

From Cost-Sensitive Classification to Tight F-measure Bounds

K. Bascol^{1,3}, R. Emonet¹, E. Fromont², A. Habrard¹, G. Metzler^{1,4}, and M. Sebban¹

1. Univ. Lyon, UJM-Saint-Etienne, Laboratoire Hubert Curien UMR CNRS 5516, F-42023, Saint-Etienne, France

2. Univ. Rennes 1, IRISA/Inria, 35042 Rennes cedex, France

3. BLUECIME inc., France and 4. BLITZ inc., France

Abstract

In an **imbalanced** setting:

→ optimizing the classical accuracy tends to predict only the majority class;

→ optimizing imbalance-proof measures (as the F-Measure) is a tough task due to its non-convexity;

⇒ approximate F-measure optimization by cost-sensitive approach.

We propose to:

- Write the difference of F-measures between the errors made by two hypotheses.
- Give an upper bound on the optimal reachable F-measure given the error made by the classifier and the used cost sensitive parameters.
- **CONE**, an algorithm to iteratively optimize the F-measure.

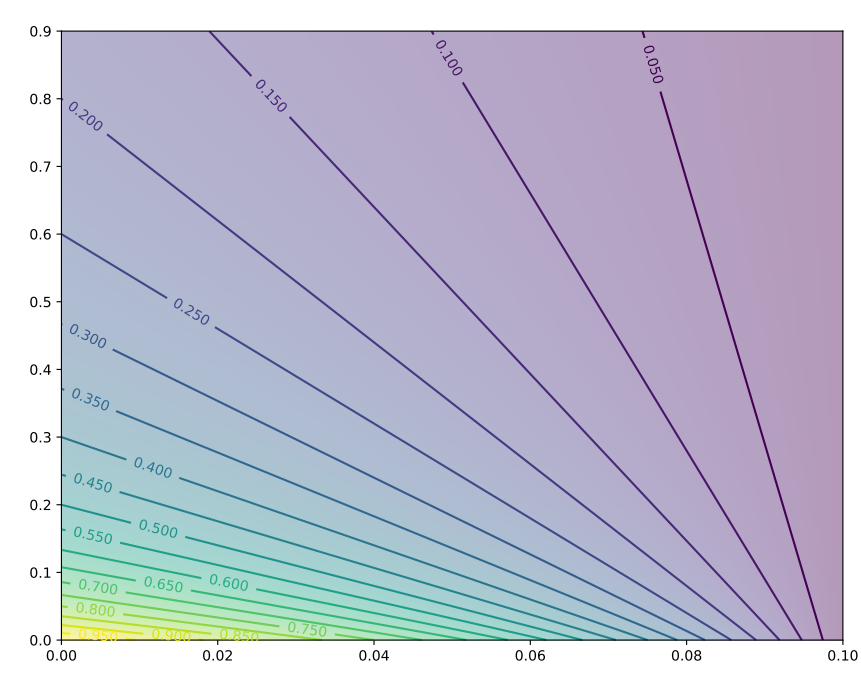
Notations and Base Result

Binary Classification

- $e = (e_1, e_2) = (FN, FP)$ the error profile obtained from h .
- A function to assign costs on each class $a(t) = (1 + \beta^2 - t, t)$.
- $F(e) = \frac{(1 + \beta^2)(P - e_1)}{(1 + \beta^2)P - e_1 + e_2}$ the associated F-measure.

Property F-measure

The level sets of the F-measure are hyperplanes: given $t \in [0, 1]$, $F(e) = t$ if and only if $\exists a, b$, two functions such that $\langle a(t), e \rangle + b(t) = 0$.



A bound on the F-measure

Step 1: impact of a change in the error profile. Given two error profiles e and e' and the previous property of the F-measure:

$$\begin{aligned} \langle a(F(e')), e - e' \rangle &= \langle a(F(e')), e \rangle + b(F(e')), \\ &= (F(e') - F(e)) \cdot ((1 + \beta^2)P_1 - e_1 + e_2), \end{aligned}$$

which leads to:

$$F(e') - F(e) = \Phi_e \cdot \langle a(F(e')), e - e' \rangle, \quad (1)$$

Step 2: bounding the difference of F-measures. Suppose that a classifier trained with $a(t)$ leads to e and $F(e)$ and consider e' obtained from an hypothetical classifier learned with $a(t')$. Then, from Eq.(1), we have:

