

High Dimensional Statistics

Yanbo Tang

Imperial College London

March. 2025

Content of this week

- LASSO MSE
- SLOPE
- LASSO sparse set recovery

Fixed versus random design

Random design

$$E \left[\frac{1}{n} \sum_{i=1}^n \| \hat{f}(x_i) - f(x_i) \|_2^2 \right]$$

$x \sim P$

Fixed design

$$\frac{1}{n} \sum_{i=1}^n \| \hat{f}(x_i) - f(x_i) \|_2^2$$
$$= E \left[\frac{1}{n} \sum_{i=1}^n \| \hat{f}(x_i) - f(x_i) \|_2^2 \right]$$

$x \sim P_n$

L^0 performance

θ is k sparse

Denote by $B_0(k)$ the ℓ_0 "ball" of \mathbb{R}^d , i.e., the set of k -sparse vectors, defined by

$$B_0(k) = \{\theta \in \mathbb{R}^d : |\theta|_0 \leq k\}.$$

Our goal is to control the MSE of $\hat{\theta}_K^{LS}$ when $K = B_0(k)$. Note that computing $\hat{\theta}_{B_0(k)}^{LS}$ is to find:

$$\hat{\theta}_{B_0(k)}^{LS} \in \operatorname{argmin}_{\theta \in B_0(k)} \|Y - \mathbb{X}\theta\|_2^2,$$

but this would require computing $\binom{d}{k}$ least squares estimators.

Best subset selection

Theorem

Fix a positive integer $k \leq d/2$. Let $K = B_0(k)$ be set of k -sparse vectors of \mathbb{R}^d , assume that $\theta^* \in B_0(k)$ and assume that the linear model with $\varepsilon \sim \text{subG}_n(\sigma^2)$. Then, for any $\delta > 0$, with probability $1 - \delta$, it holds

$$\text{MSE}(\mathbb{X}\hat{\theta}_{B_0(k)}^{\text{IS}}) \lesssim \frac{k\sigma^2}{n} \log\left(\frac{ed}{2k}\right) + \log(6) \frac{\sigma^2 k}{n} + \frac{\sigma^2}{n} \log(1/\delta).$$

Note that if we knew the exact sparsity level k_0 , then this bound would be essentially optimal.

Matching rates for LASSO?

Assumption INC(k) We say that the design matrix \mathbb{X} has incoherence k for some integer $k > 0$ if

$$\left| \frac{\mathbb{X}^T \mathbb{X}}{n} - I_d \right|_{\infty} \leq \frac{1}{32k}$$

where the $|A|_{\infty}$ denotes the largest element of A in absolute value. Equivalently,

- ① For all $j = 1, \dots, d$,

$$\left| \frac{\|\mathbb{X}_j\|_2^2}{n} - 1 \right| \leq \frac{1}{32k}.$$

- ② For all $1 \leq i, j \leq d, i \neq j$, we have

$$\frac{|\mathbb{X}_i^T \mathbb{X}_j|}{n} \leq \frac{1}{32k}.$$

Technical Lemma

Lemma

Fix a positive integer $k \leq d$ and assume that \mathbb{X} satisfies assumption $INC(k)$. Then, for any $S \in \{1, \dots, d\}$ such that $|S| \leq k$ and any $\theta \in \mathbb{R}^d$ that satisfies the cone condition

$$|\theta_{S^c}|_1 \leq 3|\theta_S|_1, \quad \text{Let } \theta = \hat{\theta} - \theta^*$$

it holds

$$|\theta|_2^2 \leq 2 \frac{|\mathbb{X}\theta|_2^2}{n}. \quad \Rightarrow \quad \|\hat{\theta} - \theta^*\|_2 \leq 2 \frac{\|\hat{\theta} - \theta^*\|_2^2}{n}$$

We will interpret the cone condition when we consider sparse recovery.

