

---

# Zero-shot Active Learning with Topological Clustering for Multiclass Classification

---

Anonymous Author(s)

Affiliation

Address

email

## Abstract

1 We present a novel approach for zero-shot active learning for multi-class classifica-  
2 tion based on a clustering technique, called ToMATo, which is guided by topological  
3 persistence. Our objective is to identify effective regions in the feature space for  
4 label querying. The labeling of examples in these regions will allow the training of  
5 efficient multi-class classification prediction functions. We have adapted ToMATo  
6 with a density aware  $\delta$ -Rips graph in order to obtain homogeneous simplicial trees.  
7 From these trees, informative simplices are identified with respect to the annotation  
8 effort, or the budget. Representative examples from each of them are labeled  
9 by an oracle and these labels are then propagated through the trees. We adapt  
10 ToMATo by computing our *persistence diagram* (PD) from a  $\delta$ -Rips graph that is  
11 estimated using a  $k$ -nearest neighbor distance matrix. This allows the application  
12 of the method to large scale scenarios. From this perspective we also propose a  
13 local density estimator from the same distance matrix. Comparisons on different  
14 benchmarks show that the proposed approach greatly improves performance with  
15 respect to a random querying strategy for label assignment that has been found  
16 outperforming state-of-the art approaches in previous works.

## 17 1 Introduction

18 In many real-life applications, the labeling of training observations for learning is costly and some-  
19 times not even realistic. For example, in web oriented applications, huge amount of observations are  
20 collected sequentially. However, there is not enough time to label these data for different purpose  
21 while unlabeled data are abundant. Different attempts have been made to reduce the annotation  
22 burden. For example, we can refer to the so many successful *semi-supervised* and *active learning*  
23 approaches that have been proposed until now [Baram et al., 2004, Settles, 2012, Chapelle et al.,  
24 2006, Amini and Usunier, 2015].

25 All of these approaches suppose that there exists a small set of labeled training data together with  
26 a large set of unlabeled examples. They also tend to identify additional informative unlabeled  
27 observations to be (pseudo-)labeled for learning. Besides, some of the proposed strategies are based  
28 on the approximation of the risk of selection, *e.g.* with iterative methods [Zhao et al., 2006, Zhu et al.,  
29 2008], or the selection of the most uncertain observations regarding some confidence measures, *e.g.*  
30 with model-driven models [Lakshminarayanan et al., 2017, Yan et al., 2011].

31 In this work, we suppose that there is no initial labeled training data and propose a new zero-shot  
32 active learning strategy based on topological persistence in order to find informative observations to  
33 be labeled for learning. More precisely, our approach is based on topological data analysis which has  
34 recently brought exciting new ideas to the machine learning community, especially in *unsupervised*  
35 *learning* with *topological clustering* [Bonis and Oudot, 2018, Cabanes et al., 2013]. Among these

36 studies, ToMATo [Chazal et al., 2013] is a mode-seeking clustering algorithm with a cluster merging  
 37 phase guided by *topological persistence*. It relies on the concept of *prominence* by computing the  
 38 PD which reflects the modes of the density, and a prominence threshold is estimated from the PD to  
 39 merge clusters and discard noise. This method is adapted and used in this work in order to detect  
 40 informative observations to be labeled.

41 The contributions of this work are the following:

- 42 • we propose a generic data driven approach for zero-shot learning based on the ToMATo  
 43 clustering technique;
- 44 • we investigate different ways to scale up computations in use in this approach and we derive  
 45 a density aware formulation of  $\delta$ -Rips graph.

46 We validate our approach by comparing it to a random querying strategy for label assignment on  
 47 different benchmarks.

## 48 2 TCAL: Topological Clustering for Active Learning

49 Let  $(\mathcal{X}, d)$  be a metric space, where  $d: \mathcal{X} \times \mathcal{X} \rightarrow [0, \infty)$  is the distance  
 50 metric. We assume that we have a set  
 51  $X := \{x_i\}_{i=1}^n \subset \mathcal{X}$  of i.i.d unlabeled  
 52 examples drawn from an arbitrary un-  
 53 known distribution  $\mathcal{D}$  over  $\mathcal{X} \times \mathcal{Y}$ , where  
 54  $\mathcal{Y} := [c]$  is a set of  $c > 1$  unknown  
 55 classes for which each example  $x_i$  has  
 56 a unique, yet unknown, label  $y_i$ . We also  
 57 have at our disposal a perfect oracle (la-  
 58 beling from the expert)  $\mathcal{O}: \mathcal{X} \rightarrow \mathcal{Y}$  over  
 59  $X$  where  $\mathbb{P}_{x_i \sim X}(\mathcal{O}(x_i) = y_i) = 1$  and  
 60 a budget  $\mathcal{B}$  which corresponds to the max-  
 61 imum number of examples the expert can  
 62 label. We describe the main steps of the  
 63 proposed method denoted by TCAL and  
 64 which is summarized in Algorithm 1.