$$\begin{aligned} F(e') - F(e) &= \Phi_e (\langle a(t'), e \rangle - \langle a(t'), e' \rangle), \\ &\leq \Phi_e (\langle a(t), e' \rangle + \varepsilon_1 - \langle a(t'), e' \rangle + (t' - t)(e_2 - e_1)), \\ &\leq \Phi_e \varepsilon_1 + \Phi_e \cdot (e_2 - e_1 - (e'_2 - e'_1))(t' - t), \end{aligned}$$

where ε_1 : *sub-optimality* of the learned classifier w.r.t. the 0-1 loss

$$\langle a(t), e \rangle \leq \varepsilon_1 + \min_{e' \in \mathcal{E}(\mathcal{H})} \langle a(t), e' \rangle$$

→ $e' = (e'_1, e'_2)$ is unknown → bound it such that $F(e') > F(e)$.

CONE: a Bound Driven Search Algorithm

Proposition. Let e be the error profile obtained with a classifier trained with the parameter t , $F(e)$ its associated F-measure value, Φ_e as defined in Eq. (1), and $\varepsilon_1 > 0$ the sub-optimality of our linear classifier.

Then for all $t' < t$:

$$F(e') \leq F(e) + \Phi_e \varepsilon_1 + \Phi_e \cdot (e_2 - e_1 - M_{max})(t' - t),$$

$$\text{where } M_{max} = \max_{\substack{e'' \in \mathcal{E}(\mathcal{H}) \\ s.t. F(e'') > F(e)}} (e''_2 - e''_1)$$

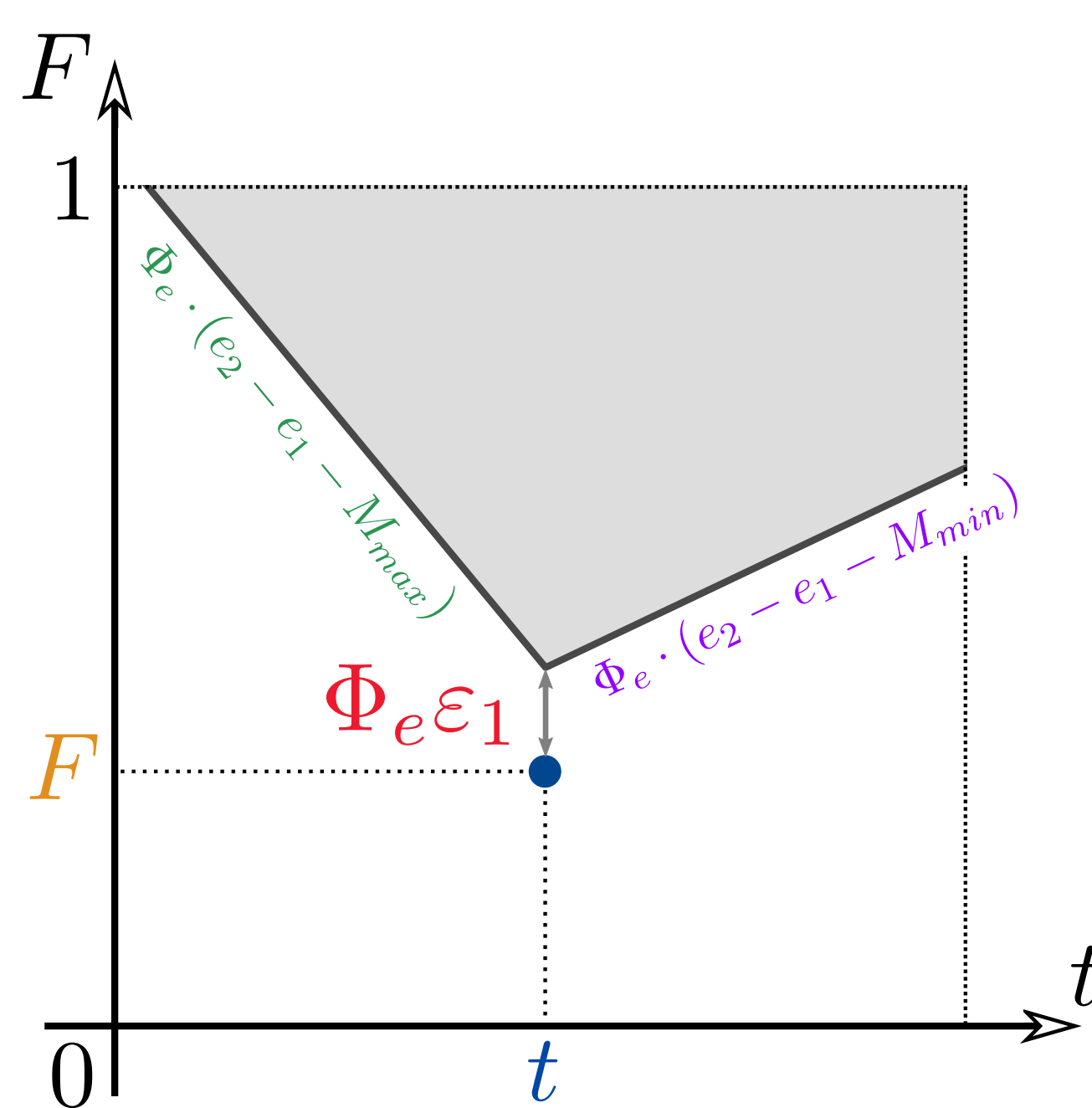
and, for all $t' > t$:

