

Variable selection:  $\ell_0$ -penalty  $\sigma^2 k \log d / n$   
 $\ell_1$ -penalty : slow rate as  $\sqrt{\frac{\log d}{n}}$

impact of dimension

$$\frac{\sigma^2 d}{n} \sim \frac{1}{M^{2/2+d}} \underset{\text{use of dim}}{\sim} \frac{1}{M^{2/d}} \frac{k-M}{k}$$

① how to pick

k "good" variables

② Can it work?

Need assumptions

variable selection

$\rightarrow d$  smaller  $\rightarrow k$

Focus on least-squares with linear models

$\ell_2$ -regularization

class 2

$\ell_1$  or  $\ell_0$

class last week

Constrained least-squares regression

$(x_i, y_i) \in X \times \mathbb{R}, i=1, \dots, n$ ,  $f_\theta(x) = \theta^\top \varphi(x)$

Design matrix  $\begin{pmatrix} \varphi(x_1)^\top \\ \vdots \\ \varphi(x_n)^\top \end{pmatrix}$

feature vector

# Design matrix

$$\text{Sign matrix } \phi \in \mathbb{R}^{n \times d} = \begin{pmatrix} \varphi(n_1) \\ \vdots \\ \varphi(n_m)^T \end{pmatrix}$$

$y \in \mathbb{R}^n$

feature vector

f | outputs  
| responses  
| labels

Empirical risk:

$$\frac{1}{n} \sum_{i=1}^n \|y_i - \phi(\hat{w})\|^2 = \frac{1}{n} \|y - \phi(\hat{w})\|_2^2$$

## Automação de performances

$$\frac{1}{m} \left| (\phi(\theta - \theta_*)) \right|^2$$

Fixed design assumption:  $x_1, \dots, x_n$  are deterministic  
 $y_1, \dots, y_n$  are random

$y_i = \phi(x_i^T \theta_*) + \epsilon_i$  model is well-specified

E Gaussian with mean  $\theta$  and covariance matrix  $\sigma^2 I$

For OLS:  $\hat{\theta} = \arg \min \frac{1}{n} \|y - \phi\theta\|_2^2$        $\hat{\theta} = (\phi^\top \phi)^{-1} \phi^\top y$

Direct computations:  $y = \phi\theta_* + \varepsilon$

Add constraints:  $\theta \in \Theta$

$$\hat{\theta} = \arg \min_{\theta \in \Theta} \frac{1}{n} \|y - \phi\theta\|_2^2 \Rightarrow \|y - \phi\hat{\theta}\|_2^2 \leq \|y - \phi\theta_*\|_2^2$$

$$\begin{aligned} y - \phi\hat{\theta} &= \phi\theta_* - \phi\hat{\theta} + \varepsilon \\ &= \phi(\theta_* - \hat{\theta}) + \varepsilon \end{aligned}$$

Optimal:  $\|\phi(\hat{\theta} - \theta_*)\|_2^2 \leq 2\|\varepsilon\|_2 \cdot \|\phi(\hat{\theta} - \theta_*)\|_2$  (Cauchy-Schwarz)

$$\begin{aligned} \|\phi(\hat{\theta} - \theta_*)\|_2 &\leq 2\|\varepsilon\|_2 \\ \frac{1}{n} \|\phi(\hat{\theta} - \theta_*)\|_2^2 &\leq \frac{4\|\varepsilon\|_2^2}{n} \quad \text{with } \mathbb{E}\|\varepsilon\|_2^2 = \sigma^2 \end{aligned}$$

Optimal 2:  $\|\phi(\hat{\theta} - \theta_*)\|_2 \leq 2\varepsilon \sqrt{\frac{\|\phi(\hat{\theta} - \theta_*)\|_2}{\|\phi(\hat{\theta} - \theta_*)\|_2}}$

$$\leq 2 \sup_{\theta \in \Theta} \varepsilon \sqrt{\frac{\|\phi(\theta - \theta_*)\|_2}{\|\phi(\theta - \theta_*)\|_2}}$$

$\mathbb{E} \|\phi(\hat{\theta} - \theta_*)\|_2^2 = \sigma^2$

~~$\|\phi(\theta_* - \hat{\theta})\|_2^2 + \|\varepsilon\|_2^2 + 2\varepsilon^\top \phi(\theta_* - \hat{\theta}) \leq \|\varepsilon\|_2^2$~~

$\|\phi(\hat{\theta} - \theta_*)\|_2^2 \leq 2\varepsilon^\top \phi(\hat{\theta} - \theta_*)$