### Algorithm 1: TCAL

**Input:**  $X := \{x_i\}_{i=1}^n$ , oracle  $\mathcal{O}$ , budget  $\mathcal{B}$  and a  
 distance metric  $d: \mathcal{X} \times \mathcal{X} \rightarrow [0, \infty)$ .

- Compute distance matrix  $M$  and density estimator  $(\tilde{f}(i))_{1 \leq i \leq n}$  with (1).
- Build graph  $G$  with the threshold distance  $\delta$ .
- Run ToMATo( $G, \tilde{f}(i), \tau$ ) to obtain the set of clusters  $(\mathcal{C}_k)_{1 \leq k \leq \mathcal{K}}$ .
- Ask an oracle to label the  $\mathcal{B}$  peaks from  $(\mathcal{C}_k)_{1 \leq k \leq \mathcal{K}}$ .
- Propagate the labels and subsample majority classes as in (2).

**Output:** Set of labeled training examples  $S$

65 For each  $x_i \in X$  we compute its distance to its  $l$ -nearest neighbors<sup>1</sup> and we obtain a sparse distance  
 66 matrix  $M = (m_{i,j}) \in \mathbb{R}^{n \times n}$  (with only  $l$  non zero values in each row) where

$$m_{i,j} = \begin{cases} d(x_i, x_j) & \text{if } x_j \text{ is one of the } l\text{-nearest neighbors of } x_i, \\ 0 & \text{elsewhere.} \end{cases}$$

67 The use of  $l$ -nearest neighbors allows us to consider large data sets, and we observe in practice no  
 68 loss in accuracy.

69 To estimate the density from  $M$ , we propose to use  $\tilde{f}$ , where, for all  $1 \leq i \leq n$ ,

$$\tilde{f}(i) = \left( \frac{1}{l} \sum_{j=1}^n m_{i,j}^2 \right)^{-1/2}. \quad (1)$$

70 A  $\delta$ -Rips graph  $G = (V, E)$  is constructed from  $\tilde{f}$  where  $V = [n]$  and

$$E = \{(i, j) \mid i \in V, j \in [k] \text{ and } 0 < m_{i,j} \leq \delta(i)\}.$$

71 Here, for all  $1 \leq i \leq n$ , the threshold distance  $\delta(i)$  is defined with respect to the density estimation  
 72 at  $x_i$  by  $\delta(i) = \delta_0 \left( \alpha - \tilde{f}(i) \right)^{1/\beta}$ , where  $\delta_0$  corresponds to the initial threshold,  $\alpha$  allows to take  
 73 into account the shift in density distribution, and  $\beta$  can be interpreted as the evolution of similarity in  
 74 different density levels.

<sup>1</sup>In practice, we use  $l = \lfloor n/10 \rfloor$ .

75 A set of  $\mathcal{K}$  clusters  $(\mathcal{C}_k)_{1 < k \leq \mathcal{K}}$  is computed by running ToMATo on the graph  $G$  with density  $\tilde{f}$ . The  
 76 prominence threshold  $\tau$  to merge clusters has to be set and its selection is discussed in Section 3.  
 77 Those clusters are sorted by decreasing size, to focus first on the largest cluster. Then, the oracle is  
 78 providing labels associated to the  $\mathcal{B}$  peaks from  $(\mathcal{C}_k)_{1 \leq k \leq \mathcal{K}}$  starting from the largest, and these labels  
 79 are propagated in their respective clusters, the result being an imbalanced set of labeled data:

$$S = \bigcup_{k=1}^{\min(\mathcal{K}, \mathcal{B})} S_k, \quad \text{where } S_k = \left\{ \left( x_l, \mathcal{O} \left( \arg \max_{i \in \mathcal{C}_k} \tilde{f}(i) \right) \right) \right\}_{l \in \mathcal{C}_k}. \quad (2)$$

80 If the hyper-parameters are carefully tuned, the set  $S$  have labels from most, if not all, of the  $c$  classes,  
 81 but with an unequal distribution. At this last step, the large classes are randomly subsampled to obtain  
 82 a balanced training set.