$$F(e') \leq F(e) + \Phi_e \varepsilon_1 + \Phi_e \cdot (e_2 - e_1 - M_{min})(t' - t),$$

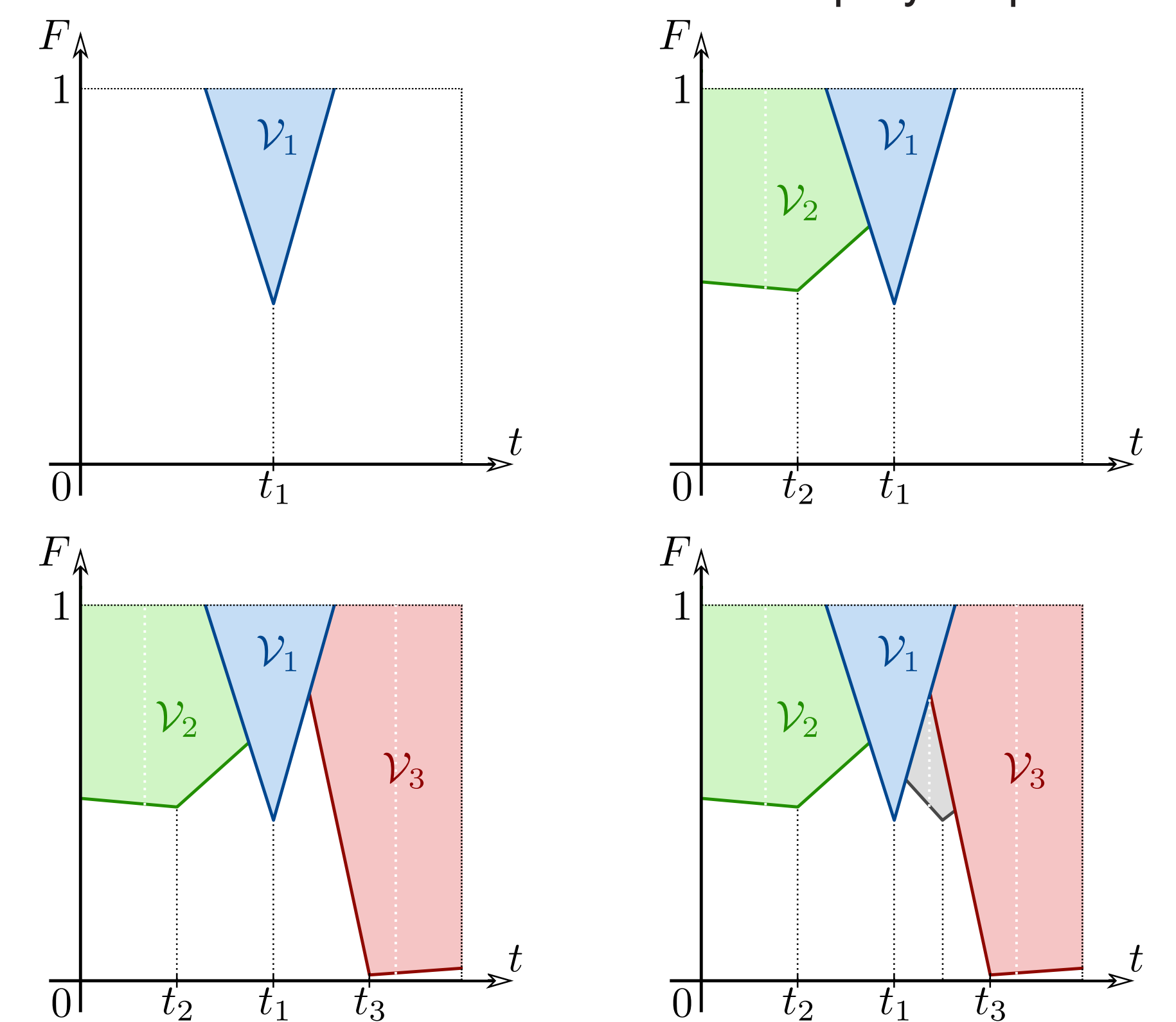
$$\text{where } M_{min} = \min_{\substack{e'' \in \mathcal{E}(\mathcal{H}) \\ s.t. F(e'') > F(e)}} (e''_2 - e''_1).$$