$$\arg \min_{\theta \in \mathbb{R}^d} \frac{\|\mathbb{X}\theta - y\|_2^2}{n} + 2\tau \|\theta\|_1$$

$$\text{just do } L_1 \text{ ball penalty } \sqrt{\frac{\log(d)}{n}}$$

Theorem

Fix $n \geq 2$. Assume that the linear model with $\varepsilon \sim \text{subG}_n(\sigma^2)$. Moreover, assume that $\|\theta^*\|_0 \leq k$ and that X satisfies assumption $\text{INC}(k)$. Then the Lasso estimator $\hat{\theta}^{\mathcal{L}}$ with regularization parameter defined by

$$2\tau = 8\sigma\sqrt{\frac{\log(2d)}{n}} + 8\sigma\sqrt{\frac{\log(1/\delta)}{n}}$$

satisfies

$$\text{MSE}(X\hat{\theta}^{\mathcal{L}}) = \frac{1}{n}\|X\hat{\theta}^{\mathcal{L}} - X\theta^*\|_2^2 \lesssim k\sigma^2 \frac{\log(2d/\delta)}{n}.$$

and

$$\|\hat{\theta}^{\mathcal{L}} - \theta^*\|_2^2 \lesssim k\sigma^2 \frac{\log(2d/\delta)}{n}.$$

with probability at least $1 - \delta$.

Proof

From definition of $\hat{\theta}_L$ (trying to show the cone condition)

$$\frac{1}{n} \|Y - X\hat{\theta}_L\| \leq \frac{1}{n} \|Y - X\theta^*\|_2^2 + 2\tau \|\theta^*\|_1 - 2\tau \|\hat{\theta}_L\|_1$$

Add $2\tau \|\hat{\theta}_L - \theta^*\|_1$ and multiply by n on both side

$$\begin{aligned} \|X\hat{\theta}_L - X\theta^*\|_2^2 + n\tau \|\hat{\theta}_L - \theta^*\|_1 &\leq 2\varepsilon^T X(\hat{\theta}_L - \theta^*) \\ &\quad + n\tau \|\hat{\theta}_L - \theta^*\|_1 \\ &\quad + 2n\tau \|\theta^*\|_1 - 2n\tau \|\hat{\theta}_L\|_1 \end{aligned} \quad (1)$$

$$\begin{aligned} \varepsilon^T X(\hat{\theta}_L - \theta^*) &\leq \|\varepsilon^T X\|_\infty \|\hat{\theta}_L - \theta^*\|_1 \\ &\leq \frac{n\tau}{2} \|\hat{\theta}_L - \theta^*\|_1 \quad \text{w.p } 1-\delta \end{aligned}$$

$$\left\{ \begin{aligned} \text{as } \text{PC } |X^T \varepsilon|_\infty \geq \varepsilon &\leq \text{PC}_{\text{chained}} |X_j^T \varepsilon| \geq \varepsilon \leq 2 \exp\left(-\frac{\varepsilon^2}{4n\sigma^2}\right) \\ \text{where we used } |X_j|_2^2 &\leq n + \frac{1}{32k} < 2n \\ \text{take } \varepsilon = n\tau &\text{ in } \end{aligned} \right.$$

Proof Cont. $\|X\hat{\theta}^L - X\theta^*\|_2^2 + n\tau \|\hat{\theta}^L - \theta^*\|$ $\|a\|_1 = \|\theta_S\|_1 + \|\theta_{S^c}\|_1$

$$\begin{aligned} &\leq 2n\tau \|\hat{\theta}^L - \theta^*\|_1 + 2n\tau \|\theta^*\|_1 - 2n\tau \|\hat{\theta}^L\|_1 \\ &= 2n\tau \|\theta_S^L - \theta^*\|_1 + 2n\tau \|\theta^*\|_1 - 2n\tau \|\theta_S^L\|_1 \\ &\leq 4n\tau \|\theta_S^L - \theta^*\|_1 \quad (2) \end{aligned}$$

