

# Boosting Algorithms for Maximizing the Soft Margin

Manfred Warmuth<sup>1</sup>, Karen Glöcer<sup>1</sup>, Gunnar Rätsch<sup>2</sup>

1: UC Santa Cruz

2: Friedrich Miescher Laboratory of the Max Planck Society

manfred,kag@soe.ucsc.edu, gunnar.raetsch@tuebingen.mpg.de

## 1. Introduction

- Boosting algorithm when there is no consistent convex combination of base hypotheses
- In  $\Theta(\frac{1}{\delta^2} \log \frac{N}{\nu})$  iterations produces a convex combination with soft margin within  $\delta$  of the maximum

Boosting protocol:

- Set of examples  $S = \langle (x_1, y_1), \dots, (x_N, y_N) \rangle$
- Maintains distribution  $d$  on examples
- At iteration  $t$ :
  - Given current distribution  $d^{t-1}$ , oracle provides hypothesis  $h_t$  of edge  $\gamma_t = d^{t-1} \cdot u^t \geq g$ , where  $u_i^t = y_i h_t(x_i)$
  - **Guarantee**  $g > 0$  not known to algorithm
  - Update distribution  $d^{t-1}$  to  $d^t$

LPBoost computes  $d^t$  by solving:

Primal		Dual	
$\min_{d, \gamma}$	$\gamma$	$\max_{w, \rho}$	$\rho + \frac{1}{\nu} \sum_{n=1}^N \xi_n$
s.t. $d \cdot u^m \leq \gamma, 1 \leq m \leq t$ , $d \in \mathcal{P}^N, d \leq \frac{1}{\nu} \mathbf{1}$ .		s.t. $\sum_{i=1}^t w_i y_i h_i + \xi_n \geq \rho, 1 \leq n \leq N$ , $w \in \mathcal{P}^M, \psi \geq \mathbf{0}$ .	
minimize maximum edge		maximize minimum soft margin	

Non-standard LPBoost formulation

- **Totally corrective**
- Capping probabilities in primal  $\leftrightarrow$  soft margin in dual

## 2. LPBoost does not have $\Omega(\log N)$ iteration bounds.

- LPBoost (Schuurmans et al) works well in practice
- No bounds have been proved for it
- In our counter examples LPBoost takes  $\Omega(N)$  iterations to achieve margin precision  $\delta \approx 1$  for separable case.
  - Forces LPBoost to concentrate its distribution on single example
  - Holds regardless of LP optimization algorithm
  - Shows need for regularization

$n \setminus m$	1	2	3	4	5
1	+1	-1 + 5 $\epsilon$	-1 + 7 $\epsilon$	-1 + 9 $\epsilon$	-1 + $\epsilon$
2	+1	-1 + 5 $\epsilon$	-1 + 7 $\epsilon$	-1 + 9 $\epsilon$	-1 + $\epsilon$
3	+1	-1 + 5 $\epsilon$	-1 + 7 $\epsilon$	-1 + 9 $\epsilon$	-1 + $\epsilon$
4	+1	-1 + 5 $\epsilon$	-1 + 7 $\epsilon$	-1 + 9 $\epsilon$	-1 + $\epsilon$
5	-1 + 2 $\epsilon$	+1	-1 + 7 $\epsilon$	-1 + 9 $\epsilon$	+1 - $\epsilon$
6	-1 + 3 $\epsilon$	-1 + 4 $\epsilon$	+1	-1 + 9 $\epsilon$	+1 - $\epsilon$
7	-1 + 3 $\epsilon$	-1 + 5 $\epsilon$	-1 + 6 $\epsilon$	+1	+1 - $\epsilon$
8	-1 + 3 $\epsilon$	-1 + 5 $\epsilon$	-1 + 7 $\epsilon$	-1 + 8 $\epsilon$	+1 - $\epsilon$

- The counter example suggests that a good algorithm should employ two tricks:
  - Cap the weight on any example
  - Spread the weight on the examples via a regularization such as the relative entropy

These two tricks used by the SoftBoost algorithm make it possible to obtain iteration bounds that grow logarithmic in  $N$ .