A Geometric Interpretation

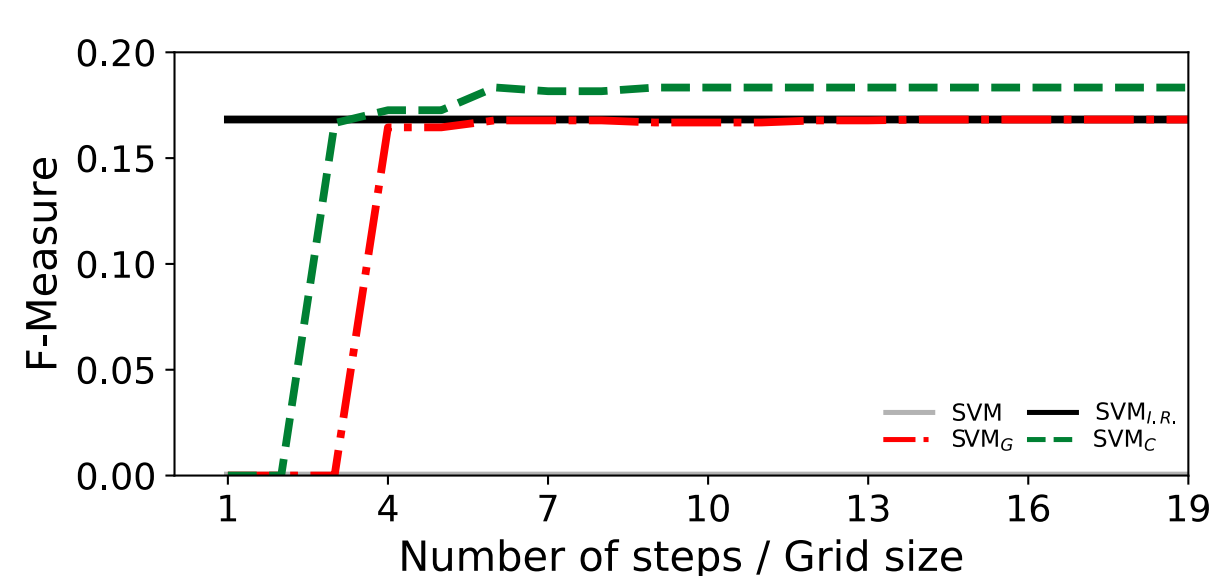
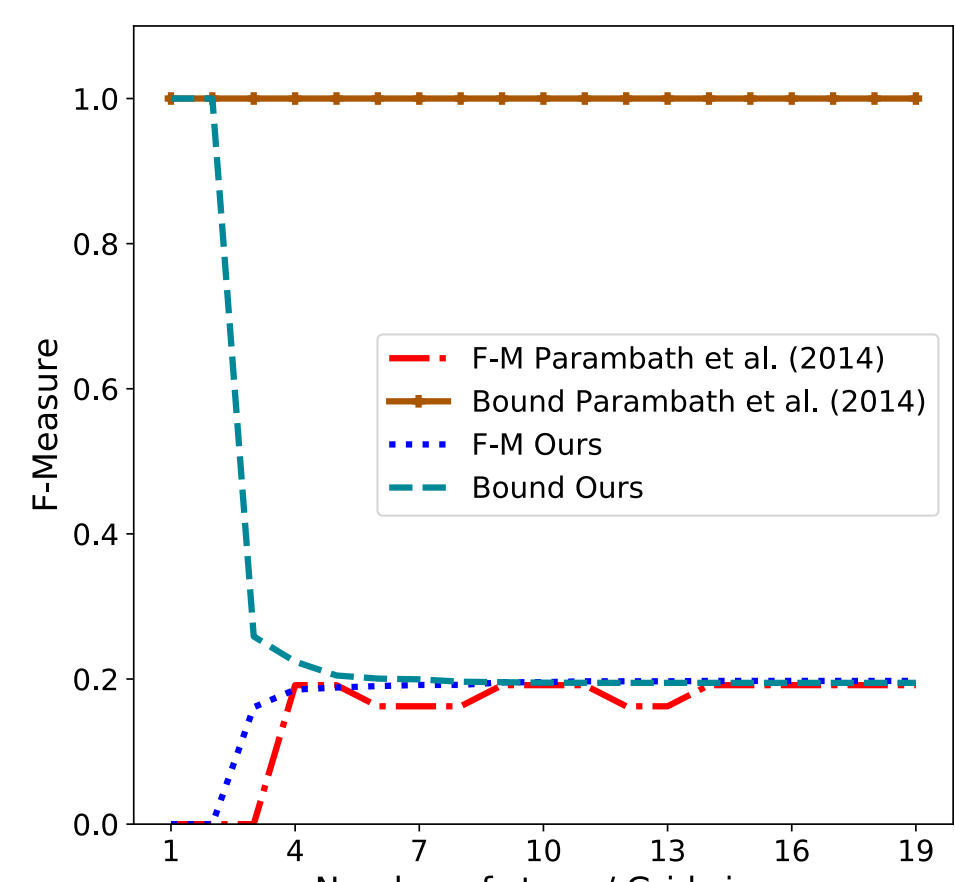