(2) implies the Cone Condition

$$\begin{aligned} &\|X\hat{\theta}^L - X\theta^*\|_2^2 + n\tau \|\hat{\theta}_{S^c}^L - \theta_{S^c}^*\|_1 + n\tau \|\theta_S^L - \theta_S^*\|_1 \\ &\leq 4n\tau \|\hat{\theta}_S^L - \theta^*\|_1 \end{aligned}$$

$$\frac{\|X\hat{\theta}^L - X\theta^*\|_2^2}{n\tau} + \|\hat{\theta}_{S^c}^L - \theta_{S^c}^*\|_1 \leq 3\|\hat{\theta}_S^L - \theta_S^*\|_1$$

this shows the Cone Condition.

$$\begin{aligned} \|\hat{\theta}_S^L - \theta_S^*\|_1 &\leq \sqrt{|S|} \|\hat{\theta}_S^L - \theta_S^*\|_2 \leq \sqrt{|S|} \|\hat{\theta}^L - \theta^*\|_2 \\ &\leq \sqrt{\frac{2K}{n}} \|X\hat{\theta}^L - X\theta^*\|_2 \end{aligned}$$

Optimal rate (MSE) - SLOPE

Definition

(*Slope estimator*). Let $\lambda = (\lambda_1, \dots, \lambda_d)$ be a non-increasing sequence of positive real numbers, $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_d > 0$. For $\theta = (\theta_1, \dots, \theta_d) \in \mathbb{R}^d$, let $(\theta_1^*, \dots, \theta_d^*)$ be a non-increasing rearrangement of the modulus of the entries, $|\theta_1|, \dots, |\theta_d|$. We define the *sorted ℓ_1 norm* of θ as

$$|\theta|_* = \sum_{j=1}^d \lambda_j |\theta_j^*|$$

or equivalently as

$$|\theta|_* = \max_{\phi \in S_d} \sum_{j=1}^d \lambda_j |\theta_{\phi(j)}|.$$

Cont.

Definition

The *Slope estimator* is then given by

$$\hat{\theta}^S \in \arg \min_{\theta \in \mathbb{R}^d} \left\{ \frac{1}{n} \|Y - X\theta\|_2^2 + 2\tau |\theta|_* \right\}$$

for a choice of tuning parameters λ and $\tau > 0$.

In what follows, we use

$$\lambda_j = \sqrt{\log(2d/j)}, \quad j = 1, \dots, d.$$

$$\leq \frac{1}{32k}$$

Theorem

Fix $n \geq 2$. Assume that the linear model holds where $\varepsilon \sim \mathcal{N}_n(0, \sigma^2 I_n)$. Moreover, assume that $|\theta^*|_0 \leq k$ and that \mathbb{X} satisfies assumption $\text{INC}(\underline{k'})$ with $k' \geq 4k \log(2de/k)$. Then the Slope estimator $\hat{\theta}^S$ with regularization parameter defined by

$$\tau = 8\sqrt{2}\sigma \sqrt{\frac{\log(1/\delta)}{n}}$$

satisfies

$$\text{MSE}(\mathbb{X}\hat{\theta}^S) = \frac{1}{n} \|\mathbb{X}\hat{\theta}^S - \mathbb{X}\theta^*\|_2^2 \lesssim \sigma^2 \frac{k \log(2d/k\delta)}{n}$$

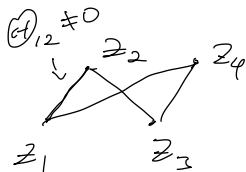
and

$$\|\hat{\theta}^S - \theta^*\|_2^2 \lesssim \sigma^2 \frac{k \log(2d/k) \log(1/\delta)}{n}.$$

with probability at least $1 - \delta$.

Why sparsity - selection of Gaussian graphical models

$$(z_1, \dots, z_d) \sim \frac{1}{\sqrt{(2\pi)^d \det(\Sigma^{-1})}} \exp\left(-\frac{1}{2} z^T \Sigma^{-1} z\right)$$



$$\Sigma^{-1}$$

$$\begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \\ 0 & 0 & 0 \end{bmatrix}$$

$$S \subseteq V := \{1, 2, \dots, d\} \quad N(S) = \{t \in V \mid \theta_{st} \neq 0\}$$

$$z_S = \langle z_{\setminus S}, \theta^* \rangle \mid z_S \sim N(0, \sigma^2) \quad z_{\setminus S} = \begin{pmatrix} z_2 \\ z_3 \\ z_4 \end{pmatrix}$$

the sparsity pattern for θ^*
indicated by Σ_S .