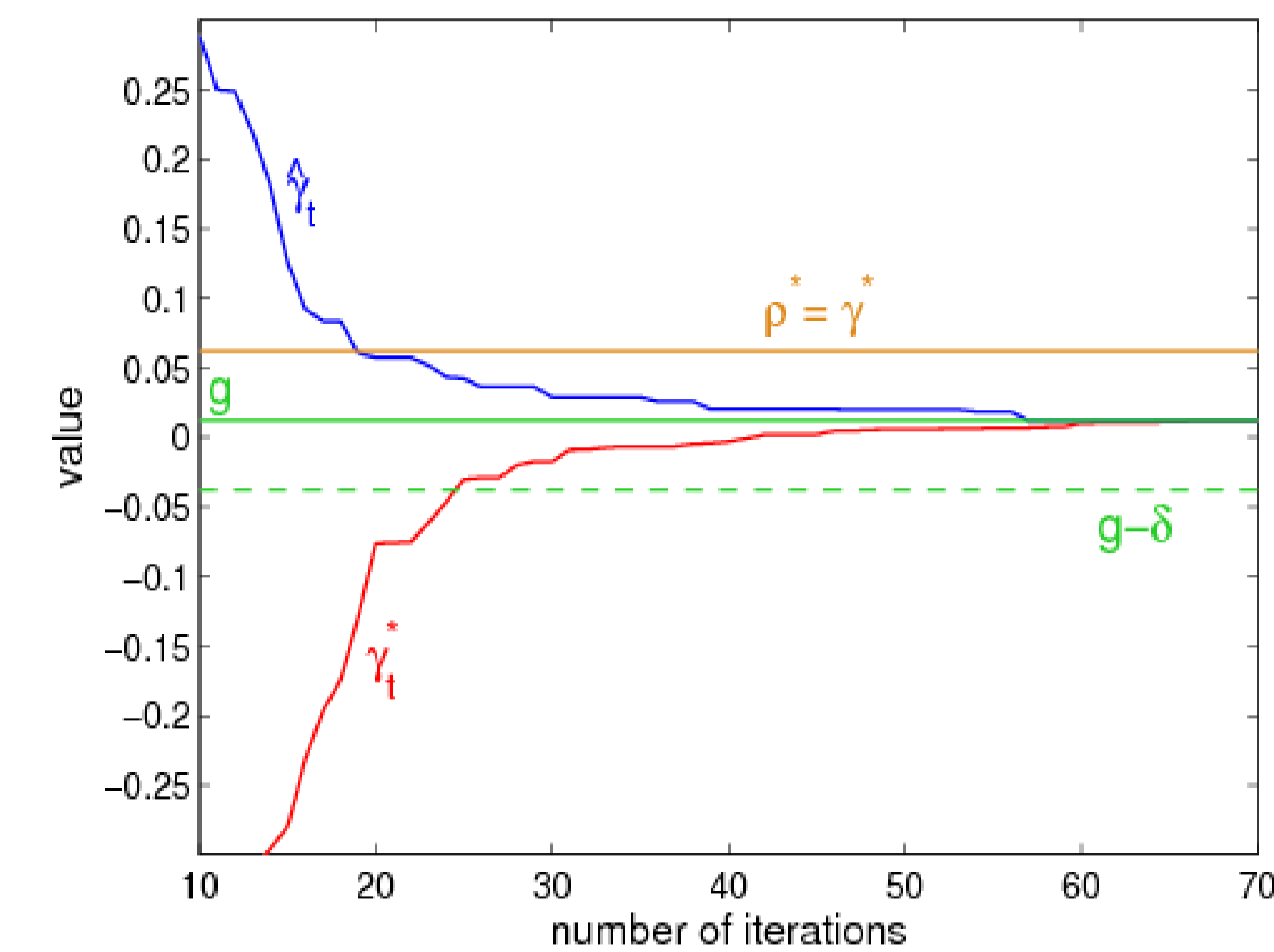
## 3. SoftBoost

- Designed for data that is not necessarily separable by convex combinations of base hypotheses
- Achieves robustness by *capping* the the weight on any example to be at most  $\frac{1}{\nu}$
- Capping the weights on the examples prevents the algorithm from focusing excessively on a few examples that it can't hope to get right
- Produces a convex combination of hypotheses whose *soft margin* is within  $\delta$  of the optimum
- SoftBoost terminates after at most  $\lceil \frac{2}{\delta^2} \ln(N/\nu) \rceil$  iterations.
- The algorithm does not need to know the **guarantee**  $g$  on the base hypotheses

