



Density Estimators for Positive-Unlabeled Learning

Teresa M. A. Basile^{2,3}, Nicola Di Mauro¹ , Floriana Esposito¹ ,
Stefano Ferilli¹, and Antonio Vergari¹

¹ Department of Computer Science, University of Bari “Aldo Moro”, Bari, Italy
`nicola.dimauro@uniba.it`

² Department of Physics, University of Bari “Aldo Moro”, Bari, Italy

³ National Institute for Nuclear Physics (INFN), Bari Division, Bari, Italy

Abstract. Positive-Unlabeled (PU) learning works by considering a set of positive samples, and a (usually larger) set of unlabeled ones. This challenging setting requires algorithms to cleverly exploit dependencies hidden in the unlabeled data in order to build models able to accurately discriminate between positive and negative samples. We propose to exploit probabilistic generative models to characterize the distribution of the positive samples, and to label as reliable negative samples those that are in the lowest density regions with respect to the positive ones. The overall framework is flexible enough to be applied to many domains by leveraging tools provided by years of research from the probabilistic generative model community. Results on several benchmark datasets show the performance and flexibility of the proposed approach.

1 Introduction

The classical supervised setting of statistical machine learning [14] aims at inducing models (classifiers) from training sets of labeled data in the form of samples (\mathbf{x}^i, y^i) i.i.d. drawn from an unknown joint probability distribution $p(\mathbf{X}, Y)$ over random variables (RVs) \mathbf{X} and Y , where Y is the *label*. For binary classification, i.e., $Y \in \{0, 1\}$, labels y^i are assumed to be modeled by a Bernoulli distribution and are associated to *positive* and *negative* samples \mathbf{x}^i .

While nowadays gathering and storing all kinds of data is easier and easier, having all these data perfectly and reliably labeled is unrealistic for several reasons, which makes classical approaches to learning classifiers inapplicable. First, the exponential rate at which data are produced contrasts the time required to produce high quality labels. Moreover, in many fields there are relatively few *labelers* effectively trained to produce reliable labels. Lastly, in many real-world domains it is sometimes unclear what should be considered as a negative sample, or the generation of negative samples is too expensive or just impossible. E.g., in process enactment, one would not waste time, money and resources to build a wrong item just for the purpose of showing how things are not to be done. Thus, the ability to learn predictive models in these scenarios may allow one to exploit the vast amount of data that are produced, saving precious time and resources.

In Positive-Unlabeled (PU) learning [8, 23], a set \mathcal{P} of positive samples, and a set \mathcal{U} of unlabeled samples—each of which may be positive or negative—are available at training time. So, discriminative information for the negative class must be found in unlabeled data. PU learning shares similarities with semi-supervised learning [26], one-class classification [28], and outlier detection [5]. Differently from the first, no negative samples are available at training time and yet it is required to learn a discriminator between the two classes, in contrast with the second. Additionally, PU learning is in opposition to the last which is usually performed in a transductive way to label unlabeled training data only.

The interest for PU learning is supported by its successful application in several domains [21, 35, 37]. PU learning approaches can be roughly grouped into *two-staged*—extracting a set of reliable negative samples (RN) from \mathcal{U} and then performing supervised learning—and *single-staged*—taking all samples in \mathcal{U} as negative. For the former, it becomes crucial to learn a metric that is able to discriminate among classes. However, each application domain needs a specific formulation for such a metric. Hence, ad hoc algorithmic solutions are often required to cope with different data representations [18, 37]. Few approaches have been proposed to deal with this categorical data [4, 18] in PU learning, due to their being quite challenging—there is no natural distance for them [19].

This work introduces *Generative Positive-Unlabeled* (GPU) learning, a novel two-staged approach to PU learning that aims to be general enough to support very different application domains¹. It estimates the marginal distribution $\mathbf{p}_{\mathcal{P}}(\mathbf{X}|Y = 1)$ of the positive samples in \mathcal{P} via a generative model, and then performs inference on such a distribution to select a set of reliable negative samples from \mathcal{U} . The modeled probability density implicitly defines a metric space among samples, and we assume negative ones to be concentrated where positive ones are less likely. Generative models such as Probabilistic Graphical Models (PGMs) [20] have been extensively studied in the literature and offer a powerful formalism to deal with complex probability distributions over continuous, categorical, or even structured data [33]. Dealing with a particular domain translates into choosing a suitable PGM from a consolidated research field. More generally, given a PGM learned as a density estimator in a certain domain, we exploit it as a negative sample extractor for partially labeled data. Albeit GPU can deal with different data representations, here we focus on categorical data, which are handled natively by PGMs. We compared GPU on real data to several PU learners that have proven to be effective on categorical data.

The paper is organized as follows: in the next section we provide a brief review of the literature about PU learning. In Sect. 3 we introduce and discuss our GPU approach, while the experimental setting and the experiment results are presented in Sect. 4. Conclusions are drawn in Sect. 5.

¹ This paper is an extended version of [2] presented at the International Workshop NFMCP held in conjunction with ECML/PKDD 2017.

2 Related Works

PU learning has attracted a great deal of attentions in machine learning and data mining research. An extensively adopted approach to PU learning is based on a negative set construction process, that first identifies reliable negative samples from the unlabeled ones and, then, directly applies traditional classification methods. Alternative methods following this paradigm differ for how they implement these two steps.

Several proposals adopt distance-based approaches to identify negative samples, as the farthest unlabeled ones from positive samples. In [34], after selecting features statistically related to positive samples, the unlabeled set is partitioned into four sets (reliable/likely/weak negative and likely positive) based on the Euclidean distance. Successively, a multi-level samples learning technique, weighted SVMs, is exploited to build a classifier. The same approach of first identifying, characterizing and discriminating features for positive samples is adopted in [18], where a particular distance function previously designed by the authors is used to determine reliable negative samples; then, distance learning is applied twice—on the positive and reliable negative samples—and the resulting distances are used for k -NN classification.

