

**APPLICATIONS OF
REGULARIZED LEAST SQUARES
TO CLASSIFICATION PROBLEMS**

Nicolò Cesa-Bianchi
Università di Milano, Italy

Contributors: Alex Conconi, Claudio Gentile, Luca Zaniboni

DESIDERATA FOR LEARNING ALGORITHMS

- Good empirical performance
- Scalability
- Versatility
- Theoretical guarantees

CONTENTS

Regularized least squares and second-order Perceptron

Mistake bounds (for individual data sequences)

Risk bounds (for probabilistic data sequences)

Application to semi-supervised learning

BINARY PATTERN CLASSIFICATION

Each data element is encoded as an attribute vector

$$\boldsymbol{x} = (x_1, \dots, x_d) \in \mathbb{R}^d \quad (\text{instance})$$

A binary label $y \in \{-1, +1\}$ classifies each instance \boldsymbol{x} according to a given semantic property

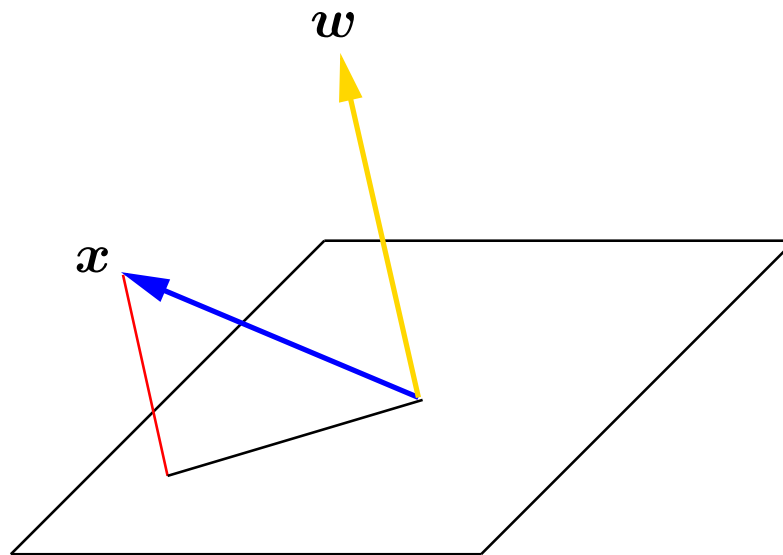
Design a rule $h : \mathbb{R}^d \rightarrow \{-1, +1\}$ for predicting labels given instances

No assumptions on the data-generating mechanism

LINEAR-THRESHOLD CLASSIFIERS

Classifier parametrized by coefficients (weights) $\mathbf{w} \in \mathbb{R}^d$

$$h_{\mathbf{w}}(\mathbf{x}) = \text{SGN}(\mathbf{w}^{\top} \mathbf{x}) \quad h_{\mathbf{w}} \text{ errs on } (\mathbf{x}, y) \quad \text{iff} \quad y \mathbf{w}^{\top} \mathbf{x} < 0$$



RLSC vs. SVM

Training set $(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_n, y_n)$ $\mathbf{x}_t \in \mathbb{R}^d$ $y_t \in \{-1, 1\}$

Margin on t -th example: $m_t(\mathbf{v}) = y_t \mathbf{v}^\top \mathbf{x}_t$

1-norm SVM classifier

(hinge loss)

$$\mathbf{w} = \operatorname{argmin}_{\mathbf{v} \in \mathbb{R}^d} \left(\sum_{t=1}^n [1 - m_t(\mathbf{v})]_+ + a \|\mathbf{v}\|^2 \right)$$

