

1 A bagging SVM to learn from positive and unlabeled  
2 examples

3 F. Mordelet<sup>a,1,\*</sup>, J-P. Vert<sup>b,c,d</sup>

4 <sup>a</sup>Duke University, LSRC building, 308 Research Drive, Durham NC 27708, USA

5 <sup>b</sup>Mines ParisTech, Centre for Computational Biology, 77300 Fontainebleau, France

6 <sup>c</sup>Institut Curie, 75005 Paris, France

7 <sup>d</sup>INSERM, U900, 75005 Paris, France

---

8 **Abstract**

We consider the problem of learning a binary classifier from a training set of positive and unlabeled examples, both in the inductive and in the transductive setting. This problem, often referred to as *PU learning*, differs from the standard supervised classification problem by the lack of negative examples in the training set. It corresponds to an ubiquitous situation in many applications such as information retrieval or gene ranking, when we have identified a set of data of interest sharing a particular property, and we wish to automatically retrieve additional data sharing the same property among a large and easily available pool of unlabeled data. We propose a new method for PU learning with a conceptually simple implementation based on bootstrap aggregating (bagging) techniques: the algorithm iteratively trains many binary classifiers to discriminate the known positive examples from

---

\*Corresponding author  
Preprint submitted to *Pattern Recognition Letters* (F. Mordelet ),  
Email addresses: fantine.mordelet@duke.edu (F. Mordelet ),  
jean-philippe.vert@mines-paristech.fr (J-P. Vert)

May 14, 2013

<sup>1</sup>Phone: +1 919 684-2124

random subsamples of the unlabeled set, and averages their predictions. We show theoretically and experimentally that the method can match and even outperform the performance of state-of-the-art methods for PU learning, particularly when the number of positive examples is limited and the fraction of negatives among the unlabeled examples is small. The proposed method can also run considerably faster than state-of-the-art methods, particularly when the set of unlabeled examples is large.

9 *Keywords:* PU learning, Bagging, SVM

---

## 10 **1. Introduction**

11 In many applications, such as information retrieval or gene ranking, one is  
12 given a finite set of data of interest sharing a particular property, and wishes  
13 to find other data sharing the same property. In information retrieval, for  
14 example, the finite set can be a user query, or a set of documents known  
15 to belong to a specific category, and the goal is to scan a large database of  
16 documents to identify new documents related to the query or belonging to  
17 the same category. In gene ranking, the query is a finite list of genes known  
18 to have a given function or to be associated to a given disease, and the  
19 goal is to identify new genes sharing the same property (Aerts et al., 2006).  
20 In fact this setting is ubiquitous in many applications where identifying a

21 data of interest is difficult or expensive, e.g., because human intervention is  
22 necessary or expensive experiments are needed, while unlabeled data can be  
23 easily collected. In such cases there is a clear opportunity to alleviate the  
24 burden and cost of interesting data identification with the help of machine  
25 learning techniques.

26 More formally, let us assign a binary label to each possible data: positive  
27 (+1) for data of interest, negative (-1) for other data. Unlabeled data  
28 are data for which we do not know whether they are interesting or not.  
29 Denoting  $\mathcal{X}$  the set of data, we assume that the “query” is a finite set of  
30 data  $\mathcal{P} = \{x_1, \dots, x_m\} \subset \mathcal{X}$  with positive labels, and we further assume that  
31 we have access to a (possibly large) set  $\mathcal{U} = \{x_{m+1}, \dots, x_n\} \subset \mathcal{X}$  of unlabeled  
32 data. Our goal is to learn, from  $\mathcal{P}$  and  $\mathcal{U}$ , a way to identify new data with  
33 positive labels, a problem often referred to as *PU learning*. More precisely  
34 we make a distinction between two flavors of PU learning:

- 35 • *Inductive PU learning*, where the goal is to learn from  $\mathcal{P}$  and  $\mathcal{U}$  a  
36 function  $f : \mathcal{X} \rightarrow \mathbb{R}$  able to associate a score or probability to be  
37 positive  $f(x)$  to any new data  $x \in \mathcal{X}$ , which may not be in the training  
38 set of unlabeled data  $\mathcal{U}$ . This may typically be the case in an image  
39 or document classification system, where a subset of the web is used as

40 unlabeled set  $\mathcal{U}$  to train the system, which must then be able to scan  
41 any new image or document out of the training set.

- 42 • *Transductive PU learning*, where the goal is to estimate a scoring func-  
43 tion  $s : \mathcal{U} \rightarrow \mathbb{R}$  from  $\mathcal{P}$  and  $\mathcal{U}$ , i.e., where we are just interested in  
44 finding positive data in the set  $\mathcal{U}$ . This is typically the case in the  
45 disease gene ranking application, where the full set of human genes is  
46 known during training and split between known disease genes  $\mathcal{P}$  and  
47 the rest of the genome  $\mathcal{U}$ . In that case we are only interested in finding  
48 new disease genes in  $\mathcal{U}$ .

49 A growing body of work has focused on PU learning recently. The fact  
50 that only positive and unlabeled examples are available prevents a priori the  
51 use of supervised classification methods, which require negative examples  
52 in the training set. A first approach to overcome the lack of negative ex-  
53 amples is to disregard unlabeled examples during training and simply learn  
54 from the positive examples, e.g., by ranking the unlabeled examples by de-  
55 creasing similarity to the mean positive example (Joachims, 1997) or using  
56 more advanced learning methods such as 1-class SVM (Schölkopf et al., 2001;  
57 Manevitz and Yousef, 2001; Vert and Vert, 2006; De Bie et al., 2007; Geurts,  
58 2011)

59       Alternatively, the problem of inductive PU learning has been studied on  
60 its own from a theoretical viewpoint (Denis et al., 2005; Scott and Blanchard,  
61 2009), and has given rise to a number of specific algorithms. Several authors  
62 have proposed two-step algorithms, heuristic in nature, which first attempt to  
63 identify negative examples in the unlabeled set, and then estimate a classifier  
64 from the positive, unlabeled and likely negative examples (Manevitz and  
65 Yousef, 2001; Liu et al., 2002; Li and Liu, 2003; Liu et al., 2003; Yu et al.,  
66 2004). Alternatively, it was observed that directly learning to discriminate  $\mathcal{P}$   
67 from  $\mathcal{U}$ , possibly after rebalancing the misclassification costs of the two classes  
68 to account for the asymmetry of the problem, leads to state-of-the-art results  
69 for inductive PU learning. This approach has been studied, with different  
70 weighting schemes, using a logistic regression or a SVM as binary classifier  
71 (Liu et al., 2003; Lee and Liu, 2003; Elkan and Noto, 2008). Inductive PU  
72 learning is also related to and has been used for novelty detection, when  
73  $\mathcal{P}$  is interpreted as “normal” data and  $\mathcal{U}$  contains an arbitrary fraction  $\pi$   
74 of negative or “novel” examples (Scott and Blanchard, 2009), or to data  
75 retrieval from a single query, when  $\mathcal{P}$  is reduced to a singleton (Shah et al.,  
76 2008).

77       Transductive PU learning is arguably easier than inductive PU learning,

78 since we know in advance the data to be screened for positive labels. Many  
79 semi-supervised methods have been proposed to tackle transductive learn-  
80 ing when both positive and negative examples are known during training,  
81 including transductive SVM (Joachims, 1999), or many graph-based meth-  
82 ods, reviewed by Chapelle et al. (2006). Comparatively little effort has been  
83 devoted to the specific transductive PU learning problem, with the notable  
84 exception of Liu et al. (2002), who call the problem *partially supervised clas-*  
85 *sification* and proposes an iterative method to solve it, and Pelckmans and  
86 Suykens (2009) who formulate the problem as a combinatorial optimization  
87 problem over a graph. Finally, Sriphaew et al. (2009) recently proposed a  
88 bagging approach which shares similarities with ours, but is more complex  
89 and was only tested on a specific application.