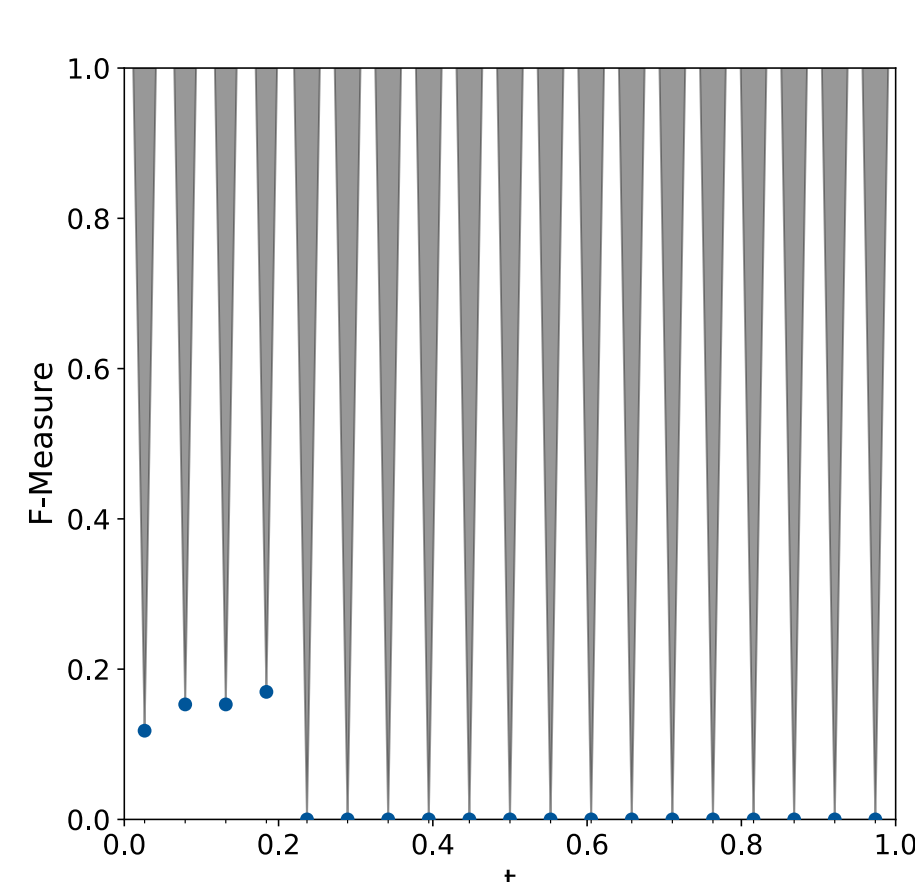
Our bound on $F(e')$ can be seen the unreachable region of F-measure in the (t, F) -space.



