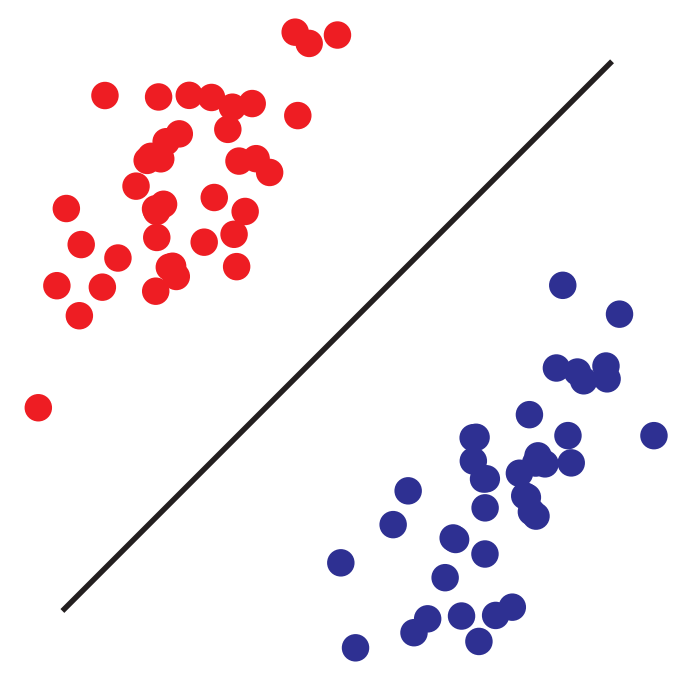


A New Perspective on Learning Linear Separators with Large $L_q L_p$ Margins

Maria-Florina Balcan and Christopher Berlind

Margins

Intuition: Learning should be easy when data is far from the decision boundary

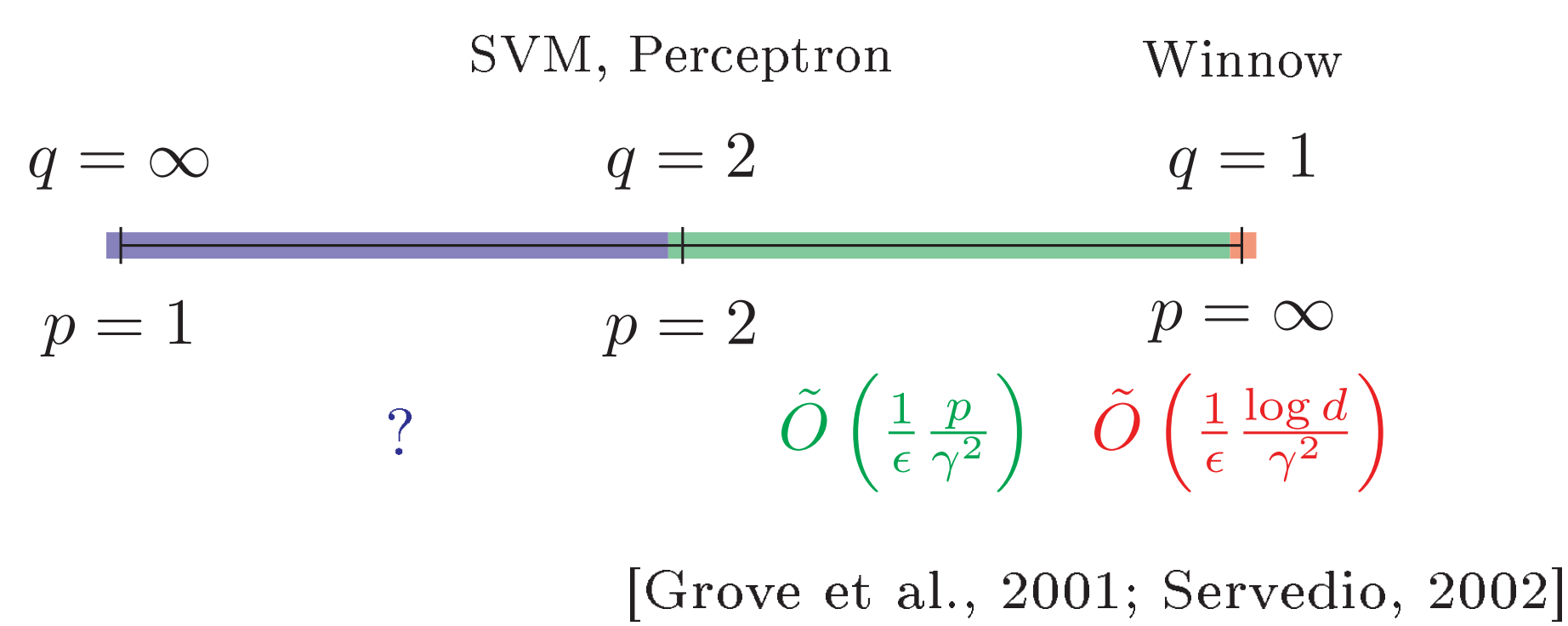


Definition: The $L_q L_p$ margin of w w.r.t. D is

$$\gamma_{q,p}(w) = \inf_{x \sim D} \frac{|w \cdot x|}{\|w\|_q \|x\|_p}$$

where $1 \leq p, q \leq \infty$ and $1/p + 1/q = 1$.