90 Several methods for PU learning, reviewed above, reduce the problem to  
91 a binary classification problem where we learn to discriminate  $\mathcal{P}$  from  $\mathcal{U}$ .  
92 This can be theoretically justified, at least asymptotically, since the ratio  
93 between the conditional distributions of positive and unlabeled examples is  
94 monotonically increasing with the ratio of positive and negative examples  
95 (Elkan and Noto, 2008; Scott and Blanchard, 2009), and has given rise to  
96 state-of-the-art methods such as biased SVM (Liu et al., 2003) or weighted

97 logistic regression (Lee and Liu, 2003). Although this reduction suggests that  
98 virtually any method for (weighted) supervised binary classification can be  
99 used to solve PU learning problems, we put forward in this paper that some  
100 methods may be more adapted than others in a non-asymptotic setting, due  
101 to the particular structure of the unlabeled class. In particular, we inves-  
102 tigate the relevance of methods based on aggregating classifiers trained on  
103 artificially perturbed training sets, in the spirit of bagging (Breiman, 1996,  
104 2001). Such methods are known to be relevant to improve the performance  
105 of unstable classifiers, a situation which, we propose, may occur particularly  
106 in PU learning. Indeed, in addition to the usual instability of learning al-  
107 gorithms confronted to a finite-size training sets, the content of a random  
108 subsample of unlabeled data in positive and negative examples is likely to  
109 strongly affect the classifier, since the contamination of  $\mathcal{U}$  with positive ex-  
110 amples makes the problem more difficult. Variations in the contamination  
111 rate of  $\mathcal{U}$  may thus have an important impact on the trained classifier, in that  
112 a higher contamination rate makes the problem harder in practice (Scott and  
113 Blanchard, 2009), a situation which bagging-like classifiers may benefit from.

114       Based on this idea, we propose a general and simple scheme for inductive  
115 PU learning, akin to an asymmetric form of bagging for supervised binary

116 classification. The method, which we call *bagging SVM*, consists in aggregat-  
117 ing classifiers trained to discriminate  $\mathcal{P}$  from a small random subsample of  $\mathcal{U}$ ,  
118 where the size of the random sample plays a specific role. This method can  
119 naturally be adapted to the transductive PU learning framework. We demon-  
120 strate on simulated and real data that bagging SVM performs at least as well  
121 as existing methods for PU learning, while being often faster in particular  
122 when  $|\mathcal{P}| \ll |\mathcal{U}|$ .

123 This paper is organized as follows. We present and study theoretically  
124 the bagging SVM for inductive PU learning in Section 2.1. Its extension to  
125 transductive PU learning is considered in Section 2.2. Experimental results  
126 are presented in Section 3, followed by a discussion in Section 4.

## 127 **2. Methods**

### 128 *2.1. Bagging for inductive PU learning*

129 Our starting point to learn a classifier in the PU learning setting is the  
130 observation that learning to discriminate positive from unlabeled samples is  
131 a good proxy to our objective, which is to discriminate positive from negative  
132 samples. Even though the unlabeled set is contaminated by hidden positive  
133 examples, it is generally admitted that its distribution contains some infor-



134 mation which should be exploited. That is for instance, the foundation of  
135 semi-supervised methods.

136 Indeed, let us assume for example that positive and negative examples  
137 are randomly generated by class-conditional distributions  $\mathbb{P}_+$  and  $\mathbb{P}_-$  with  
138 densities  $h_+$  and  $h_-$ . If we model unlabeled examples as randomly sampled  
139 from  $\mathbb{P}_+$  with probability  $\gamma$  and from  $\mathbb{P}_-$  with probability  $1 - \gamma$ , then the  
140 distribution of unlabeled has a density

$$h_u = \gamma h_+ + (1 - \gamma) h_- . \quad (1)$$

141 Now notice that

$$\frac{h_u(x)}{h_+(x)} = \gamma + (1 - \gamma) \frac{h_-(x)}{h_+(x)} , \quad (2)$$

142 showing that the ratio between the conditional distributions of positive and  
143 unlabeled examples is monotonically increasing with the ratio of positive  
144 and negative examples (Elkan and Noto, 2008; Scott and Blanchard, 2009).  
145 Hence any estimator of the conditional probability of positive vs. unlabeled  
146 data should in theory also be applicable to discriminate positive from neg-  
147 ative examples. This is the case for example of logistic regression or some  
148 forms of SVM (Steinwart, 2003; Bartlett and Tewari, 2007). In practice it  
149 seems useful to train classifiers to discriminate  $\mathcal{P}$  from  $\mathcal{U}$  by penalizing more  
150 false negative than false positive errors, in order to account for the fact that

151 positive examples are known to be positive, while unlabeled examples are  
152 known to contain hidden positives. Using soft margin SVM while giving  
153 high weights to false negative errors and low weights to false positive er-  
154 rors leads to the biased SVM approach described by Liu et al. (2003), while  
155 the same strategy using a logistic regression leads to the weighted logistic  
156 regression approach of Lee and Liu (2003). Both methods, tested on text  
157 categorization benchmarks, were shown to be very efficient in practice, and  
158 in particular outperformed all approaches based on heuristic identifications  
159 of true negatives in  $\mathcal{U}$ .

160     Among the many methods for supervised binary classification which could  
161 be used to discriminate  $\mathcal{P}$  from  $\mathcal{U}$ , bootstrap aggregating or “bagging” is an  
162 interesting candidate (Breiman, 1996). The idea of bagging is to estimate a  
163 series of classifiers on datasets obtained by perturbing the original training  
164 set through bootstrap resampling, and to combine these classifiers by some  
165 aggregation technique. The method is conceptually simple, can be applied in  
166 many settings, and works very well in practice (Breiman, 2001; Hastie et al.,  
167 2001). Bagging generally improves the performance of individual classifiers  
168 when they are not too correlated to each other, which happens in particular  
169 when the classifier is highly sensitive to small perturbations of the training

170 set. For example, Breiman (2001) showed that the difference between the  
171 expected mean square error (MSE) of a classifier trained on a single bootstrap  
172 sample and the MSE of the aggregated predictor increases with the variance  
173 of the classifier.

174 We propose that, by nature, PU learning problems have a particular  
175 structure that leads to instability of classifiers, which can be advantageously  
176 exploited by a bagging-like procedure which we now describe. Intuitively,  
177 an important source of instability in PU learning situations is the empirical  
178 contamination  $\hat{\gamma}$  of  $\mathcal{U}$  with positive examples, i.e., the percentage of positive  
179 examples in  $\mathcal{U}$  which on average equals  $\gamma$  in (1). If by chance  $\mathcal{U}$  is mostly  
180 made of negative examples, i.e., has low contamination by positive examples,  
181 then we will probably estimate a better classifier than if it contains mostly  
182 positive examples, i.e., has high contamination. Moreover, we can expect  
183 the classifiers in these different scenarios to be little correlated, since intu-  
184 itively they estimate different log-ratios of conditional distribution. Hence,  
185 in addition to the “normal” instability of a classifier trained on a finite-size  
186 sample, which is exploited by bagging in general, we can expect an increased  
187 instability in PU learning due to the sensitivity of the classifier to the em-  
188 pirical contamination  $\hat{\gamma}$  of  $\mathcal{U}$  in positive examples. In order to exploit this

189 sensitivity in a bagging-like procedure, we propose to randomly subsample  
190  $\mathcal{U}$  and train classifiers to discriminate  $\mathcal{P}$  from each subsample, before ag-  
191 gregating the classifiers. By subsampling  $\mathcal{U}$ , we hope to vary in particular  
192 the empirical contamination between samples. This will induce a variety of  
193 situations, some lucky (small contamination), some less lucky (large contam-  
194 ination), which eventually will induce a large variability in the classifiers that  
195 the aggregation procedure can then exploit.

196 In opposition to classical bagging, the size  $K$  of the samples generated  
197 from  $\mathcal{U}$  may play an important role to balance the accuracy against the  
198 stability of individual classifiers. On the one hand, larger subsamples should  
199 lead on average to better classifiers, since any classification method generally  
200 improves on average when more training points are available. On the other  
201 hand, the empirical contamination varies more for smaller subsamples.