$\ell_0$ -penalty /  $\ell_0$ -constraint:  $\|\theta\|_0 = \# \text{ of non-zero components of } \theta$

$$\mathcal{C} = \{\theta \in \mathbb{R}^d; \|\theta\|_0 \leq h\}$$

$$\hat{\theta}^T \varphi(x) = \sum_{j=1}^d \hat{\theta}_j \varphi(x)_j$$

$$\hat{\theta} = \arg \min_{\|\theta\|_0 \leq h} \|y - \varphi\theta\|_2^2$$

$$\|\theta\|_0 \leq h$$

$$= \arg \min_{\|A\| \leq h}$$

$$\begin{aligned} & \min_{\theta_{AC}=0} \|y - \varphi\theta\|_2^2 \\ & \theta_{AC}=0 \end{aligned}$$

function of  $A$   
is mon-increasing

$$\|\theta\|_0 = \sum_{j=1}^d |\theta_j|^0 / "0^0 = 0"$$

can restrict  
to  $|A| \leq h$

Algorithms:

- ① exhaustive search a sets  $|A| \leq h$

Support of  $\theta$   
There are  $\binom{d}{h}$  such sets

$$= \frac{d!}{h!(d-h)!} \leq d^h$$

$$\leq \left(\frac{ed}{h}\right)^h \Rightarrow \text{see back}$$

② greedy: efficient but not optimal

$$\begin{aligned} \|\phi(\hat{\theta} - \theta_\infty)\|_2 &\leq 2\epsilon^T \left( \frac{\phi(\hat{\theta} - \theta_\infty)}{\|\phi(\hat{\theta} - \theta_\infty)\|_2} \right) \\ &\leq 2 \sup_{\|\theta\|_2 \leq h} \epsilon^T \left( \frac{\phi(\theta - \theta_\infty)}{\|\phi(\theta - \theta_\infty)\|_2} \right) \end{aligned}$$

$$\|\theta_\infty\| \leq h$$

$$\|\theta - \theta_\infty\|_2 \leq 2h$$

$$\leq 2 \sup_{\|\theta - \theta_\infty\|_2 \leq 2h} \epsilon^T \frac{\phi(\theta - \theta_\infty)}{\|\phi(\theta - \theta_\infty)\|_2}$$

$$\leq 2 \sup_{|B| \leq 2h} \sup_{(\theta - \theta_\infty)_B = 0} \epsilon^T \frac{\phi(\theta - \theta_\infty)}{\|\phi(\theta - \theta_\infty)\|_2}$$

$$\boxed{\epsilon^T \frac{\phi(\theta - \theta_\infty)}{\|\phi(\theta - \theta_\infty)\|_2}}$$

$$\begin{aligned} &\leq 2 \sup_{|B| \leq 2h} \sup_{\substack{\|\beta\|_2 = 1 \\ \beta \in \text{Im } \phi_B}} \epsilon^T \beta \\ &= \|\phi_B \epsilon\|_2 \end{aligned}$$

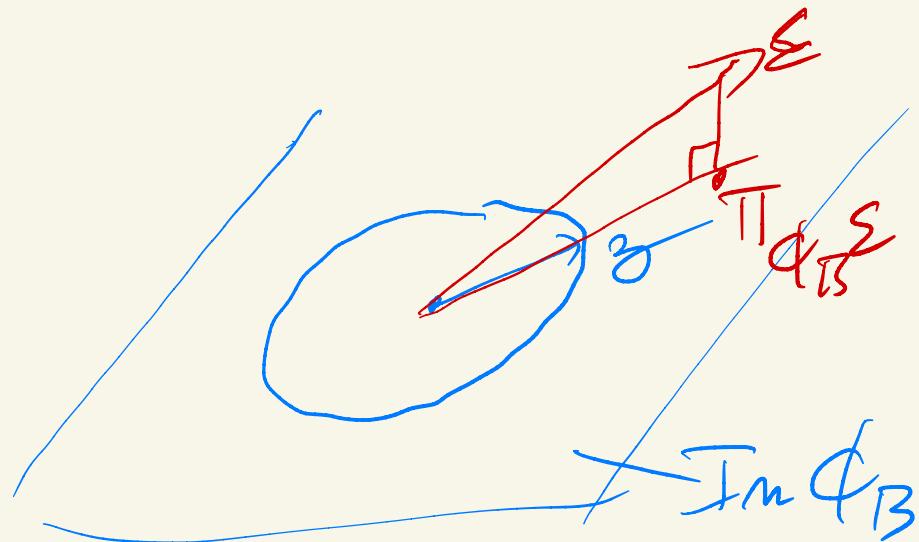
$$M \begin{pmatrix} d \\ \beta \\ \vdots \\ \beta \end{pmatrix} = \begin{pmatrix} d \\ \beta \\ \vdots \\ \beta \end{pmatrix}$$

