

# Learning Maximum Excluding Ellipsoids in Unbalanced Scenarios with Theoretical Guarantees

G. Metzler<sup>1,3</sup>, X. Badiche<sup>3</sup>, B. Belkasmi<sup>3</sup>, S. Canu<sup>2</sup>, E. Fromont<sup>1</sup>,  
A. Habrard<sup>1</sup> and M. Sebban<sup>1</sup>

1. Univ. Lyon, Univ. St-Etienne F-42000, UMR CNRS 5516, Laboratoire Hubert-Curien  
2. LITIS EA 4108, Univ. Rouen 76800 St-Etienne du Rouvray  
3. BLITZ BUSINESS SCEB, 38090 Villefontaine

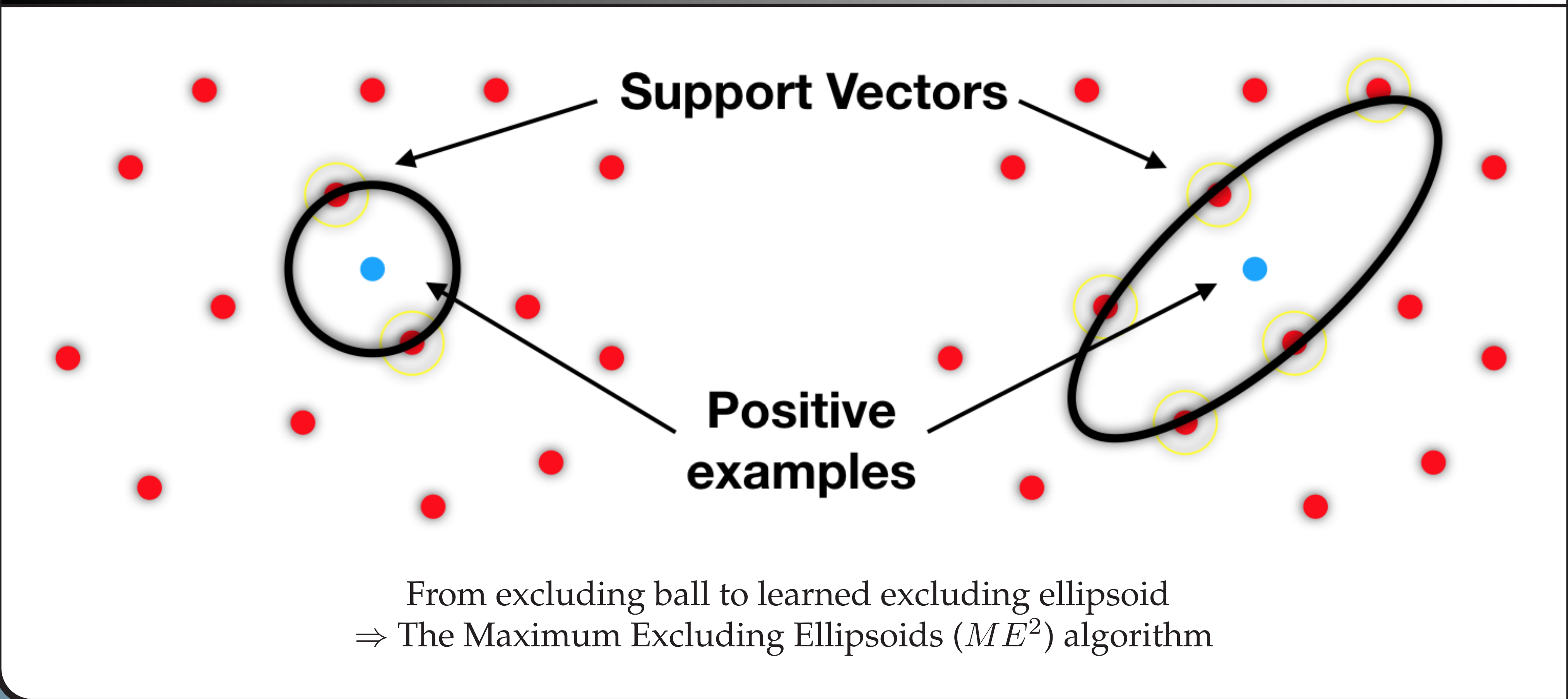
## ABSTRACT

We propose a theoretically-founded method to learn maximum excluding balls in the context of unbalanced binary classification. The objective is to learn a set of local balls centered at the minority class examples which exclude the examples of the majority class. Our contribution is twofold:

- 1) we address this problem from a metric learning point of view [2],
- 2) we derive generalization guarantees on the radius of the sphere and the learned metric using the uniform stability framework [3]

Our experimental evaluation on classic benchmarks shows the effectiveness of our approach.

