# CS395T Data Mining
# Project report
# One-class SVM formulations for Multiple Instance learning

Sudheendra Vijayanarasimhan

## 1    Introduction

Multiple Instance learning (MIL) considers a particular form of weak supervision in which the learner is given a set of *positive* bags and *negative* bags. *Positive* bags are sets of instances containing atleast one positive example and *negative* bags are sets of instances all of which are negative. A number of binary SVM based solutions have been proposed to this problem like the Normalized Set Kernel of Gartner et. al, 2002 ([1]) which represents the bag as the sum of all its instances normalized by its 1 or 2-norm and the sparse MIL (sMIL) technique of Razvan and Mooney, 2007 ([2]) which improves upon NSK by considering a weaker balancing constraint. In this project I plan to look at equivalent formulations for a one-class SVM and empirically evaluate if ignoring the negative bags in the formulation is detrimental to the solution found.

## 2    Related Work

A number of 2-class SVM based formulations have been looked at in the literature. The following are a few relevant MIL SVM formulations

- Normalized Set Kernel (NSK)

  In the Normalized Set Kernel of Gartner et. al, 2002 ([1]) a bag is represented as the sum of all its instances, normalized by its 1 or 2-norm. The resulting representation is then trained using a standard SVM. The formulation for NSK is as follows

- sparse MIL (sMIL)

  The sparse MIL formulation of Razvan and Mooney, 2007 ([2]) considers the equation for the positive bags in the formulation of NSK as a balancing constraint. The balancing constraint of NSK is too strong since it assumes that all the instances in a positive bag are positive. Since this is

problematic when the positive bag is particularly sparse in positive examples they consider the constraint that expresses that *at least* one instance from the bag is positive. The formulation for the same is as follows

- MI-SVM

minimize: $\frac{1}{2}||w||^2 + \frac{C}{|\mathcal{X}_n|}\sum_{x \in \mathcal{X}_n} \xi_x + \frac{C}{|\mathcal{X}_p|}\sum_{X \in \mathcal{X}_p} \xi_X$ (1)

subject to: $max_{x \in X} \quad y(w\phi(x) + b) \geq 1 - \xi_X, \forall X \in \mathcal{X}_p, X_n$

(2)

This is the maximum bag margin formulation of Andrews et. al (2003) [3]. The associated heuristic algorithm starts by training a standard SVM. This is followed by relabeling of instances in positive bags using the decision hyperplane. If a positive bag contains no instances that are positive according to this hyperplane then the instance with the maximum value of the decision function is relabeled as positive and the SVM is retrained on this relabeled data. This is continued till there are no more labels to be changed.

- A regularization framework for Multiple-Instance learning In the method proposed by Cheung and Kwok (2006) [4] a loss function is introduced between the label of a bag and the label of the most positive instance in the bag and SVM is formulated by including this loss function in the objective. This is based on relaxing the idea of mi-SVM that the label of a positive bag is equal to the label of its most positive instance. But the objective function is no longer convex because of the max function used for the loss and therefore they directly formulate the dual problem instead of the primal and solve it using CCCP([5]).

Positive bags are easily constructed in all of the above cases. For example, in image retrieval each returned set of images can be considered as a positive bag, segmentation where each image is a positive bag containing atleast one valid segmentation. On the other hand it is not clear on how to choose negative bags in these cases and they are typically constructed from examples that are known to be non-positive.

Ray and Craven (2005) ([6]) observe that the nature of the negative instances in the positive bags may be different from the nature of the negative instances in the negative bags. If this is the case then one would be dealing with 3 different distributions which might or might not be separable using the single hyperplane found by all of these methods. And so it appears that most of these methods work well only when the negatives in the positive bags are similarly distributed to the negatives in the negative bags ([2]).

Therefore even though we've the freedom of constructing the negative bags from any set of instances that is not positive the most gains are obtained when these are sampled from a distribution similar to that of the negatives in the

positive bags. This might not be possible in some cases like image retrieval where no information is known about the distribution of the noisy images in each retrieved set.

In such situations completely ignoring the negative bags in the formulation and considering only the positive bags and using clustering techniques or a one-class SVM might be fruitful. In the following work we will formulate a one-class SVM for the MI problem and compare it with the standard one-class SVM on an image dataset. We will also compare the one-class SVM solutions with the 2-class methods outlined above to study the effect of ignoring the negative bags.