### 83 3 Parameter selection

84 **Graph parameters** Finding the right graph representation that explains class similarity across  
 85 different data collections is a universal problem for graph based methods. The threshold distance  
 86 function  $\delta$  for the graph  $G$  has three parameters  $\delta_0, \beta, \alpha$ , describing that the similarity between  
 87 data points is not uniform overall density levels, especially for multi-class data sets. A radius that  
 88 gives good representatives of dense classes will capture less information in low density classes.  
 89 Inversely, larger radius will capture representatives in low density levels but also diffuse noise in  
 90 high density regions. In practice, we consider hyperparameters  $(\delta_0, \beta, \alpha)$  such that the threshold  
 91 function has a decreasing behavior in the region  $[\min \tilde{f}, \max \tilde{f}] \times [\min M_{avg}, \max M_{avg}]$  where  
 92  $M_{avg} = \left\{ \sum_{j=1}^n m_{i,j} / l \mid 1 \leq i \leq n \right\}$ . We notice that on image data sets (Coil20, MNIST and  
 93 Statlog), the same hyperparameters are providing good performance.

94 **Prominence threshold** The parameter  $\tau$  in ToMATo is used to filter out topological noise and to  
 95 distinguish between relevant peaks from subsidiary peaks, coming from parents in  $\tilde{f}$ . It makes a  
 96 trade-off between purity and cluster size. More precisely, increasing  $\tau$  will merge more clusters (thus  
 97 reduce the number  $\mathcal{K}$  of clusters) and diffuse more noise, so that  $(\mathcal{C}_k)_{1 < k \leq \mathcal{K}}$  contains larger clusters  
 98 with less purity with respect to their true labels. We follow a similar but more conservative procedure  
 99 to Chazal et al. [2013]: we minimize the diffusion of the noise and keep reasonable cluster sizes.  
 100 To do so, we sort the prominent peaks given by  $\tilde{f}$  by decreasing order in  $V$ , and we fix  $\tau$  to be the  
 101 more stable value after a significant gap in the distribution. In practice, we compute the variance on a  
 102 sliding window of size  $\lfloor |V|/10 \rfloor$ , we select the window that has the lowest variance after the window  
 103 of maximum variance, and we fix  $\tau$  to its median value.

## 104 4 Experiments

105 **Data sets** We conduct experiments on benchmark data sets for classification problems also often  
 106 used in active learning: MNIST [LeCun et al., 1998], COIL-20 [Yang et al., 2011], Isolet [Fanty and  
 107 Cole, 1991] and sensorless drive diagnosis SDD [Paschke et al., 2013] as well as two imbalanced  
 108 data sets including Protein [Higuera et al., 2015] and Statlog [King et al., 2000]. Table 1 presents  
 109 statistics of the data sets on the four first columns.

110 **Baseline** Following results from Siméoni et al. [2019], we only use random labeling strategy, as it  
 111 outperforms many recent strategies in active learning with small budget scenarios.

112 A simple linear support vector machine with stochastic gradient descent is used for the classifier for  
 113 both methods, with default parameters and a single epoch, since our objective is to show the gap in  
 114 learning performance of any hypothesis class with respect to the training set in an online setting. For  
 115 the metric distance for TCAL, we use the euclidean distance overall data sets.

116 We run the two procedures for several budgets  $\mathcal{B}$  (less than 0.5% of the sample size) and compare  
 117 the evolution. As the expert labels only one observation per cluster, in our experiments the budget is  
 118 upper bounded by the number of clusters detected by ToMATo. 20 random splits are considered, with  
 119 70% of the data in the training set and 30% in the test set.

Table 1: Average classification accuracy (in %) and standard deviation over 20 random splits for different budgets  $\mathcal{B}$ . The third column corresponds to the dimension of the feature space  $\mathbb{R}^p$ .

