### *l*<sub>1</sub>-Regularized Linear Regression: Persistence and Oracle Inequalities

Peter Bartlett EECS and Statistics UC Berkeley

slides at http://www.stat.berkeley.edu/~bartlett

Joint work with Shahar Mendelson and Joe Neeman.

◆□ ▶ ◆ ■ ▶ ◆ ■ ▶ ◆ ■ ● ● ● ●

- Random pair:  $(X, Y) \sim P$ , in  $\mathbb{R}^d \times \mathbb{R}$ .
- *n* independent samples drawn from *P*:  $(X_1, Y_1), \ldots, (X_n, Y_n)$ .
- Find  $\beta$  so linear function  $\langle X, \beta \rangle$  has small risk,

$$\mathcal{P}\ell_{eta} = \mathcal{P}\left(\langle X, eta 
angle - Y
ight)^2.$$

Here,  $\ell_{\beta}(X, Y) = (\langle X, \beta \rangle - Y)^2$  is the quadratic loss of the linear prediction.

・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・

## ℓ<sub>1</sub>-regularized linear regression

- Random pair:  $(X, Y) \sim P$ , in  $\mathbb{R}^d \times \mathbb{R}$ .
- *n* independent samples drawn from *P*:  $(X_1, Y_1), \ldots, (X_n, Y_n)$ .
- Find  $\beta$  so linear function  $\langle X, \beta \rangle$  has small risk,

$$P\ell_{\beta} = P(\langle X, \beta \rangle - Y)^2.$$

**Example.**  $\ell_1$ -regularized least squares:

$$\hat{\beta} = \arg\min_{\beta \in \mathbb{R}^d} |P_n \ell_\beta + \rho_n ||\beta||_{\ell_1^d} ,$$
  
where  $P_n \ell_\beta = \frac{1}{n} \stackrel{\not\sim}{\underset{i=1}{\longrightarrow}} (\langle X_i, \beta \rangle - Y_i)^2 , \text{ and } ||\beta||_{\ell_1^d} = \stackrel{\not\sim}{\underset{j=1}{\longrightarrow}} |\beta_j|.$ 

(日) (日) (日) (日) (日) (日) (日)

**Example.**  $\ell_1$ -regularized least squares:

◆□▶ ◆□▶ ▲□▶ ▲□▶ ■ ののの

- Tends to select sparse solutions (few non-zero components β<sub>j</sub>).
- Useful, for example, if  $d \gg n$ .

**Example.**  $\ell_1$ -regularized least squares:

$$\hat{\beta} = \arg\min_{\beta \in \mathbb{R}^d} P_n \ell_\beta + \rho_n \|\beta\|_{\ell_1^d} ,$$

**Example.**  $\ell_1$ -constrained least squares:

$$\hat{\beta} = \arg \min_{\|\beta\|_{\ell_1^{\alpha}} \le b_n} P_n \ell_{\beta}.$$

◆□▶ ◆□▶ ▲□▶ ▲□▶ ■ ののの

[Recall:  $\ell_{\beta}(X, Y) = (\langle X, \beta \rangle - Y)^2$ .]

**Example.**  $\ell_1$ -regularized least squares:

$$\hat{\beta} = \arg\min_{\beta \in \mathbb{R}^d} \ \mathbf{P}_n \ell_\beta + \rho_n \|\beta\|_{\ell_1^d} \ ,$$

**Example.**  $\ell_1$ -constrained least squares:

$$\hat{\beta} = \arg\min_{\|\beta\|_{\ell_1^d} \le b_n} P_n \ell_{\beta}.$$

Some questions:

► Prediction: Does give accurate forecasts? e.g., How does Pl<sub>Â</sub> compare with Pl<sub>β\*</sub>?

Here, 
$$\beta^* = \arg \min^{\bigcap} \mathcal{P}\ell_{\beta} : \|\beta\|_{\ell_1^d} \leq b_n^{\bigcirc}$$
.