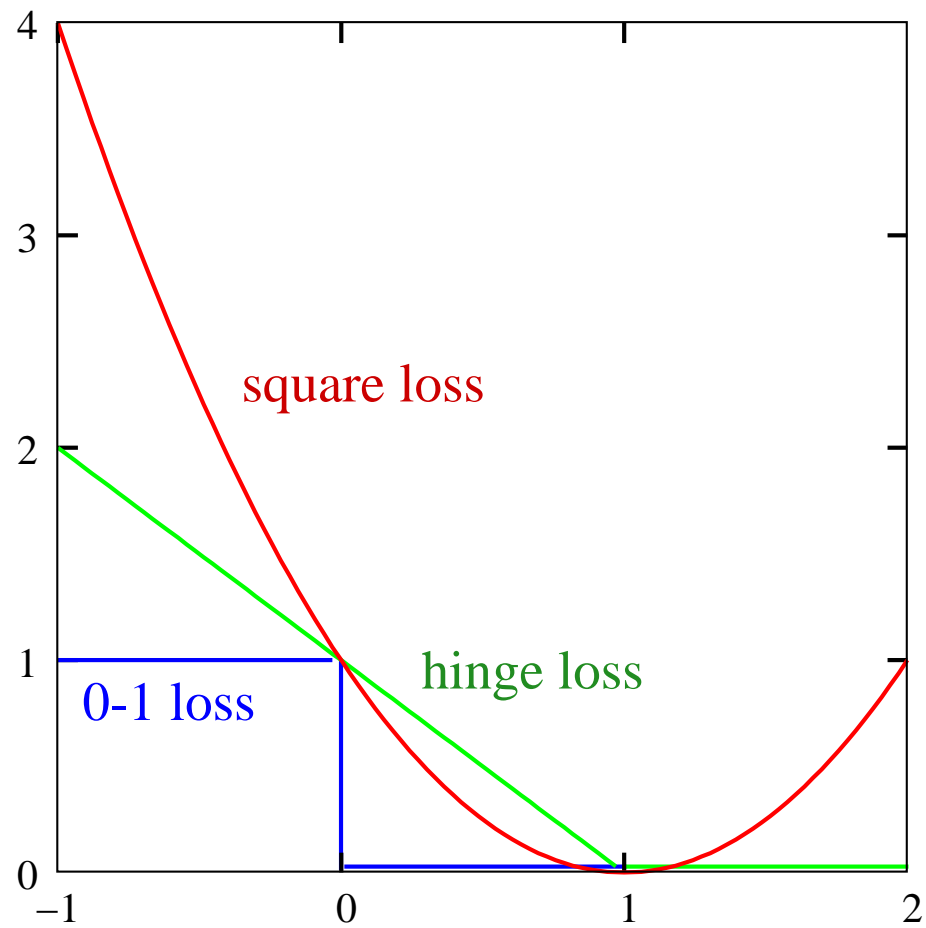
Regularized Least Squares for Classification

(square loss)

$$\mathbf{w} = \operatorname{argmin}_{\mathbf{v} \in \mathbb{R}^d} \left(\sum_{t=1}^n (1 - m_t(\mathbf{v}))^2 + a \|\mathbf{v}\|^2 \right)$$

2-norm SVM classifier $\left([1 - m_t(\mathbf{v})]_+ \right)^2$

LOSSES AS FUNCTIONS OF MARGIN



RLSC IS SOLVED ANALYTICALLY

$$\mathbf{w} = \operatorname{argmin}_{\mathbf{v} \in \mathbb{R}^d} \left(\sum_{t=1}^n (1 - m_t(\mathbf{v}))^2 + a \|\mathbf{v}\|^2 \right)$$

Closed form solution $\mathbf{w} = (aI + S S^\top)^{-1} S \mathbf{y}$

- I = identity matrix
- $S = [\mathbf{x}_1 \dots \mathbf{x}_n]$ matrix of instances
- $\mathbf{y} = (y_1, \dots, y_n)$ vector of labels

Solution is not sparse...

ENFORCING SPARSITY IN RLSC

Second-order Perceptron

Initial classifier $\mathbf{w} = (0, \dots, 0)$

For $t = 1, 2, \dots$

1. get next example (\mathbf{x}_t, y_t)
2. if $\text{SGN}(\mathbf{w}^\top \mathbf{x}_t) \neq y_t$ then perform **RLSC-update** of \mathbf{w} :

$\mathbf{S} = [\mathbf{x}_{s_1}, \dots, \mathbf{x}_{s_m}, \mathbf{x}_t]$ matrix of **mistaken** instances

$\mathbf{y} = (y_{s_1}, \dots, y_{s_m}, y_t)$ vector of **mistaken** labels

$$\mathbf{w} = (aI + \mathbf{S}\mathbf{S}^\top)^{-1} \mathbf{S}\mathbf{y}$$

REMARKS

Algorithms works in dual variables (kernels OK!)