To formalize a bit more this line of thought, let us denote by  $\hat{\gamma}$  the true contamination rate in  $\mathcal{U}$ , that is, the true proportion of positive examples hidden in  $\mathcal{U}$ . Whenever a bootstrap sample  $\mathcal{U}_t$  of size  $K$  is drawn from  $\mathcal{U}$ , its empirical number of positive examples is a binomial random variable  $\sim B(K, \hat{\gamma})$ , leading to a contamination rate  $\hat{\gamma}_t$  with mean and variance:

$$\mathbb{E}(\hat{\gamma}_t) = \hat{\gamma} \text{ and } \mathbb{V}(\hat{\gamma}_t) = \frac{1}{K} \hat{\gamma}(1 - \hat{\gamma}).$$

The intuition that less contaminated samples allow to estimate better classifiers can be formalized in many different ways. Here we follow the analysis of Scott and Blanchard (2009) who study the inductive PU learning framework and consider the setting where we want to learn a function  $f : \mathcal{X} \rightarrow \mathbb{R}$  which has a small probability of predicting negative examples as positives  $R_-(f) = \mathbb{P}_-(f(X) > 0)$  for a probability of predicting positive examples as negatives  $R_+(f) = \mathbb{P}_+(f(X) < 0)$  bounded by a fixed level  $\alpha > 0$ . For any level  $\alpha > 0$ , denoting by  $R_{-, \alpha}^* = \inf_{f: R_+(f) \leq \alpha} R_-(f)$  the smallest possible risk, Scott and Blanchard (2009, Theorem 2) show that the excess risk of a function  $\hat{f}_t$  trained to discriminate a set  $\mathcal{P}$  of size  $P$  of positive examples from a subsample  $\mathcal{U}_t$  of size  $K$  and contamination  $\hat{\gamma}_t$  is upper bounded with large probability  $\delta$  as follows:

$$L_{-, \alpha}(\hat{f}_t) = R_-(\hat{f}_t) - R_{-, \alpha}^* \leq \frac{\varepsilon_P + \varepsilon_K}{1 - \hat{\gamma}_t},$$

202 where  $\varepsilon_i$  is an upper bound on the excess risk due to a finite sample of size  $i$ ,  
 203 typically proportional to  $i^{-1/2}$  for a classifier trained to discriminate  $\mathcal{P}$  from  
 204  $\mathcal{U}_t$  by empirical risk minimization on a finite set (Scott and Nowak, 2005).  
 205 This leads to an upper bound of the excess risk of the form:

$$L_{-, \alpha}(\hat{f}_t) \leq c \frac{P^{-1/2} + K^{-1/2}}{1 - \hat{\gamma}_t}, \quad (3)$$

206 where  $c$  is a constant. Details relating constant  $c$  and the probability  $\delta$  can be  
 207 found in Scott and Blanchard (2009, Theorem 2). Equation 3 shows that the  
 208 quality of the estimator increases when the size of the unlabeled sample  $K$  in-  
 209 creases and its contamination  $\hat{\gamma}_t$  decreases, as expected. When we aggregate  
 210 different classifiers  $\hat{f}_t$  trained on subsamples with varying contamination  $\hat{\gamma}_t$ ,  
 211 we can expect that the excess risk of the aggregated classifier reaches the per-  
 212 formance of individual classifiers with smaller-than-average contamination,  
 213 typically with contamination  $\hat{\gamma} - c_2\sqrt{\hat{\gamma}(1-\hat{\gamma})/K}$  where  $c_2 > 0$  is a constant  
 214 (independent of  $c$ ). Plugging this estimate into (3), we obtain that the excess  
 215 risk of the aggregated classifier is upper bounded with large probability by

$$c \frac{P^{-1/2} + K^{-1/2}}{1 - \hat{\gamma} + c_2\sqrt{\hat{\gamma}(1-\hat{\gamma})}K^{-1/2}}. \quad (4)$$

216 Now we see that, when  $P > \hat{\gamma}c_2^2/(1-\hat{\gamma})$ , the upper bound (4) is a decreasing  
 217 function of  $K$  and there is no apparent gain in subsampling and aggregating  
 218 with  $K < N$ . On the other hand, when  $P < \hat{\gamma}c_2^2/(1-\hat{\gamma})$ , the upper bound  
 219 (4) is an *increasing* function of  $K$ , suggesting that choosing  $K < N$  may  
 220 lead to more accurate classifiers. In words, when  $P$  is not too large, it may  
 221 be better to subsample the set of unlabeled examples with  $K < N$  samples  
 222 and aggregated the resulting classifiers, because the gain in performance due  
 223 to the stochastic decrease in contamination in a fraction of subsamples can

224 be exploited by aggregation and outperforms the loss due the fact that each  
225 classifier is trained on a smaller data set.

226 In summary, the method we propose for PU learning is presented in Al-  
227 gorithm 1. It creates a series of classifiers trained to discriminate  $\mathcal{P}$  from  
228 random subsamples of  $\mathcal{U}$ . The output of each of these classifiers is a function  
229  $f_t$  that assigns a prediction score to any example. The score function  $f$  of  
230 the final aggregated classifier is simply defined as the average score of the  
individual classifiers. We call it bagging SVM when the classifier used to

---

**Algorithm 1** Inductive bagging PU learning

---

INPUT :  $\mathcal{P}, \mathcal{U}, K = \text{size of bootstrap samples}, T = \text{number of bootstraps}$

OUTPUT : a function  $f : \mathcal{X} \rightarrow \mathbb{R}$

**for**  $t = 1$  to  $T$  **do**

    Draw a subsample  $\mathcal{U}_t$  of size  $K$  from  $\mathcal{U}$ .

    Train a classifier  $f_t$  to discriminate  $\mathcal{P}$  against  $\mathcal{U}_t$ .

**end for**

Return

$$f = \frac{1}{T} \sum_{t=1}^T f_t$$

---

231

232 discriminate  $\mathcal{P}$  from a random subsample of  $\mathcal{U}$  is a biased SVM. It is akin to  
233 bagging to learn to discriminate  $\mathcal{P}$  from  $\mathcal{U}$ , with two important specificities.

234 First, only  $\mathcal{U}$  is subsampled. This is to account for the fact that elements in  
235  $\mathcal{P}$  are known to be positive, and moreover that the number of positive exam-  
236 ples is often limited. Second, the size of subsamples is a parameter  $K$  whose  
237 effect needs to be studied. If an optimal value exists, then this parameter  
238 may need to be adjusted.

239 The number  $T$  of bootstrap samples is also a user-defined parameter.  
240 Intuitively, the larger  $T$  the better, although we observed empirically little  
241 improvement for  $T$  larger than 100 (see Section 3, Figure 3). Finally, al-  
242 though we propose to aggregate the  $T$  classifiers by a simple average, other  
243 aggregation rules could easily be used. On preliminary experiments on sim-  
244 ulated and real data, we did not observed significant differences between the  
245 simple average and majority voting, another popular aggregation method.

## 246 *2.2. Bagging SVM for transductive PU learning*

247 We now consider the situation where the goal is only to assign a score to  
248 the elements of  $\mathcal{U}$  reflecting our confidence that these elements belong to the  
249 positive class. Liu et al. (2002) have studied the same problem which they  
250 call “partially supervised classification”. Their proposed technique combines  
251 Naive Bayes classification and the Expectation-Maximization algorithm to  
252 iteratively produce classifiers. The training scores of these classifiers are



253 then directly used to rank  $\mathcal{U}$ . Following this approach, a straightforward  
254 solution to the transductive PU learning problem is to train any classifier  
255 to discriminate between  $\mathcal{P}$  and  $\mathcal{U}$  and to use this classifier to assign a score  
256 to the unlabeled data that were used to train it. Using SVMs this amounts  
257 to using the biased SVM training scores. We will subsequently denote this  
258 approach by transductive biased SVM.

259       However, one may argue that assigning a score to an unlabeled example  
260 that has been used as negative training example is problematic. In partic-  
261 ular, if the classifier fits too tightly to the training data, a false negative  
262  $x_i$  will hardly be given a high training score when used as a negative. In  
263 a related situation in the context of semi-supervised learning, Zhang et al.  
264 (2009) showed for example that unlabeled examples used as negative training  
265 examples tend to have underestimated scores when an SVM is trained with  
266 the classical hinge loss. More generally, most theoretical consistency prop-  
267 erties of machine learning algorithms justify predictions on samples outside  
268 of the training set, raising questions on the use of all unlabeled samples as  
269 negative training samples at the same time.