$$\|\beta\|_2 = 1$$

$$\beta \in \text{Im } \phi_B$$

Submatrix of  $\phi$   
with columns indexed  
by  $B$

$$2 \sup_{|\beta| \leq 2h} \left\{ \sup_{\gamma \text{ s.t. } \|\beta\|_2 = 1, \gamma \in \text{Im } \phi_B} \varepsilon^\top \beta \right\} = \|\Pi \phi_B \varepsilon\|_2$$



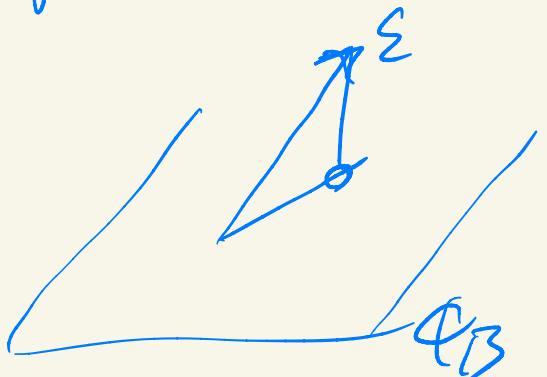
$$\begin{aligned} \varepsilon^\top \beta &= (\Pi \phi_B \varepsilon + \varepsilon - \Pi \phi_B \varepsilon)^\top \beta \\ &= \beta^\top \Pi \phi_B \varepsilon + 0 \leq \|\beta\|_2 \|\Pi \phi_B \varepsilon\|_2 \quad \text{by Cauchy-Schwarz} \end{aligned}$$

Summary

$$\|\phi(\hat{\theta} - \theta_*)\|_2 \leq 2 \sup_{|\beta| \leq 2h} \|\Pi \phi_B \varepsilon\|_2$$

$\varepsilon$  is Gaussian with covariance matrix  $\sigma^2 I$   
 $\Pi \phi_B \varepsilon$  is Gaussian

If  $\varepsilon$  is Gaussian with mean zero and covariance matrix  $\sigma^2 I$



$\Pi_{Q_B} \varepsilon$  is a Gaussian in dimension  $|B|$  and covariance matrix  $\sigma^2 I$

WLOG,  $Q_B = \text{span}(e_1, \dots, e_{|B|})$

canonical basis of  $\mathbb{R}^d$

$$\mathbb{E} \| \Pi_{Q_B} \varepsilon \|^2 = \sigma^2 |B|$$

$$\mathbb{E} \sup_{|B|=2h} \| \Pi_{Q_B} \varepsilon \|^2$$

$\boxed{|B|=2h}}$

$$= \| \beta_B \|_2^2$$

$m = \binom{d}{2h}$

$$\mathbb{E} e^{S|B|/2}$$

$\boxed{|B|/2}$   $R$

$$= \frac{1}{(1 - 2\sigma^2 S)}$$

Lemma 1:  $z_1, \dots, z_m$  random variables not necessarily independent

$$\begin{aligned} \mathbb{E} \max\{z_1, \dots, z_m\} &\leq \frac{1}{s} \log (\mathbb{E} e^{s \max\{z_1, \dots, z_m\}}) \\ &\leq \frac{1}{s} \log (\mathbb{E} \max\{e^{sz_1}, \dots, e^{sz_m}\}) \quad (\text{exp is increasing}) \end{aligned}$$

$$\mathbb{E} Y \leq \frac{1}{s} \log (\mathbb{E} e^{sY})$$

$t \mapsto e^{st}$  convex

$$\mathbb{E}(e^{sY}) \geq e^{s\mathbb{E} Y}$$

Jensen's

$$\mathbb{E} \max\{z_1, \dots, z_m\} \leq \frac{1}{s} \log (\mathbb{E} e^{sz_1} + \dots + \mathbb{E} e^{sz_m})$$