After having theoretically shown that, under appropriate conditions, \mathcal{P} and \mathcal{U} provide sufficient information for learning, in [23] PU learning is posed as a constrained optimization problem. In such a setting, the set of reliable negative samples is selected by using a Naive Bayes (NB) classifier and EM. To the extreme, all the unlabeled samples are treated as negative samples in the NB classifier initially learned and successively used to extract the set of reliable negatives from unlabeled data [22]. The dataset so obtained is finally exploited to learn a classifier using SVM. The obtained augmented set, as representative of negatives samples, is exploited, along with positive ones, to compute the parameters of the NB classifier devoted to reliable negative samples identification. Finally, an EM-based algorithm is exploited to learn the predictive model.

A different policy is the weighted-based approach on unlabeled data exploited in [12]. The study shows that a classifier trained on positive and unlabeled samples is able to predict probabilities that differ by only a constant factor from the true conditional probabilities produced by a model trained on fully labeled positive and negative samples, provided that the labeled positive samples are chosen completely at random from all positive samples. This result is used in two different ways: learning from \mathcal{P} versus \mathcal{U} with adjustment of output probabilities finally assigned to unlabeled samples, and learning from \mathcal{P} and \mathcal{U} after double weighting of \mathcal{U} . The basic learning algorithm for each method is an SVM with a linear kernel whose outputs are post-processed into calibrated probabilities by fitting a one-dimensional logistic regression function.

Naive Bayes is the classifier extensively adopted for categorical data in the four methods proposed in [4], namely (Average) Positive Naive Bayes ((A)PNB), based on Naive Bayes, and (Average) Positive TAN ((A)PTAN), two variants of the Tree Augmented Naive Bayes model [13] able to deal with positive and unlabeled samples. The difference lies in the way the prior probability for the negative

class is estimated. For PNB and PTAN this probability is derived directly from the whole set of unlabeled samples while for APNB and APTAN the uncertainty is modeled by a Beta distribution.

The above survey shows that many works on PU learning [4, 18, 22, 23] have adopted the text categorization perspective, which is quite peculiar. Indeed, features are intrinsically categorical, there is a huge number of features compared to other settings, the representation of samples is very sparse, and there is a heavy impact of text pre-processing in setting up the classification problem. Others have faced biomedical problems [12, 34], where it is typical that databases specify which genes or proteins are related to some specific consequence, but this does not mean that all the others are unrelated to that consequence and, on the contrary, there is a strong interest in identifying which ones actually are [12].

As previously pointed out, although PU learning shares similarities to outlier detection [5] and to one-class classification [28], it shows difference from these settings in both the goal to fulfill and in the training set exploited, even in the case of probability density estimation techniques are used as solving strategies [15, 27, 29, 32]. Indeed, both methods aim at learning a model able to reject the new incoming samples using positive training data only. They do not require to learn a discriminator between the two classes and, hence, no effort to learn a model for the negative class is done. Intuitively, this type of approach is inferior because it ignores useful information that is present in the unlabeled samples.

3 Methodology

Let RVs be denoted by upper-case letters, e.g., X , and their values as the corresponding lower-case letters, e.g., $x \sim X$. We denote sets of RVs as \mathbf{X} , and their combined values as \mathbf{x} . When we refer to a joint probability distribution $\mathbf{p}(\mathbf{X})$ over RVs \mathbf{X} , we are either considering the joint probability density function for continuous RVs, or the probability mass function for discrete RVs, or a hybrid combination of both in hybrid domains [20, 33]. To denote a finite domain of a discrete RV X_j we introduce the following notation $\text{Val}(X_j) = \{x_j^k\}_{k=1}^K$. If \mathcal{D} is a set of samples over RVs \mathbf{X} , we indicate with $\mathbf{p}_{\mathcal{D}}(\mathbf{X})$ the real (unknown) probability distribution that generated the data, while if \mathcal{M} indicates a generative model, $\mathbf{p}_{\mathcal{M}}(\mathbf{X})$ refers to the probability distribution estimated by such a model on finite sample sets. Disambiguation is provided by context. Generally one wants the estimate $\mathbf{p}_{\mathcal{M}}(\mathbf{X})$ to be as close as possible to $\mathbf{p}_{\mathcal{D}}(\mathbf{X})$. A common way to measure this closeness is via the *log-likelihood* function [20], or one of its variants, defined as:

$$\ell_{\mathcal{D}}(\mathcal{M}) = \sum_{\mathbf{x}^i \in \mathcal{D}} \log \mathbf{p}_{\mathcal{M}}(\mathbf{x}^i).$$

In the classical PU learning setting, one has a training set $\mathcal{D} = \mathcal{P} \cup \mathcal{U}$ i.i.d. from $\mathbf{p}(\mathbf{X}, Y)$. Samples in \mathcal{P} are provided with a known positive class label, i.e., $\mathcal{P} = \{(\mathbf{x}^i, 1)\}_{i=1}^{m_{\mathcal{P}}} \sim \mathbf{p}_{\mathcal{P}}(\mathbf{X}|Y = 1)$. On the other hand, class information, i.e., labels, is not provided for samples in \mathcal{U} , i.e., $\mathcal{U} = \{\mathbf{x}^i\}_{i=1}^{m_{\mathcal{U}}} \sim \mathbf{p}_{\mathcal{U}}(\mathbf{X})$, where $\mathbf{p}_{\mathcal{U}}(\mathbf{X})$ is the marginal probability distribution w.r.t. $\mathbf{p}_{\mathcal{U}}(\mathbf{X}, Y)$. Let \mathcal{D}_0 (resp. \mathcal{D}_1)

denote the subset of all negative (resp. positive) samples in \mathcal{D} . The aim of PU learning is to build a discriminator model $f : \mathbf{X} \rightarrow Y$ from \mathcal{D} in order to make accurate predictions about the labels of unseen test data samples. Following [12], we assume that samples in \mathcal{P} are *selected completely at random* from all positive samples in \mathcal{D} , i.e., $\mathbf{p}_{\mathcal{P}}(\mathbf{X}|Y = 1) = \mathbf{p}_{\mathcal{D}}(\mathbf{X}|Y = 1)$.

3.1 Generative Models for PU Learning

Our proposed approach, Generative PU learning (GPU), falls in the category of two-staged methods for PU learning. First it extracts a set of reliable negative samples \mathcal{N} from \mathcal{U} , then \mathcal{N} is employed to perform supervised learning. In the following we detail our contribution to the first step, discussing possible approaches for the second one.