270       Alternatively, the inductive bagging PU learning lends itself particularly  
271 well to the transductive setting, through the procedure described in Algo-

272 rithm 2. Each time a random subsample  $\mathcal{U}_t$  of  $\mathcal{U}$  is generated, a classifier is  
273 trained to discriminate  $\mathcal{P}$  from  $\mathcal{U}_t$ , and used to assign a predictive score to  
274 any element of  $\mathcal{U} \setminus \mathcal{U}_t$ . At the end the score of any element  $x \in \mathcal{U}$  is obtained  
275 by aggregating the predictions of the classifiers trained on subsamples that  
276 did not contain  $x$  (the counter  $n(x)$  simply counts the number of such classi-  
277 fiers). As such, no point of  $\mathcal{U}$  is used simultaneously to train a classifier and  
278 to test it. In practice, it is useful to ensure that we average the predictions  
279 over a sufficient number of classifiers. Typically, if we wish to average over  
280  $n$  scores, we need to choose  $T$  such that  $T(1 - \frac{K}{|\mathcal{U}|}) \gg n$

### 281 3. Results

282 In this section we investigate the empirical behavior of our bagging algo-  
283 rithm on one simulated dataset (Section 3.1) and two real applications: text  
284 retrieval with the 20 newsgroup benchmark (Section 3.2), and reconstruc-  
285 tion of gene regulatory networks (Section 3.3). We compare the new bagging  
286 SVM to the state-of-the-art biased SVM, and also add in the comparison for  
287 real data two one-class approaches, namely, ranking unlabeled examples by  
288 decreasing mean similarity to the positive examples (called *Baseline* below),  
289 and the one-class SVM (Schölkopf et al., 2001). The biased SVM consists

---

**Algorithm 2** Transductive bagging PU learning

---

INPUT :  $\mathcal{P}, \mathcal{U}, K = \text{size of bootstrap samples}, T = \text{number of bootstraps}$

OUTPUT : a score  $s : \mathcal{U} \rightarrow \mathbb{R}$

Initialize  $\forall x \in \mathcal{U}, n(x) \leftarrow 0, f(x) \leftarrow 0$

**for**  $t = 1$  to  $T$  **do**

    Draw a bootstrap sample  $\mathcal{U}_t$  of size  $K$  in  $\mathcal{U}$ .

    Train a classifier  $f_t$  to discriminate  $\mathcal{P}$  against  $\mathcal{U}_t$ .

    For any  $x \in \mathcal{U} \setminus \mathcal{U}_t$ , update:

$$f(x) \leftarrow f(x) + f_t(x),$$

$$n(x) \leftarrow n(x) + 1.$$

**end for**

Return  $s(x) = f(x)/n(x)$  for  $x \in \mathcal{U}$

---

290 in training a soft margin SVM to discriminate between  $\mathcal{P}$  and  $\mathcal{U}$ . Both bag-  
 291 ging and biased methods involve an SVM with asymmetric penalties  $C_+$  and  
 292  $C_-$  for the positive and negative class, respectively. By default we always  
 293 set them to ensure that the total penalty is equal for the two classes, i.e.,  
 294  $C_+n_+ = C_-n_-$ , where  $n_+$  and  $n_-$  are the number of positive and negative  
 295 examples fed to the SVM, and optimized the single parameter  $C = C_+ + C_-$   
 296 over a grid. We checked on all experiments that this choice was never signif-  
 297 icantly outperformed by another penalty ratio  $C_+/C_-$ .

298 All methods were implemented in MATLAB, using the LIBSVM software  
 299 (Chang and Lin, 2011) to train one- and two-class SVM. All experiments were  
 300 run under Linux on a machine with two 4-core Intel Xeon 3.16GHz processors  
 301 and 16Gb of RAM.

### 302 3.1. Simulated data

303 A first series of experiments were conducted on simulated data to compare  
 304 our bagging procedure to the biased approach in an inductive setting. We  
 305 consider the simple situation where the positive examples are generated from  
 306 an isotropic Gaussian distribution in  $\mathbb{R}^p$  :  $\mathcal{P} \sim \mathbb{P}_+ = \mathcal{N}(0_p, \sigma * I_p)$ , with  
 307  $p = 50$  and  $\sigma = 0.6$ , while the negative examples are generated from another  
 308 Gaussian distribution with same isotropic covariance and a different mean,

309 of norm 1. We replicate the following iteration 50 times for different values  
310 of  $\gamma$  :

- 311 • Draw a sample  $\mathcal{P}$  of 5 positives examples, and a sample  $\mathcal{U}$  of 50 unlabeled examples from  $\gamma * \mathbb{P}_+ + (1 - \gamma) * \mathbb{P}_-$ .
- 312
- 313 • Train respectively the biased and bagging logit (with 200 bootstraps)<sup>2</sup>.
- 314 • Compare their performance on a test set of 1000 examples containing  
315 50% positives.

316 For  $K$ , we tested equally spaced values between 1 and 50, and we varied  
317  $\gamma$  on the interval  $[0; 0.8]$ . The performance is measured by computing the  
318 area under the Receiving Operator Characteristic curve (AUC) on the inde-  
319 pendent test set. Figure 1 (left) shows the performance of bagging logit for  
320 different levels of contamination of  $\mathcal{U}$ , as a function of  $K$ , the size of the ran-  
321 dom samples. The uppermost curve thus corresponds to  $\gamma = 0$ , i.e., the case  
322 where all unlabeled data are negative, while the bottom curve corresponds  
323 to  $\gamma = 0.8$ , i.e., the case where 80% of unlabeled data are positive. Note that  
324  $K = 50$  corresponds to classical bagging on the biased logit classifier, i.e., to

---

<sup>2</sup>The bagging logit corresponds to the procedure described above, when the classifier is a logistic regression. This is the same for the biased logit, see also (Lee and Liu, 2003)

325 the case where all unlabeled examples are used to train the classifier.

326 We observe that in the classical setting of supervised binary classifica-  
327 tion where  $\mathcal{U}$  is not contaminated by positive samples ( $\gamma = 0$ ), the bagging  
328 procedure does not improve performance, whatever the size of the bootstrap  
329 samples. On the other hand, as contamination increases, we observe an over-  
330 all decrease of the performance, confirming that the classification problem  
331 becomes more difficult when contamination increases. In addition, the bag-  
332 ging logit always succeeds in reaching at least the same performance for a  
333 value of  $K$  below 50, even for high rates of contamination. Figure 1 (right)  
334 shows the evolution of AUC as  $\gamma$  increases, for both methods. For the bag-  
335 ging logit we report the AUC reached for the best  $K$  value. We see that  
336 bagging logit slightly outperforms the biased logit method.

337 To further illustrate the assumption that motivated bagging SVM, namely  
338 that decreasing  $K$  would decrease the average performance of single classifiers  
339 but would increase their variance due to the variations in contamination, we  
340 take a closer look at the successive classifiers learnt when training Algorithm  
341 1. Each classifier corresponds to a random bootstrap subsample  $\mathcal{U}_t$ . We  
342 show in Figure 2 a scatter plot of the AUC of these individual classifiers as a  
343 function of  $\hat{\gamma}$ , the empirical contamination of the bootstrap sample  $\mathcal{U}_t$ , for two