HS > C

Lemma 2: if  $y$  is Gaussian with covariance matrix  $\sigma^2 I$

$y \in \mathbb{R}^h$

then

$$\left[ \mathbb{E} e^{s \|y\|_2^2} \right]$$

$$= \mathbb{E} \prod_{i=1}^h e^{sy_i^2} = \prod_{i=1}^h \mathbb{E} e^{sy_i^2} = \left( \mathbb{E} e^{sg_i^2} \right)^h$$

sy independence

$$= \left( \frac{1}{\sqrt{2\pi\sigma^2}} \right) \left( e^{sy_1^2} e^{-\frac{1}{2} \frac{y_1^2}{\sigma^2}} dy_1 \right)^h$$

$$= \left( \frac{1}{\sqrt{4\sigma^2 - 2s}} \right)^d$$

$$= \left( \frac{1}{\sqrt{1 - 2s\sigma^2}} \right)^{\frac{h}{2}}$$

$$s < \frac{1}{2\sigma^2}$$

$$\mathbb{E} \sup_{\|B\|=2h} \|\Pi \Phi_B \varepsilon\|_2^2$$

$\| \Phi_B \varepsilon \|_2^2$

$\| \Phi_B \varepsilon \|_2^2$

$m = \binom{d}{2h}$

$$\mathbb{E} e^{S \|\Phi_B \varepsilon\|_2^2}$$

$\left( \frac{1}{1 - 2\sigma^2 S} \right) h$

0.32
0.50

$$\mathbb{E} \sup_{\|B\|=2h} \|\Pi \Phi_B \varepsilon\|_2^2 \leq \frac{1}{S} \log \left[ m \frac{1}{(-2\sigma^2 S)^h} \right]$$

$S = \frac{1}{4\sigma^2}$

$$\leq 4\sigma^2 \log [m^{2h}]$$

$$\leq 4\sigma^2 \left[ \log \binom{d}{2h} + h \log 2 \right]$$

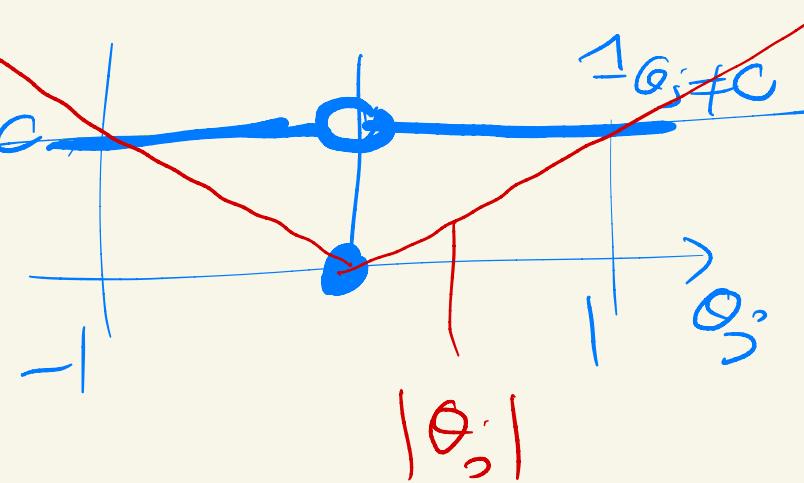
*price of adaptivity*

$$\leq \log \left( \frac{(ed)^{2h}}{2h} \right) = 2h \log \left( \frac{ed}{h} \right)$$

