A probabilistic approach to binary classification
Other elements to take into account
Numerica results and comparisons
Conclusion

# Probabilistic Approach to One Class SVM

V. Leclère (ENPC), L. El Ghaoui, E. Graves

MODE 2016
March 25, 2016

A probabilistic approach to binary classification
Other elements to take into account
Numerica results and comparisons
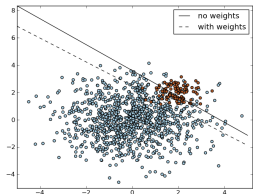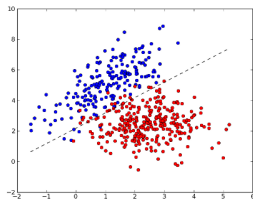Conclusion

## Presentation Outline

1. **A probabilistic approach to binary classification**
   - Problem formulation
   - Toward the SVM formulation

2. **Other elements to take into account**
   - Some other ingredients
   - Large scale problem
   - Kernelization

3. **Numerica results and comparisons**
   - Some recalls
   - Small experiment : protein classification
   - Large experiment : text classification

4. **Conclusion**

A probabilistic approach to binary classification
Other elements to take into account
Numerica results and comparisons
Conclusion

Problem formulation
Toward the SVM formulation

# Contents

A probabilistic approach to binary classification
Other elements to take into account
Numerica results and comparisons
Conclusion

Problem formulation
Toward the SVM formulation

## Binary classification

- Consider that you have a set of training data points $\{x_i\}_{i \in I}$.
- Each data is either labeled as a positive or negative point, through $y_i \in \{-1, +1\}$.
- We are looking for a linear classifier $(w, b)$ such that
  - $w^T x_i - b \geq 0$, if $y_i = 1$
  - $w^T x_i - b \leq 0$, if $y_i = -1$
- We are intereseted in imbalanced classification where there are a lot more negative points than positive points.

A probabilistic approach to binary classification
Other elements to take into account
Numerica results and comparisons
Conclusion

Problem formulation
Toward the SVM formulation

## Probabilistic formulation

- We represent the negative points as a random variable **x**.
- We want a classifier that truly classify each positive point, and minimize the probability of a false negative, i.e.

$$\max_{\mathbf{w}, b} \quad \mathbb{P}(\mathbf{w}^\top \mathbf{x} - b \leq 0),$$

$$s.t. \quad \mathbf{w}^\top \mathbf{x}_i - b \geq 0, \quad \forall i \in I^+.$$

- However, as we do not know the probability distribution of the negative points, we consider a robust approach where only the mean x̄ and covariance Σ are known.

$$\max_{\mathbf{w}, b} \quad \inf_{\mathbf{x} \sim (\bar{\mathbf{x}}, \Sigma)} \quad \mathbb{P}(\mathbf{w}^\top \mathbf{x} - b \leq 0),$$

$$s.t. \quad \mathbf{w}^\top \mathbf{x}_i - b \geq 0, \quad \forall i \in I^+.$$

A probabilistic approach to binary classification
Other elements to take into account
Numerica results and comparisons
Conclusion

Problem formulation
Toward the SVM formulation

## Probabilistic formulation

- We represent the negative points as a random variable $\mathbf{x}$.
- We want a classifier that truly classify each positive point, and minimize the probability of a false negative, i.e.

$$\max_{\mathbf{w},b} \qquad \mathbb{P}(\mathbf{w}^\top \mathbf{x} - b \leq 0),$$
$$s.t. \qquad \mathbf{w}^\top \mathbf{x}_i - b \geq 0, \qquad \forall i \in I^+.$$

- However, as we do not know the probability distribution of the negative points, we consider a robust approach where only the mean $\bar{\mathbf{x}}$ and covariance $\Sigma$ are known.

$$\max_{\mathbf{w},b} \qquad \inf_{\mathbf{x}\sim(\bar{\mathbf{x}},\Sigma)} \quad \mathbb{P}(\mathbf{w}^\top \mathbf{x} - b \leq 0),$$
$$s.t. \qquad \mathbf{w}^\top \mathbf{x}_i - b \geq 0, \qquad \forall i \in I^+.$$

A probabilistic approach to binary classification
Other elements to take into account
Numerica results and comparisons
Conclusion