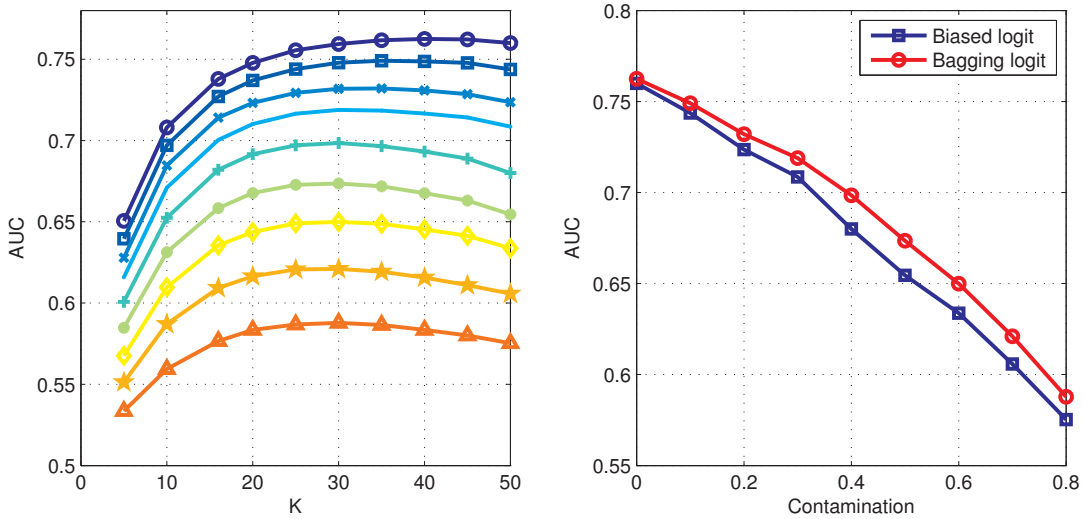


Figure 1: Results on simulated data. *Left:* AUC of the bagging logit as a function of  $K$ , the size of the bootstrap samples, on simulated data. Each curve, from top to bottom, corresponds to a contamination level  $\gamma \in \{0; 0.1; 0.2; \dots; 0.8\}$ . *Right* Performance of two methods as a function of  $\gamma$ , the contamination level, on simulated data. The performance of bagging logit was taken at the optimal  $K$  value.

344 values of  $K$  (10 and 40). Here the mean contamination was set to  $\gamma = 0.2$ .  
 345 Obviously, the variations of  $\hat{\gamma}$  are much larger for  $K = 10$  (between 0 and 0.5)  
 346 than for  $K = 40$  (between 0.1 and 0.25). The correlation coefficient between  
 347  $\hat{\gamma}$  and the performance (reported above each plot) is strongly negative, in  
 348 particular for smaller  $K$ . It is quite clear that less contaminated subsamples  
 349 tend to yield better classifiers, and that the variation in the contamination is  
 350 an important factor to increase the variance between individual predictors,

which aggregation can benefit from.

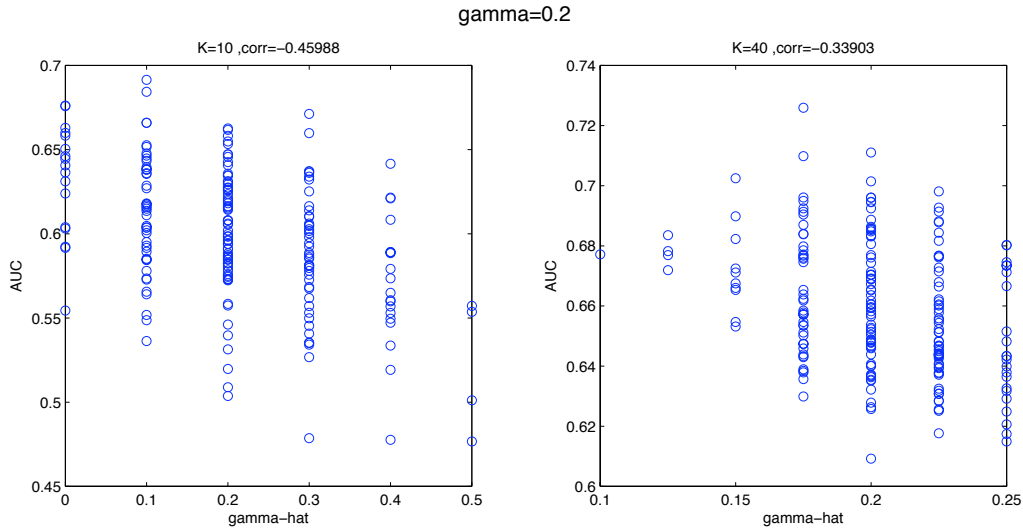


Figure 2: Distribution of AUC and  $\hat{\gamma}$  over the 500 iterations of one bootstrap loop on the simulated dataset,  $\gamma = 0.2$ .

351

### 352 3.2. *Newsgroup dataset*

353 The 20 Newsgroup benchmark is widely used to test PU learning methods.

354 The version we used is a collection of 11293 articles partitioned into 20 sub-

355 sets of roughly the same size (around 500)<sup>3</sup>, corresponding to post articles of

356 related interest. For each newsgroup, the positive class consists of those  $\sim 500$

---

<sup>3</sup>We used the Matlab pre-processed version available at <http://renatocorrea.googlepages.com/ng2011293x8165itrn.mat>



357 articles known to be relevant, while the negative class is made of the remain-  
358 der. After pre-processing, each article is represented by a 8165-dimensional  
359 feature vector, based on word counts, using the TFIDF representation over  
360 a dictionary of 8165 words (Joachims, 1997).

361 To simulate a PU learning problem, we applied the following strategy.  
362 For a given newsgroup, we created a set  $\mathcal{P}$  of known positive examples by  
363 randomly selecting a given number of positive examples, while  $\mathcal{U}$  contains  
364 the non-selected positive examples and all negative examples. We varied  
365 the size  $N_P$  of  $\mathcal{P}$  in  $\{5, 10, 20, 50, 100, 200, 300\}$  to investigate the influence  
366 of the number of known positive examples. For each newsgroup and each  
367 value of  $N_P$ , we train all 4 methods described above (bagging SVM, biased  
368 SVM, baseline, one-class SVM) and rank the samples in  $\mathcal{U}$  by decreasing score  
369 (transductive setting). We then compute the AUC, and average this measure  
370 over 10 replicates of each newsgroup and each value of  $N_P$ . For bagging and  
371 biased SVM, we varied the  $C$  parameter over the grid  $\{e^{-12}, e^{-10}, \dots, e^2\}$ ,  
372 while we vary parameter  $\nu$  in  $\{0.1, 0.2, \dots, 0.9\}$  for 1-class SVM. We only  
373 used the linear kernel.

374 We first investigated the influence of  $T$ . Figure 3 shows, for the first  
375 newsgroup (alt.atheism), the performance reached as a function of  $T$ , for

376 different settings in  $N_P$  and  $K$ . As expected we observe that in general the  
377 performance increases with  $T$ , but quickly reaches a plateau beyond which  
378 additional bootstraps do not improve performance. Overall the smaller  $K$ ,  
379 the larger  $T$  must be to reach the plateau. From these preliminary results  
380 we set  $T = 35$  for  $K \leq 20$ , and  $T = 10$  for  $K > 30$ , and kept it fix for the  
381 rest of the experiments. To further clarify the benefits of bagging, we show  
382 in Figure 4 the performance of the bagging SVM versus the performance of a  
383 SVM trained on a single bootstrap sample ( $T = 1$ ), for different values of  $K$   
384 and a fixed number of positives  $N_P = 10$ . We observe that, for  $K$  below 200,  
385 aggregating classifiers over several bootstrap subsamples is clearly beneficial,  
386 while for larger values of  $K$  it does not really help. This is coherent with the  
387 observation that SVM usually rarely benefit from bagging: here the benefits  
388 come from our particular bagging scheme. Interestingly, we see that very  
389 good performance is reached even for small values of  $K$  with the bagging.

