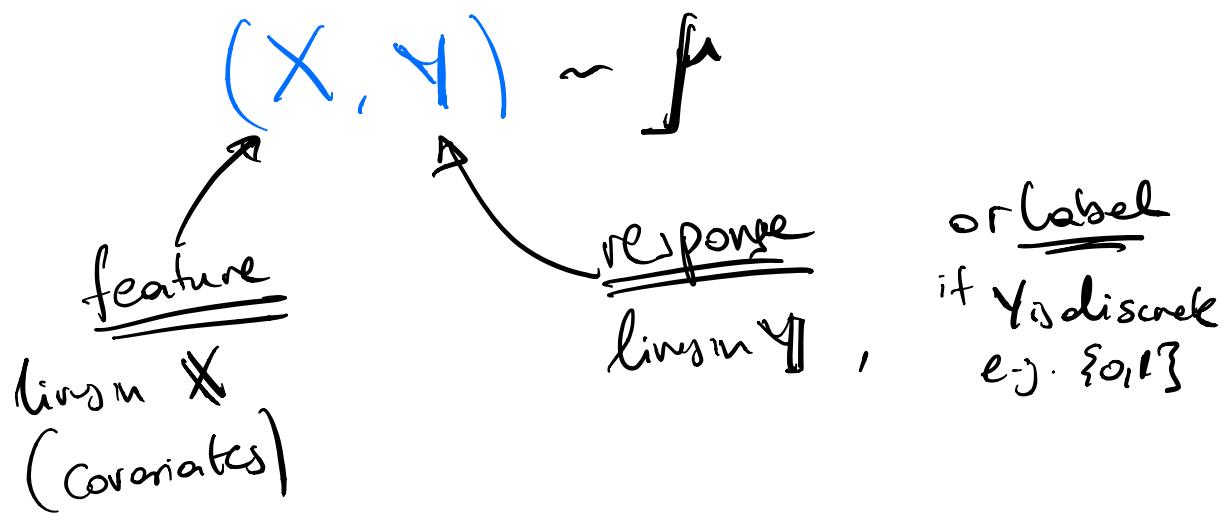


## Definitions :

The setting of the problem is the following

we have data of the form

$\{(x_i, y_i) ; i=1, \dots, n\}$ , three observations of



GOAL : "learn from a data  
how to predict  $y$  from  $x_0$ "

In some ways the best predictor of  
 $y$  given  $x$  is the  
conditional expectation  
 $E[Y | X=x]$ .

e.g. for it we want a function  $f$  that  
minimizes

$$\mathbb{E}[(Y - f(X))^2]$$

then

$$\begin{aligned} & \mathbb{E}[(Y - f(X))^2] \\ &= \mathbb{E}[(Y - \mathbb{E}[Y|X] + \mathbb{E}[Y|X] - f(X))^2] \\ &= \mathbb{E}[(Y - \mathbb{E}[Y|X])^2] + \mathbb{E}\{(\mathbb{E}[Y|X] - f(X))^2\} \\ &\quad + 2 \mathbb{E}\{(Y - \mathbb{E}[Y|X])(\mathbb{E}[Y|X] - f(X))\} \end{aligned}$$

But  $\mathbb{E}\{(Y - \mathbb{E}[Y|X])(\mathbb{E}[Y|X] - f(X))\}$

$$\begin{aligned} &= \mathbb{E}\left\{\mathbb{E}\left[(Y - \mathbb{E}(Y|X)) \times (\mathbb{E}[Y|X] - f(X)) \mid X\right]\right\} \\ &= \mathbb{E}\left\{(\mathbb{E}[Y|X] - f(X)) \times \mathbb{E}\{Y - \mathbb{E}[Y|X] \mid X\}\right\} \end{aligned}$$

(Since  $\mathbb{E}[Y|X], f(X)$  are m'ble wrt  $\sigma(X)$ )

$$= 0 \quad \text{since } \mathbb{E}[Y - \mathbb{E}[Y|X] \mid X] = \mathbb{E}[Y|X] - \mathbb{E}[Y|X]$$

Thus

$$\mathbb{E}[(Y - f(X))^2]$$

$$= \mathbb{E}[(Y - \mathbb{E}[Y|X])^2] + \mathbb{E}[(\mathbb{E}[Y|X] - f(x))^2]$$

$$\geq \mathbb{E}[(Y - \mathbb{E}[Y|X])^2].$$

The function

$\mathbb{E}[Y|X=x]$  is called the  
Regression function of  $Y$  on  $X$

In general, for any loss/risk  
we'd be done if we knew the  
Law of  $Y|X=x$  for any  $x$ .