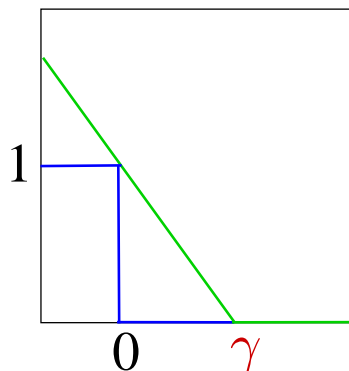
Final classifier only depends on mistaken examples:

$$\text{mistakes} \approx \text{support vectors}$$

Incremental update of $(I + S S^\top)^{-1}$ in time $\Theta(\text{mistakes}^2)$ when working with dual variables

Strong theoretical guarantees...

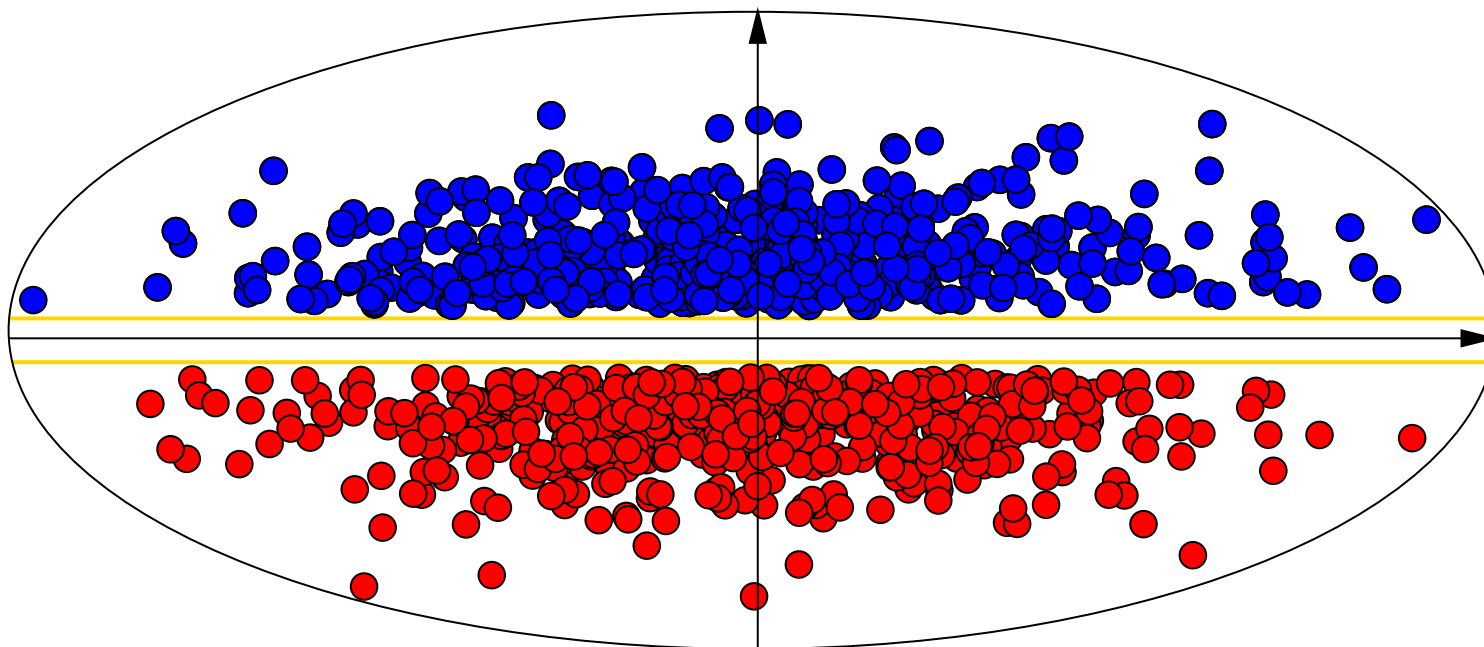
BOUND ON UPDATES (MISTAKES)



$$\inf_{\gamma > 0} \min_{\|\mathbf{u}\|=1} \frac{1}{\gamma} \left(\sum_{t=1}^n [\gamma - m_t(\mathbf{u})]_+ + \sqrt{(1 + \lambda(\mathbf{u})) \sum_{i=1}^d \ln(1 + \lambda_i)} \right)$$

- $\lambda_1, \dots, \lambda_d$ are the eigenvalues of $S S^\top$ $\lambda_{\min} \leq \lambda(\mathbf{u}) \leq \lambda_{\max}$
- Bound holds on any individual data sequence!