Problem formulation
Toward the SVM formulation

# Contents

A probabilistic approach to binary classification
Other elements to take into account
Numerica results and comparisons
Conclusion

Problem formulation
Toward the SVM formulation

## Toward an SVM-like representation 1/3

- The condition

$$\inf_{\mathbf{x} \sim (\bar{\mathbf{x}}, \boldsymbol{\Sigma})} \mathbb{P}\left(\mathbf{x}^\top \mathbf{w} - b \leq 0\right) \geq \alpha,$$

holds if and only if

$$b - \bar{\mathbf{x}}^\top \mathbf{w} \geq \kappa(\alpha)\sqrt{\mathbf{w}^\top \boldsymbol{\Sigma} \mathbf{w}},$$

where $\kappa(\alpha) = \sqrt{\frac{\alpha}{1-\alpha}}$.

- Hence, our problem reads

$$\max_{x, w, b} \quad \alpha$$

$$\text{s.t.} \quad b - \bar{\mathbf{x}}^\top \mathbf{w} \geq \kappa(\alpha)\sqrt{\mathbf{w}^\top \boldsymbol{\Sigma} \mathbf{w}},$$

$$\mathbf{x}_i^\top \mathbf{w} - b \geq 0, \qquad \forall i \in I^+.$$

A probabilistic approach to binary classification
Other elements to take into account
Numerica results and comparisons
Conclusion

Problem formulation
Toward the SVM formulation

## Toward an SVM-like representation 1/3

- The condition

$$\inf_{\mathbf{x} \sim (\bar{\mathbf{x}}, \boldsymbol{\Sigma})} \mathbb{P}\left(\mathbf{x}^\top \mathbf{w} - b \leq 0\right) \geq \alpha,$$

holds if and only if

$$b - \bar{\mathbf{x}}^\top \mathbf{w} \geq \kappa(\alpha) \sqrt{\mathbf{w}^\top \boldsymbol{\Sigma} \mathbf{w}},$$

where $\kappa(\alpha) = \sqrt{\frac{\alpha}{1-\alpha}}$.

- Hence, our problem reads

$$\begin{aligned}
\max_{\alpha, \mathbf{w}, b} \quad & \alpha \\
\text{s.t.} \quad & b - \bar{\mathbf{x}}^\top \mathbf{w} \geq \kappa(\alpha) \sqrt{\mathbf{w}^\top \boldsymbol{\Sigma} \mathbf{w}}, \\
& \mathbf{x}_i^\top \mathbf{w} - b \geq 0, \qquad\qquad \forall i \in I^+.
\end{aligned}$$

A probabilistic approach to binary classification
Other elements to take into account
Numerica results and comparisons
Conclusion

Problem formulation
Toward the SVM formulation

## Toward an SVM-like representation 2/3

- $\kappa : \alpha \mapsto \sqrt{\frac{\alpha}{1-\alpha}}$ is increasing on $[0, 1[$, this problem is equivalent to the program:

$$\max_{\kappa, \mathbf{w}, b} \quad \kappa$$
$$\text{s.t.} \quad b - \bar{\mathbf{x}}^\top \mathbf{w} \geq \kappa \sqrt{\mathbf{w}^\top \Sigma \mathbf{w}},$$
$$\mathbf{x}_i^\top \mathbf{w} - b \geq 0, \qquad \forall i \in I^+.$$

- Note that $w \neq 0$, hence we can impose $\kappa\sqrt{w^\top \Sigma w} = 1$.
- Leading to

$$\max_{w, b} \quad \frac{1}{\sqrt{w^\top \Sigma w}}$$
$$\text{s.t.} \quad b - \bar{x}^\top w \geq 1,$$
$$x_i^\top w - b \geq 0, \qquad \forall i \in I^+.$$

A probabilistic approach to binary classification
Other elements to take into account
Numerica results and comparisons
Conclusion

Problem formulation
Toward the SVM formulation

## Toward an SVM-like representation 2/3

- $\kappa : \alpha \mapsto \sqrt{\frac{\alpha}{1-\alpha}}$ is increasing on $[0, 1[$, this problem is equivalent to the program:

$$\max_{\kappa, \mathbf{w}, b} \quad \kappa$$
$$\text{s.t.} \quad b - \bar{\mathbf{x}}^\top \mathbf{w} \geq \kappa \sqrt{\mathbf{w}^\top \Sigma \mathbf{w}},$$
$$\mathbf{x}_i^\top \mathbf{w} - b \geq 0, \qquad \forall i \in I^+.$$

- Note that $w \neq 0$, hence we can impose $\kappa \sqrt{\mathbf{w}^\top \Sigma \mathbf{w}} = 1$.

