 2013.
- [10] Tomoki Tokuda, Junichiro Yoshimoto, Yu Shimizu, Go Okada, Masahiro Takamura, Yasumasa Okamoto, Shigeto Yamawaki, and Kenji Doya. Multiple co-clustering based on nonparametric mixture models with heterogeneous marginal distributions. *PloS one*, 12(10):e0186566, 2017.
- [11] Jayaram Sethuraman. A constructive definition of dirichlet priors. *Statistica sinica*, pages 639–650, 1994.
- [12] Andrew Gelman, John B Carlin, Hal S Stern, David B Dunson, Aki Vehtari, and Donald B Rubin. *Bayesian data analysis*. CRC press, 2013.
- [13] Mike West. *Hyperparameter estimation in Dirichlet process mixture models*. Duke University ISDS Discussion Paper# 92-A03, 1992.
- [14] Edward Meeds and Sam Roweis. Nonparametric bayesian biclustering. Technical report, Citeseer, 2007.