WHERE SECOND ORDER PROVABLY WINS



Mistakes of first-order Perceptron: $\leq \lambda_{\max}/\text{margin}^2$

Mistakes of second-order Perceptron: $\leq \lambda_{\min}/\text{margin}^2$

FROM MISTAKE BOUNDS TO RISK BOUNDS

Assume data (\mathbf{x}, y) are drawn i.i.d. from an unknown distribution on $\mathbb{R}^d \times \{-1, 1\}$ (statistical learning model)

Statistical risk of classifier h is $\text{RISK}(h) = \mathbb{P}(h(\mathbf{x}) \neq y)$

Risk bound for typical second-order classifier

$$\text{risk at most } \frac{M_n}{n} + \sqrt{\frac{2}{n} \ln \frac{1}{\delta}} \quad \text{with probability at least } 1 - \delta$$

M_n is number of updates (mistakes) made by 2nd order Perceptron on training set

HOW TO PICK A GOOD CLASSIFIER

h_0, h_1, \dots, h_m classifiers generated during a run

1. test each h_t on $(\mathbf{x}_{t+1}, y_{t+1}), (\mathbf{x}_{t+2}, y_{t+2}), \dots$
2. pick h_t minimizing a **penalized risk estimate**

Risk bound

$$\text{risk at most } \frac{M_n}{n} + 5\sqrt{\frac{1}{n} \ln \frac{2(n+1)}{\delta}} \quad \text{with probability at least } 1 - \delta$$

PREDICTION, LEARNING, AND GAMES

Repeated game theory

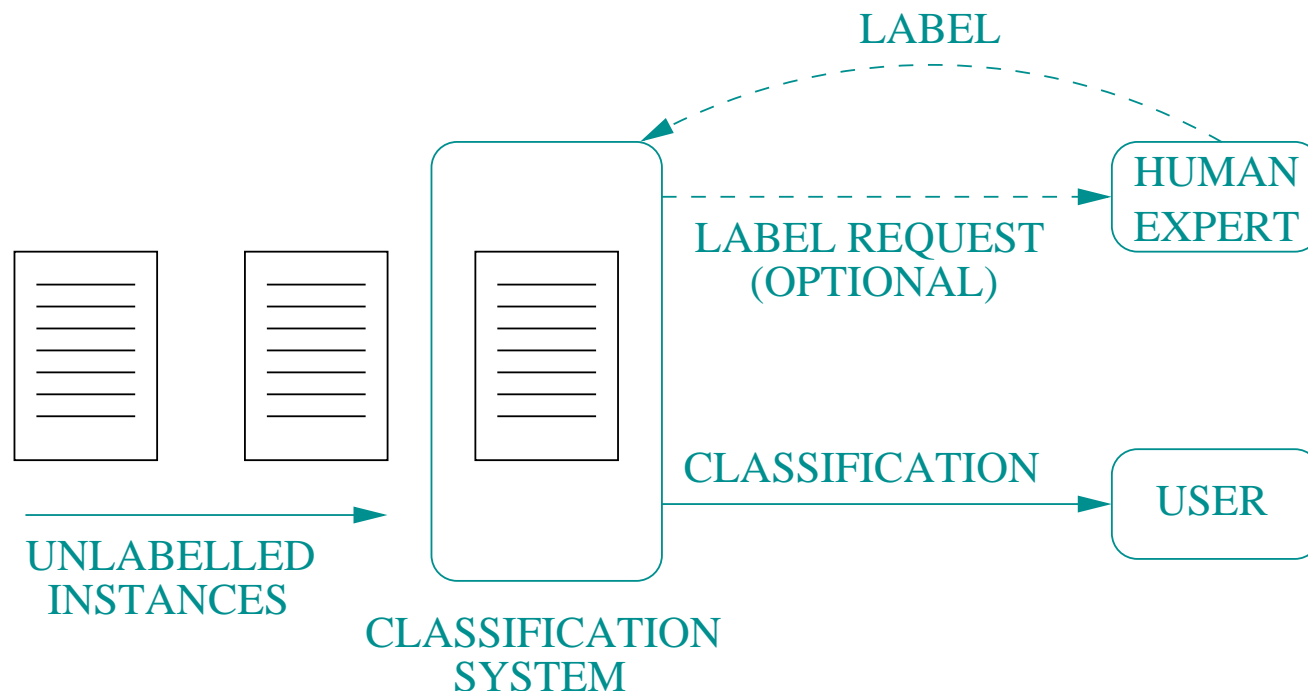
[Hannan, 1957]