◆□▶ ◆□▶ ▲□▶ ▲□▶ ■ ののの

**Example.**  $\ell_1$ -regularized least squares:

$$\hat{\beta} = \arg\min_{\beta \in \mathbb{R}^d} \ \mathbf{P}_n \ell_\beta + \rho_n \|\beta\|_{\ell_1^d} \ ,$$

**Example.**  $\ell_1$ -constrained least squares:

$$\hat{\beta} = \arg\min_{\|\beta\|_{\ell_1^d} \le b_n} P_n \ell_{\beta}.$$

Some questions:

- ► Does  $\hat{\beta}$  give accurate forecasts? e.g.,  $P\ell_{\hat{\beta}}$  versus  $P\ell_{\beta^*} = \min P\ell_{\beta} : \|\beta\|_{\ell^d_1} \le b_n$ ?
- **Estimation:** Under assumptions on *P*, is  $\hat{\beta} \approx$  correct?
- Sparsity Pattern Estimation: Under assumptions on P, are the non-zeros of  $\hat{\beta}$  correct?

# Outline of Talk

- 1. For  $\ell_1$ -constrained least squares, bounds on  $P\ell_{\beta} P\ell_{\beta^*}$ .
  - ► **Persistence:** (Greenshtein and Ritov, 2004) For what  $d_n, b_n \to \infty$  does  $P\ell_{\hat{\beta}} - P\ell_{\beta^*} \to 0$ ?
  - ► Convex Aggregation: (Tsybakov, 2003) For b = 1 (convex combinations of dictionary functions), what is rate of  $P\ell_{\beta} - P\ell_{\beta*}$ ?

(日) (日) (日) (日) (日) (日) (日)

- 2. For  $\ell_1$ -regularized least squares, oracle inequalities.
- 3. Proof ideas.

**Key Issue**:  $\ell_{\beta}$  is unbounded, so some key tools (e.g., concentration inequalities) cannot immediately be applied.

- For (X, Y) bounded, ℓ<sub>β</sub> can be bounded using ||β||<sub>ℓ<sup>d</sup><sub>1</sub></sub>, but this gives loose prediction bounds.
- ► We use chaining to show that metric structures of ℓ<sub>1</sub>-constrained linear functions under P<sub>n</sub> and P are similar.

(日) (日) (日) (日) (日) (日) (日)

Main Results: Excess Risk

For  $\ell_1$ -constrained least squares,

$$\hat{\beta} = \arg\min_{\|\beta\|_{\ell_1^q} \le b} P_n \ell_{\beta},$$

if *X* and *Y* have suitable tail behaviour then, with probability  $1 - \delta$ ,

$$P\ell_{\hat{\beta}} - P\ell_{\beta^*} \leq \frac{c \log^{\alpha}(nd)}{\delta^2} \min - \frac{b^2}{n} + \frac{d}{n}, \frac{b}{\sqrt{n}} - 1 + \frac{b}{\sqrt{n}}$$

.

- Small *d* regime: d/n.
- Large *d* regime:  $b/\sqrt{n}$ .

# Main Results: Excess Risk

For  $\ell_1$ -constrained least squares, with probability  $1 - \delta$ ,

$$P\ell_{\hat{\beta}} - P\ell_{\beta^*} \leq \frac{c \log^{\alpha}(nd)}{\delta^2} \min - \frac{b^2}{n} + \frac{d}{n}, \frac{b}{\sqrt{n}} - 1 + \frac{b}{\sqrt{n}}$$

Conditions:

- 1.  $PY^2$  is bounded by a constant.
- 2.  $\|X\|_{\infty}$  bounded a.s.,
  - X log concave and  $\max_{j} \|\langle X, e_{j} \rangle \|_{L_{2}} \leq c$ , or

◆□▶ ◆□▶ ▲□▶ ▲□▶ ■ ののの

X log concave and isotropic.