## 3 A one-class SVM approach to MIL

A one-class SVM is a function $f$ that takes the value $+1$ in a "small" region capturing most of the data points and -1 elsewhere. One-class SVMs are typically used for novelty detection where the task is to say whether a new example is unlike any one of training examples. One-class SVMs have also been applied to the task of unsupervised learning for character regognition ([7]).

The MI problem could be solved using the one-class SVM by simply considering all instances in positive bags as unlabeled data and then estimating a function that returns $+1$ in a "small" region that should correspond to the true positives. The function is found by mapping the data into feature space corresponding to a kernel and then separating them from the origin with maximum margin. This corresponds to the following quadratic problem

$$\text{minimize:} \quad \tfrac{1}{2}||w||^2 + \tfrac{1}{\nu l}\sum_i \xi_i - \rho \tag{3}$$
$$\text{subject to:} \quad (w.\phi(x)) \geq \rho - \xi_i, \xi_i \geq 0. \tag{4}$$

But the above formulation does not respect the MI constraint which states that positive bags should contain atleast one positive instance. From Ray and Craven (2005) [6] it is clear that even though ignoring the bag constraint and solving the standard supervised problem produces results comparable to MI methods in most datasets, when the bags are very sparse MI methods invariably perform better.

As seen in Section 2 the MI constraint can be captured in a number of ways. Of these the idea of Andrews et. al [3] is closest to capturing the MI constraint since it states that the maximum value of $y$ within a bag should be greater than $+1$ if the bag is positive. But even though the max function is convex it is not smooth and so the standard quadratic optimization techniques cannot be applied. Therefore we will apply the technique used in [4]. The one-class MI-svm formulation is given in Figure 1

Here $l$ denotes the total number of instances within all bags while $n$ denotes the total number of bags. The penalty for bags and instances have been seperated out because the bag constraint is a stronger constraint as we want atleast

$$\text{minimize:} \quad \tfrac{1}{2}||w||^2 + \tfrac{1}{\nu l}\sum_i \xi_i + \tfrac{1}{\nu n}\sum_i \Xi_i - \rho \qquad (5)$$

$$\text{subject to:} \quad (w.\phi(x)) \geq \rho - \xi_i, \xi_i \geq 0.$$

$$max_{x \in X} \quad (w\phi(x)) \geq \rho - \Xi_X, \forall X \in \mathcal{X}_p$$

$$(6)$$

Figure 1: one-class MI-svm formulation

one instance to be positive in each bag. On the other hand individual instances might not be all positive and therefore the penalty should be lower.

Using multipliers $\alpha_i, \beta_i \geq 0$, we introduce a Lagrangian

$$
L(w,\xi,\Xi,\rho,\alpha,\beta) = \quad \tfrac{1}{2}||w||^2 + \tfrac{1}{\nu l}\sum_i \xi_i + \tfrac{1}{\nu n}\sum_i \Xi_i - \rho - \sum_i \alpha_i((w.\phi(x_i)) - \rho + \xi_i)
$$
$$
- \sum_i \alpha_i (max_{x \in X_i}(w.\phi(x)) - \rho + \Xi_i) - \sum_i \beta_i \xi_i - \sum_i \beta_i \Xi_i
$$
$$(7)$$

Now, because of the presence of the $max$ function the lagragian is not differential. But by using the sub-gradient of the $max$ function and setting the derivatives to zero we can obtain the dual of the above problem. For the point-wise maximum function $h(x) = max_{1 \leq i \leq p} h_i(x)$ its subdifferential at $x_0$ is the convex hull of the union of subgradients of "active" functions at $x_0$. Function $h_i$ is said to be active if $h_i = max_{1 \leq i \leq p} h_i(x)$. Introducing variables $a_{ij}$ to denote whether $(w.\phi(x_j))$ is active or not in the max function yields the solution

$$w = \sum_i \alpha_i \phi(x_i) + \sum_i \alpha_i (\sum_j a_{ij}\phi(x_j))$$

$$\alpha_i \leq \frac{1}{\nu l}, when\ i\ is\ an\ instance$$

$$\alpha_i \leq \frac{1}{\nu n}, when\ i\ is\ a\ bag$$

$$\sum_i \alpha_i = 1.$$

$$\sum_j a_{ij} = 1, \forall i$$

$$(8)$$

Similar to [4] we initialize $a_{ij}^{(0)} = 1/n_i$ for all bags, where $n_i$ denotes the number of instances in the bag and the $a_{ij}$s are updated as $a_{ij} = 0, if x_j$ is not active in the max function and $a_{ij} = 1/n_a$ if it is. Here $n_a$ denotes the number of active instances.

