

Recall our setup:

- Estimator

$$\hat{\theta} \in \underset{\theta \in \mathbb{R}^d}{\operatorname{argmin}} \left\{ \mathcal{L}(\theta; \{z_i\}_{i=1}^n) + \lambda R(\theta) \right\};$$

$$z_i = (x_i, y_i); \quad y = X\theta^* + w. \quad \rightarrow \text{norm}$$

-  $m \subseteq \bar{m} \subseteq \mathbb{R}^d$ : model subspace,

$\bar{m}^\perp$ : perturbation subspace,

Decomposability:  $R(\theta + \gamma) = R(\theta) + R(\gamma) \quad \forall \theta \in m, \gamma \in \bar{m}^\perp$

Proposition: Suppose that  $\mathcal{L}$  is smooth and convex,

$\lambda \geq 2 R^*(\nabla \mathcal{L}(\theta^*))$  and  $R$  is decomposable wrt  $(m, \bar{m}^\perp)$ .

$\rightarrow$  dual norm:  $R^*(\theta) = \max_{\eta: R(\eta) \leq 1} \eta^T \theta$  e.g.:  $R(\theta) = \|\theta\|_2, R^*(\theta) = \max_{\|\eta\|_2 \leq 1} \eta^T \theta = \|\theta\|_2$

Then,

$$\hat{\Delta} = \hat{\theta} - \theta^* \in \mathcal{C} \triangleq \left\{ \Delta \in \mathbb{R}^d: R(\Delta_{\bar{m}^\perp}) \leq 3R(\Delta_{\bar{m}}) + 4R(\theta_{\bar{m}^\perp}^*) \right\}$$

Remarks:

1) Consider  $\mathcal{L}(\theta, \{z_i\}_{i=1}^n) = \frac{1}{2} \|y - X\theta\|_2^2$ . Then,  $\nabla \mathcal{L}(\theta^*) = -X^T(y - X\theta^*) = -X^T w$ . Hence, if  $R(\theta) = \|\theta\|_1$ , then  $R^*(\theta) = \|\theta\|_\infty = \max_i |\theta_i|$

and we require  $\lambda \geq 2 \|X^T w\|_\infty$ . The last quantity can be estimated if we assume probabilistic models on  $X$  and  $w$ .

2) If  $\theta^* \in m$ , then  $R(\theta_{\bar{m}^\perp}^*) = 0$  and  $\mathcal{C}$  is a cone: If  $\Delta \in \mathcal{C}$ , then  $t\Delta \in \mathcal{C} \quad \forall t > 0$ .

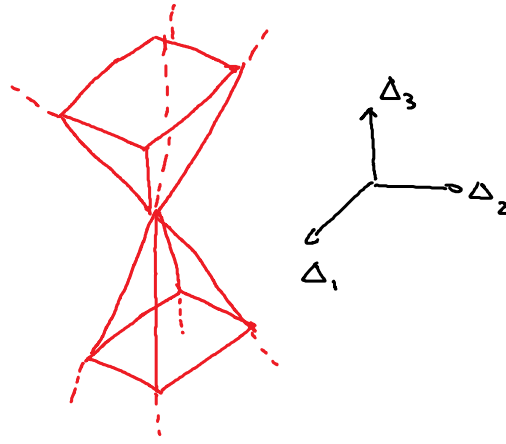
Example:  $d=3$ , 1-sparse vector  $\theta = \text{support set} = \{3\}$

$$m(\mathcal{S}) = \{ \theta \in \mathbb{R}^3 : \theta_1 = \theta_2 = 0 \} = \{ \theta \in \mathbb{R}^3 : \theta_i = 0 \quad \forall i \notin \mathcal{S} \}$$