Application: Persistence

Consider  $\ell_1$ -constrained least squares,

$$\hat{\beta} = \arg\min_{\|\beta\|_{\ell_1^d} \le b} P_n \ell_{\beta}.$$

Suppose that  $PY^2$  is bounded by a constant and tails of *X* decay nicely (e.g.,  $||X||_{\infty}$  bounded a.s. or *X* log concave and isotropic).

Then for increasing  $d_n$  and

$$b_n = o \quad \frac{\sqrt{n}}{\log^{3/2} n \log^{3/2} (nd_n)}$$
,

1

(日) (日) (日) (日) (日) (日) (日)

 $\ell_1$ -constrained least squares is persistent (i.e.,  $P\ell_{\hat{\beta}} - P\ell_{\beta^*} \rightarrow 0$ ).

Application: Persistence

If  $PY^2$  is bounded and tails of X decay nicely, then  $\ell_1$ -constrained least squares is persistent provided that  $d_n$  is increasing and

$$b_n = o \quad \frac{\sqrt{n}}{\log^{3/2} n \log^{3/2} (nd_n)}$$

Previous Results (Greenshtein and Ritov, 2004):

- 1.  $b_n = \omega(n^{1/2} / \log^{1/2} n)$  implies empirical minimization is not persistent for Gaussian (*X*, *Y*).
- 2.  $b_n = o(n^{1/2} / \log^{1/2} n)$  implies empirical minimization is persistent for Gaussian (*X*, *Y*).
- 3.  $b_n = o(n^{1/4} / \log^{1/4} n)$  implies empirical minimization is persistent under tail conditions on (X, Y).

# **Application: Convex Aggregation**

Consider b = 1, so that the  $\ell_1$ -ball of radius b is the convex hull of a dictionary of d functions (the components of X). Tsybakov (2003) showed that, for any aggregation scheme  $\hat{\beta}$ , the rate of convex aggregation satisfies

$$P\ell_{\hat{\beta}} - P\ell_{\beta^*} = \Omega \quad \min \quad \frac{d}{n}, \stackrel{\Gamma}{\longrightarrow} \frac{\log d}{n}$$

For bounded, isotropic distributions, our result implies that this rate can be achieved, up to log factors, by least squares over the convex hull of the dictionary.

Previous positive results (Tsybakov, 2003; Bunea, Tsybakov and Wegkamp, 2006) involved complicated estimators.

# Outline of Talk

- 1. For  $\ell_1$ -constrained least squares, bounds on  $P\ell_{\hat{\beta}} P\ell_{\beta^*}$ .
  - Persistence:

For what  $d_n, b_n \to \infty$  does  $P\ell_{\hat{\beta}} - P\ell_{\beta^*} \to 0$ ?

Convex Aggregation:
 For b = 1 (convex combinations of dictionary functions), what is rate of Pl<sub>β</sub> - Pl<sub>β\*</sub>?

(日) (日) (日) (日) (日) (日) (日)

- 2. For  $\ell_1$ -regularized least squares, oracle inequalities.
- 3. Proof ideas.

### **Proof Ideas:** 1. $\epsilon$ -equivalence of *P* and *P<sub>n</sub>* structures

#### Define

$$egin{aligned} G_\lambda = & rac{\lambda}{m{P}(\ell_eta-\ell_{eta^*})}(\ell_eta-\ell_{eta^*}):m{P}(\ell_eta-\ell_{eta^*})\geq\lambda \end{aligned}$$

#### Then:

E sup<sub>*g*∈*G*<sub>λ</sub> |*P*<sub>*n*</sub>*g* − *Pg*| is small ⇒ with high probability, for all β with  $P(\ell_\beta - \ell_{\beta^*}) \ge \lambda$ ,</sub>

$$(1-\epsilon)P(\ell_{\beta}-\ell_{\beta^*}) \leq P_n(\ell_{\beta}-\ell_{\beta^*}) \leq (1+\epsilon)P(\ell_{\beta}-\ell_{\beta^*})$$

(日) (日) (日) (日) (日) (日) (日)

 $\Rightarrow P(\ell_{\hat{\beta}} - \ell_{\beta^*}) \leq \lambda$ , where  $\hat{\beta} = \arg \min_{\beta} P_n \ell_{\beta}$ .