- Leading to

$$\max_{w, b} \quad \frac{1}{\sqrt{\mathbf{w}^\top \Sigma \mathbf{w}}}$$
$$\text{s.t.} \quad b - \bar{\mathbf{x}}^\top w \geq 1,$$
$$\mathbf{x}_i^\top w - b \geq 0, \qquad \forall i \in I^+.$$

A probabilistic approach to binary classification
Other elements to take into account
Numerica results and comparisons
Conclusion

Problem formulation
Toward the SVM formulation

## Toward an SVM-like representation 2/3

- $\kappa : \alpha \mapsto \sqrt{\frac{\alpha}{1-\alpha}}$ is increasing on $[0, 1[$, this problem is equivalent to the program:

$$\max_{\kappa, \mathbf{w}, b} \quad \kappa$$
$$\text{s.t.} \quad b - \bar{\mathbf{x}}^\top \mathbf{w} \geq \kappa \sqrt{\mathbf{w}^\top \Sigma \mathbf{w}},$$
$$\mathbf{x}_i^\top \mathbf{w} - b \geq 0, \qquad \forall i \in I^+.$$

- Note that $w \neq 0$, hence we can impose $\kappa \sqrt{\mathbf{w}^\top \Sigma \mathbf{w}} = 1$.
- Leading to

$$\max_{\mathbf{w}, b} \quad \frac{1}{\sqrt{\mathbf{w}^\top \Sigma \mathbf{w}}}$$
$$\text{s.t.} \quad b - \bar{\mathbf{x}}^\top \mathbf{w} \geq 1,$$
$$\mathbf{x}_i^\top \mathbf{w} - b \geq 0, \qquad \forall i \in I^+.$$

A probabilistic approach to binary classification
Other elements to take into account
Numerica results and comparisons
Conclusion

Problem formulation
Toward the SVM formulation

## Toward an SVM-like representation 3/3

- Finally, since the function $x \mapsto 1/\sqrt{x}$ is decreasing on $\mathbb{R}_+^*$, we obtain the equivalent program:

$$\min_{\mathbf{w}} \quad \mathbf{w}^\top \Sigma \mathbf{w}$$
$$\text{s.t.} \quad b - \bar{\mathbf{x}}^\top \mathbf{w} \geq 1,$$
$$\mathbf{x}_i^\top \mathbf{w} - b \geq 0, \qquad \forall i \in I^+.$$

- Which is a Support Vector Machine formulation with two differences :
  - instead of minimizing the $\ell_2$-norm of $\mathbf{w}$, we minimize the Mahalanobis norm corresponding to the covariance matrix of the negative class distribution,
  - the negative class contains only $\bar{\mathbf{x}}$.

A probabilistic approach to binary classification
Other elements to take into account
Numerica results and comparisons
Conclusion

Problem formulation
Toward the SVM formulation

## Link to the One-Class SVM

- In the previous formulation we see that the optimal $b$ can be derived as

$$b^\sharp = \min_{i \in I^+} x_i^T w^\sharp = 1 + \bar{\mathbf{x}}^\top \mathbf{w}^\sharp$$

- Hence, we have the formulation

$$\min_{\mathbf{w}} \quad \mathbf{w}^\top \Sigma \mathbf{w}$$
$$\text{s.t.} \quad (\mathbf{x}_i - \bar{\mathbf{x}})^\top \mathbf{w} \geq 1, \qquad \forall i \in I^+.$$

- This is almost a one-class SVM :
  - we separate in the Mahalanobis norm,
  - we separate from the mean of the negative class instead of the origin.
- It is equivalent to apply classical one class SVM to preprocessed positive datapoints:

$$\tilde{x}_i = \Sigma^{-1/2}(x_i - \bar{x}).$$

A probabilistic approach to binary classification
Other elements to take into account
Numerica results and comparisons
Conclusion