390 Figure 5 shows the mean AUC averaged over the 10 folds and the 20  
391 newsgroups for bagging SVM as a function of  $K$ , and compares it to that  
392 of the biased SVM. More precisely, each point on the curve corresponds to  
393 the performance averaged over the 20 Newsgroups after choosing a posteriori  
394 the best  $C$  parameter for each newsgroup. This is equivalent to comparing

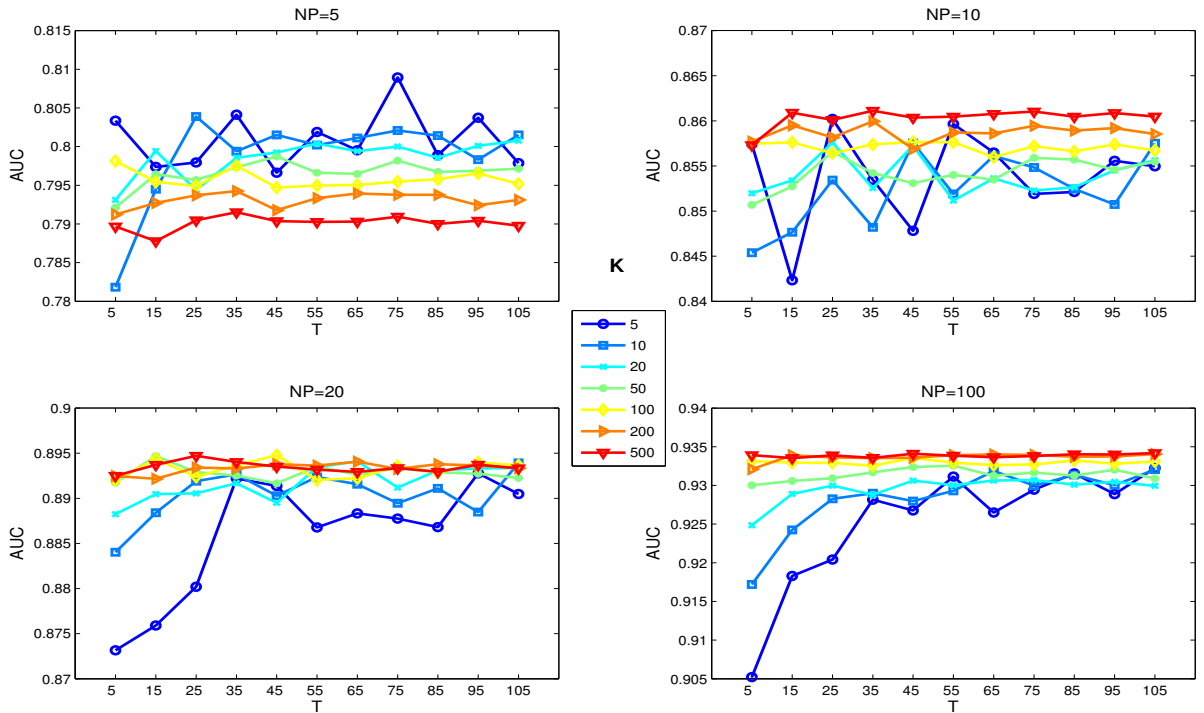


Figure 3: Performance on newsgroup 1 (alt.atheism) as a function of the number of bootstraps  $T$ , for different values of  $N_P$  and  $K$ .

395 optimal cases for both methods. Contrary to what we observed on simulated  
 396 data, we observe that  $K$  has in general very little influence on the perfor-  
 397 mance. The AUC of the bagging SVM is similar to that of the biased SVM  
 398 for most values of  $K$ , although for  $N_P$  larger than 50, a slight advantage can  
 399 be observed for the biased SVM over bagging SVM when  $K$  is too small.  
 400 We conclude that in practice, parameter  $K$  may not need to be finely tuned  
 401 since it does not always have a big impact on the performance. In all cases,

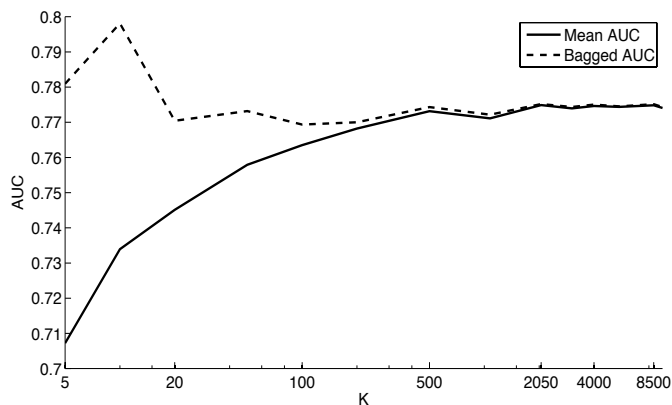


Figure 4: Performance on one newgroup of bagging SVM (*bagged AUC*) vs a SVM trained on a single bootstrap sample (*mean AUC*), for different values of  $K$ .

402  $K = N_P$  seems to be a safe choice for the bagging SVM.

403 Finally, Figure 6 shows the average AUC over the 20 newsgroups for all  
 404 four methods, as a function of  $N_P$ . Overall all methods are very similar,  
 405 with the Baseline slightly below the others. In details, the bagging SVM  
 406 curve dominates all other methods for  $N_P \geq 20$ , while the 1-class SVM is  
 407 the one which dominates for smaller values of  $N_P$ . Although the differences  
 408 in performance are small, the bagging SVM outperforms the biased SVM  
 409 significantly for  $N_P > 20$  according to a Wilcoxon paired sample test (at 5%  
 410 confidence). For small values of  $N_P$  however, no significant difference can be  
 411 proven in either way between bagging SVM and 1-class SVM, which remains  
 412 a very competitive method.

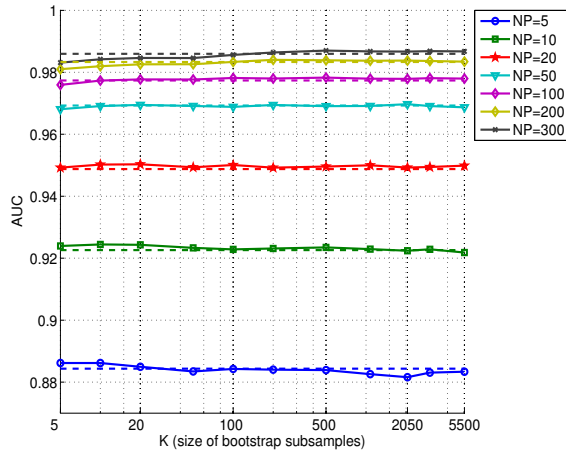


Figure 5: Macro averaged performance of the bagging SVM as a function of  $K$ . The dashed horizontal lines show the AUC level of the biased SVM. The curves are plotted for different values of  $N_P$ , the size of the positive set.

### 413 3.3. *E. coli* dataset : inference of transcriptional regulatory network

414 In this section we test the different PU learning strategies on the problem  
 415 of inferring the transcription regulatory network of the bacteria *Escherichia*  
 416 *coli* from gene expression data. The problem is, given a transcription fac-  
 417 tor (TF), to predict which genes it regulates. Following Mordelet and Vert  
 418 (2008), we can formulate this problem as transductive PU learning by start-  
 419 ing from known regulated genes (considered positive examples), and looking  
 420 for additional regulated genes in the bacteria’s genome.

421 To represent the genes, we use a compendium of microarray expression  
 422 profiles provided by Faith et al. (2008), in which 4345 genes of the *E. Coli*

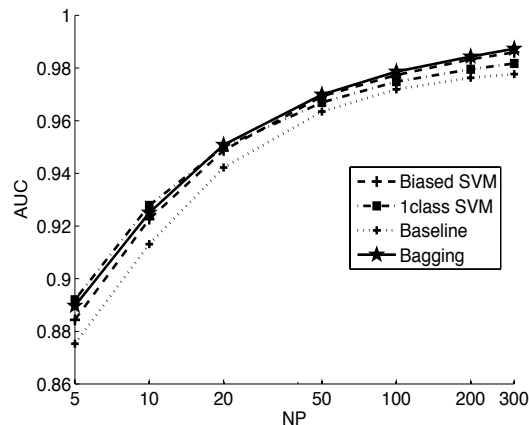


Figure 6: Performance of the baseline method, the 1-class SVM, the biased SVM and the newly proposed bagging SVM methods on the 20 Newsgroups dataset. Each curve shows how the mean AUC varies with the number of positive training examples  $N_P$ . For each value of  $N_P$ , the performance of bagging SVM is computed at the optimal value for  $K$ , as shown in Figure 5.

423 genome are represented by vectors in 445 dimensions, corresponding to their  
 424 expression level in 445 different experiments. We extracted the list of known  
 425 regulated genes for each TF from the RegulonDB (Salgado et al., 2006).