As usual in statistical machine learning, we assume $\mathbf{p}_{\mathcal{D}}$ to be modeled as a mixture of probability distributions for the positive and negative class, i.e.,

$$\mathbf{p}_{\mathcal{D}} = \sum_{y \in \{0,1\}} \mathbf{p}(Y = y)\mathbf{p}(\mathbf{X}|Y = y) = w_{\mathcal{D}_0}\mathbf{p}_{\mathcal{D}_0}(\mathbf{X}) + w_{\mathcal{D}_1}\mathbf{p}_{\mathcal{D}_1}(\mathbf{X}),$$

where $w_{\mathcal{D}_0}$ (resp. $w_{\mathcal{D}_1}$) denotes the marginal probabilities of the label w.r.t the negative (resp. positive) class and $\mathbf{p}_{\mathcal{D}_0}(\mathbf{X})$ (resp. $\mathbf{p}_{\mathcal{D}_1}(\mathbf{X})$) denotes the conditional probability of a sample w.r.t the negative (resp. positive) class. As already said, as it is common practice in PU learning [12], we assume that the positive samples in \mathcal{P} are highly representative for all positive samples in \mathcal{D}_1 . As an additional assumption, we consider the distribution generating \mathcal{D}_0 and \mathcal{D}_1 to be fairly *distinguishable* [1]. That is, we assume that high density regions of $\mathbf{p}_{\mathcal{D}_0}$ correspond to low density regions of $\mathbf{p}_{\mathcal{D}_1}$ and *vice versa*. While this assumption might seem too strict for real data, in practice, it is commonly adopted when performing unsupervised clustering (e.g., Gaussian densities must be separable in EM and K -means). As future research, we plan to investigate how to adapt GPU learning to more complex learning settings.

The high level idea behind our approach is the following. By correctly modeling the probability distribution of positive samples over RVs \mathbf{X} , one can model discriminative patterns among samples in the form of *probabilistic dependencies* among their RVs. If this is done accurately, then a metric space is implicitly defined, associating low probability regions to negative samples and high probability ones to positive samples. Similar ideas have also been successfully investigated in applications for anomalous or outlier training samples [27, 32]. Algorithm 1 illustrates the general schema of our proposed GPU approach. In order to estimate $\mathbf{p}_{\mathcal{P}}$ we fit a generative model, \mathcal{G} , over the RVs \mathbf{X} of the positive training set (line 3). We discuss the choice of such an estimator in Sect. 3.2. After that, we derive an empirical estimation of the less dense (i.e., less likely) regions by computing the point-wise log-likelihood of \mathcal{G} over the samples in \mathcal{U} . Based on this information we build a set of reliable negative samples, denoted as \mathcal{N} (line 7), to be exploited in the second stage of PU learning. As already stated, such a schema is general enough to be adapted to different data domains

Algorithm 1. LearnGPU(\mathcal{P}, \mathcal{U})

-
- 1: **Input:** a set $\mathcal{P} = \{(\mathbf{x}^i, 1)\}_{i=1}^{m_{\mathcal{P}}}$ of positive samples, and a set $\mathcal{U} = \{\mathbf{x}^i\}_{i=1}^{m_{\mathcal{U}}}$ of unlabeled samples over RVs $\mathbf{X} \cup \{Y\}$, with $\text{Val}(Y) = \{0, 1\}$.
 - 2: **Output:** a trained discriminative model learned on positive samples \mathcal{P} and reliable negative samples \mathcal{N} extracted from \mathcal{U}
 - 3: $\mathcal{G} \leftarrow \text{learnGenerativeModel}(\mathcal{P}, \mathbf{X})$ \triangleright learn a generative model \mathcal{G} from \mathcal{P}
 - 4: $\mathcal{L} \leftarrow \{\log p_{\mathcal{G}}(\mathbf{x}^i) | \mathbf{x}^i \in \mathcal{U}\}$
 - 5: $\mathcal{N} \leftarrow \text{reliableNegativeSamples}(\mathcal{L}, \mathcal{P}, \mathcal{U})$
 - 6: $f \leftarrow \text{fitClassifier}(\mathcal{P}, \mathcal{N})$
 - 7: **return** f
-

by leveraging different density estimators. Moreover, by specifying algorithmic variants to build \mathcal{N} and the final discriminator f , one can improve its robustness and accuracy. We discuss such extensions in the following sections.

3.2 Bayesian Networks and Mixtures of Trees

A question arises on which generative model to employ. The main challenge in learning generative models is balancing the *representation expressiveness* of the learned models against the *cost of learning and performing inference* on them.

Probabilistic Graphical Models (PGMs), like Bayesian Networks (BNs) and Markov Networks (MNs), are able to model highly complex probability distributions and have been successfully employed as density estimators. However, exact inference with them is generally *intractable*. Since our GPU learning schema only requires the computation of complete evidence queries, employing BNs in GPU would lead to tractable inference to build \mathcal{N} .

Nevertheless, learning a complex model could still pose a challenge on very large datasets. Guaranteeing exact and tractable inference, a series of *tractable probabilistic models* (TPMs) have been recently proposed: either by restricting the expressiveness of PGMs by bounding their treewidth, or by exploiting local structures in a distribution. The limited expressive capabilities of TPMs, like mixtures of Bayesian trees (MT) [25] and Cutset Networks [9–11], or their ability to compile a high treewidth network into a deep probabilistic architecture, like Sum-Product Networks [31], allow for more efficient learning schemes.

In this work we evaluate GPU by exploiting both BNs and MTs to investigate how the model expressiveness affects the estimation of $p_{\mathcal{P}}$ and therefore ultimately the accuracy of the learned discriminator (see Sect. 4). In the following we briefly review both models.