Some other ingredients
Large scale problem
Kernelization

# Contents

A probabilistic approach to binary classification
Other elements to take into account
Numerica results and comparisons
Conclusion

Some other ingredients
Large scale problem
Kernelization

## Non-separability

- Our formulation require that the mean of the negative class is not in the convex hull of the positive points.
- To relax this strong assumption we add slack variable $\xi_i$, which are penalized.
- The formulation reads

$$
\min_{\mathbf{w}, b, \xi \geq 0} \quad \mathbf{w}^\top \Sigma \mathbf{w} + \frac{1}{\nu |I^+|} \sum_{i \in I^+} \xi_i
$$
$$
\text{s.t.} \quad b - \bar{\mathbf{x}}^\top \mathbf{w} \geq 1,
$$
$$
\mathbf{x}_i^\top \mathbf{w} - b + \xi_i \geq 0, \qquad \forall i \in I^+.
$$

A probabilistic approach to binary classification
Other elements to take into account
Numerica results and comparisons
Conclusion

Some other ingredients
Large scale problem
Kernelization

## Sparsity

- We might want the optimal classifier $w$ to be sparse to be able to interpret it, and to regularize the solution.
- A classical way of asking for Sparsity consists in adding a penalization of the $L_1$ norm of the solution $w$.
- Leading to

$$\min_{\mathbf{w}, b, \xi \geq 0} \quad \mathbf{w}^\top \Sigma \mathbf{w} + \frac{1}{\nu |I^+|} \sum_{i \in I^+} \xi_i + \eta \|w\|_1$$

$$\text{s.t.} \quad b - \bar{\mathbf{x}}^\top \mathbf{w} \geq 1,$$

$$\mathbf{x}_i^\top \mathbf{w} - b + \xi_i \geq 0, \qquad \forall i \in I^+.$$

A probabilistic approach to binary classification
Other elements to take into account
Numerica results and comparisons
Conclusion

Some other ingredients
Large scale problem
Kernelization

# Contents

A probabilistic approach to binary classification
Other elements to take into account
Numerica results and comparisons
Conclusion

Some other ingredients
Large scale problem
Kernelization

## Factor model of covariance matrix

- To reduce the computation cost and to reduce the effect of noise in the data we look to a factor model of the covariance matrix.

- More precisely we look for matrices $D$ and $F$ such that

$$\Sigma \approx D + FF^T,$$

where $D$ is diagonal and $F$ is $n \times k$ with $k << n$.
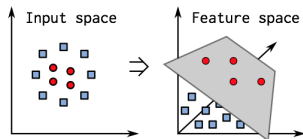
- This approximation can be obtained by use of svd decomposition of the centered dataset.

- Note that if the data is sparse the svd decomposition of the centered data can be done efficiently.

A probabilistic approach to binary classification
Other elements to take into account
Numerica results and comparisons
Conclusion

Some other ingredients
Large scale problem
Kernelization

## Contents

A probabilistic approach to binary classification
**Other elements to take into account**
Numerica results and comparisons
Conclusion

Some other ingredients
Large scale problem
**Kernelization**

# Kernelization

- It is classical for different application to train svm with kernel.

- Basically a kernel approach consists in looking for non linear separation of the datapoints by considering a linear separation in a bigger space:
  - consider a feature function $\varphi$ such that for any $(x, y)$ we have $\varphi(x)^T \varphi(y) = K(x, y)$;
  - apply linear SVM to the points $\varphi(x_i)$.



- Practically it only requires to replace the scalar product $x_i^T x_j$ by another symmetric function $K(x_i, x_j)$ satisfying some properties.

A probabilistic approach to binary classification
Other elements to take into account
Numerica results and comparisons
Conclusion

Some other ingredients
Large scale problem
Kernelization

## Kernelization

- We have seen that our approach is equivalent to applying one-class SVM to some preprocessed data. So why not applying kernel to this preprocessed data ?

- We don't want to preprocess the datapoints (consume time and lose sparsity).

- However applying classical kernels (polynomial, RBF and sigmoidal) applied to the preprocessed data are equivalent to customized kernels applied to the original positive data.