But we only have access  
to  $(X_1, Y_1), \dots, (X_n, Y_n) =: S_n$   
data

So our goal is to use  $S_n$  to  
construct a function  $f_n$  that approximates

$$f_n = \mathbb{E}[Y|X=n], \text{ say in } L_2$$

that is we want

$$\hat{f}_n : \mathbb{X} \rightarrow \mathbb{Y}$$

s. that

$$\begin{aligned} R(f_n) &:= \mathbb{E}[(Y - \hat{f}_n(x))^2] \\ &= \mathbb{E}\left[\int (y - \hat{f}_{S_n}(x))^2 f_n(dx, dy)\right] \end{aligned}$$

is minimized (in some sense wrt  $S_n$ )

**Note!** In  $\int l(y; \hat{f}_{S_n}(x)) f_n(dx, dy)$

the integral is wrt the distribution of  $(X, Y)$

whereas the dependence on the

sample  $f_n$  is still there.

i.e. if  $\hat{f}_{S_n}$  depends non-trivially

on  $S_n$ , the quantity

$$R(\hat{f}_n) = \|f - \hat{f}_n\|_2^2 \text{ is } \underline{\text{RANDOM}}$$

Note 2:

We proved earlier that

$$\begin{aligned} R(\hat{f}_n) &= \mathbb{E}[(Y - \hat{f}_n)^2] \\ &= \mathbb{E}[(Y - f(x))^2] + \mathbb{E}[(f - \hat{f}_n)^2(x)] \\ &:= \mathbb{E}[Y|X] \end{aligned}$$

That is, for any "estimator"  $\hat{f}_n$ , the risk can be decomposed as

$$R(\hat{f}_n) = \underbrace{R(f)}_{\text{approximation error}} + \underbrace{\|f - \hat{f}\|_2^2}_{\text{estimation error}}$$

Approximation error is related to the richness of the function class we're

allowed our estimator to come from.

e.g.  $f(x) = E[Y | X=x]$  is the optimal

when we solve

$$\inf \left\{ \|Y - f\| : Y \text{ is } x\text{-mable} \in L^2 \right\}$$

+

the optimal in the class  $R(f)$ .

---

This error is not related to the statistical procedure so we will ignore it & focus on the

"estimation error"

which answers the question

Q: how efficiently can we learn  $f_{(\text{optimal})}$  from observations  $S_n$ ?

---

Note 3: bounds on  $\mathbb{E}[\|f_n - f\|_2^2]$  will

be given in expectation with high prob.

e.g. if  $Y$  is a label then one may want  $\mathbb{P}[Y \neq h(x)]$ ; the prob. of misclassification.

We have not so far discussed the question of how to construct the estimator  $\hat{f}_n$  from the data.

One principle that is often used goes

i) that of

EMPIRICAL  
RISK  
MINIMISATION

E  
R  
M

That is: the optimal solution of our problem

is given as an optimizer of an expectation

$$f^*(x) = \mathbb{E}[Y|x=x] = \inf_{f \in L^2} \mathbb{E}[(Y-f(x))^2].$$

We don't have access to the full measure for generating  $(X, Y)$  but we have samples from it, so we can approximate it

with the empirical measure

$$\hat{f}_n(dz) = \frac{1}{n} \sum_{i=1}^n \delta_{(X_i, Y_i)}(dz)$$

So

we may try

$$L_f(x, y) = (y - f(x))^2$$

.. .. .. ..

$$\hat{f}_n = \arg \inf_{f \in \mathcal{F}} f_n(L_f) = \inf_{f \in \mathcal{F}} \frac{1}{n} \sum_{i=1}^n (y_i - f(x_i))^2$$

$$\arg \inf_{f \in \mathcal{F}} f_n(L_f) = f.$$

Note: in certain cases the function class

considered (e.g. neural nets) may be rich enough to include  $g$  s.t.  $y_i = g(x_i)$   $\forall i=1 \dots n$  for all reasonable sample sizes  $n$ .  $\rightarrow$  interpolate

In principle, this may lead to overfitting.

(i) memorizing the data w/o extracting any information about  $f$ .

Deep neural nets tend to interpolate w/o overfitting, but no-one has a clear picture why or how this happens.

To avoid this we may add a

regularizer, e.g.