BNs are a PGM encoding a probability distribution by means of a directed acyclic graph and a set of weights, where nodes correspond to RVs and edges to dependencies among RVs. Given a set of n RVs \mathbf{X} , for each variable $X_i \in \mathbf{X}$, Pa_i denotes the set of parents of node X_i in the DAG. The structure of the BN \mathcal{G} , induces a factorization of the joint distribution into local factors, that is:

$$p_{\mathcal{G}}(\mathbf{X}) = \prod_{i=1}^n p(X_i | \text{Pa}_i).$$

Learning a BN corresponds to learning both the structure and the conditional probability distribution corresponding to each local factor from the data. Classical structure learning algorithms search in the space of BNs guided by a scoring function. On the other hand, parameter learning is obtained by maximum likelihood estimation.

Concerning mixtures of generative models, a very competitive density estimation algorithm is MT [25]. MT learns a mixture model \mathcal{M} whose distribution factorizes according to

$$p_{\mathcal{M}}(\mathbf{X}) = \sum_{i=1}^k \lambda_i p_{\mathcal{T}_i}(\mathbf{X}),$$

where the distributions $p_{\mathcal{T}_i}$, learned using the Chow-Liu algorithm [6], are the mixture components and $\lambda_i \geq 0$, with $\sum_{i=1}^k \lambda_i = 1$ are their coefficients. The Chow-Liu algorithm learns BNs with lower treewidth (i.e., nodes have at most one parent in the network), thus leading to efficient learning and inference time. In [25] the best components and weights are found as (local) likelihood maxima by using EM, with k fixed in advance.

3.3 Reliable Negative Sample Elicitation

After learning a generative model \mathcal{G} , the density estimation information \mathcal{G} provides can be exploited in several ways. The most straightforward one would be to impose a threshold hyperparameter θ such that each sample in \mathcal{U} whose log-likelihood $\log p_{\mathcal{G}}$ falls under θ can be added to \mathcal{N} . However, determining the best value for θ would require to perform additional hyperparameter optimization. To alleviate this issue we propose to implicitly compute it by building \mathcal{N} to comprise the $m_{\mathcal{P}} = |\mathcal{P}|$ samples in \mathcal{U} with the lowest log-likelihood score according to \mathcal{G} . In such a way we ensure that the resulting labeled set $\mathcal{P} \cup \mathcal{N}$ is balanced w.r.t. the positive and negative class. The risk of including positive samples into $\mathcal{P} \cup \mathcal{N}$ can be mitigated by adopting a robust classifier in the following supervised step, whose generalization ability on test data may also additionally benefit from the regularization capability of mis-specifying some sample labels. Lastly, we note how density information in the form of the finite set log-likelihoods can be directly incorporated into the construction of the classifier over $\mathcal{P} \cup \mathcal{N}$.

While we employ the likelihoods to select the most reliable negative samples from \mathcal{U} , they could also be used to select the most reliable positive samples instead. Adopting such a strategy, GPU can be turned into an iterative schema in which at each iteration \mathcal{P} is augmented with the samples belonging to the most dense regions. After a stopping criterion is met, \mathcal{N} can be built by collecting all the samples in \mathcal{U} not added to \mathcal{P} .

3.4 Supervised Classification Step

In principle, every supervised classifier can be employed in GPU after the set \mathcal{N} is constructed. In the empirical evaluation we provide in Sect. 4 we adopt the

regular implementation of Support Vector Machines (SVMs). Nevertheless, we now discuss other interesting variants for GPU learning. First, if one builds \mathcal{N} to be unbalanced w.r.t. \mathcal{P} , it would be possible to adopt the more robust variant of *biased SVMs* [17]. Alternatively, if one focuses on iteratively augmenting the set \mathcal{P} only with GPU, then 1-class SVMs [28] could be employed to derive a max-margin hypersphere for the positive class.

Additionally, the likelihoods associated to samples in \mathcal{U} could be interpreted as sample confidence weights. Approaches like that of [36] could be adopted to learn a *weighted classifier* over $\mathcal{P} \cup \mathcal{U}$ without the need to build \mathcal{N} either.

Lastly, our probabilistic generative approach for the first stage can be plugged in an *unsupervised clustering* approach for the second stage, as done with the EM algorithm in [23]. A principled end-to-end probabilistic formulation would allow estimating both $\mathbf{p}_{\mathcal{D}_0}$ and $\mathbf{p}_{\mathcal{D}_1}$ iteratively and jointly.

4 Experiments

In this section we empirically evaluate the proposed GPU approach, applying it to categorical data. We are interested in this kind of data because they are challenging for classical metric based approaches. Since there is no general consensus on how to build a metric to evaluate categorical data, ad-hoc solutions have been adopted on a domain-wise perspective [19], and only recently PU learning schemes have been devised for it [18]. On the other hand, PGMs have been extensively investigated for categorical data and estimating a probability distribution over discrete RVs is a consolidated practice for extracting new representations in a domain-agnostic unsupervised way [3, 16, 30]. As stated in the previous sections, adapting GPU to other domains reduces to selecting an appropriate generative toolbox from the probabilistic model literature. Specifically, we aim at answering the following research questions: **(Q1)** how does GPU compare to state-of-the-art PU learning approaches? **(Q2)** how does the quantity of available positive samples affect GPU learning? **(Q3)** how much does the choice of a generative model in estimating $\mathbf{p}_{\mathcal{P}}$ affect GPU’s performance?

4.1 Experimental Setting

We took 10 datasets publicly available on the UCI machine learning repository², derived 3 experimental settings for each, and ran 10-fold cross validations exactly as in [18]³. The three settings were generated by putting in \mathcal{P} 30%, 40%, and 50% labeled samples of the positive class respectively, and in \mathcal{U} the remaining positive samples plus all the negative ones. When the dataset does not describe a binary classification problem, the two heavily populated classes were considered. In our experiments, all numerical attributes were discretized into 10 equal-width bins. Detailed dataset statistics are reported in Table 1.

² <http://archive.ics.uci.edu/ml/>.

³ The datasets and settings used in [18] were kindly provided by Dino Ienco.

Table 1. Dataset statistics. #pos and #unl denote the number of positive and unlabeled samples respectively.