The Margin Spectrum:



Contributions

- Sample complexity bound covering the entire spectrum of margins
- Sufficient condition on data under which large margins lead to fast learning
- Upper and lower bounds for a family of problems showing a concrete advantage for $p = 1$
- Experimental confirmation that the theoretical results are relevant in practice

Discussion

- Important to consider entire margin spectrum
- Performance depends on both γ and $\|\mathbf{X}\|_{2,p}$
- Non-realizable case: L_q -norm regularization
- Relative sparsity of data and weight vector

Open questions:

- Algorithms that adaptively choose optimal p
- Generalization to multiple kernel learning
- Use $\|\mathbf{X}\|_{2,p}$ to aid feature selection

Experiments

$L_q L_p$ SVM: Given a set of n training examples, we can efficiently solve the convex program

$$\begin{aligned} \min_w \quad & \|w\|_q + C \sum_{i=1}^n \xi_i \\ \text{s.t.} \quad & \frac{y^i (w \cdot x^i)}{\|x^i\|_p} \geq 1 - \xi_i, \quad 1 \leq i \leq n. \end{aligned}$$

Equivalent to minimizing the hinge loss of an L_p -normalized data set using L_q -norm regularization.

Empirical advantage for $p = 1$ on synthetic data and on several real data sets from the UCI repository.

Generalization Bound

Theorem 1. Let $\|X\|_p = \left(\sup_{x \sim D} \|x\|_p \right)$ and

$$\|\mathbf{X}\|_{2,p} := \left(\sum_{i=1}^d \left(\sum_{j=1}^n |x_i^j|^2 \right)^{p/2} \right)^{1/p}$$

If there are constants $C = C(d, p)$ and $0 \leq \alpha < 1$ such that $\|\mathbf{X}\|_{2,p} \leq C n^\alpha \|X\|_p$ for any data set from D , then

$$\tilde{O} \left(\frac{1}{\epsilon} \left(\frac{C \sqrt{p}}{\gamma_{q,p}} \right)^{\frac{1}{1-\alpha}} \right)$$

samples suffices to achieve error ϵ for any $1 \leq p < \infty$.

Proof summary:

- Novel bound on fat-shattering dimension
- Use Khintchine inequality and bound on $\|\mathbf{X}\|_{2,p}$
- Apply standard generalization error bound

Example 1: Unhelpful Margins

Basis vectors with $w^* \in \{-1, 1\}^d$.

$$\begin{aligned} w^* : & \quad + - + - + - + - + - + - + - + - + - + - \\ D : & \quad 000\mathbf{1}0000000000000000 \quad + \\ & \quad 000000000000\mathbf{1}000000 \quad - \\ & \quad 0000000000\mathbf{1}00000000 \quad - \\ & \quad \vdots \end{aligned}$$

p	1	2	∞
$\gamma_{q,p}(w^*)$	1	$1/\sqrt{d}$	$1/d$
$\ \mathbf{X}\ _{2,p}$	\sqrt{dn}	\sqrt{n}	$\sqrt{n/d}$
s.c.	$\tilde{O}(d/\epsilon)$	$\tilde{O}(d/\epsilon)$	$\tilde{O}(d/\epsilon)$

Also have lower bound of $\tilde{\Omega}(d)$.