Prediction with expert advice

[Littlestone and Warmuth, 1989; Vovk, 1990]

LEARNING	PREDICTION	GAMES
FEATURES	EXPERTS	PURE STRATEGIES
LINEAR CLASSIFIER	WEIGHTED FORECASTER	MIXED STRATEGY
PERCEPTRON TH.	EXPERTS TH.	HANNAN'S TH.
[1962]	[1993]	[1957]

SEMI-SUPERVISED LEARNING



Goal: Minimize classification mistakes and rate of requested labels

SUFFICIENT QUERY RATE

Measure performance against best classifier in a finite set (**experts model**)

Sampling data at random with rate $(\log n)(\log \log n)$ guarantees
convergence to best expert

A RLSC-BASED SEMI-SUPERVISED ALGORITHM

- Classify next instance \mathbf{x}_t with $\text{SGN}(\mathbf{w}^\top \mathbf{x}_t)$
- Compute $\Delta_t = \mathbf{w}^\top \mathbf{x}_t / \|\mathbf{x}_t\|$
- Query label y_t with probability

$$\frac{1}{1 + |\Delta_t|/C + \Theta(\Delta_t^2)}$$

- If label queried and $\text{SGN}(\Delta_t) \neq y_t$,
then perform RLSC-update of \mathbf{w}