$$\leq 2h + \boxed{2h \log \frac{d}{h}}$$

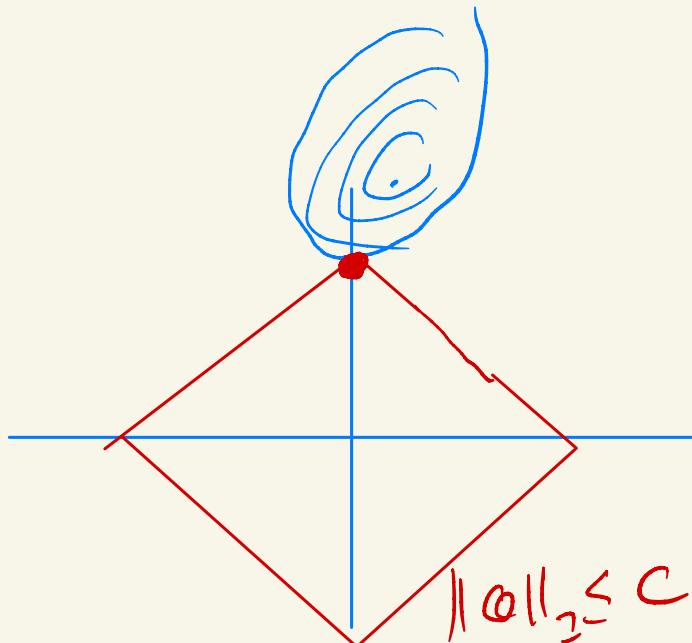
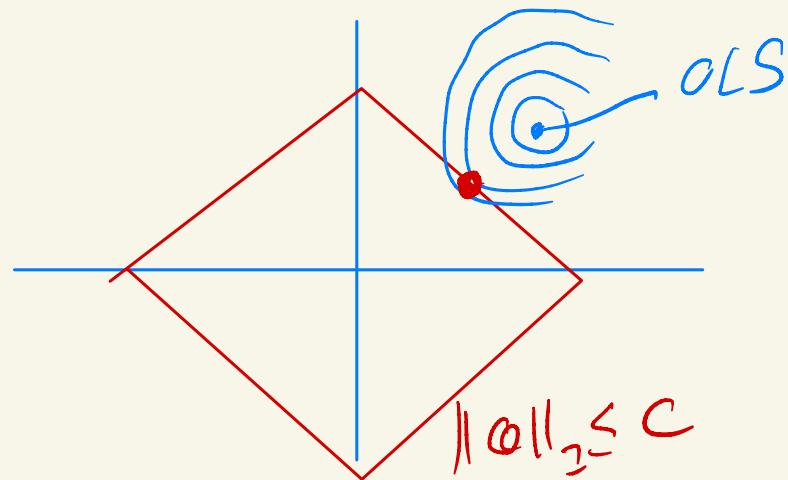
From  $\ell_0$  to  $\ell_1$ :  $\|\theta\|_0 = \sum_{j=1}^d 1_{\theta_j \neq 0}$

$\|\theta\|_1$  is the convex envelope of  $\|\theta\|_0$  on  $\mathbb{E}^{(1,1)^d}$



constraint or penalize by  $\ell_1$ -norm

$$\min \|(\phi\theta - y)\|^2 \text{ s.t. } \|\theta\|_1 \leq c$$



Lasso :  $\min_{\theta} \frac{1}{2n} \|y - \phi\theta\|_2^2 + \lambda \|\theta\|_1$

$$+ \lambda \sum_j |\theta_j|$$

## ① Algorithms

"y-trick" :  $|\theta_j| = \inf_{q_j > 0} \frac{1}{2} \frac{\theta_j^2}{q_j} + \frac{1}{2} q_j$  y-trick

with  $q_j^* = |\theta_j|$

$$\min_{\theta \in \mathbb{R}^d} \min_{y \in \mathbb{R}_+^d} \frac{1}{2n} \|y - \phi\theta\|_2^2 + \frac{\lambda}{2} \sum_j \frac{\theta_j^2}{q_j} + \frac{\lambda}{2} \sum_j q_j$$

With  $y$  fixed = quadratic problem in  $\theta \Rightarrow$  closed form  
 With  $\theta$  fixed  $\Rightarrow q_j^* = |\theta_j|$

$\Rightarrow$  alternating minimization

for each  $y \Rightarrow$  proximal methods  $\Rightarrow$  see blog post

(2) Analysis

- slow rate with "nb" assumptions  $\sqrt{\frac{\log n}{n}}$
- fast rate with strong assumptions  $\frac{1}{\sqrt{n}}$

$\sigma^2 \text{ hlgd}$   $\Rightarrow$  see back

Consider  $\hat{\theta} = \arg \min_{\theta} \frac{1}{2n} \|y - \Phi\theta\|_2^2 + \lambda \mathcal{R}(\theta)$

$$y = \Phi\theta_* + \varepsilon$$

$$\frac{1}{2n} \|y - \Phi\hat{\theta}\|_2^2 + \lambda \mathcal{R}(\hat{\theta}) \leq \frac{1}{2n} \|y - \Phi\theta_*\|_2^2 + \lambda \mathcal{R}(\theta_*)$$

$$\frac{1}{2n} \left[ \|\Phi(\theta_* - \hat{\theta})\|_2^2 + \|\varepsilon\|_2^2 + 2\varepsilon^\top \Phi(\theta_* - \hat{\theta}) \right] + \lambda \mathcal{R}(\hat{\theta}) \leq \frac{1}{2n} \|\varepsilon\|_2^2 + \lambda \mathcal{R}(\theta_*)$$

$$\begin{aligned} \|\Phi(\hat{\theta} - \theta_*)\|_2^2 &\leq 2\Omega^*(\Phi^\top \varepsilon) \Omega(\hat{\theta} - \theta_*) + 2n\lambda\Omega(\theta_*) - 2n\lambda\Omega(\hat{\theta}) \\ &\leq n\lambda\Omega(\hat{\theta} - \theta_*) + 2n\lambda\Omega(\theta_*) - 2n\lambda\Omega(\hat{\theta}) \\ &\leq n\lambda\Omega(\hat{\theta}) + n\lambda\Omega(\theta_*) + 2n\lambda\Omega(\theta_*) - 2n\lambda\Omega(\hat{\theta}) \leq 3n\lambda\Omega(\theta_*) - n\lambda\Omega(\hat{\theta}). \end{aligned}$$