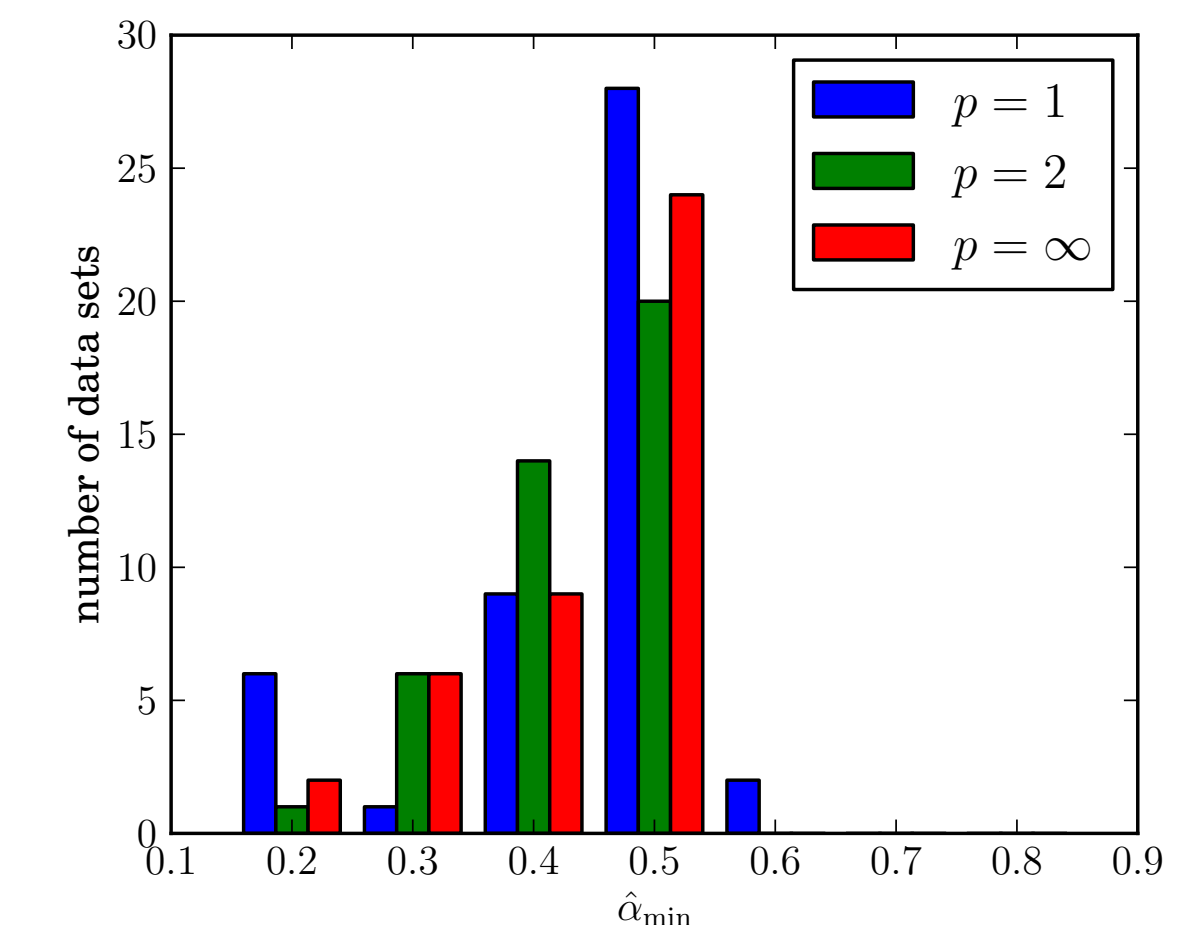
Bounding the $L_{2,p}$ -norm

Theoretically:

When $C = 1$, can always use $\alpha = 1/2$ when $p \geq 2$, but may need as much as $\alpha = 1/p$ when $p < 2$.

In reality:

For almost all data sets tested, we can bound $\|\mathbf{X}\|_{2,p}$ with $C = 1$ and $\alpha \leq 1/2$, regardless of p .



Histogram of α estimates on 47 real data sets

Example 2: Helpful Margins

Divide the d coordinates evenly into k blocks.

$$\begin{aligned} w^* : & \quad + + + + + - - - - - + + + + + \\ D : & \quad \mathbf{111111}0000000000000000 \quad + \\ & \quad 000000000000\mathbf{111111} \quad + \\ & \quad 000000\mathbf{111111}000000 \quad - \\ & \quad \vdots \end{aligned}$$

p	1	2	∞
$\gamma_{q,p}(w^*)$	1	$1/\sqrt{k}$	$1/k$
$\ \mathbf{X}\ _{2,p}$	\sqrt{kn}	\sqrt{n}	$\sqrt{n/k}$
s.c.	$\tilde{O}(k/\epsilon)$	$\tilde{O}(k/\epsilon)$	$\tilde{O}(k/\epsilon)$

Significant improvement when $k = o(d)$.

Making the Case for $L_\infty L_1$ Margins

Divide the d coordinates evenly into k blocks.

Distribution D randomly picks a block and either

- sets to 1 a single variable in the block or
- sets to 1 exactly $d/(2k)$ variables in the block.

Target w^* maximizes $L_\infty L_1$ margin.

$$\begin{aligned} w^* : & \quad + + + + + - - - - - + + + + + \\ D : & \quad 00\mathbf{1}0000000000000000 \quad + \\ & \quad 0000000\mathbf{11}00\mathbf{1}000000 \quad - \\ & \quad 000000000000\mathbf{11}00\mathbf{1} \quad + \\ & \quad 000000\mathbf{1}000000000000 \quad - \\ & \quad \vdots \end{aligned}$$

Theorem 2. If $k = O(d^{1/4})$ and $\epsilon = \Omega(d^{-1/4})$ in the above learning setting, then any algorithm restricted to using the large-margin class

$$W_p = \{w \in \mathbb{R}^d : \gamma_{q,p}(w) \geq \gamma_{q,p}(w^*)\}$$

for a fixed p has sample complexity

$$\begin{aligned} p = 1 : & \quad \tilde{O}(\sqrt{d}) \\ p > 1 : & \quad \tilde{\Omega}(d). \end{aligned}$$

Proof summary:

- W_1 : ++++++----- (unanimous on each block)
- W_2 : +++-+-+----- (some dissenters allowed)
- W_∞ : +++-+-+----- (any $w \in \{-1, 1\}^d$)
- Bound covering number for each W_p
- Apply distribution-specific s.c. bounds

Synthetic and Real Data Results

Synthetic data (top):
Blocks, Gaussian

Real data (bottom):
Fertility, SPECTF, CNAE-9