Dataset	#attributes	% pos						#test
		30		40		50		
		#pos	#unl	#pos	#unl	#pos	#unl	
Audiology	69	15	79	20	74	26	68	11
Breast-cancer	9	54	203	72	185	91	166	29
Chess	36	451	2425	601	2275	751	2125	320
Dermatology	34	30	136	40	126	50	116	19
Hepatitis	19	9	130	12	127	15	124	16
Lymph	18	17	111	22	106	28	100	14
Nursery	8	1166	6562	1555	6173	1944	5784	859
Pima	8	135	556	180	511	225	466	77
Soybean	35	25	140	33	132	42	123	18
Vote	16	72	319	96	295	120	271	44

We evaluate GPU by employing either BNs (GPU^{BN}) or MTs (GPU^{MT}) as generative models (see Sect. 3.2). BNs are learnt using the R package `bnlearn`⁴ (release 4.1.1). To learn their structure we employed the simple score-based hill-climbing algorithm. In order to avoid overfitting of the network to the positive samples, the following $K2$ scoring function [7] was adopted⁵:

$$\text{score}_{K2}(\mathcal{G} : \mathcal{P}) = \log p(\mathcal{G}) + \sum_{i=1}^n \sum_{j=1}^{q_i} \left(\log \left(\frac{(r_i - 1)!}{(N_{ij} + r_i - 1)!} \right) + \sum_{k=1}^{r_i} \log(N_{ijk}!) \right),$$

where $p(\mathcal{G})$ represents the prior probability of the network G over the n RVs X_i , r_i is the number of states of variable X_i , q_i is the number of possible configurations of the parent set Pa_i , N_{ijk} is the number of instances in the data where variable X_i takes value x_{ik} and the set of variables Pa_i takes value w_{ij} , N_{ij} is the number of instances in the data where the variables in Pa_i take their j -th configuration w_{ij} . Concerning parameter estimation, we set the imaginary sample size to 1. MTs are learnt using the `Libra` [24] toolkit⁶ (version 1.1.2). We imposed the number of components to be 10.

As the classifier for the supervised second stage, we adopt the commonly used SVMs⁷ with an RBF kernel as implemented in `scikit-learn`⁸. The penalty

⁴ <http://www.bnlearn.com/>.

⁵ The same set of experiments have been conducted using the likelihood as scoring function, leading to overfitted models with an overall result worst than that obtained using the $K2$ score.

⁶ <http://libra.cs.uoregon.edu/>.

⁷ For this stage only, categorical data is one-hot encoded.

⁸ <http://scikit-learn.org/>.

parameter C and the kernel coefficient γ have been optimized with a cross validation on the following grid $C \in \{0.001, 0.01, 0.1, 1, 10, 100, 1000\}$ and $\gamma \in \{0.001, 0.01, 0.1, 1, 10, 100, 1000\}$.

We compared GPU with Positive Naive Bayes (PNB), Average Positive Naive Bayes (APNB), Positive TAN (PTAN), Average Positive TAN (APTAN) [4] and Pulce [18] with $k = 7$. See Sect. 2 for a description of these methods.

Source code, in Python and R, of the proposed approach and scripts to reproduce the results are available at <https://github.com/nicoladimauro/GPU>.

4.2 Results and Discussion

Performance on the test set was evaluated using the F1-score measure of the accuracy, defined as $F = 2PR/(P+R) = 2tp/(2tp+fp+fn)$, where P and R are, respectively, the precision and the recall obtained by the algorithm, and tp , fp and fn are, respectively, the true positive, false positive and false negative samples. Since the number of positive samples is much larger than that of negative ones, as in [18] we directed the computation of P , R , and F1-score to the negative samples, differently from their classical setting, i.e., as $F = 2tn/(2tn+fp+fn)$. In particular, since no information for the negative class is provided, correctly predicting negative samples should be somehow harder than focusing on the positive counterparts.

Overall results are reported in Tables 2 and 3. We may note that PTAN and APTAN never won against the other approaches, while the two GPU approaches won 73.3% of the times (53.3% of the times GPU^{BN} alone), and each GPU approach won more times than any competitor (GPU^{BN} more than doubled the number of wins of each competitor). The worst-performing dataset for GPU approaches, and the only one where they perform neatly worse than all other competitors, is ‘hepatitis’. This may indicate that for such a dataset the distributions of the negative and positive class are hard to estimate as very different densities. Concerning question Q1, therefore, we can say that both GPU^{BN} and GPU^{MT} are competitive to the current state-of-the-art for categorical data. On datasets on which GPU^{BN} does not win in all settings, it still performs comparably or better on settings with larger \mathcal{P} sets. Overall, increasing the size of \mathcal{P} improves the models’ accuracy in a consistent way. At the same time, on datasets where both GPU approaches are competitive, they improve over other methods even with only 30% positive samples available (Q2). Lastly, we observe that while GPU^{BN} generally outperforms GPU^{MT}, the latter is still comparable to Pulce (see average ranks, Table 2) and overall more accurate than all other methods. To answer question Q3, we can state that the greater expressiveness of BNs, allowing better modeling the probability distribution of the positive class, is fairly relevant for achieving better performances. Nevertheless, note that for both GPU^{BN} and GPU^{MT} we employed out-of-the-box PGMs and did not invest too much time optimizing the hyperparameters for their structure and weight learning algorithms. It is left for future work to explore how increasing a model complexity can degrade its performance, that is when too accurate probability distribution estimates can lead to overfitting.

Table 2. F1-score results over the 30 samples, comparing GPU^{BN} and GPU^{MT} against the competitors Pulce, PNB, APNB, PTAN, APTAN. The second column indicates the percentage of positive samples in \mathcal{P} .