A probabilistic approach to binary classification
Other elements to take into account
Numerica results and comparisons
Conclusion

Some recalls
Small experiment : protein classification
Large experiment : text classification

# Contents

A probabilistic approach to binary classification
Other elements to take into account
Numerica results and comparisons
Conclusion

Some recalls
Small experiment : protein classification
Large experiment : text classification

## Other method available

Standard binary classification on an imbalanced subset gives poor results for the small class. In particular most classifier have a tendency to always classify points as element of the big class.

Here are two classical ways of dealing with this problem

- Subsampling: randomly selectionning a subsample of the negative class and consider it as the whole negative class.
- Differential costs: assign different penalty $C$ to each class, penalizing more any misclassification of the small class.

A probabilistic approach to binary classification
Other elements to take into account
Numerica results and comparisons
Conclusion

Some recalls
Small experiment : protein classification
Large experiment : text classification

## Area Under the Curve Metric

- linear classifier: $w^T x_i + b$.
- Parameter $b$ depend on your willingness to have false positive or false negative.
- The ROC curve is a curve in the true positive / false positive plane.
- Interpretation : take randomly a positive and a negative point. AUC is the probability of finding the positive with classifier $w$.



Comparing ROC Curves

A probabilistic approach to binary classification
Other elements to take into account
Numerica results and comparisons
Conclusion

Some recalls
Small experiment : protein classification
Large experiment : text classification

# Contents

A probabilistic approach to binary classification
Other elements to take into account
Numerica results and comparisons
Conclusion

Some recalls
Small experiment : protein classification
Large experiment : text classification

## Dataset presentation

| Dataset | # positive | # negative | ratio |
|---------|-----------|-----------|-------|
| PHOSS   | 613       | 10,798    | 17    |
| PHOST   | 140       | 9,051     | 64    |
| PHOSY   | 136       | 5,103     | 37    |
| CAM     | 942       | 17,974    | 19    |

Table : Basic statistics about the different datasets.

Available here :
www.informatics.indiana.edu/predrag/publications.htm

A probabilistic approach to binary classification
Other elements to take into account
Numerical results and comparisons
Conclusion

Some recalls
Small experiment : protein classification
Large experiment : text classification

## Results comparison

|  | This work | Cost-sensitive | Sampling |
|---|---|---|---|
| PhosS | $77.2^{\dagger} \pm 0.7$ | $76.8 \pm 0.8$ | $74.3 \pm 1.1$ |
| PhosT | $77.4^{\dagger} \pm 1.7$ | $73.0 \pm 2.0$ | $72.0 \pm 1.5$ |
| PhosY | $76.2^{\dagger} \pm 1.5$ | $72.8 \pm 1.7$ | $70.1 \pm 2.1$ |
| CaM | $78.2 \pm 0.5$ | $78.1 \pm 0.5$ | $75.3 \pm 0.4$ |

Table : Areas under the ROC curve (with confidence intervals), averaged over twenty experiments. $^{\dagger}$ indicates that our method is significantly better than the two others, (with $p$-value $p < 0.01$).

A probabilistic approach to binary classification
Other elements to take into account
Numerica results and comparisons
Conclusion

Some recalls
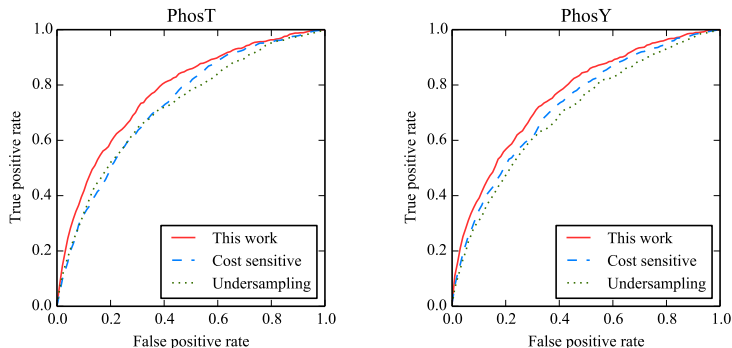Small experiment : protein classification
Large experiment : text classification

Figure : ROC curves, averaged over twenty experiments, on the PHOST and PHOSY datasets.