### Algorithm 1: SoftBoost

- Input:**  $S = \langle (x_1, y_1), \dots, (x_N, y_N) \rangle$ , desired accuracy  $\delta$ , and capping parameter  $\nu \in [1, N]$ .
- Initialize:**  $d_n^0$  to the uniform distribution
- Do for**  $t = 1, \dots$ 
  - Train classifier on  $d^{t-1}$  and  $\{u_1, \dots, u_{t-1}\}$  and obtain hypothesis  $h^t$ . Set  $u_n^t = h^t(x_n)y_n$ .
  - Calculate the edge  $\gamma_t$  of  $h^t$ :  $\gamma_t = d^t \cdot u^t$
  - Set  $\hat{\gamma}_t = (\min_{m=1 \dots t} \gamma_m) - \delta$
  - Set  $\gamma^*$  = solution to the primal linear programming problem.
  - Update
 
$$d^{t+1} = \operatorname{argmin}_d \sum_{n=1}^N d_n \log \frac{d_n}{d_n^t}$$
 s.t.  $d \cdot u^m \leq \hat{\gamma}_t - \delta$ , for  $1 \leq m \leq t$   
 $\sum_n d_n = 1, d \leq \frac{1}{\nu} \mathbf{1}$ .
- If**  $\hat{\gamma}_t - \gamma^* \leq \delta$  then  $T = t - 1$  and break
- Output:**  $f_w(x) = \sum_{m=1}^T w_m h^m(x)$ , where the coefficients  $w_m$  maximize the soft margin over the hypothesis set  $\{h^1, \dots, h^t\}$  using the LP problem.

