Constrained Clustering Challenges and Frontiers

Dr. Germán González Almagro



Categorical Info



Relational Info



- Semi-Supervised Clustering
- **Constrained Clustering**
- Hybrid Models for CC
- New Paradigms in CC
 - **Challenges and Frontiers**

Knowledge Discovery in Databases (KDD)



KDD: a set of stages that enable the identification of valuable patterns and relationships within the data. **Our goal is to extract meaningful insights from the data.**

Classic Machine Learning Paradigms

Supervised Learning: the goal is to build a classifier or regressor that, trained with a set of examples and their corresponding output values, can **predict the value of unseen inputs**.



Classic Machine Learning Paradigms

Unsupervised Learning: only the set of examples is available, and no output value is provided. The goal here is to **discover some underlying structure in the data.**



Semi-Supervised Learning (SSL)

Semi-Supervised Learning: tries to combine the benefits of supervised and unsupervised learning by making use of both labeled and unlabeled data, or other kinds of expert knowledge.



Assumptions

Smoothness Assumption.

Two instances that are close in the input space should have the same label.

Low-density Assumption.

Decision boundaries should preferably pass through low-density spatial regions.





Assumptions

Manifold Assumption

Instances in a high-dimensional input space are usually gathered along lower-dimensional structures.



Cluster Assumption

Instances which belong to the same cluster also belong to the same class.



Dichotomy

Inductive Methods: aim to build a classifier capable of outputting a label for any instance in the input space. The predictions for unseen instances are independent. Inductive methods in SSL are categorized as semi-supervised classification.



Transductive Methods: their predictions are limited to the data used during the training phase. Transductive methods do not have separated training and testing phases. Transductive methods in SSL are categorized as semi-supervised clustering.



Semi-Supervised Learning Semi-Supervised Clustering **Constrained Clustering** Hybrid Models for CC New Paradigms in CC **Challenges and Frontiers**

Semi-Supervised Learning: tries to combine the benefits of supervised and unsupervised learning by making use of both labeled and unlabeled data, or **other kinds of expert knowledge.**



Semi-Supervised Clustering: in addition to the unlabeled dataset, background knowledge is given to perform clustering. When the background knowledge is given in the form of constraints, the resulting clustering paradigm is known as Constrained Clustering (CC).



Background Knowledge can be found in many forms, which can be categorized as follows.

Background Knowledge



Background Knowledge



Background Knowledge



Other

Constraints

Background Knowledge



Background Knowledge



Constraints

Background Knowledge



Background Knowledge



Equivalencies



Must-link & Cannot-link



Instance-level Pairwise Constrained Clustering (CC)

The goal of constrained clustering is to **find a partition of the dataset that ideally meets all constraints** in the union of both constraint sets.



The Infeasibility

The infeasibility refers to the number of constraints broken by a given partition.



 ${\sf Infs}(C,CS)=0$ ${\sf Infs}(C,CS)=2$

$$\operatorname{Infs}(C, CS) = \sum_{C_{=}(x_i, x_j) \in CS} \mathbb{1} \left[\left[l_i^C \neq l_j^C \right] \right] + \sum_{C_{\neq}(x_i, x_j) \in CS} \mathbb{1} \left[\left[l_i^C = l_j^C \right] \right]_{4/51}$$

The Feasibility Problem

How do constraints affect the complexity of clustering? Intuitively, the clustering problem goes from its classic formulation "find the best partition for a given dataset" to its constrained form "find the best partition for a given dataset satisfying all constraints in the constraint set".

In **partitional CC** the number of clusters is bounded and fixed.

In **hierarchical CC** the full dendrogram is available.





The Complexity of the Feasibility Problem

- **Partitional CC is NP-Complete** It can be reduced from the Graph K-Colorability problem.
- **Hierarchical CC is NP-Complete** It can be reduced from the One-in-three 3SAT with positive literals problem.

Constraints	Partitional CC	Hierarchical CC	Dead Ends?
ML	Р	Р	No
CL	NP-complete	NP-complete	Yes
ML and CL	NP-complete	NP-complete	Yes

Knowing that a feasible solution exists does not help us find it.

Hard Constraints VS. Soft Constraints

Hard Constraints



large Infs(C,CS)=0?

Soft Constraints



;
$$\mathsf{Infs}(C,CS)=0$$
 ?; $\mathsf{Infs}(C,CS)=1$?

Soft constraints are **resilient to noise**, and allow for **flexibility** in the cost/objective function and in the optimization procedure.



Taxonomic Tree



- 2 Major Categories: Constrained Partitional & Constrained DML.
- **17 Subcategories:** Constrained K-Means, Ensemble CC, Active CC, etc.
- 29 Final Categories.
- 307 methods reviewed.