## ILLUSTRATION



## METHOD

We learn a PSD matrix  $M$  of a Mahalanobis-distance defined as follows:

$$\|x - c\|_M^2 = (x - c)^T M (x - c).$$

$$\begin{aligned} \min_{R, M, \xi} \quad & \frac{1}{n} \sum_{i=1}^n \xi_i + \mu(B - R)^2 + \lambda \|M - I\|_F^2 \\ \text{s.t.} \quad & \|x_i - c\|_M^2 \geq R - \xi_i, \quad \forall i = 1, \dots, n, \\ & \xi_i \geq 0, \\ & B \geq R \geq 0, \end{aligned} \quad (1)$$

where,

- $B$  is the bound of the radius,
- $\mu$  controls the size of the ellipsoid,
- $\lambda$  controls the distortion w.r.t. an Euclidean ball.

## DUAL FORM

$$\begin{aligned} \min_{\alpha, \beta, \delta} \quad & \alpha^T \left( \frac{1}{4\lambda} G' + \frac{1}{4\mu} \mathbb{1}_{d \times d} \right) \alpha + \frac{\beta^2}{4\mu} + \frac{\delta^2}{4\mu} + \\ & \alpha^T \left( \text{diag}(G) - \left( B + \frac{\beta}{2\mu} - \frac{\delta}{2\mu} \right) \mathbb{1}_d \right) \\ & + \beta \left( B - \frac{\delta}{2\mu} \right), \\ \text{s.t.} \quad & 0 \leq \alpha_i \leq \frac{1}{n}, \quad \forall i = 1, \dots, n, \\ & \beta, \delta \geq 0, \end{aligned}$$

The Dual Problem gives an explicit expression of both the **Radius** and the **Metric**:

$$\begin{aligned} R &= \frac{\beta - \delta + 2\mu B - \sum_{i=1}^n \alpha_i}{2\mu}, \\ M &= I + \frac{1}{2\lambda} \sum_{i=1}^n \alpha_i (x_i - c)(x_i - c)^T. \end{aligned}$$

This expression directly shows that  $M$  is PSD.

## THEORETICAL BOUND

**Definition:** A learning algorithm has a uniform stability in  $\frac{\gamma}{n}$  w.r.t. a loss function  $\ell$  and a parameter set  $\theta$ , with  $\gamma$  a positive constant if:

$$\forall S, \forall i, 1 \leq i \leq n, \sup_x |\ell(\theta_S, x) - \ell(\theta_{S^i}, x)| \leq \frac{\gamma}{n}.$$

**Theorem:** Let  $\delta > 0$  and  $n > 1$ . There exists a constant  $\kappa > 0$ , such that with probability at least  $1 - \delta$  over the random draw over  $S$ , we have for any  $(M, R)$  solution of Problem (1):