426 For each TF, we ran a double 3-fold cross validation with an internal loop  
 427 on each training set to select parameter  $C$  of the SVM (or  $\nu$  for the 1-class  
 428 SVM). To make this possible, we restrict ourselves to 31 TFs with at least  
 429 8 known regulated genes. Following Mordelet and Vert (2008), we normalize  
 430 the expression data to unit norm, use a Gaussian RBF kernel with  $\sigma = 8$ , and

431 perform a particular cross-validation scheme to ensure that operons are not  
 432 split between folds. Finally, following our previous results on simulated data  
 433 and the newsgroup benchmark, we test two variants of bagging SVM, setting  
 434  $K$  successively to  $N_P$  and  $5 * N_P$ . These choices are denoted respectively by  
 435 *bagging1 SVM* and *bagging5 SVM*.

436 Figure 7 shows the average precision/recall curves of all methods tested.  
 437 Overall we observe that all three PU learning methods give significantly bet-  
 438 ter results than the two methods which use only positive examples (Wilcoxon  
 439 paired sample test at 5% significance level). No significant difference was  
 440 found between the three PU learning methods. This confirms again that for  
 441 different values of  $K$  bagging SVM matches the performance of biased SVM.

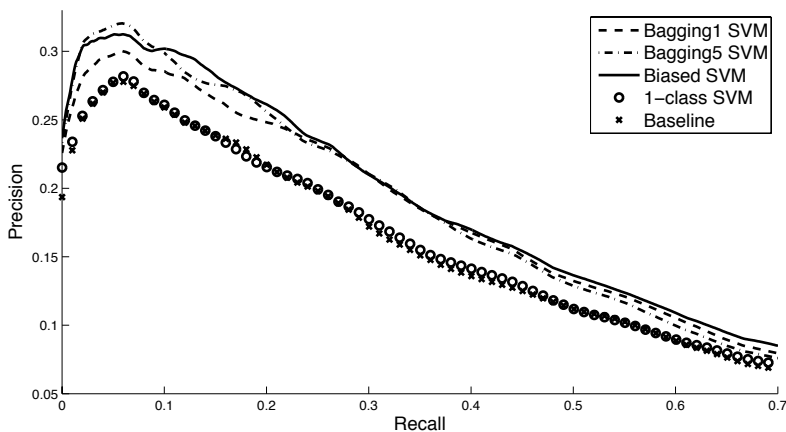


Figure 7: Precision-recall curves to compare the performance between the bagging1 SVM, the bagging5 SVM, the biased SVM, the 1-class SVM and the baseline method.

#### 442 4. Discussion

443 The main contribution of this work is to propose a new method, bagging  
444 SVM, both for inductive and transductive PU learning, and to assess in detail  
445 its performance and the influence of various parameters on simulated and real  
446 data.

447 The motivation behind bagging SVM was to exploit an intrinsic feature of  
448 PU learning to benefit from classifier aggregation through a random subsam-  
449 ple strategy. Indeed, by randomly sampling  $K$  examples from the unlabeled  
450 examples, we can expect various contamination rates, which in turn can lead  
451 to very different single classifiers (good ones when there is little contamina-  
452 tion, worse ones when contamination is high). Aggregating these classifiers  
453 can in turn benefit from the variations between them. This suggests that  
454  $K$  may play an important role in the final performance of bagging SVM,  
455 since it controls the trade-off between the mean and variance of individual  
456 classifiers. While we showed on simulated data that this is indeed the case,  
457 and that there can be some optimum  $K$  to reach the best final accuracy,  
458 the two experiments on real data did not show any strong influence of  $K$   
459 and suggested that  $K = N_P$  may be a safe default choice. This is a good  
460 news since it does not increase the number of parameters to optimize for



461 the bagging SVM and leads to balanced training sets that most classification  
462 algorithms can easily handle. Regarding parameter  $C$  optimization, our ex-  
463 periments on the Newsgroup dataset were designed so as to compare optimal  
464 cases and therefore did not include any parameter selection strategy. Hence  
465 they were rather intended as a proof of concept, to show that if one is able  
466 to successfully select optimal parameters  $C$ , then one would be able to reach  
467 same performance with the bagging SVM scheme as with the biased SVM  
468 method. However note that in practice, parameter optimization is a crucial  
469 step which may be carried out using cross validation, as was done on the *E.*  
470 *coli* dataset.

471 The comparison between different methods is mitigated. While bagging  
472 SVM outperforms biased SVM on simulated data, they are not significantly  
473 different on the two experiments with real data. Interestingly, while these PU  
474 learning methods were significantly better than two methods that learned  
475 from positive examples only on the gene regulatory network example, the  
476 1-class SVM behaved very well on the 20 newsgroup benchmark, even out-  
477 performing the PU learning methods when less than 10 training examples  
478 were provided. Many previous works, including Liu et al. (2003) and Yu  
479 et al. (2004) discard 1-class SVMs for showing a bad performance in terms

480 of accuracy, while Manevitz and Yousef (2001) report the lack of robustness  
481 of this method arguing that it has proved very sensitive to changes of pa-  
482 rameters. Our results suggest that there are cases where it remains very  
483 competitive, and that PU learning may not always be a better strategy than  
484 simply learning from positives.

485 Finally, the main advantage of bagging SVM over biased SVM is that it  
486 greatly alleviates the computation burden, in particular when there are far  
487 more unlabeled than positive examples. Indeed, a typical algorithm, such  
488 as an SVM, trained on  $N$  samples, has time complexity proportional to  $N^\beta$ ,  
489 with  $\beta$  between 2 and 3. Therefore, biased SVM has complexity proportional  
490 to  $(P + U)^\beta$  while bagging SVM's complexity is proportional  $T * (P + K)^\beta$ .  
491 With the default choice  $K = P$  ratio of CPU time to train the biased SVM  
492 vs the bagging SVM can therefore be expected to be  $((P + U)/(2P))^\beta / T$ .  
493 Then we conclude that bagging SVM should be faster than biased SVM as  
494 soon as  $U/P > 2T^{1/\beta} - 1$ . For example, taking  $T = 35$  and  $\beta = 3$ , bagging  
495 SVM should be faster than biased SVM as soon as  $U/P > 6$ , a situation  
496 very often encountered in practice where the ratio  $U/P$  is more likely to be  
497 several orders of magnitude larger. In the two real datasets, this was always  
498 the case. Table 1 reports CPU time in seconds and performance measure for

499 training bagging SVM on the first fold of newsgroup 1 with  $C$  fixed at its  
 500 best value a posteriori and  $N_P = 10$ .

Table 1: CPU time (in s) and performance measures (AUC: Area under the ROC curve and AUP: Area under the Precision/recall curve) for different settings of  $T$  and  $K$  for bagging SVM.

Bagging		CPU			AUC-AUP		
		K=10	K=50	K=200	K=10	K=50	K=200
T	35	13	39	91	0.921-0.531	0.917-0.524	0.902-0.518
	50	18	54	127	0.920-0.539	0.914-0.522	0.904-0.522
	200	72	170	473	0.918-0.539	0.910-0.528	0.904-0.511

501 In comparison, the biased SVM's CPU time is 227s for  $AUC = 0.932$   
 502 and  $AUP = 0.491$ . This confirms that for reasonable values of  $T$  and  $K$ ,  
 503 the bagging SVM is much faster than the biased SVM for a comparable  
 504 performance.

## 505 5. Conclusion

506 We have presented an original approach to the problem of learning from  
 507 positive and unlabeled examples. Our approach uses a bagging-like strategy  
 508 to exploit the availability of the numerous unlabeled examples. Extensive

509 experiments on simulated and real data have allowed us to assess the sen-  
510 sitivity of the algorithm to its parameter and to compare its performance  
511 to existing methods. We have provided safe choices of the parameters, thus  
512 reducing the number of parameters to optimize and making our algorithm  
513 simple to implement and to apply. We have shown that our bagging SVM  
514 outperforms existing approaches on simulated data. On real data, the results  
515 were more mitigated, but the bagging SVM remained competitive being ei-  
516 ther the dominant method or performing equally (no significance difference  
517 was found). Finally, our method greatly improves over the state-of-art biased  
518 SVM in terms of computation time.