A probabilistic approach to binary classification
Other elements to take into account
Numerica results and comparisons
Conclusion

Some recalls
Small experiment : protein classification
Large experiment : text classification

# Contents

A probabilistic approach to binary classification
Other elements to take into account
Numerica results and comparisons
Conclusion

Some recalls
Small experiment : protein classification
Large experiment : text classification

## The reuter dataset

- About 200'000 documents, with 50'000 features, classified in 40 classes.
- Available on Liblinear website.
- We select one-class to be positive all the other are negative.

A probabilistic approach to binary classification
Other elements to take into account
Numerica results and comparisons
Conclusion

Some recalls
Small experiment : protein classification
Large experiment : text classification

## Reuters results                                                          I

| Topic | This work | Cost-sensitive | Sampling |
|-------|-----------|----------------|----------|
| 2 | $89.7 \pm 1.0$ | $89.9 \pm 1.4$ | $87.7 \pm 1.2$ |
| 9 | $96.1 \pm 0.7$ | $96.3 \pm 0.8$ | $94.1 \pm 1.3$ |
| 25 | $95.1 \pm 0.8$ | $94.3 \pm 1.6$ | $93.7 \pm 1.2$ |
| 33 | $96.0 \pm 0.4$ | $95.7 \pm 0.6$ | $93.9 \pm 0.7$ |
| 59 | $96.1 \pm 0.4$ | $95.9 \pm 1.4$ | $95.0 \pm 0.6$ |
| 84 | $96.9 \pm 0.8$ | $96.4 \pm 1.5$ | $96.3 \pm 0.9$ |

Table : Areas under the ROC curve (with confidence intervals), averaged over ten experiments. Differences between our moment-based imbalanced binary classifier and subsampling results are statistically significant (with $p$-value $p < 0.01$).

A probabilistic approach to binary classification
Other elements to take into account
Numerica results and comparisons
Conclusion

Some recalls
Small experiment : protein classification
Large experiment : text classification

## Reuters results                                                           II

| Topic | This work | Cost-sensitive | Speed-up |
|-------|-----------|----------------|----------|
| 2     | 33        | 1088           | $33\times$ |
| 9     | 49        | 1451           | $29\times$ |
| 25    | 56        | 1211           | $21\times$ |
| 33    | 74        | 1788           | $24\times$ |
| 59    | 62        | 1299           | $21\times$ |
| 84    | 56        | 2056           | $36\times$ |

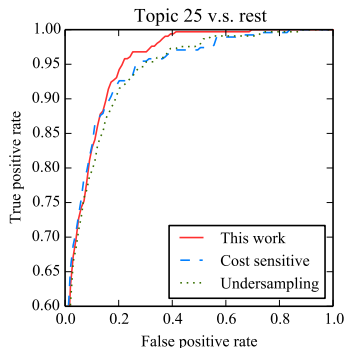Table : Computational times, in milliseconds, required to solve one problem, averaged over ten experiments.
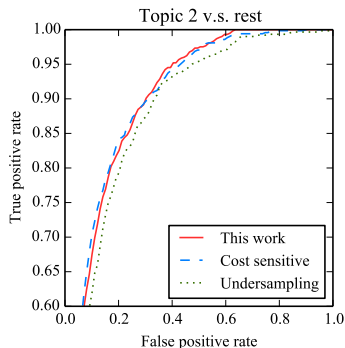
A probabilistic approach to binary classification
Other elements to take into account
Numerica results and comparisons
Conclusion

Some recalls
Small experiment : protein classification
Large experiment : text classification



Figure : ROC curves, averaged over ten experiments, on the REUTERS RCV1 dataset.

A probabilistic approach to binary classification
Other elements to take into account
Numerica results and comparisons
**Conclusion**

## Conclusion

- We give a theoretical interpretation for the one-class SVM method.
- We show how to adapt the one-class SVM in the case where we have first and second order information over the negative class.
- We apply this approach to imbalanced classification with good results both in precision and in computational speed.
- We apply this approach to large-scale imbalanced classification with significative speed improvement.

A probabilistic approach to binary classification
Other elements to take into account
Numerica results and comparisons
Conclusion

# The end

Thank you for your attention !