1°:  $\theta^* = (0, 0, x)$  correctly specified  $m$

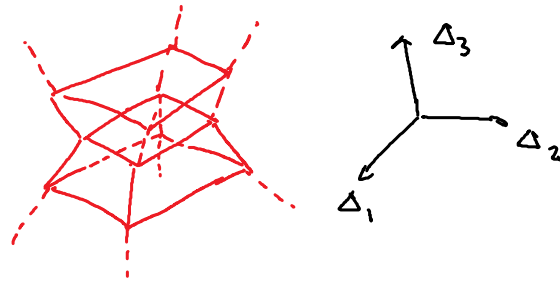
$$R(\theta) = \|\theta\|_1, \quad m = \bar{m}; \quad m^\perp = \{ \theta \in \mathbb{R}^3 : \theta_3 = 0 \}$$

$$\therefore \mathcal{C} = \{ \Delta \in \mathbb{R}^3 : |\Delta_1| + |\Delta_2| \leq 3|\Delta_3| \}$$



2°:  $\theta^* = (0, x, xx)$  incorrectly specified  $m$

$$\mathcal{C} = \{ \Delta \in \mathbb{R}^3 : |\Delta_1| + |\Delta_2| \leq 3|\Delta_3| + \underbrace{4|\theta_2^*|}_{\text{fixed}} \}$$



Proof of Proposition:

Define

$$\mathcal{D}(\Delta) = \mathcal{L}(\theta^* + \Delta) - \mathcal{L}(\theta^*) + \lambda(\mathcal{R}(\theta^* + \Delta) - \mathcal{R}(\theta^*))$$

(recall: we are minimizing  $\mathcal{L}(\theta) + \lambda\mathcal{R}(\theta)$ ). Note that with

$\hat{\Delta} = \hat{\theta} - \theta^*$ , we have  $\mathcal{D}(\hat{\Delta}) \leq 0$  by optimality of  $\hat{\theta}$

Claim 1:  $\mathcal{R}(\theta^* + \Delta) - \mathcal{R}(\theta^*) \geq \mathcal{R}(\Delta_{\bar{m}^\perp}) - \mathcal{R}(\Delta_{\bar{m}}) - 2\mathcal{R}(\theta_{\eta^\perp}^*)$

Claim 2: If  $\lambda \geq 2\mathcal{R}^*(\nabla\mathcal{L}(\theta^*))$  and  $\mathcal{L}$  is convex, then

$$\mathcal{L}(\theta^* + \Delta) - \mathcal{L}(\theta^*) \geq -\frac{\lambda}{2} [\mathcal{R}(\Delta_{\bar{m}}) + \mathcal{R}(\Delta_{\bar{m}^\perp})]$$

With these claims:

$$0 \geq \mathcal{D}(\hat{\Delta}) \geq \lambda [\mathcal{R}(\hat{\Delta}_{\bar{m}^\perp}) - \mathcal{R}(\hat{\Delta}_{\bar{m}}) - 2\mathcal{R}(\theta_{\eta^\perp}^*)] \leftarrow \text{Claim 1}$$

$$-\frac{\lambda}{2} [R(\hat{\Delta}_{\bar{m}}) + R(\hat{\Delta}_{\bar{m}^\perp})] \quad \leftarrow \text{claim 2}$$

$$= \frac{\lambda}{2} [R(\hat{\Delta}_{\bar{m}^\perp}) - 3R(\hat{\Delta}_{\bar{m}}) - 4R(\Theta_{\eta^\perp}^*)]$$

Proof of Claim 1: We want to prove

$$R(\Theta^* + \Delta) - R(\Theta^*) \geq R(\Delta_{\bar{m}^\perp}) - R(\Delta_{\bar{m}}) - \underline{2R(\Theta_{\eta^\perp}^*)}$$