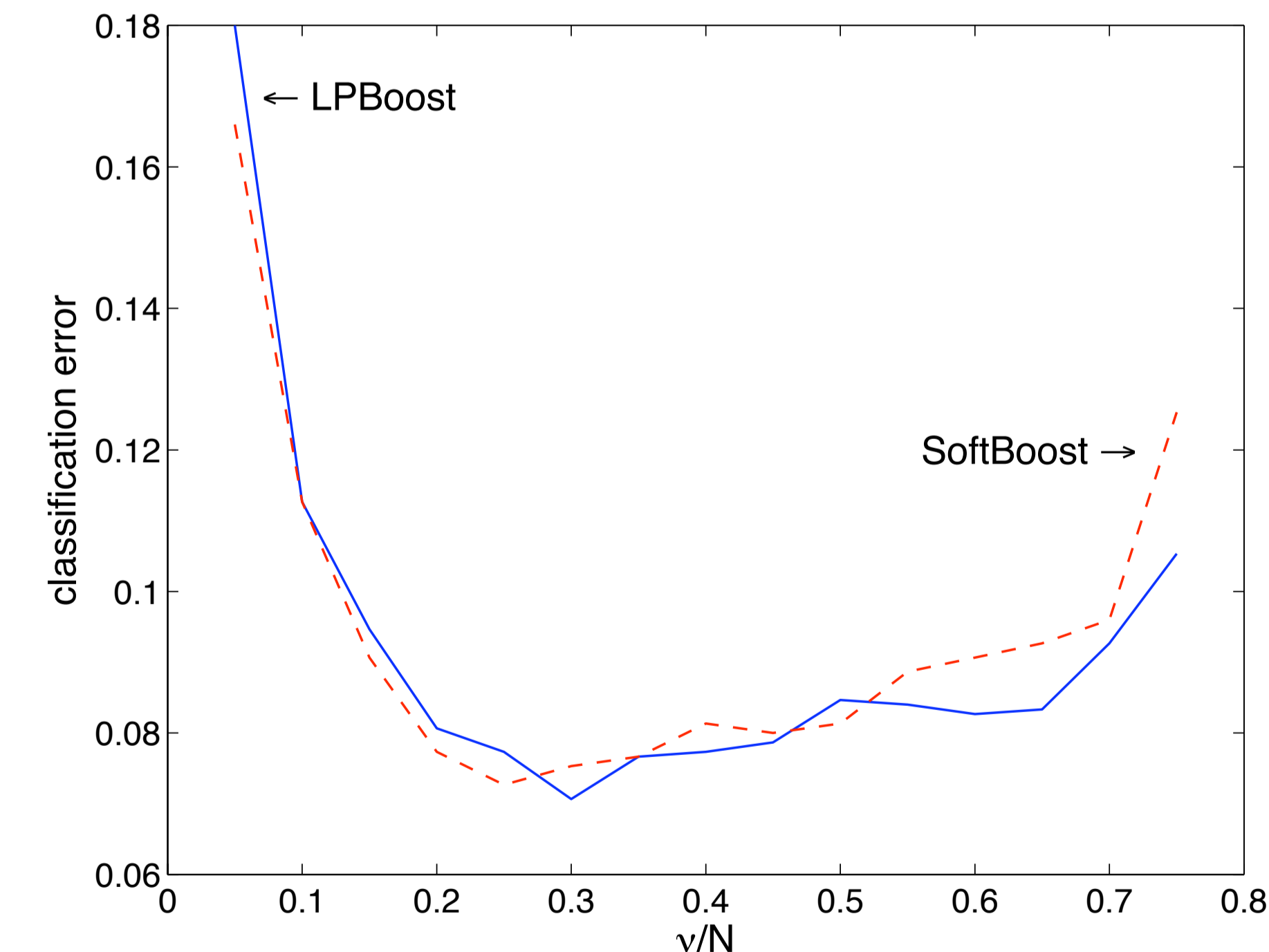
Illustration of Stopping Criterion



- $\hat{\gamma}_t := (\min_{m=1 \dots t} \gamma_m) - \delta$  is nonincreasing and  $\hat{\gamma}_t \geq g$
- $\gamma^*$  (solution to (1)) is non decreasing and  $\gamma^* \leq g$
- Algorithm terminates when they are sufficiently close together

## 4. Experimental Results

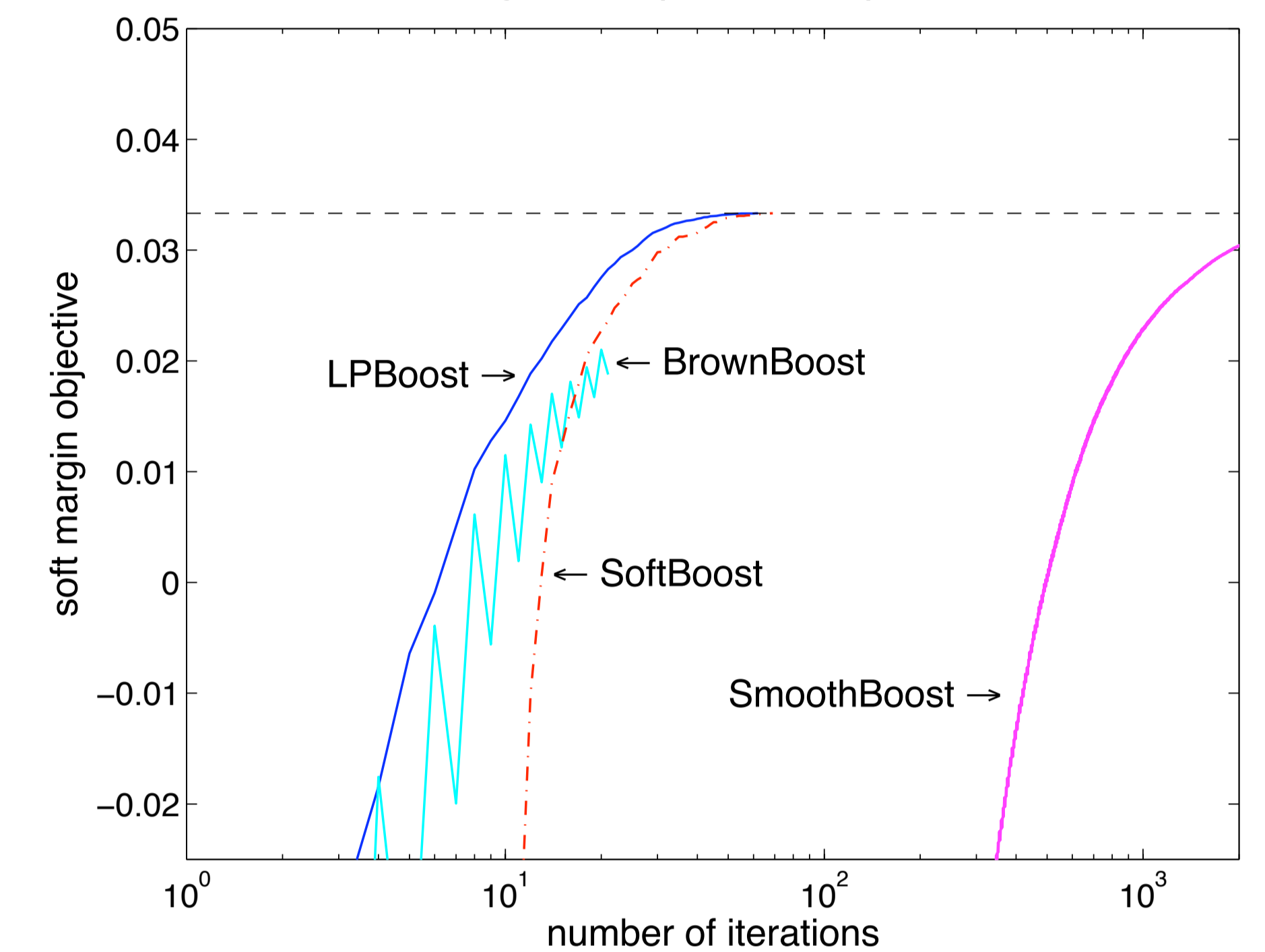
Generalization Performance of SoftBoost and LPBoost