## Proof Ideas: 2. Symmetrization, subgaussian tails

<ロ>

# Proof Ideas: 3. Chaining

For a subgaussian process  $\{Z_t\}$  indexed by a metric space (T, d), and for  $t_0 \in T$ ,

$$\mathsf{E}\sup_{t\in\mathcal{T}}|Z_t-Z_{t_0}|\leq c\mathcal{D}(\mathcal{T},d)=c\int_{0}^{Z}\operatorname{diam}(\mathcal{T},d)\operatorname{p}\overline{\log N(\epsilon,\mathcal{T},d)}\,d\epsilon,$$

▲□▶ ▲□▶ ▲□▶ ▲□▶ □ のQで

where  $N(\epsilon, T, d)$  is the  $\epsilon$  covering number of T.

# Proof Ideas: 4. Bounding the Entropy Integral

It suffices to calculate the entropy integral  $\mathcal{D}(\sqrt{\lambda}D \cap 2bB_1^d, d)$ . We can approximate this by

 $\mathcal{D}(\sqrt{\lambda} D \cap 2\textit{bB}_1^{\textit{d}}, \textit{d}) \leq \min \ \mathcal{D}(\sqrt{\lambda} D, \textit{d}), \mathcal{D}(2\textit{b}B_1^{\textit{d}}, \textit{d}) \ .$ 

This leads to:

$$P\ell_{\hat{\beta}} - P\ell_{\beta^*} \leq rac{c\log^{lpha}(nd)}{\delta^2}\min - rac{b^2}{n} + rac{d}{n}, rac{b}{\sqrt{n}} - 1 + rac{b}{\sqrt{n}}$$

◆□▶ ◆□▶ ▲□▶ ▲□▶ ■ ののの

# Proof Ideas: 5. Oracle Inequalities

We get an isomorphic condition on  $\{\ell_{\beta} - \ell_{\beta^*}\}$ ,

$$\frac{1}{2}\boldsymbol{P}_{n}(\ell_{\beta}-\ell_{\beta^{*}})-\epsilon_{n}\leq \boldsymbol{P}(\ell_{\beta}-\ell_{\beta^{*}})\leq 2\boldsymbol{P}_{n}(\ell_{\beta}-\ell_{\beta^{*}})+\epsilon_{n},$$

and this implies that  $\hat{\beta} = \arg \min_{\beta} (P_n \ell_{\beta} + c \epsilon_n)$  has

$$P\ell_{eta} \leq \inf_{eta} P\ell_{eta} + c'\epsilon_n$$
 .

This leads to oracle inequality: For  $\ell_1$ -regularized least squares,

$$\hat{\beta} = \arg\min_{\beta} P_n \ell_{\beta} + \rho_n \|\beta\|_{\ell_1^{d_n}},$$

with probability at least 1 - o(1),

$$P\ell_{\hat{\beta}} \leq \inf_{\beta} P\ell_{\beta} + c\rho_n 1 + \|\beta\|_{\ell_1^{d_n}}$$

Outline of Talk

1. For  $\ell_1$ -constrained least squares,

$$m{P}\ell_{\hat{eta}} - m{P}\ell_{eta^*} \leq rac{c\log^lpha(nd)}{\delta^2} {
m min} \quad rac{b^2}{n} + rac{d}{n}, rac{b}{\sqrt{n}} \quad 1 + rac{b}{\sqrt{n}}$$

# Persistence: If b<sub>n</sub> = õ(√n), then Pℓ<sub>β̂</sub> - Pℓ<sub>β\*</sub> → 0. Convex Aggregation: Empirical risk minimization gives optimal rate (up to log factors): Õ min(d/n, <sup>D</sup> log d/n).

・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・

- 2. For  $\ell_1$ -regularized least squares, oracle inequalities.
- 3. Proof ideas: subgaussian Rademacher process.