We compute  $\begin{matrix} \Theta^* \\ \hline \end{matrix} = \begin{matrix} \Theta_m^* \\ \hline \end{matrix} + \begin{matrix} \Theta_{\eta^\perp}^* \\ \hline \end{matrix}$

$$R(\Theta^* + \Delta) = R(\Theta_m^* + \Theta_{\eta^\perp}^* + \Delta_{\bar{m}} + \Delta_{\bar{m}^\perp})$$

$$\geq R(\Theta_m^* + \Delta_{\bar{m}^\perp}) - R(\Theta_{\eta^\perp}^* + \Delta_{\bar{m}})$$

$$\geq R(\Theta_m^* + \Delta_{\bar{m}^\perp}) - R(\Theta_{\eta^\perp}^*) - R(\Delta_{\bar{m}})$$

} triangle inequality

$$= \underline{R(\Theta_m^*)} + R(\Delta_{\bar{m}^\perp}) - \underline{R(\Theta_{\eta^\perp}^*)} - R(\Delta_{\bar{m}}) \quad \text{decomposability}$$

On the other hand,

$$R(\Theta^*) \leq \underline{R(\Theta_m^*)} + \underline{R(\Theta_{\eta^\perp}^*)} \quad \text{triangle inequality}$$

Summary: We know  $\hat{\Delta} = \hat{\Theta} - \Theta^* \in \mathcal{C}$ .

To bound the error, it suffices to show that  $\mathcal{L}$  is not "too flat" on  $\mathcal{C}$ . This motivates the following

definition:

Definition: (Restricted Strong Convexity (RSC))

We say that  $\mathcal{L}$  is RSC on  $\mathcal{C}$  if there exists

a constant  $\kappa > 0$  and a function  $\tau(\cdot)$  s.t.

$$\underline{\mathcal{L}(\Theta^* + \Delta)} \geq \underline{\mathcal{L}(\Theta^*)} + \nabla \mathcal{L}(\Theta^*)^\top \Delta + \kappa \|\Delta\|_2^2 - \underline{\tau(\Theta^*)} \quad \forall \Delta \in \mathcal{C}.$$

Statistical tolerance "restricted"

Example: Consider  $\mathcal{L}(\theta) = \frac{1}{2} \|y - X\theta\|_2^2$ . Then,

$$\mathcal{L}(\theta^* + \Delta) - \mathcal{L}(\theta^*) - \nabla \mathcal{L}(\theta^*)^T \Delta = \frac{1}{2} \|X\Delta\|_2^2 \quad (\text{check})$$

Then, the RSC property means

$$\frac{1}{2} \|X\Delta\|_2^2 \geq \kappa \|\Delta\|_2^2 - \tau^2(\theta^*)$$

Next: With the above, we have

Theorem: Under the assumptions of the Proposition and the assumption that  $\mathcal{L}$  is RSC on  $\mathcal{C}$ , we have

$$\|\hat{\theta} - \theta^*\|_2^2 \leq \underbrace{\frac{9\lambda^2}{4\kappa^2} \Psi^2(\bar{m})}_{\text{tightness measure}} + \underbrace{\frac{2}{\kappa} \left[ \tau^2(\theta^*) + 2\lambda R(\theta_{\eta^*}^*) \right]}_{\text{model specification error}},$$

where

$$\Psi(\mathcal{M}) = \sup_{u \in \mathcal{M} \setminus \{0\}} \frac{R(u)}{\|u\|_2}$$

can be regarded as the Lipschitz constant of  $R$  over the subspace  $\mathcal{M}$ .

Example: Consider  $R(\theta) = \|\theta\|_2$ . For  $\theta \in \mathbb{R}^d$ ,

$$\|\theta\|_2 \leq \sqrt{d} \|\theta\|_2 \Rightarrow \Psi(\mathbb{R}^d) = \sqrt{d}.$$

However, consider  $\mathcal{M} = \{ \theta \in \mathbb{R}^d : \theta_j = 0 \ \forall j \notin \mathcal{S} \}$  with  $|\mathcal{S}| = s \ll d$ . Then, one can show

$$\Psi(\mathcal{M}) = \sqrt{s}$$