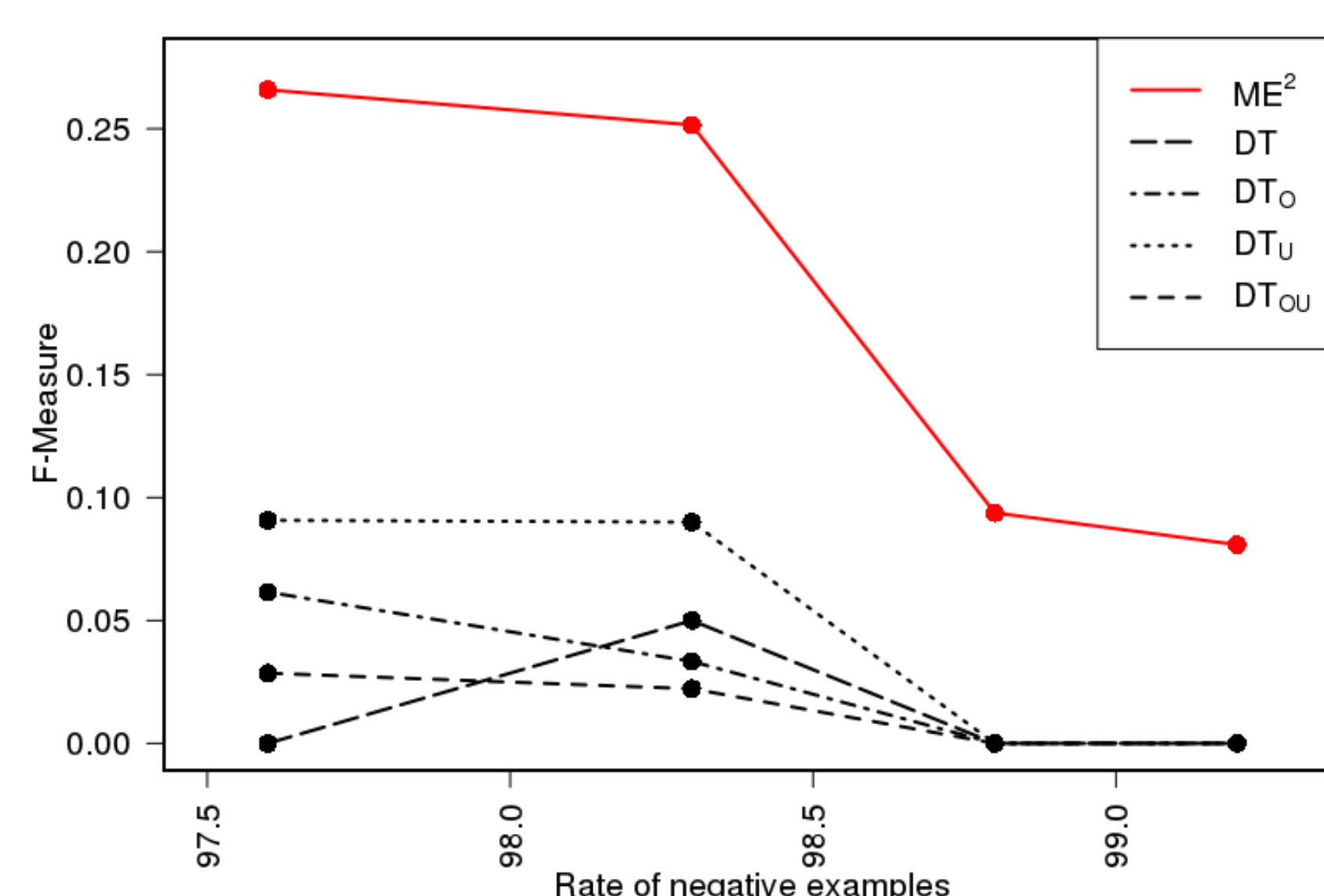
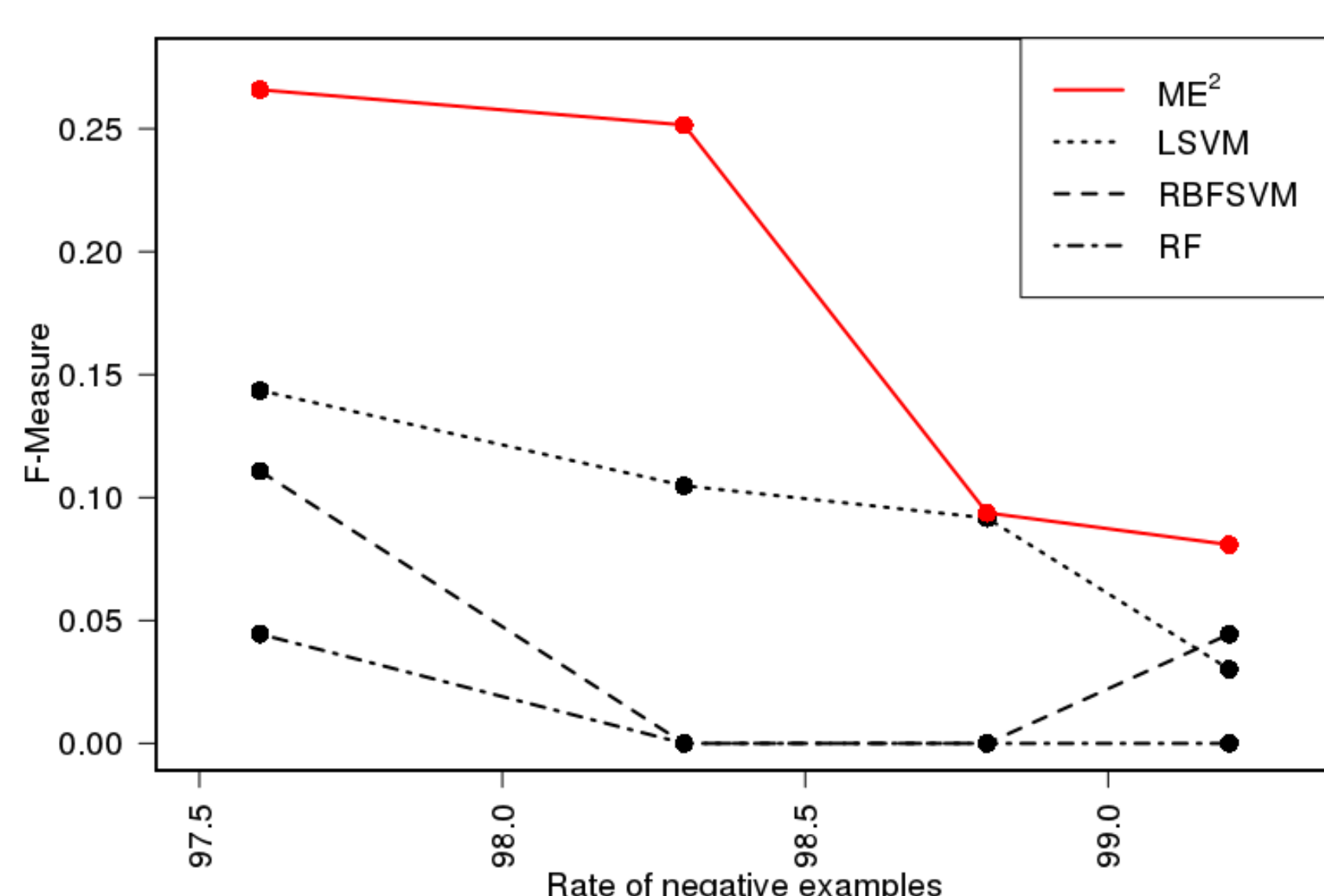
$$\begin{aligned} L(M, R) &\leq \hat{L}_S(M, R) + \frac{4 \max(1, 4B^2)}{n\kappa \min(\mu, \lambda)} \\ &+ \left( \frac{8 \max(1, 4B^2)}{\kappa \min(\mu, \lambda)} + B \right) \sqrt{\frac{\ln 1/\delta}{2n}} \\ &+ \left( 4B^2 \sqrt{\frac{\mu B^2}{\lambda} + d} \right) \sqrt{\frac{\ln 1/\delta}{2n}} \end{aligned}$$