Lemma:  $\mathcal{R}^*(u) = \sup_{\mathcal{R}(\beta) \leq 1} \beta^\top u \Rightarrow$  dual norm

$$\text{Dual norm} : \Sigma^*(u) = \sup_{\|z\|_1 \leq 1} z^T u$$

$$\Sigma: l_2 = \Sigma^* = l_2$$

$$\Sigma: l_1: \Sigma^* = l_\infty$$

$$\Sigma: l_\infty: \Sigma^* = l_1$$

$$\Sigma: l_p, \Sigma^* = l_q$$

$$\frac{1}{p} + \frac{1}{q} = 1$$

$$\text{Property} = z^T u \leq \Sigma(z) \Sigma^*(u)$$

Proof

$$z^T u = \frac{z^T u}{\Sigma(z)} \cdot \Sigma(z) \leq \Sigma(z) \sup_{\|z'\|_1 \leq 1} z'^T u$$

$$\text{"Cauchy-Schwarz"} \\ \Sigma(z) \Sigma^*(u) = z^T u$$

$$\begin{aligned} \sup_{\|z\|_1 \leq 1} z^T u &= \sum_{j=1}^d z_j u_j \\ &\leq \sum_j |z_j| \max_j |u_j| \\ &\leq \max_j |u_j| \\ &= \|u\|_\infty \end{aligned}$$

Hölder-Lp

Proposition: if  $\mathcal{R}^*(\Phi^\top \Sigma) \leq \frac{n\lambda}{2}$  then  $\underbrace{\frac{1}{2} \|\Phi(\hat{\theta} - \theta_*)\|_2^2}_{\text{in performance}} \leq 3\lambda \mathcal{R}(\theta_*)$

Proof: Model:  $y = \Phi \theta_* + \varepsilon$

$$\frac{1}{2n} \|y - \Phi \hat{\theta}\|_2^2 + \lambda \mathcal{R}(\hat{\theta}) \leq \frac{1}{2n} \|y - \Phi \theta_*\|_2^2 + \lambda \mathcal{R}(\theta_*)$$

$$\frac{1}{2n} \left[ \|\Phi(\theta_* - \hat{\theta})\|_2^2 + \|\varepsilon\|_2^2 + 2\varepsilon^\top \Phi(\theta_* - \hat{\theta}) \right] + \lambda \mathcal{R}(\hat{\theta}) \leq \frac{1}{2n} \|\varepsilon\|_2^2 + \lambda \mathcal{R}(\theta_*)$$

$$\begin{aligned} \|\Phi(\hat{\theta} - \theta_*)\|_2^2 &\leq 2\Omega^*(\Phi^\top \varepsilon)\Omega(\hat{\theta} - \theta_*) + 2n\lambda\Omega(\theta_*) - 2n\lambda\Omega(\hat{\theta}) \\ &\leq n\lambda\Omega(\hat{\theta} - \theta_*) + 2n\lambda\Omega(\theta_*) - 2n\lambda\Omega(\hat{\theta}) \\ &\leq n\lambda\Omega(\hat{\theta}) + n\lambda\Omega(\theta_*) + 2n\lambda\Omega(\theta_*) - 2n\lambda\Omega(\hat{\theta}) \leq 3n\lambda\Omega(\theta_*) - n\lambda\Omega(\hat{\theta}). \end{aligned}$$

Using  $\varepsilon^\top \Phi(\hat{\theta} - \theta_*) \leq \mathcal{R}(\hat{\theta} - \theta_*) \mathcal{R}^*(\Phi^\top \Sigma)$

Proposition: if  $\|\boldsymbol{\varphi}^*(\boldsymbol{\varphi}^\top \boldsymbol{\varepsilon})\| \leq \frac{n\epsilon}{2}$  then  $\frac{1}{n} \|\boldsymbol{\varphi}(\hat{\boldsymbol{\theta}} - \boldsymbol{\theta}_*)\|_2^2 \leq 3 \lambda \sigma^2(\boldsymbol{\theta}_*)$