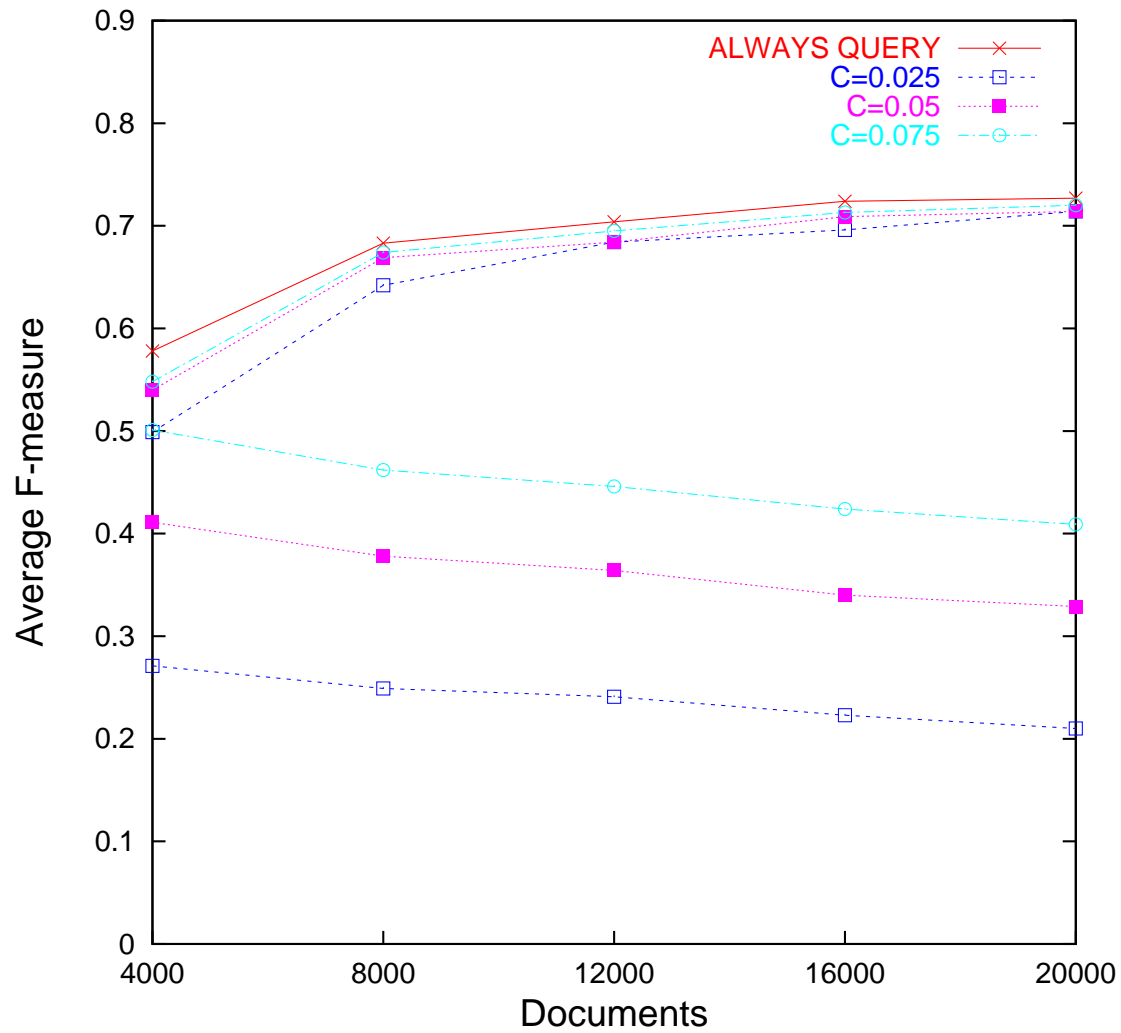
EXPECTED BOUNDS ON MISTAKES

Note: Only mistakes on queried labels trigger updates

$$\inf_{\gamma > 0} \min_{\|\mathbf{u}\|=1} \frac{1}{\gamma} \left(\sum_{t=1}^n [\gamma - m_t(\mathbf{u})]_+ + \sqrt{(1 + \mathbb{E}\Lambda(\mathbf{u})) \sum_{i=1}^d \mathbb{E} \ln(1 + \Lambda_i)} \right)$$

- Bound holds on any data sequence (proper tuning of C needed)
- Apparently, no loss of performance due to semi-supervision!

REUTERS — 50 CATEGORIES



A PROBABILISTIC MODEL FOR LABELS

Labels $y_t \in \{-1, +1\}$ drawn from parametric distribution

$$\mathbb{P}(y_t = 1 \mid \mathbf{x}_t) = \frac{1 + \mathbf{u}^\top \mathbf{x}_t}{2} \quad \|\mathbf{u}\| = 1$$

If \mathbf{w} is solution of RLSC problem, then

$$\mathbb{E}[\mathbf{w}^\top \mathbf{x}_t] = \mathbf{u}^\top \mathbf{x}_t + \text{smaller order terms}$$

We conjecture a Chernoff bound for $|\mathbf{w}^\top \mathbf{x}_t - \mathbf{u}^\top \mathbf{x}_t|$