Dataset	%	GPU ^{BN}	GPU ^{MT}	Pulce	PNB	APNB	PTAN	APTAN
Audiology	30	0.839	0.902	0.745	0.68	0.7	0.66	0.66
Audiology	40	0.879	0.804	0.846	0.75	0.74	0.71	0.66
Audiology	50	0.980	0.991	0.899	0.80	0.80	0.78	0.71
Breast-cancer	30	0.475	0.450	0.534	0.40	0.39	0.43	0.43
Breast-cancer	40	0.483	0.513	0.438	0.42	0.40	0.43	0.45
Breast-cancer	50	0.517	0.535	0.443	0.42	0.41	0.44	0.44
Chess	30	0.689	0.663	0.696	0.58	0.64	0.59	0.64
Chess	40	0.691	0.665	0.688	0.58	0.64	0.60	0.64
Chess	50	0.773	0.650	0.655	0.58	0.64	0.60	0.64
Dermatology	30	1.000	0.834	0.992	0.57	0.57	0.57	0.56
Dermatology	40	1.000	0.836	0.992	0.57	0.58	0.57	0.57
Dermatology	50	0.992	0.951	0.992	0.59	0.60	0.57	0.58
Hepatitis	30	0.822	0.665	0.843	0.87	0.87	0.85	0.86
Hepatitis	40	0.778	0.654	0.873	0.88	0.88	0.85	0.85
Hepatitis	50	0.764	0.742	0.855	0.88	0.88	0.86	0.85
Lymph	30	0.827	0.782	0.851	0.84	0.85	0.79	0.84
Lymph	40	0.825	0.795	0.827	0.84	0.83	0.79	0.81
Lymph	50	0.824	0.835	0.814	0.86	0.87	0.81	0.82
Nursery	30	0.809	0.761	0.739	0.65	0.65	0.56	0.50
Nursery	40	1.000	0.762	0.773	0.69	0.69	0.61	0.56
Nursery	50	0.960	0.779	0.807	0.69	0.70	0.74	0.44
Pima	30	0.588	0.576	0.532	0.49	0.50	0.50	0.50
Pima	40	0.568	0.593	0.547	0.49	0.50	0.50	0.51
Pima	50	0.609	0.605	0.528	0.49	0.51	0.50	0.52
Soybean	30	0.893	0.766	0.738	0.81	0.86	0.80	0.81
Soybean	40	0.883	0.852	0.767	0.86	0.86	0.84	0.83
Soybean	50	0.890	0.923	0.823	0.92	0.92	0.88	0.86
Vote	30	0.826	0.799	0.679	0.62	0.62	0.56	0.55
Vote	40	0.850	0.790	0.800	0.71	0.71	0.58	0.54
Vote	50	0.844	0.829	0.829	0.77	0.77	0.61	0.56
# wins		16	6	3	5	6	0	0
Avg. F1-score		0.796	0.743	0.751	0.677	0.686	0.653	0.643
30%		0.777	0.720	0.735	0.651	0.665	0.631	0.635
40%		0.796	0.727	0.755	0.679	0.683	0.648	0.642
50%		0.815	0.784	0.764	0.700	0.710	0.679	0.652
Avg. ranking		2.16	3.23	3.2	4.57	3.95	5.47	5.42

Table 3. Number of wins/ties among all the methods, the average of the number of wins for each method and its ranking in parenthesis.

	GPU ^{BN}	GPU ^{MT}	Pulce	PNB	APNB	PTAN	APTAN	Avg.
GPU ^{BN}	—	24/0	23/1	23/0	23/0	27/0	25/1	24.17 (1)
GPU ^{MT}	6/0	—	13/0	22/0	22/0	25/0	24/0	18.67 (3)
Pulce	6/0	17/0	—	22/0	22/0	25/0	24/0	19.33 (2)
PNB	7/0	8/0	8/0	—	5/12	18/2	18/3	10.67 (5)
APNB	7/0	8/0	8/0	13/12	—	23/3	21/4	13.33 (4)
PTAN	3/0	5/0	5/0	10/2	4/3	—	12/6	6.50 (7)
APTAN	4/0	6/0	6/0	9/3	5/4	12/6	—	7.00 (6)

Table 4. Number of positive samples incorrectly predicted as negative ones in the negative sample elicitation phase. In parenthesis the average number of errors for each fold over the cardinality of both \mathcal{P} and \mathcal{N} .

Dataset	30%		40%		50%	
Audiology	0.17	(2.6/15)	0.11	(2.3/20)	0.11	(2.8/26)
Breast-cancer	0.48	(25.7/54)	0.46	(33.3/72)	0.43	(38.8/91)
Chess	0.33	(148.6/451)	0.27	(162.4/601)	0.21	(159.7/751)
Dermatology	0.00	(0.0/30)	0.00	(0.0/40)	0.00	(0.0/50)
Hepatitis	0.19	(1.7/9)	0.00	(0.0/12)	0.11	(1.7/15)
Lymph	0.16	(2.8/17)	0.15	(3.4/22)	0.14	(3.8/28)
Nursery	0.00	(0.0/1166)	0.00	(0.0/1555)	0.03	(50.4/1944)
Pima	0.33	(44.4/135)	0.30	(53.6/180)	0.28	(64.0/225)
Soybean	0.14	(3.4/25)	0.15	(4.9/33)	0.10	(4.4/42)
Vote	0.30	(22.9/72)	0.22	(21.0/96)	0.23	(27.9/120)

As already said, in the reliable negative sample elicitation phase, the number of negative samples to be included in the set \mathcal{N} was set to the same number of positive samples available in \mathcal{P} , i.e., $|\mathcal{N}| = |\mathcal{P}|$. However, the generated set \mathcal{N} may contain some true positive samples that have been incorrectly predicted as negative ones. In order to quantify the accuracy of the proposed approach, Table 4 reports, for each dataset, the number of errors occurred in the negative elicitation step, when GPU^{BN} has been used as a density estimator. As we can see, the percentage of errors decreases as the percentage of positive samples in \mathcal{P} increases. For datasets like ‘breast-cancer’, the number of errors reaches 50% thus confirming the low accuracy in terms of F1-score obtained on this dataset. Anyway, also other competitors are not able to properly separate positive and negative samples on this dataset.

A concluding experiment has been done by varying the number of negative samples in \mathcal{N} as a percentage of the number of samples in \mathcal{P} , i.e., by setting

Table 5. Detailed F1-score results over the 30 samples obtained by GPU^{BN} by varying the percentage of the reliable negative samples added to \mathcal{N} . Last column reports Pulce results for comparison.