"performance"

Consequences:  $\ell_1$ -norm  
Band with high probability

$$\begin{aligned} P\left(\|\boldsymbol{\varphi}^*(\boldsymbol{\varphi}^\top \boldsymbol{\varepsilon})\| \geq \frac{n\epsilon}{2}\right) &= P\left(\|\boldsymbol{\varphi}^\top \boldsymbol{\varepsilon}\|_\infty \geq \frac{n\epsilon}{2}\right) \\ &= P\left(\sup_{j \in \{1, \dots, d\}} |(\boldsymbol{\varphi}^\top \boldsymbol{\varepsilon})_j| \geq \frac{n\epsilon}{2}\right) \\ &\stackrel{\text{union band}}{\rightarrow} \leq \sum_{j=1}^d P(|(\boldsymbol{\varphi}^\top \boldsymbol{\varepsilon})_j| \geq \frac{n\epsilon}{2}) \end{aligned}$$

If  $\boldsymbol{\varepsilon}$  is Gaussian  
with mean zero  
and covariance

matrix  $\boldsymbol{\Sigma}^\top \boldsymbol{\Sigma}$ ,

$(\boldsymbol{\varphi}^\top \boldsymbol{\varepsilon})_j = \boldsymbol{\varphi}_j^\top \boldsymbol{\varepsilon} \in \mathbb{R}$  is Gaussian  
with mean  $c$   
and variance  $\boldsymbol{\varphi}_j^\top \boldsymbol{\varphi}_j \approx 1/n$ .

Assumption:  $\frac{1}{n} \sum_{i=1}^n \boldsymbol{\varphi}_i^\top \boldsymbol{\varphi}_j = \frac{1}{n} \sum_{i=1}^n \boldsymbol{\varphi}_i^\top \boldsymbol{\varphi}_j = \frac{1}{n} \sum_{i=1}^n \boldsymbol{\varphi}_j^\top \boldsymbol{\varphi}_i = \frac{1}{n}$

Lemma: if  $z \sim \text{Gaussian}$  with mean 0 and variance  $\sigma^2$

$$P(|z| \geq t) \leq 2e^{-t^2/2\sigma^2}$$

$$L = \frac{1}{\sqrt{2\pi\sigma^2}} \int_{|z|=t} e^{-z^2/2\sigma^2} dz$$

$$= \frac{1}{\sqrt{2\pi}} \int_{|z'|=\frac{t}{\sigma}} e^{-z'^2/2} dz'$$

$$\leq 2e^{-t^2/2\sigma^2}$$

$$\text{with } z' = \frac{z}{\sigma}$$

$$P(z \geq t) = P(e^{tz} \geq e^{zt})$$

rigor

$$\leq \frac{\mathbb{E} e^{tz}}{e^{zt}}$$

$$\boxed{\mathbb{E} e^{tz} = e^{\frac{1}{2}\sigma^2 t^2}}$$

$$P(|(\phi^\top \xi)_j| \geq \frac{n\delta}{2}) \leq 2 \exp\left(-\left(\frac{n^2\delta^2}{4}\right)\frac{1}{2n}\right)$$

gaussian  
with variance  $n$

$$= 2 \exp\left(-\frac{n\delta^2}{8}\right)$$

$$P\left(\sum_i (\phi^\top \xi)_{ij} \geq \frac{n\delta}{2}\right) \leq 2d \exp\left(-\frac{n\delta^2}{8}\right)$$

Consequence: for the task with probability  $-n\delta/8$

$$\geq 1 - 2de$$

$$\frac{1}{n} \|\phi(\hat{\theta} - \theta_0)\|^2 \leq 3\lambda \|\theta_0\|_1 \quad \frac{\delta}{8}$$

$$\leq 3 \|\theta_0\|_1 \sqrt{\frac{8}{n} \left( \log 2d + \log \frac{1}{\delta} \right)}$$