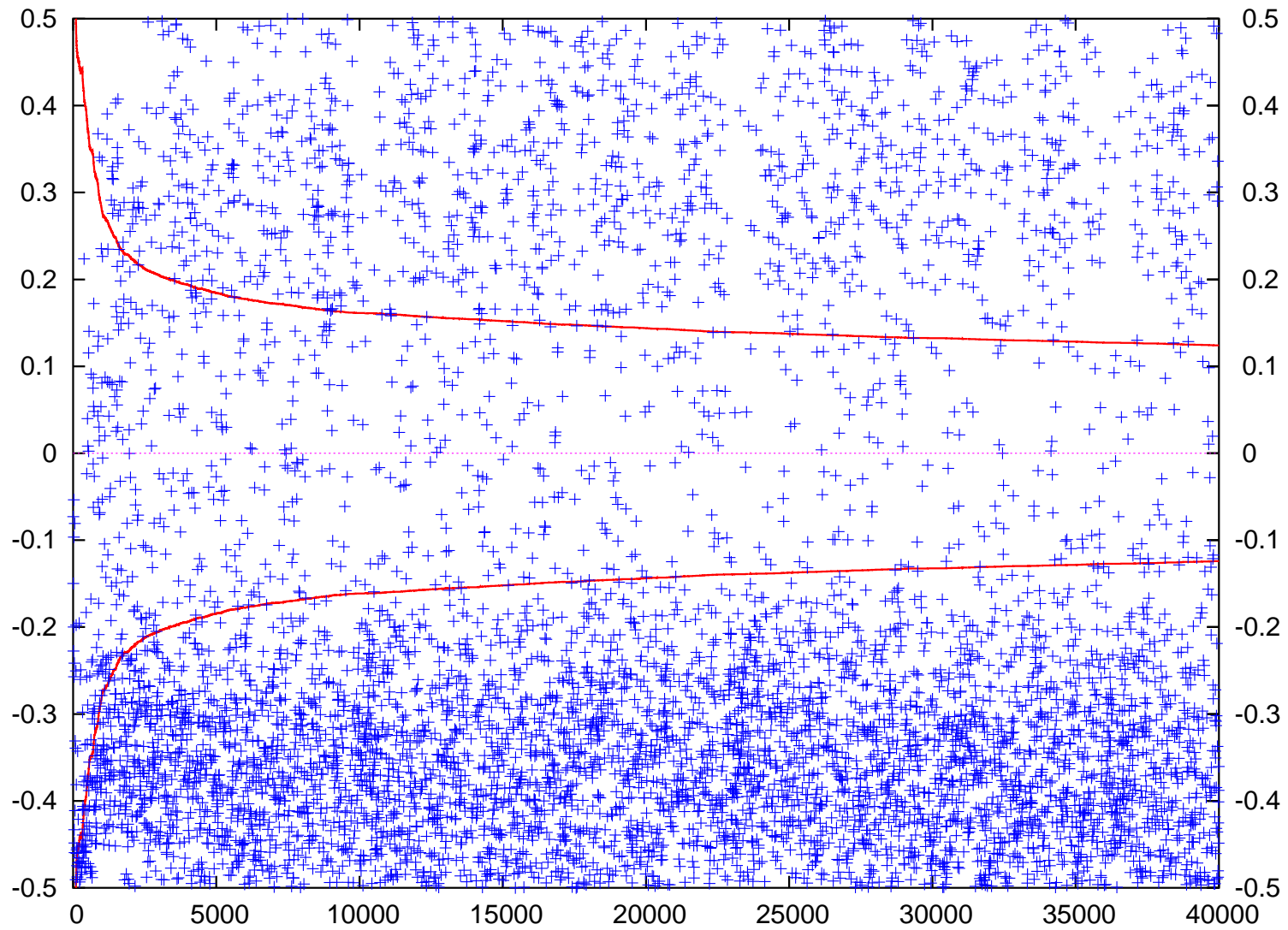
A SEMI-SUPERVISED DETERMINISTIC ALGORITHM

- Classify next instance \mathbf{x}_t with $\text{SGN}(\mathbf{w}^\top \mathbf{x}_t)$
- Compute $\Delta_t = \mathbf{w}^\top \mathbf{x}_t / \|\mathbf{x}_t\|$
- If $|\Delta_t| \leq \sqrt{\frac{4 \ln t}{N_t}}$ then query label y_t of \mathbf{x}_t
- If label queried then perform RLSC-update of \mathbf{w}