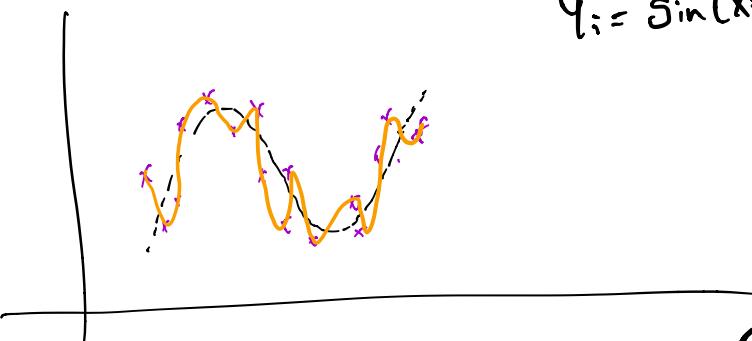
solve instead

$$\inf_g \{ R_n(g) + \text{Penalty}(g) \}$$

where the penalty may penalise overly complicated answers.

Why?

e.g. suppose



$$y_i = \sin(x_i) + \varepsilon_i$$

the orange curve achieves 0 ERM  
ie it interpolates  $\hat{f}(x_i) = y_i$  for all  $i = 1, \dots, n$

BUT: the real generating function is  $\sin(\cdot)$

so potentially here, if we are thinking of using trigonometric functions

we could perhaps use a penalty

of the form  $\sum_{k=0}^{\infty} (1+k^2) |\hat{c}_k|^2$ , where

$\hat{c}_k$  is the coefficient of  $\exp(ikx)$   
to penalise "high energy" solutions.

In fact this penalty is equivalent to

$H^1$  (Sobolev) norm so it restricts the norm of the derivatives

$$\text{i.e. } \|\hat{f}\|^2 + \|f'\|^2 \leq \sum (1+n^n) |\hat{f}_n|^2 \\ \leq C_0 (\|f\|^2 + \|f'\|^2)$$

So it indeed acts as a regularizer.

BUT: regularization may also play a different role,  
i.e. it may in some way be used to

choose among many solutions

that would o.w. look identical in terms of  
our loss function, by imposing restrictions  
such as sparsity, etc.

This will be crucial in high-dimensional problems  
where typically there are many more parameters  
than observations

$$P \gg n$$

Params                      sample size

## HIGH-DIMENSIONAL REGRESSION

### & LASSO

#### Observation Model

$$y = X\theta^* + \varepsilon$$

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

observation  
vector

$$X = \begin{bmatrix} x_{11} & x_{12} & \cdots & x_{1d} \\ x_{21} & x_{22} & \cdots & x_{2d} \\ \vdots & & & \\ x_{n1} & x_{n2} & \cdots & x_{nd} \end{bmatrix} \in \mathbb{R}^{n \times d}$$

design  
matrix

$\varepsilon \in \mathbb{R}^n$  noise vector.

$n = \#$  observations       $d = \#$  degrees of freedom?  
dimension?

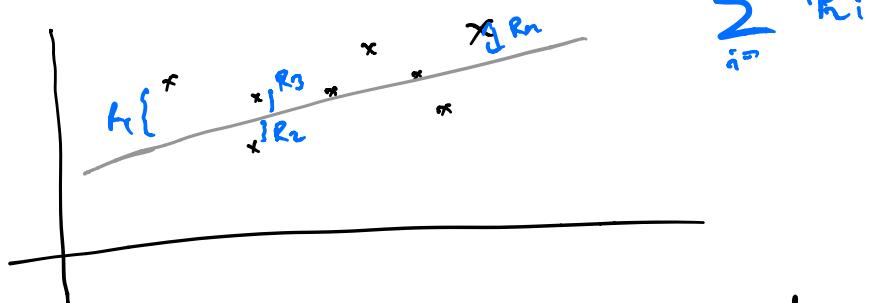
fixed design:  $X$  is deterministic

random design:  $X$  is random

WARM-UP: Let's consider first the classical  
setting i.e.  $d \ll n$ .

The problem is often written as OLS

$$\hat{\theta}^{LS} := \underset{\theta}{\operatorname{arg\,min}} \|y - X\theta\|_2^2$$
$$= \underset{\theta}{\operatorname{arg\,min}} \sum_{i=1}^n \left( y_i - \sum_{j=1}^m x_{ij} \theta_j \right)^2$$



The quadratic loss function we used is convex  
so it is not difficult to verify that  
the first order condition for optimality

$$X^T X \hat{\theta}^{LS} = X^T y$$

Assume for now that

" $X$  has full rank"  
i.e. its columns are linearly independent.