| Category | AUROC - training | | AUROC - test | |
|---|---|---|---|---|
| | one-class | one-class MI-svm | one-class | one-class MI-svm |
| ajaxorange | 64.22 | **65.40** | 63.72 | **64.14** |
| apple | 50.64 | **50.96** | 49.70 | 49.62 |
| banana | 63.16 | **64.82** | 61.24 | **63.20** |
| bluescrunge | 49.00 | **52.44** | 47.86 | **50.50** |
| candlewithholder | 76.14 | **76.62** | 77.22 | 77.20 |
| cardboardbox | 77.62 | **78.84** | 75.80 | **77.00** |
| checkeredscarf | 78.32 | 78.06 | 78.62 | 78.42 |
| cokecan | 76.78 | 74.32 | 76.00 | 72.98 |
| dataminingbook | 81.36 | 80.84 | 77.16 | 76.34 |
| dirtyrunningshoe | 76.54 | **77.70** | 71.86 | **72.60** |
| dirtyworkgloves | 74.82 | **75.90** | 77.78 | **78.74** |
| fabricsoftenerbox | 83.62 | **83.76** | 82.42 | 82.24 |
| feltflowerrug | 60.56 | 57.00 | 56.66 | 52.18 |
| glazedwoodpot | 48.52 | **48.84** | 44.08 | 44.06 |
| goldmedal | 49.66 | **55.78** | 49.86 | **55.24** |
| greenteabox | 64.36 | **64.92** | 63.76 | **64.28** |
| juliespot | 50.86 | **52.00** | 47.68 | **48.58** |
| largespoon | 69.94 | **72.86** | 75.74 | **77.72** |
| rapbook | 69.62 | **71.32** | 72.22 | **73.50** |
| smileyfacedoll | 56.50 | 56.20 | 59.80 | **59.98** |
| spritecan | 68.92 | **69.04** | 67.82 | 67.22 |
| stripednotebook | 86.06 | **87.84** | 82.72 | 84.36 |
| translucentbowl | 43.06 | **44.18** | 44.48 | **46.24** |
| wd40can | 69.78 | 66.02 | 65.68 | 60.64 |
| woodrollingpin | 78.62 | **79.20** | 76.94 | 76.50 |
| Average | 66.75 | **67.39** | 65.87 | **66.14** |

Table 1: Area under the ROC curve for different categories in the SIVAL dataset on both training and test data averaged over 5 random trials. Numbers highlighted in bold area cases where adding the MI constraint improves the area under the ROC. We can clearly see that there is an improvement for majority of the categories even though overall average is only slightly larger.

# 4    Experiments and Results

The one-class MI-svm formulation was tested on the image segmentation domain using the SIVAL dataset. The dataset contains segmented images of various objects in different scenes. A positive instance is a segment containing the object, while all others are negative. An image (bag) is labeled positive if it contains the object.

The classification task is to say whether a given *segment* is positive or not. Since we do not use the labels on positive instances in either the standard one-class SVM or the one-class MI-svm both the results on the training data and the test data are equally relevant here. The dataset was randomly split into 5 runs each containing a training set of 20 images ($\tilde{2}0$*30 segments) and a test set of 40 images. All results are averaged over the 5 runs.

Kernel and other parameters were optimized separately and the same value was used for both the standard SVM and the MI-svm since only the comparison on the same set of parameters would be relevant. All results are for a quadratic kernel with a coefficient of $5 * 10^{-6}$ and $\nu = 0.9$.

Table 1 shows the Area under the ROC curve for the classification task of predicting whether a *segment* is positive or not. The first two columns are on the training data while the last two are on the test data. Numbers highlighted in bold refer to cases where the MI method did better than the standard one-class method. The MI method performs better than the standard one-class method in a majority of the categories with a maximum improvement of 3.44 in the case of *bluescrunge* on the training set. But the overall averges differ by less than 1 point and it cannot be considered as definitive that the MI method is better than the one-class. Figure 2 shows a whisker and box plot of the difference in the area under the ROC for the two methods. We see that there are a number of negative outliers which could be responsible for the low average improvement.