Basis pursuit

Before we begin with the sparse recovery of the noisy version of the problem, let us think about the deterministic version of the problem. Suppose that you are asked to solve the following system of equations for θ :

$$\begin{array}{l} \text{if } Q \gg n \\ \text{impossible} \end{array} \quad \begin{array}{l} X\theta = Y. \quad X \in \mathbb{R}^{n \times d} \\ Y \in \mathbb{R}^n \end{array}$$

if additionally $\|\theta\|_0 \leq K$
is known to be a soln.

$$\min_{\theta \in \mathbb{R}^d} \|\theta\|_0 : X\theta = Y$$

$$\min_{\theta \in \mathbb{R}^d} \|\theta\|_1 : X\theta = Y$$

Restricted nullspace

$$\text{null}(X) = \{ \Delta \in \mathbb{R}^d : X\Delta = 0 \}$$

$$\forall \Delta \in \text{null}(X) \quad X\theta^* = y \Rightarrow X(\theta^* + \Delta) = y$$

$$\stackrel{0}{=} T(\theta^*) = \{ \Delta \in \mathbb{R}^d \mid \|\theta^* + \epsilon \Delta\|_1 \leq \|\theta^*\|_1 \text{ for some } \epsilon > 0 \}$$

$$\begin{aligned} \|\theta^*\|_1 &\geq \|\theta^* + \epsilon \Delta\|_1 = \|\theta^* + \epsilon \Delta_S\|_1 + \|\epsilon \Delta_{S^c}\|_1 \\ &\geq \|\theta^*\|_1 - \epsilon \|\Delta_S\|_1 + \epsilon \|\Delta_{S^c}\|_1 \\ \Rightarrow \|\Delta_{S^c}\|_1 &\leq \|\Delta_S\|_1 \end{aligned}$$

Restricted nullspace illustration

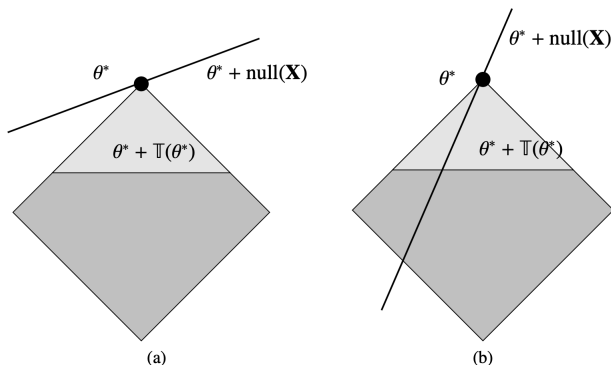


Figure 7.2 Geometry of the tangent cone and restricted nullspace property in $d = 2$ dimensions. (a) The favorable case in which the set $\theta^* + \text{null}(\mathbf{X})$ intersects the tangent cone only at θ^* . (b) The unfavorable setting in which the set $\theta^* + \text{null}(\mathbf{X})$ passes directly through the tangent cone.

Image taken from High-Dimensional Statistics by Martin Wainwright.

Restricted null space

Let:

$$C(S) := \{\Delta \in \mathbb{R}^d : \|\Delta_{S^c}\| \leq \|\Delta_S\|_1\}$$

Definition

Restricted nullspace The matrix X satisfies the restricted nullspace property with respect to S if $C(S) \cap \text{null}(X) = \{0\}$.