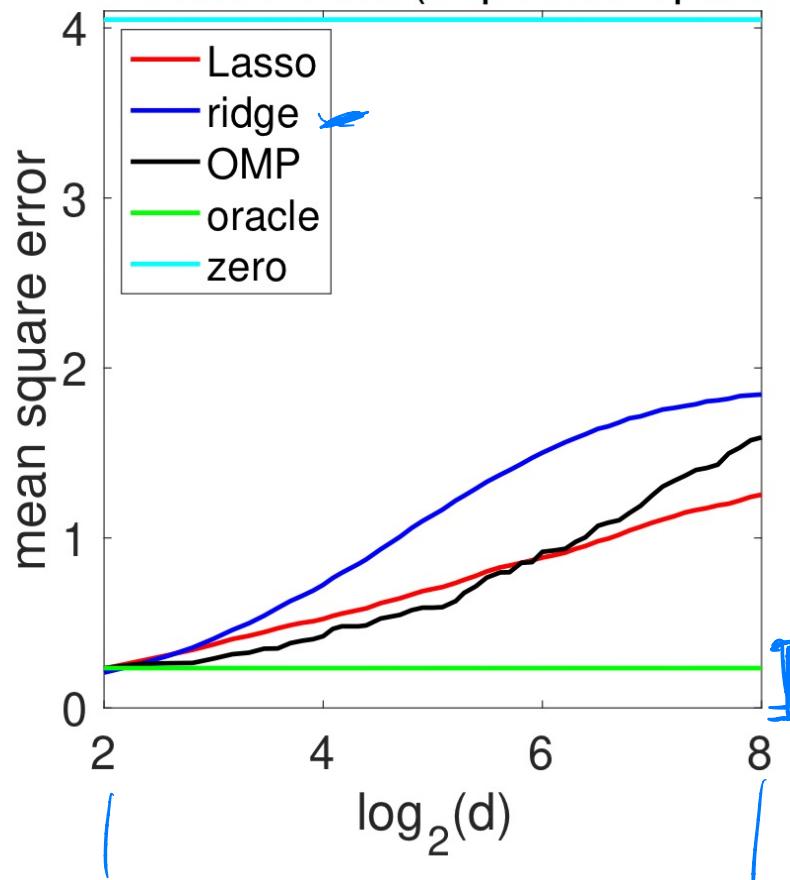
$\downarrow \log d / n$

$$\frac{n\delta^2}{8} = \log \frac{2d}{\delta}$$

$$\delta = \sqrt{\frac{8}{n} \log \frac{2d}{\delta}}$$

least-squares :  $y = \phi \theta_0 + \varepsilon$  only non-zero components in the first 4 component

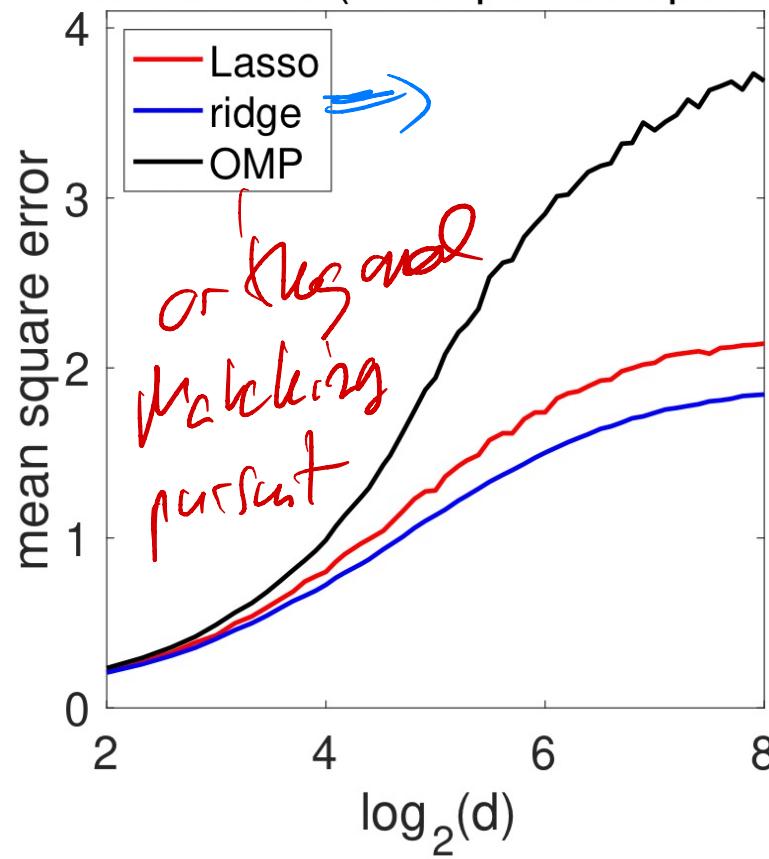
Non-rotated data (expected sparsity)



$$d=4$$

$$d=28$$

Rotated data (no expected sparsity)



$$\phi \leftarrow \phi R^{-\text{rotation}}$$

$$\phi \theta = \phi R R^\top \theta$$

$$\|\alpha\|_2 \\ \|R\alpha\|_2$$