## RESULTS

Algo.	Data	Abalone 10.7%	Abalone17 2.5%	Yeast6 2.4%	Abalone20 1.4%	Abalone19 0.76%
RF (10 trees)		0.67	0.20	0.04	0.00	0.00
$DT$ (simple version)		<b>0.71</b>	0.00	0.00	0.00	0.00
$DT_O$ (oversampling)		0.67	0.35	0.09	0.018	0.02
$DT_U$ (undersampling)		0.69	0.33	0.09	0.18	0.00
$DT_{OU}$ (both)		0.62	0.31	0.17	0.15	0.04
LSVM		0.62	0.29	0.15	<b>0.21</b>	0.04
RBFSVM		0.63	0.17	0.09	0.00	0.00
$ME^2$		0.62	<b>0.37</b>	<b>0.26</b>	<b>0.21</b>	<b>0.04</b>

Comparison in terms of  $F$ -Measure with some state of the art algorithms.

Datasets are ordered w.r.t. to an increasing imbalance ratio. The same global weight is given to both classes to learn the LSVM and RBF-SVM.



Effectiveness of  $ME^2$  when the disequilibrium increases. We can note that some of the state of the art methods lead to a null  $F$ -Measure.

## CONCLUSION

The  $ME^2$  algorithm presents the following advantages:

- Captures non linearity via local linear models
- Theoretically founded (uniform stability)
- Models can be learned in parallel
- Promising results in unbalanced scenarios

**Perspective:** study the link between the stability of the ellipsoids and the generalization error of a Nearest Neighbor algorithm.

## REFERENCES

### References

- [1] D. M. J. Tax and R. P. W. Duin, Support vector data description *Machine Learning Journal* (2013), vol.5, 287-364.
- [2] A. Bellet, A. Habrard and M. Sebban, A survey on metric learning for feature vectors and structured data, *arXiv*, (2013)
- [3] O. Bousquet and A. Elisseeff, Stability and generalization, *Journal of Machine Learning Research*, (2002) 499-526.