Unique solution - basis pursuit

Theorem

The following two properties are equivalent:

- (a) For any vector $\theta^* \in \mathbb{R}^d$ with support S , the basis pursuit program (7.9) applied with $y = X\theta^*$ has unique solution $\hat{\theta} = \theta^*$.
- (b) The matrix X satisfies the restricted nullspace property with respect to S .

b \Rightarrow a $\hat{\theta}$ and θ^* are solutions

$$(1) \quad \|\hat{\theta}\|_1 \leq \|\theta^*\|_1 \quad \text{define} \quad \tilde{\Delta} := \hat{\theta} - \theta^*$$

$$\begin{aligned} \|\theta_S^*\|_1 &= \|\theta^*\|_1 \geq \|\theta^* + \tilde{\Delta}\|_1 = \|\theta_S^* + \tilde{\Delta}_S\|_1 + \|\tilde{\Delta}_{S^c}\|_1 \\ &\geq \|\theta_S^*\|_1 - \|\tilde{\Delta}_S\|_1 + \|\tilde{\Delta}_{S^c}\|_1 \\ \Rightarrow \|\tilde{\Delta}_S\|_1 &\geq \|\tilde{\Delta}_{S^c}\|_1 \end{aligned}$$

Proof $a \Rightarrow b$ \forall sparse pattern.
 $\text{null}(X) \setminus \text{SOS} \cap C(S)$ is empty.

$$\forall \theta \in \text{null}(X) \setminus \text{SOS} \quad \text{Consider } \min_{\beta \in \mathbb{R}^d} \|\beta\|_1 : X\beta = X \begin{bmatrix} \theta_S^* \\ 0 \end{bmatrix}$$

the unique solution is $\hat{\beta} = \begin{bmatrix} \theta_S^* \\ 0 \end{bmatrix}$.

as $X\theta^* = 0 \Rightarrow \begin{bmatrix} 0 \\ \theta_{S^c}^* \end{bmatrix}^T$ is also a solution

$$X \begin{pmatrix} \theta_S^* \\ \theta_{S^c}^* \end{pmatrix} = 0 \quad \therefore X \begin{pmatrix} \theta_S^* \\ 0 \end{pmatrix} = -X \begin{pmatrix} 0 \\ \theta_{S^c}^* \end{pmatrix}$$

but $\begin{bmatrix} \theta_S^* \\ 0 \end{bmatrix}^T$ is the unique minimizer

$$\therefore \|\theta_S^*\|_1 < \|\theta_{S^c}^*\|_1$$

$$\therefore \theta^* \notin C(S)$$

Restricted isometry

$$\sum_n^1 = \frac{x^T x}{n}$$

Proposition

If

$$\left| \frac{\mathbb{X}^T \mathbb{X}}{n} - I_d \right|_{\infty} \leq \frac{1}{3k},$$

then the restricted null space condition holds for all subsets of cardinality at most k .

RIP

Definition

For an integer k , $X \in \mathbb{R}^{n \times d}$ satisfies a restricted isometry property of order k with constant $\delta_k(X) > 0$ if

$$\left\| \frac{X_S^\top X_S}{n} - I_k \right\|_{op} \leq \delta_k(X)$$

for all subsets of size at most k .

Proposition

If the RIP constant of order $2k$ is bounded as $\delta_{2k} < 1/3$, then the restricted null space condition holds for any subset S of cardinality $|S| \leq k$.

Assumptions

Assumption

The smallest eigenvalue of the sample covariance submatrix indexed by S is bounded below:

$$\gamma_{\min} \left(\frac{\mathbf{X}_S^T \mathbf{X}_S}{n} \right) \geq c_{\min} > 0.$$

Assumption

There exists some $\alpha \in [0, 1)$ such that mutual coherence

$$\max_{j \in S_c} \|(\mathbf{X}_S^T \mathbf{X}_S)^{-1} \mathbf{X}_S^T \mathbf{X}_j\|_1 \leq \alpha.$$

want to predict x_j using \mathbf{X}_S
by $\hat{w} = \arg\min \|x_j - \mathbf{X}_S w\|_2^2 = (\mathbf{X}_S^T \mathbf{X}_S)^{-1} \mathbf{X}_S^T x_j.$

$$\lambda_{-1} \lambda_n \quad \Pi_{S^\perp} = [I - X_S(X_S^T X_S)^{-1} X_S^T]$$

Theorem

Consider an S -sparse linear regression model for which the design matrix satisfies Assumptions 2 and 3. Then for any choice of regularization parameter such that

$$\lambda_n \geq \frac{2}{1 - \alpha} \left\| X_S^T \Pi_{S^\perp}(\mathbf{X}) \frac{\epsilon}{n} \right\|_\infty, \quad (1)$$

the Lasso program has the following properties:

- (a) **Uniqueness:** There is a unique optimal solution $\hat{\theta}$.
- (b) **No false inclusion:** This solution has its support set \hat{S} contained within the true support set S .
- (c) **ℓ_∞ -bounds:** The error $\hat{\theta} - \theta^*$ satisfies

$$\|\hat{\theta}_S - \theta_S^*\|_\infty \leq \underbrace{\left\| \left(\frac{X_S^T X_S}{n} \right)^{-1} X_S^T \frac{\epsilon}{n} \right\|_\infty}_{B(\lambda_n, \mathbf{X})} + \left\| \left(\frac{X_S^T X_S}{n} \right)^{-1} \right\|_\infty \lambda_n, \quad (2)$$

where $\|A\|_\infty = \max_{i=1, \dots, S} \sum_j |A_{i,j}|$ is the matrix ℓ_∞ -norm.

- (d) **No false exclusion:** The Lasso includes all indices $i \in S$ such that $|\theta_i^*| > B(\lambda_n; \mathbf{X})$, and hence is variable selection consistent if $\min_{i \in S} |\theta_i^*| > B(\lambda_n; \mathbf{X})$.

Corollary

For a S -sparse linear model based on a noise vector ϵ with zero-mean i.i.d. σ -sub-Gaussian entries, and a deterministic design matrix X that satisfies Assumptions 2 and 3, as well as the C -column normalization condition $\max_{j=1,\dots,d} \|X_j\|_2/\sqrt{n} \leq C$. Suppose that we solve the Lasso program with regularization parameter

$$\lambda_n = \frac{2C\sigma}{1-\alpha} \left\{ \sqrt{\frac{2\log(d-k)}{n}} + \delta \right\}$$

for some $\delta > 0$. Then the optimal solution $\hat{\theta}$ is unique with its support contained within S , and satisfies the ℓ_∞ -error bound

$$\|\hat{\theta}_S - \theta_S^*\|_\infty \leq \frac{\sigma}{\sqrt{c_{\min}}} \left(\sqrt{\frac{2\log s}{n}} + \delta \right) + \left\| \left(\frac{X_S^T X_S}{n} \right)^{-1} \right\|_\infty \lambda_n, \quad (3)$$

all with probability at least $1 - 4e^{-n\delta^2/2}$.

Proof

Sub-gradients

Primal Dual Witness

Definition

Primal–dual witness (PDW) construction:

- 1 Set $\hat{\theta}_{S^c} = 0$.
- 2 Determine $(\hat{\theta}_S, \hat{z}_S) \in \mathbb{R}^s \times \mathbb{R}^s$ by solving the *oracle subproblem*

$$\hat{\theta}_S \in \arg \min_{\theta_S \in \mathbb{R}^s} \left\{ \underbrace{\frac{1}{2n} \|y - X_S \theta_S\|_2^2 + \lambda_n \|\theta_S\|_1}_{=: f(\theta_S)} \right\}, \quad (4)$$

and then choosing $\hat{z}_S \in \partial \|\hat{\theta}_S\|_1$ such that $\nabla f(\theta_S)|_{\theta_S = \hat{\theta}_S} + \lambda_n \hat{z}_S = 0$.

- 3 Solve for $\hat{z}_{S^c} \in \mathbb{R}^{d-s}$ via the zero-subgradient equation, and check whether or not the *strict dual feasibility* condition $\|\hat{z}_{S^c}\|_\infty < 1$ holds.

Witness for LASSO

Lemma

If the lower eigenvalue condition holds, then success of the PDW construction implies that the vector $(\hat{\theta}_S, 0) \in \mathbb{R}^d$ is the unique optimal solution of the Lasso.

Proof

Proof of main theorem