Dataset	$n$	$p$	$c$	Budget $\mathcal{B}$	Random	TCAL
Protein	1080	77	8	5	$18.70 \pm 5.54$	<b><math>32.29 \pm 4.47</math></b>
				10	$20.35 \pm 8.38$	<b><math>38.56 \pm 4.60</math></b>
COIL-20	1440	1024	20	10	$23.00 \pm 3.24$	<b><math>41.57 \pm 3.48</math></b>
				50	$57.28 \pm 5.90$	<b><math>82.74 \pm 3.43</math></b>
				100	$75.44 \pm 4.55$	<b><math>95.08 \pm 1.59</math></b>
Isolet	6238	617	26	10	$13.22 \pm 2.53$	<b><math>33.29 \pm 1.88</math></b>
				50	$26.28 \pm 2.81$	<b><math>55.09 \pm 2.81</math></b>
				100	$43.69 \pm 4.18$	<b><math>63.13 \pm 4.05</math></b>
Statlog	6435	36	6	10	$31.32 \pm 12.83$	<b><math>65.60 \pm 2.64</math></b>
				20	$32.16 \pm 14.61$	<b><math>66.28 \pm 1.97</math></b>
SDD	58.5k	48	11	100	$33.84 \pm 6.55$	<b><math>46.92 \pm 3.01</math></b>
				250	$38.57 \pm 6.50$	<b><math>47.54 \pm 4.05</math></b>
				500	$43.13 \pm 5.44$	<b><math>49.37 \pm 3.15</math></b>
MNIST	70k	784	10	100	$68.22 \pm 2.96$	<b><math>81.00 \pm 1.01</math></b>
				700	$82.91 \pm 1.05$	<b><math>87.29 \pm 0.50</math></b>
				1000	$83.45 \pm 1.24$	<b><math>88.12 \pm 0.54</math></b>
				1400	$84.03 \pm 0.71$	<b><math>88.75 \pm 0.46</math></b>

120 **Results** Table 1 presents the results of our approach and of the random labeling strategy over all  
 121 data sets. In all cases, TCAL provides significantly better results than the random strategy. For  
 122 imbalanced data sets (Protein and Statlog), the random strategy is affected by the imbalance in class  
 123 distribution and performs badly, whereas our method has benefit from the clustering step and performs  
 124 the best. We remark that even with very few labels, if the labels are given for observations particularly  
 125 discriminant, the performance are very high, *e.g.* for COIL-20 with 100 labels over the 20 classes,  
 126 where the accuracy of TCAL is 95%, while the performance for the random strategy is 75%.

127 **Highlight: the effect of the prominence**  
 128 **threshold** We investigate the effect of  
 129 the prominence threshold  $\tau$  on the classi-  
 130 fier performance. Figure 1 shows the aver-  
 131 age accuracy curve of 20 random splits  
 132 on MNIST data set for 100 points of  $\tau$  uni-  
 133 formly in  $[0, 10^{-2}]$  with a budget of 100,  
 134 where  $\hat{\tau}$  is estimated following the proce-  
 135 dure in Section 3. First, we remark that  
 136 whatever the value for  $\tau$  (small enough),  
 137 we get good performance compared with  
 138 the random labeling strategy. In addition,  
 139 the strategy of largest gap for  $\tau$  estimation  
 140 in ToMATo fails in this case, where it gives  
 141 a value around 0.02. Finally, note that our  
 142 estimation is very close to the best one.

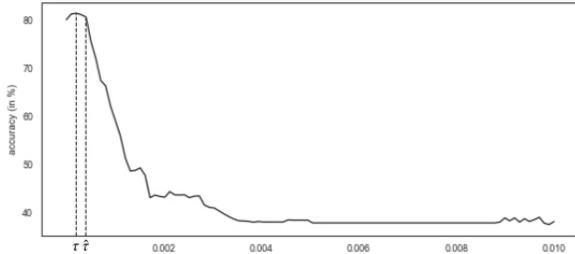


Figure 1: Average classification accuracy (in %) of 20 random splits on MNIST data set with different values for the prominence threshold. Optimal value  $\tau$  is shown with dashed line whereas  $\hat{\tau}$  is computed as explained in Section 3.

## 143 5 Conclusion

144 We propose a data driven method for zero-shot active learning in multiclass classification problems  
 145 with topological clustering. Our empirical study validates this method on different benchmark data  
 146 sets. This work is, to our knowledge, the first significant step to use topological data analysis to detect  
 147 relevant observations for zero-shot active learning. Challenging open questions are left, as the use of  
 148 semi-supervised model to conclude the analysis (instead of a supervised classifier) and theoretical  
 149 results that guarantee good performance in active learning.