Dataset	GPU ^{BN}										Pulce
	Negative percentage										
	%	60%	70%	80%	90%	100%	110%	120%	130%	140%	
Audiology	30	0.830	0.832	0.812	0.834	0.839	0.627	0.857	0.857	0.880	0.745
Audiology	40	0.896	0.949	0.940	0.958	0.879	0.938	0.940	0.923	0.931	0.846
Audiology	50	1.000	0.991	0.989	0.971	0.980	0.989	0.966	0.957	0.942	0.899
Breast-cancer	30	0.478	0.422	0.426	0.410	0.475	0.485	0.497	0.464	0.450	0.534
Breast-cancer	40	0.421	0.416	0.409	0.470	0.483	0.494	0.492	0.500	0.452	0.438
Breast-cancer	50	0.457	0.486	0.454	0.508	0.517	0.522	0.497	0.481	0.476	0.443
Chess	30	0.604	0.629	0.644	0.668	0.689	0.693	0.693	0.700	0.700	0.696
Chess	40	0.601	0.637	0.671	0.675	0.691	0.703	0.724	0.731	0.740	0.688
Chess	50	0.623	0.645	0.701	0.739	0.773	0.777	0.788	0.786	0.783	0.655
Dermatology	30	0.992	0.992	0.992	0.992	1.000	0.992	0.992	0.992	0.992	0.992
Dermatology	40	0.992	0.992	0.992	0.992	1.000	1.000	1.000	1.000	1.000	0.992
Dermatology	50	0.992	0.992	0.992	0.992	0.992	0.992	1.000	1.000	0.993	0.992
Hepatitis	30	0.336	0.620	0.605	0.513	0.822	0.822	0.862	0.842	0.876	0.843
Hepatitis	40	0.562	0.611	0.678	0.741	0.778	0.825	0.823	0.871	0.839	0.873
Hepatitis	50	0.679	0.684	0.730	0.695	0.764	0.871	0.859	0.898	0.891	0.855
Lymph	30	0.792	0.810	0.791	0.793	0.827	0.782	0.800	0.829	0.794	0.851
Lymph	40	0.757	0.772	0.781	0.842	0.825	0.823	0.838	0.819	0.849	0.827
Lymph	50	0.721	0.792	0.814	0.841	0.824	0.833	0.816	0.823	0.792	0.814
Nursery	30	0.799	0.810	0.816	0.819	0.809	0.810	0.817	0.818	0.809	0.739
Nursery	40	1.000	1.000	1.000	1.000	1.000	1.000	0.832	0.832	0.832	0.773
Nursery	50	1.000	1.000	1.000	1.000	0.960	1.000	1.000	1.000	1.000	0.807
Pima	30	0.481	0.515	0.535	0.545	0.588	0.543	0.562	0.565	0.567	0.532
Pima	40	0.504	0.509	0.531	0.574	0.568	0.600	0.590	0.607	0.620	0.547
Pima	50	0.519	0.532	0.552	0.603	0.609	0.620	0.613	0.623	0.610	0.528
Soybean	30	0.802	0.834	0.816	0.902	0.893	0.914	0.902	0.915	0.929	0.738
Soybean	40	0.830	0.846	0.893	0.902	0.883	0.915	0.924	0.925	0.926	0.767
Soybean	50	0.819	0.886	0.864	0.854	0.890	0.922	0.916	0.935	0.938	0.823
Vote	30	0.819	0.841	0.838	0.832	0.826	0.798	0.795	0.801	0.779	0.679
Vote	40	0.883	0.850	0.861	0.864	0.850	0.848	0.831	0.814	0.806	0.800
Vote	50	0.892	0.918	0.888	0.854	0.844	0.845	0.846	0.791	0.824	0.829
# wins	12	14	16	20	22	21	25	23	24		
Avg. F1-score	0.736	0.760	0.767	0.779	0.796	0.799	0.802	0.803	0.801	0.751	

$|\mathcal{N}| = \alpha|\mathcal{P}|$. Table 5 reports the results adopting GPU^{BN} as a density estimator for $\alpha \in \{0.6, 0.7, 0.8, 0.9, 1.0, 1.1, 1.2, 1.3, 1.4\}$. It is possible to see that for $\alpha = 1.2$ the number of wins of the proposed approach over Pulce increases to 25.

5 Conclusions

In Positive-Unlabeled (PU) learning only positive samples are labeled at training time. PU learning requires algorithms to cleverly exploit dependencies hidden in the data in order to build models able to discriminate between positive and negative samples. In this paper, we proposed to exploit probabilistic generative models for PU learning by characterizing the density distribution for the positive class. The overall GPU framework is flexible enough to be applied on many domains by leveraging tools provided by PGMs. Results on several benchmark datasets empirically confirmed the validity of our new proposed approach.

References

1. Balasubramanian, V.: MDL, Bayesian inference, and the geometry of the space of probability distributions. In: Grünwald, P.D., Myung, I.J., Pitt, M.A. (eds.) *Advances in Minimum Description Length: Theory and Applications*, pp. 81–98. MIT Press, Cambridge (2005)
2. Basile, T., Mauro, N.D., Esposito, F., Ferilli, S., Vergari, A.: Generative probabilistic models for positive-unlabeled learning. In: *Workshop on NFMCP Held with ECML/PKDD* (2017)
3. Bengio, Y., Courville, A.C., Vincent, P.: Unsupervised feature learning and deep learning: a review and new perspectives. *CoRR* abs/1206.5538 (2012)
4. Calvo, B., Larraaga, P., Lozano, J.A.: Learning bayesian classifiers from positive and unlabeled examples. *Pattern Recogn. Lett.* **28**(16), 2375–2384 (2007)
5. Chandola, V., Banerjee, A., Kumar, V.: Anomaly detection: a survey. *ACM Comput. Surv.* **41**(3), 15:1–15:58 (2009)
6. Chow, C., Liu, C.: Approximating discrete probability distributions with dependence trees. *IEEE Trans. Inf. Theor.* **14**, 462–467 (1968)
7. Cooper, G.F., Herskovits, E.: A Bayesian method for the induction of probabilistic networks from data. *Mach. Learn.* **9**(4), 309–347 (1992)
8. De Comité, F., Denis, F., Gilleron, R., Letouzey, F.: Positive and unlabeled examples help learning. In: Watanabe, O., Yokomori, T. (eds.) *ALT 1999. LNCS (LNAI)*, vol. 1720, pp. 219–230. Springer, Heidelberg (1999). https://doi.org/10.1007/3-540-46769-6_18
9. Di Mauro, N., Vergari, A., Basile, T.M.A., Esposito, F.: Fast and accurate density estimation with extremely randomized cutset networks. In: Ceci, M., Hollmén, J., Todorovski, L., Vens, C., Džeroski, S. (eds.) *ECML PKDD 2017. LNCS (LNAI)*, vol. 10534, pp. 203–219. Springer, Cham (2017). https://doi.org/10.1007/978-3-319-71249-9_13
10. Di Mauro, N., Vergari, A., Basile, T.M.A.: Learning Bayesian random cutset forests. In: Esposito, F., Pivert, O., Hacid, M.-S., Raś, Z.W., Ferilli, S. (eds.) *ISMIS 2015. LNCS (LNAI)*, vol. 9384, pp. 122–132. Springer, Cham (2015). https://doi.org/10.1007/978-3-319-25252-0_13
11. Di Mauro, N., Vergari, A., Esposito, F.: Learning accurate cutset networks by exploiting decomposability. In: Gavanelli, M., Lamma, E., Riguzzi, F. (eds.) *AI*IA 2015. LNCS (LNAI)*, vol. 9336, pp. 221–232. Springer, Cham (2015). https://doi.org/10.1007/978-3-319-24309-2_17
12. Elkan, C., Noto, K.: Learning classifiers from only positive and unlabeled data. In: *KDD*, pp. 213–220 (2008)