## 519 **References**

- 520 Aerts, S., Lambrechts, D., Maity, S., Van Loo, P., Coessens, B., De Smet, F.,  
521 Tranchevent, L.-C., De Moor, B., Marynen, P., Hassan, B., Carmeliet, P.,  
522 Moreau, Y., May 2006. Gene prioritization through genomic data fusion.  
523 Nat. Biotechnol. 24 (5), 537–544.
- 524 Bartlett, P. L., Tewari, A., 2007. Sparseness vs estimating conditional prob-  
525 abilities: Some asymptotic results. J. Mach. Learn. Res. 8, 775–790.
- 526 Breiman, L., 1996. Bagging predictors. Mach. Learn. 24 (2), 123–140.

- 527 Breiman, L., 2001. Random forests. *Mach. Learn.* 45 (1), 5–32.
- 528 Chang, C.-C., Lin, C.-J., 2011. LIBSVM: A library for support vector ma-  
529 chines. *ACM Transactions on Intelligent Systems and Technology* 2 (3),  
530 27:1–27:27.
- 531 Chapelle, O., Schölkopf, B., Zien, A., 2006. *Semi-Supervised Learning*. MIT  
532 Press, Cambridge, MA.
- 533 De Bie, T., Tranchevent, L.-C., van Oeffelen, L. M. M., Moreau, Y., Jul 2007.  
534 Kernel-based data fusion for gene prioritization. *Bioinformatics* 23 (13),  
535 i125–i132.
- 536 Denis, F., Gilleron, R., Letouzey, F., 2005. Learning from positive and unlabeled  
537 examples. *Theor. Comput. Sci.* 348 (1), 70–83.
- 538 Elkan, C., Noto, K., 2008. Learning classifiers from only positive and unlabeled  
539 data. In: *KDD '08: Proceeding of the 14th ACM SIGKDD international conference on Knowledge discovery and data mining*. ACM, New  
540 York, NY, USA, pp. 213–220.
- 542 Faith, J., Driscoll, M., Fusaro, V., Cosgrove, E., Hayete, B., Juhn, F.,  
543 Schneider, S., Gardner, T., Jan 2008. Many microbe microarrays database:

544 uniformly normalized affymetrix compendia with structured experimental  
545 metadata. *Nucleic Acids Res.* 36 (Database issue), D866–D870.

546 Geurts, P., 2011. Learning from positive and unlabeled examples by enforcing  
547 statistical significance. *J. Mach. Learn. Res. - Proceedings Track 15*, 305–  
548 314.

549 Hastie, T., Tibshirani, R., Friedman, J., 2001. *The elements of statistical*  
550 *learning: data mining, inference, and prediction*. Springer.

551 Joachims, T., 1997. A probabilistic analysis of the Rocchio algorithm with  
552 TFIDF for text categorization. In: Fisher, D. H. (Ed.), *ICML '97: Pro-*  
553 *ceedings of the Fourteenth International Conference on Machine Learning*.  
554 Morgan Kaufmann Publishers Inc., Nashville, Tennessee, USA, pp. 143–  
555 151.

556 Joachims, T., 1999. Transductive inference for text classification using sup-  
557 port vector machines. In: Fürnkranz, J. and Kubat, M. (Eds.), *ICML*  
558 *'99: Proceedings of the Sixteenth International Conference on Machine*  
559 *Learning*. Morgan Kaufmann Publishers Inc., San Francisco, CA, USA,  
560 pp. 200–209.

561 Lee, W. S., Liu, B., 2003. Learning with positive and unlabeled examples us-

562 ing weighted logistic regression. In: Fawcett, T., Mishra, N. (Eds.), ICML  
563 2003: Proceedings of the Twentieth International Conference on Machine  
564 Learning. AAAI Press, pp. 448–455.

565 Li, X., Liu, B., 2003. Learning to classify texts using positive and unlabeled  
566 data. In: Gottlob, G. and Walsh, T. (Eds.), IJCAI’03: Proceedings of  
567 the 18th International Joint Conference on Artificial Intelligence. Morgan  
568 Kaufmann Publishers Inc., San Francisco, CA, USA, pp. 587–592.

569 Liu, B., Dai, Y., Li, X., Lee, W. S., Yu, P. S., 2003. Building text classifiers  
570 using positive and unlabeled examples. In: Wu, X. and Tuzhilin, A. and  
571 Shavlik, J. (Eds), ICDM ’03: Proceedings of the Third IEEE International  
572 Conference on Data Mining. IEEE Computer Society, Washington, DC,  
573 USA, pp. 179–186.

574 Liu, B., Lee, W. S., Yu, P. S., Li, X., 2002. Partially supervised classifi-  
575 cation of text documents. In: Sammut, C. and Hoffmann, A. G. (Eds.),  
576 ICML’02: Proceedings of the Nineteenth International Conference on Ma-  
577 chine Learning. Morgan Kaufmann Publishers Inc., San Francisco, CA,  
578 USA, pp. 387–394.

- 579 Manevitz, L. M., Yousef, M., 2001. One-class SVMs for document classifica-  
580 tion. *J. Mach. Learn. Res.* 2, 139–154.
- 581 Mordelet, F., Vert, J.-P., 2008. SIRENE: Supervised inference of regulatory  
582 networks. *Bioinformatics* 24 (16), i76–i82.
- 583 Pelckmans, K., Suykens, J., 2009. Transductively learning from positive ex-  
584 amples only. In: *ESANN 2009: Proceedings of the 17th European Sym-  
585 posium on Artificial Neural Networks*.
- 586 Salgado, H., Gama-Castro, S., Peralta-Gil, M., Díaz-Peredo, E., Sánchez-  
587 Solano, F., Santos-Zavaleta, A., Martínez-Flores, I., Jiménez-Jacinto,  
588 V., Bonavides-Martínez, C., Segura-Salazar, J., Martínez-Antonio, A.,  
589 Collado-Vides, J., Jan 2006. RegulonDB (version 5.0): *Escherichia coli*  
590 K-12 transcriptional regulatory network, operon organization, and growth  
591 conditions. *Nucleic Acids Res.* 34 (Database issue), D394–D397.
- 592 Schölkopf, B., Platt, J. C., Shawe-Taylor, J., Smola, A. J., Williamson, R. C.,  
593 2001. Estimating the support of a high-dimensional distributions. *Neural  
594 Comput.* 13, 1443–1471.
- 595 Scott, C., Nowak, R., 2005. A Neyman-Pearson approach to statistical learn-  
596 ing. *IEEE Trans. Inf. Theory* 51 (11), 3806–3819.



- 597 Scott, C., Blanchard, G., 2009. Novelty detection: Unlabeled data definitely  
598 help. In: van Dyk, D., Welling, M. (Eds.), AISTATS '09 Proceedings of the  
599 Twelfth International Conference on Artificial Intelligence and Statistics.  
600 Vol. 5. JMLR: W&CP 5, Clearwater Beach, Florida, pp. 464–471.
- 601 Shah, A. R., Oehmen, C. S., Webb-Robertson, B.-J., Mar 2008. SVM-  
602 HUSTLE—an iterative semi-supervised machine learning approach for pair-  
603 wise protein remote homology detection. *Bioinformatics* 24 (6), 783–790.
- 604 Sriphaew, K., Takamura, H., Okumura, M., 2009. Cool blog classification  
605 from positive and unlabeled examples. In: Theeramunkong, T. and Ki-  
606 jsirikul, B. and Cercone, C. and Ho, T-B. (Eds.), PAKDD 2009: Pro-  
607 ceedings of the 13th Pacific-Asia Conference on Advances in Knowledge  
608 Discovery and Data Mining, Springer-Verlag, Berlin, Heidelberg, pp. 62–  
609 73.
- 610 Steinwart, I., 2003. Sparseness of Support Vector Machines. *J. Mach. Learn.*  
611 *Res.* 4, 1071–1105.
- 612 Vert, R., Vert, J.-P., 2006. Consistency and convergence rates of one-class  
613 SVMs and related algorithms. *J. Mach. Learn. Res.* 7, 817–854.
- 614 Yu, H., Han, J., Chang, K. C.-C., 2004. PEBL: Web page classification

615 without negative examples. IEEE Trans. Knowl. Data Eng. 16 (1), 70–  
616 81.

617 Zhang, K., Tsang, I., Kwok, J., 2009. Maximum margin clustering made  
618 practical. IEEE Trans. Neural Netw 20 (4), 583–596.