## 150 References

- 151 Yoram Baram, Ran El-Yaniv, and Kobi Luz. Online choice of active learning algorithms. *J. Mach.*  
152 *Learn. Res.*, 5:255–291, 2004.
- 153 Burr Settles. *Active Learning*. Morgan & Claypool Publishers, 2012.
- 154 Olivier Chapelle, Bernhard Schölkopf, and Alexander Zien. *Semi-Supervised Learning*. MIT Press,  
155 2006.
- 156 Massih-Reza Amini and Nicolas Usunier. *Learning with Partially Labeled and Interdependent Data*.  
157 Springer, 2015.
- 158 Yue Zhao, Ciwen Xu, and Yongcun Cao. Research on query-by-committee method of active learning  
159 and application. In Xue Li, Osmar R. Zaiane, and Zhanhuai Li, editors, *Advanced Data Mining*  
160 *and Applications*, pages 985–991, Berlin, Heidelberg, 2006. Springer Berlin Heidelberg.
- 161 Jingbo Zhu, Huizhen Wang, Tianshun Yao, and Benjamin K Tsou. Active learning with sampling by  
162 uncertainty and density for word sense disambiguation and text classification. In *Proceedings of*  
163 *the 22nd International Conference on Computational Linguistics (Coling 2008)*, pages 1137–1144,  
164 Manchester, UK, 2008.
- 165 Balaji Lakshminarayanan, Alexander Pritzel, and Charles Blundell. Simple and scalable predictive  
166 uncertainty estimation using deep ensembles. In *Advances in Neural Information Processing*  
167 *Systems 30*, pages 6402–6413, 2017.
- 168 Yan Yan, Romer Rosales, Glenn Fung, and Jennifer G. Dy. Active learning from crowds. In *Proceed-*  
169 *ings of the 28th International Conference on International Conference on Machine Learning*, page  
170 1161–1168, 2011.
- 171 Thomas Bonis and Steve Oudot. A fuzzy clustering algorithm for the mode-seeking framework.  
172 *Pattern Recognition Letters*, 102:37 – 43, 2018.
- 173 Guénaél Cabanes, Younès Bennani, Renaud Destenay, and André Hardy. A new topological clustering  
174 algorithm for interval data. *Pattern Recognition*, 46(11):3030 – 3039, 2013.
- 175 Frédéric Chazal, Leonidas J. Guibas, Steve Y. Oudot, and Primoz Skraba. Persistence-based clustering  
176 in riemannian manifolds. *J. ACM*, 60(6), 2013.
- 177 Yann LeCun, Léon Bottou, Yoshua Bengio, and Patrick Haffner. Gradient-based learning applied to  
178 document recognition. *Proceedings of the IEEE*, 86(11):2278–2324, 1998.
- 179 Jianwei Yang, Zirun Chen, Wen-Sheng Chen, and Yunjie Chen. Robust affine invariant descriptors.  
180 *Mathematical Problems in Engineering*, 2011.
- 181 Mark Fanty and Ronald Cole. Spoken letter recognition. In R. P. Lippmann, J. E. Moody, and  
182 D. S. Touretzky, editors, *Advances in Neural Information Processing Systems 3*, pages 220–226.  
183 Morgan-Kaufmann, 1991.
- 184 Fabian Paschke, Christian Bayer, Martyna Bator, Uwe Mönks, Alexander Dicks, Olaf Enge-  
185 Rosenblatt, and Volker Lohweg. Sensorlose Zustandsüberwachung an Synchronmotoren. In  
186 F. Hoffmann, editor, *23. Workshop Computational Intelligence 2013. Proceedings*, volume 46,  
187 pages 211–225. KIT Scientific Publishing, 2013.
- 188 Clara Higuera, Katheleen J. Gardiner, and Krzysztof J. Cios. Self-organizing feature maps identify  
189 proteins critical to learning in a mouse model of down syndrome. *PLOS ONE*, 10(6):1–28, 2015.
- 190 Ross D. King, Chen-Chieh Feng, and Alistair A.P. Sutherland. Statlog: Comparison of classification  
191 algorithms on large real-world problems. *Applied Artificial Intelligence*, 9, 2000.
- 192 Oriane Siméoni, Mateusz Budnik, Yannis Avrithis, and Guillaume Gravier. Rethinking deep active  
193 learning: Using unlabeled data at model training. 2019.