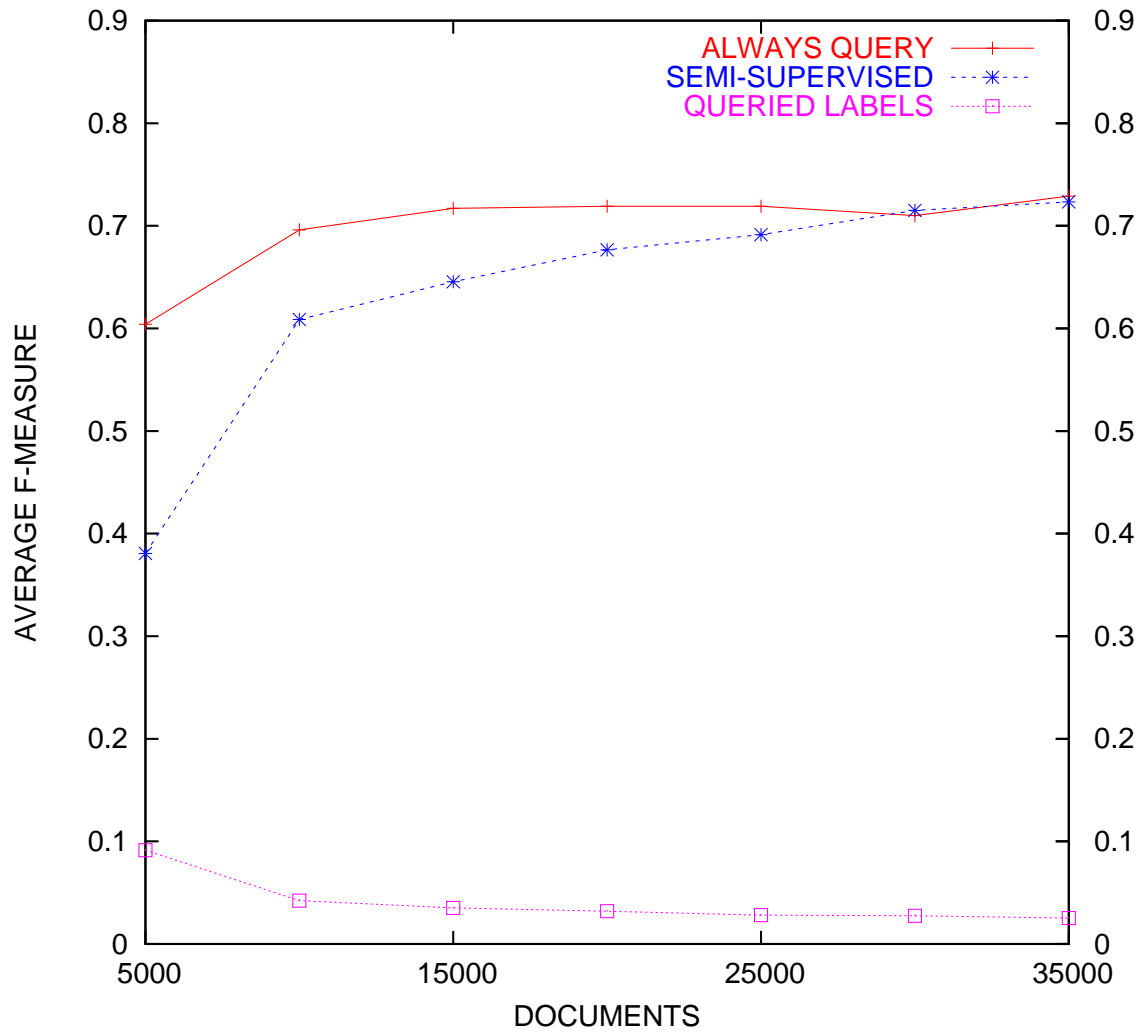
N_t = number of labels sampled so far

Schedule $\sqrt{(4 \ln t)/N_t}$ derived from conjectured Chernoff bound

MARGIN PLOT



REUTERS — 50 CATEGORIES



CONCLUSIONS

- Second-order Perceptron (incremental RLSC) is a practically reasonable algorithm with strong theoretical guarantees
- Trading margin with probability of mistake may save lots of labels in semi-supervised learning
- Applications of RLSC-based algorithms to [filtering](#) and [hierarchical classification](#)

REFERENCES

R. Rifkin, G. Yeo, and T. Poggio. [Regularized least squares classification](#). In *Advances in Learning Theory: Methods, Model and Applications*. IOS Press, 2003.

N. Cesa-Bianchi, A. Conconi and C. Gentile. [A second-order Perceptron algorithm](#). *SIAM J. on Computing*. To appear.

N. Cesa-Bianchi, A. Conconi and C. Gentile. [On the generalization ability of on-line learning algorithms](#). *IEEE Transactions on Information Theory*. To appear.

N. Cesa-Bianchi, A. Conconi and C. Gentile. [Learning probabilistic linear-threshold classifiers via selective sampling](#). COLT 2003, pages 373–386. Springer, 2003.

REFERENCES (CONT.)

N. Cesa-Bianchi, C. Gentile, and L. Zaniboni. [Worst-case analysis of selective sampling for linear-threshold algorithms](#). NIPS 17. MIT Press, to appear.

N. Cesa-Bianchi, A. Conconi, and C. Gentile. [Margin-based algorithms for information filtering](#). NIPS 15. MIT Press, 2003.

N. Cesa-Bianchi, A. Conconi and C. Gentile. [Regret bounds for hierarchical classification with linear-threshold functions](#). COLT 2004, pages 93–108. Springer, 2003.

N. Cesa-Bianchi, C. Gentile, A. Tironi, and L. Zaniboni. [Incremental Algorithms for Hierarchical Classification](#). NIPS 17. MIT Press, to appear.