13. Friedman, N., Geiger, D., Goldszmidt, M.: Bayesian network classifiers. *Mach. Learn.* **29**(2–3), 131–163 (1997)
14. Hastie, T., Tibshirani, R., Friedman, J.: *The Elements of Statistical Learning*. Springer, New York (2009). <https://doi.org/10.1007/978-0-387-84858-7>
15. Hempstalk, K., Frank, E., Witten, I.H.: One-class classification by combining density and class probability estimation. In: Daelemans, W., Goethals, B., Morik, K. (eds.) *ECML PKDD 2008. LNCS (LNAI)*, vol. 5211, pp. 505–519. Springer, Heidelberg (2008). https://doi.org/10.1007/978-3-540-87479-9_51
16. Hinton, G.E., Salakhutdinov, R.R.: Reducing the dimensionality of data with neural networks. *Science* **313**, 504–507 (2006)
17. Hoi, C.H., Chan, C.H., Huang, K., Lyu, M.R., King, I.: Biased support vector machine for relevance feedback in image retrieval. In: *IJCNN*, pp. 3189–3194 (2004)
18. Ienco, D., Pensa, R.G.: Positive and unlabeled learning in categorical data. *Neurocomputing* **196**, 113–124 (2016)
19. Ienco, D., Pensa, R.G., Meo, R.: From context to distance: learning dissimilarity for categorical data clustering. *TKDD* **6**(1), 1:1–1:25 (2012)
20. Koller, D., Friedman, N.: *Probabilistic Graphical Models: Principles and Techniques*. MIT Press, Cambridge (2009)
21. Li, H., Chen, Z., Liu, B., Wei, X., Shao, J.: Spotting fake reviews via collective positive-unlabeled learning. In: *ICDM*, pp. 899–904 (2014)
22. Liu, B., Dai, Y., Li, X., Lee, W.S., Yu, P.S.: Building text classifiers using positive and unlabeled examples. In: *ICDM*, pp. 179–188 (2003)
23. Liu, B., Lee, W.S., Yu, P.S., Li, X.: Partially supervised classification of text documents. In: *ICML*, pp. 387–394 (2002)
24. Lowd, D., Rooshenas, A.: The Libra toolkit for probabilistic models. *CoRR* abs/1504.00110 (2015)
25. Meila, M., Jordan, M.I.: Learning with mixtures of trees. *JMLR* **1**, 1–48 (2000)
26. du Plessis, M.C., Sugiyama, M.: Semi-supervised learning of class balance under class-prior change by distribution matching. *Neural Netw.* **50**, 110–119 (2014)
27. Riahi, F., Schulte, O., Li, Q.: A proposal for statistical outlier detection in relational structures. In: *SRAI AAAI Workshop* (2014)
28. Schölkopf, B., Platt, J.C., Shawe-Taylor, J.C., Smola, A.J., Williamson, R.C.: Estimating the support of a high-dimensional distribution. *Neural Comput.* **13**(7), 1443–1471 (2001)
29. Tax, D.M.J., Duin, R.P.W.: Support vector data description. *Mach. Learn.* **54**(1), 45–66 (2004)
30. Vergari, A., Di Mauro, N., Esposito, F.: Visualizing and understanding sum-product networks. *CoRR* abs/1608.08266 (2016)
31. Vergari, A., Di Mauro, N., Esposito, F.: Simplifying, regularizing and strengthening sum-product network structure learning. In: Appice, A., Rodrigues, P.P., Santos Costa, V., Gama, J., Jorge, A., Soares, C. (eds.) *ECML PKDD 2015. LNCS (LNAI)*, vol. 9285, pp. 343–358. Springer, Cham (2015). https://doi.org/10.1007/978-3-319-23525-7_21
32. Xu, J., Shelton, C.R.: Intrusion detection using continuous time Bayesian networks. *J. Artif. Int. Res.* **39**(1), 745–774 (2010)
33. Yang, E., Baker, Y., Ravikumar, P., Allen, G., Liu, Z.: Mixed graphical models via exponential families. In: *AISTATS*, pp. 1042–1050 (2014)
34. Yang, P., Li, X.L., Mei, J.P., Kwok, C.K., Ng, S.K.: Positive-unlabeled learning for disease gene identification. *Bioinformatics* **28**, 2640–2647 (2012)
35. Zhao, Y., Kong, X., Philip, S.Y.: Positive and unlabeled learning for graph classification. In: *ICDM*, pp. 962–971 (2011)

36. Zhou, J., Pan, S., Mao, Q., Tsang, I.: Multi-view positive and unlabeled learning. In: ACML, pp. 555–570 (2012)
37. Zhou, K., Gui-Rong, X., Yang, Q., Yu, Y.: Learning with positive and unlabeled examples using topic-sensitive PLSA. TKDE **22**(1), 46–58 (2010)