# 5    Conclusion

A one-class SVM was formulated for the MI problem and it was compared with a standard one-class SVM for an image dataset. Preliminary results suggest that the MI method might be slightly better than the standard one-class SVM but this needs to be tested on more datasets to come to a strong conclusion. Also the above methods need to be compared with existing 2-class SVMs for MIL to analyse if ignoring negative bags can help simplify the problem.

# References

[1] Gartner, T., Flach, P., Kowalczyk, A. & Smola, A. (2002). Multiple Instance kernels. *In Proceedings of the 19th International Conference on Machine Learning*
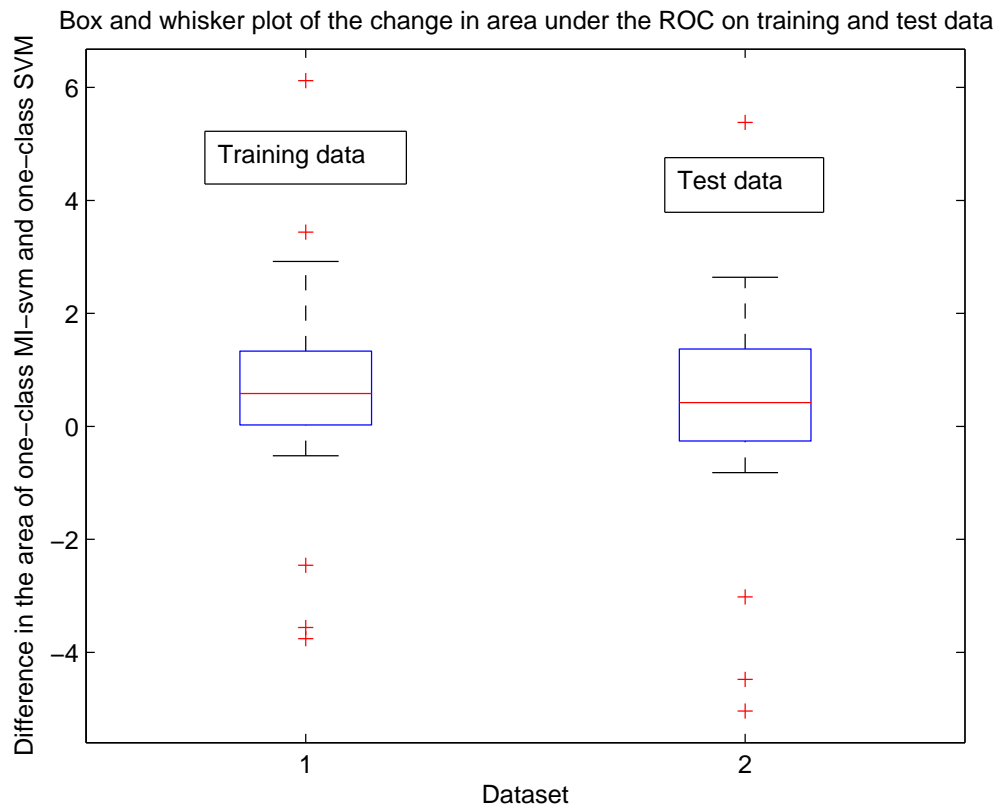
Figure 2: A box and whisker plot of the difference in the area under ROC for the two methods. The mean and lower quartile are above zero for both test and training data implying that the MI constraint does improve the area under the ROC for most categories

[2] Razvan Bunescu, Ray Mooney (2007). Multiple Instance learning for sparse positive bags.

[3] Andrews, S., Tsochantaridis, I., & Hofmann, T. (2003). Support vector machines for multipleinstance learning. *Advances in Neural Information Processing Systems 15. Cambridge, MA: MIT Press.*

[4] Pak-Ming Cheung and James T. Kwok (2006) A regularization framework for multiple-instance learning. *ICML '06: Proceedings of the 23rd international conference on Machine learning*

[5] Yuille, A., & Rangarajan, A. (2003). The concaveconvex procedure. *Neural Computation, 15, 915*

[6] Ray, S., and Craven, M. (2005). Supervised versus multiple instance learning: An empirical comparison. *Proceedings of 22nd International Conference on Machine Learning (ICML-2005) (pp. 697Bonn, Germany.*

[7] B. Sch and o Platt and J. Shawe-Taylor and A. Smola and R. Williamson. Estimating the support of a high-dimensional distribution. Technical Report 99-87, Microsoft Research, 1999.