Taxonomic Tree



Some categories are exclusive to CC and are proposed for the first time in this taxonomization.

Taxonomic Tree



Other categories are exclusive to CC and existed before this taxonomization.

Taxonomic Tree



Other categories are exclusive to CC and existed before this taxonomization.

Experimental elements - Datasets



- Classification datasets are used as benchmarks, using the labels as an oracle to generate the constraint sets.
- Most studies use between 1 and 10 datasets in their experiments.
- The number of datasets used in experiments shows an increasing tendency.

Experimental elements – Methods



Used Not used $(149 \sim 48.5\%)$ (158 ~ 51.5%)

Proportion of CC methods

used in experiments

Experimental elements – Methods





Experimental elements - Validity Indices



Proportion of CC studies using statistical tests









3SHACC – First Stage

Run WLSI over X, Iterate over PRun $\mathcal A$ over X γ W, $C_{=}$ and $C_{
eq}$ to to compute times to build get metric constraints partitions matrix Pweights matrix Wmatrix M

3SHACC – Second Stage



The similarity computation is based on the reconstruction coefficient, which computes similarities in terms of how much an instance is explained by others.

3SHACC – Third Stage





Types of Background Knowledge





- There is an order relationship between classes.
- Our goal is to minimize the number of misclassifications regarding the class order.
- The costs of misclassifications are different for every class.

Monotonicity in MCDA

MCDA introduces the concept of preference. The **preference** quantifies the addition of differences between the features of two instances, **limited to the features in which one of them is strictly better than the other**.



PCKM-Mono – Objective Function

An instances is assigned to a cluster if **its centroid is the most similar to it in terms of preference**, taking the number of violated constraints into account.

$$J_{PCKMM} = \frac{1}{K} \sum_{k=1}^{K} \sum_{x_i \in c_k} \left[(r(x_i, \mu_k) - r(\mu_k, x_i)) + \sum_{(x_i, x_j) \in C_{=}} \mathbb{1} \left[l_i \neq l_j \right] + \sum_{(x_i, x_j) \in C_{\neq}} \mathbb{1} \left[l_i = l_j \right]$$

PCKM-Mono – Objective Function

An instances is assigned to a cluster if **its centroid is the most similar to it in terms of preference**, taking the number of violated constraints into account.

$$J_{PCKMM} = \frac{1}{K} \sum_{k=1}^{K} \sum_{x_i \in c_k} |(r(x_i, \mu_k) - r(\mu_k, x_i)| + \sum_{(x_i, x_j) \in C_{\neq}} \mathbb{1}[[l_i \neq l_j]] + \sum_{(x_i, x_j) \in C_{\neq}} \mathbb{1}[[l_i = l_j]]$$

PCKM-Mono - Expectation-Minimization

An instances is assigned to a cluster if its centroid is the most similar to it in terms of preference, taking the number of violated constraints into account.

$$\begin{aligned} \arg\min_{h} \left(|\sum_{j=1}^{u} (x_{[i,j]} - \mu_{[h,j]})| + \right. \\ \sum_{x_{j}: (x_{i},x_{j}) \in C_{=}} \mathbb{I}[l(c_{h}) \neq l_{j}] + \sum_{x_{j}: (x_{i},x_{j}) \in C_{\neq}} \mathbb{I}[l(c_{h}) = l_{j}] \right) \end{aligned}$$

The **centroids must be neutral** in terms of preference with respect to all instances in the cluster.

$$\mu_i = \frac{1}{|c_i|} \sum_{x_i \in c_i} x_i$$

Case of Study



Shanghai Ranking of World Universities

Clustering problem with both monotonicity constraints and instance-level pairwise constraints.

Features monotonicity faults in 7% of its instances.

Case of Study



47/51

Semi-Supervised Learning Semi-Supervised Clustering **Constrained Clustering** Hybrid Models for CC New Paradigms in CC **Challenges and Frontiers**

Future Work

- **Creation of a free CC library**. The creation of an open-access library specialized in CC methods would greatly stimulate research in the area.
- **Constraint-based preprocessing**. We argue that constraints can be beneficial in preprocessing procedures.
- **Preprocessing the constraint set**. As the constraint set can be considered a dataset itself, it can suffer from the same imperfections as traditional datasets, such as: missing values, noise or redundancies.
- New combinations of types of background knowledge. Another potential research direction is to investigate how to automatically identify the best combination of background knowledge for a given problem, hence keeping low human effort and cost.

Thanks for your attention



Applications

Proportion of proposals by field of application