- Generalization performance of SoftBoost (red) and LPBoost (blue) for different values of  $\nu$
- The data is a synthetic data set with 10% label noise in the training set

- If  $\nu$  is too small, the algorithm concentrates on a very few, presumably wrongly labeled examples and does not generalize well

Convergence Speed Comparison



- SoftBoost starts more slowly than LPBoost
- Both converge to within  $\delta$  of **guarantee**  $g$  in approximately the same number of iterations
- BrownBoost, which was designed to deal with noisy data but is not a smooth boosting algorithm, does not maximize the soft margin
- SmoothBoost, the best of previously existing smooth boosting algorithms, converges much more slowly and does not achieve the optimal soft margin

	AdaBoost	LPBoost	SoftBoost	BrownBoost
Banana	13.3 ± 0.7	11.1 ± 0.6	11.1 ± 0.5	12.9 ± 0.7
B.Cancer	32.1 ± 3.8	27.8 ± 4.3	28.0 ± 4.5	30.2 ± 3.9
Diabetes	27.9 ± 1.5	24.4 ± 1.7	24.4 ± 1.7	27.2 ± 1.6
German	26.9 ± 1.9	24.6 ± 2.1	24.7 ± 2.1	24.8 ± 1.9
Heart	20.1 ± 2.7	18.4 ± 3.0	18.2 ± 2.7	20.0 ± 2.8
Ringnorm	1.9 ± 0.3*	1.9 ± 0.2	1.8 ± 0.2	1.9 ± 0.2
F.Solar	36.1 ± 1.5	35.7 ± 1.6	35.5 ± 1.4	36.1 ± 1.4
Thyroid	4.4 ± 1.9*	4.9 ± 1.9	4.9 ± 1.9	4.6 ± 2.10
Titanic	22.8 ± 1.0	22.8 ± 1.0	23.0 ± 0.8	22.8 ± 0.8
Waveform	10.5 ± 0.4	10.1 ± 0.5	9.8 ± 0.5	10.4 ± 0.4

- Generalization error estimates and standard deviations for ten UCI benchmark data sets
- SoftBoost and LPBoost outperform AdaBoost and BrownBoost on most data sets
- SoftBoost and LPBoost perform similarly