An illustration of CONE step by step



Practical Evaluation of Theoretical Guarantees



On Abalone12, in function of the number of **CONE** steps and [1] grid size:
left: on train set, evolution of the F-Measure and of the considered bound
right: on test set, evolution of the F-Measure

Dataset	SVM	SVM _{I.R.}	SVM _G	SVM _C	SVM _T	LR ^T	LR _{I.R.} ^T	LR _G ^T	LR _B
Adult	62.5 (0.2)	64.9 (0.3)	66.4 (0.1)	66.5 (0.1)	66.4 (0.1)	66.5 (0.1)	66.5 (0.1)	66.5 (0.1)	66.6 (0.1)
Abalone10	0.0 (0.0)	30.9 (1.2)	32.4 (1.3)	32.2 (0.8)	31.8 (1.9)	30.8 (2.2)	30.7 (1.9)	30.7 (1.9)	31.6 (0.6)
Satimage	0.0 (0.0)	23.4 (4.3)	20.4 (5.3)	20.6 (5.6)	30.9 (2.0)	21.2 (11.1)	28.6 (1.9)	28.6 (1.9)	21.4 (4.6)
IJCNN	44.5 (0.4)	53.3 (0.4)	61.6 (0.6)	61.6 (0.6)	62.6 (0.4)	59.4 (0.5)	56.5 (0.3)	56.5 (0.3)	59.2 (0.3)
Abalone12	0.0 (0.0)	16.8 (2.7)	16.8 (4.2)	18.3 (3.3)	16.3 (3.0)	15.5 (3.1)	17.0 (3.3)	17.0 (3.3)	17.7 (3.7)
Pageblocks	48.1 (5.8)	39.6 (4.7)	66.4 (3.2)	62.8 (3.9)	67.6 (4.0)	59.2 (8.1)	55.9 (6.4)	55.9 (6.4)	55.7 (5.7)
Yeast	0.0 (0.0)	29.4 (2.9)	38.6 (7.1)	39.0 (7.5)	35.4 (15.6)	37.4 (10.1)	39.9 (6.5)	27.6 (6.8)	27.6 (6.8)
Wine	0.0 (0.0)	15.6 (5.2)	20.0 (6.4)	22.7 (6.0)	19.3 (7.9)	21.5 (3.7)	25.2 (4.5)	25.2 (4.5)	18.3 (7.2)
Average	19.4 (0.8)	34.2 (2.8)	40.3 (3.5)	40.5 (3.5)	41.3 (4.4)	38.9 (5.2)	40.0 (3.1)	38.5 (3.2)	37.3 (3.6)

F-Measure for Logistic Regression (LR) [3] and SVM algorithms (averaged over 5 experiments).

"I.R." I.R. based class costs; "G": Reproduction of [1]; "C": CONE "T" thresholded predictions; "B": Bisection [2]

Examples of runs on Abalone12
top: wrapper from [1]; bottom:
CONE;
both with SVM classifier ($C = 1$).

Perspectives

- Extend our study to more complex class of hypotheses (non linear hypotheses such as neural networks).
- Prove the convergence of our algorithm.
- Work on the notion of sub-optimality to improve the bound.
- Work on a generalization bound on the F-measure.

Acknowledgements



References

- [1] S. P. Parambath, N. Usunier, and Y. Grandvalet, *Optimizing F-Measures by Cost-Sensitive Classification*, **NIPS** 2014.
- [2] H. Narasimhan, H. Ramaswamy, A. Saha, and S. Agarwal, *Consistent Multiclass Algorithms for Complex Performance Measures*, **ICML** 2015.
- [3] O. Koyejo, N. Natarajan, P. Ravikumar, and I. Dhillon, *Consistent Binary Classification with Generalized Performance Metrics*, **NIPS** 2014.