Claim: then  $X^T X$  is invertible

Pf of claim: Consider  $Xv$  for some  $v \in \mathbb{R}^d$

then  $Xv = \begin{bmatrix} x_{11} & \cdots & x_{1d} \\ \vdots & & \vdots \\ x_{n1} & \cdots & x_{nd} \end{bmatrix} \begin{bmatrix} v_1 \\ \vdots \\ v_d \end{bmatrix}$

$$= \begin{pmatrix} x_{11}v_1 + x_{12}v_2 + \cdots + x_{1d}v_d \\ \vdots \\ x_{n1}v_1 + x_{n2}v_2 + \cdots + x_{nd}v_d \end{pmatrix}$$

$$= v_1 \begin{bmatrix} x_{11} \\ \vdots \\ x_{n1} \end{bmatrix} + v_2 \begin{bmatrix} x_{12} \\ \vdots \\ x_{n2} \end{bmatrix} + \cdots + v_d \begin{bmatrix} x_{1d} \\ \vdots \\ x_{nd} \end{bmatrix}$$

$\Rightarrow$  linear combination of column vectors of  $X$   
By assumption columns of  $X$  are lin-independent  
 $Xv = 0 \Leftrightarrow v \in \text{Null}(X)$

Thus  $v^T X^T X v > 0$  if  $v \neq 0 \Rightarrow$   
 $X^T X$  "strictly" positive definite.

$\Rightarrow X^T X$  invertible

□

Thus if  $\text{rank}(X) = d$

$$\hat{\theta}^{\text{LS}} = (X^T X)^{-1} X^T y.$$

But obviously  $\text{rank}(X) \leq \min\{d, n\}$  which may only happen if  $d < n$ .

What if  $d \geq n$ :

we still require  $X^T \hat{\theta}^{\text{LS}} = X^T y$

But now all we can observe is

$$X \hat{\theta}^{\text{LS}}$$

$\nearrow$        $\uparrow$        $\nwarrow$

$(n \times d)$        $(d \times 1)$        $(n \times 1)$  an  $n$ -dimensional vector.

That is now  $\text{rank}(X) < d$   
so  $\text{null}(X) \neq \emptyset \Rightarrow$  if  $q \in \text{null}(X)$

$$X \hat{\theta} = X(\hat{\theta} + q).$$

So the system is underdetermined &  
we cannot hope to obtain a unique sol<sup>n</sup>  
→ instead we get a subspace of solutions

At this stage we may consider regularizing the problem, that is perhaps we indeed consider

$$\min_{\theta} \|\theta\|_2^2, \text{ s.t. } X^T X \theta = X^T g$$

find the minimum norm solution

It is not hard to see that the above always admits a unique solution,

think of the subspace of solutions  $X^T X \theta = X^T g$  & its distance from the origin.

In fact the solution above is denoted by

$(X^T X)^+$  is known as the (it coincides with usual inverse if invertible)

Moore-Penrose Inverse.

historical note?

So we can write the solution to

$$\hat{\theta}_{LS} = (X^T X)^+ X^T g.$$

We would now like to see how well this procedure performs in terms of dimension.

First we need to decide how to measure performance here there is an important decision to be made

it really depends on the purpose of the analysis, i.e.

CASE 1: we are interested in **prediction**  
 that is we want to estimate the value  
 $X\theta^*$  as well as possible from the noisy  
 observations  $y = X\theta^* + \varepsilon$ . (also known as sample prediction). vs out-of-sample

CASE 2: we are interested in the actual value  $\theta^*$ .  
 In case (1) one may consider e.g. the mean squared error, i.e.

$$\text{MSE}(X\hat{\theta}^{\text{LS}}) = \frac{1}{n} \|X\theta^* - X\hat{\theta}^{\text{LS}}\|_2^2 \\ = (\theta^* - \hat{\theta}^{\text{LS}})^T \frac{X^T X}{n} (\theta^* - \hat{\theta}^{\text{LS}})$$

whereas  
 in the second case we may consider e.g.  
 $\|\theta^* - \hat{\theta}^{\text{LS}}\|_2^2$ .

Let's consider prediction error for now  
 & more specifically let's focus on the MSE.

Notice that since  $\hat{\theta}^{\text{LS}}$  solves  
 $\min_{\theta} \|y - X\theta\|_2^2$  s.t.  $\theta$  minimizes  $\|y - X\theta\|_2^2$

we have by definition that

$\hat{\theta}^L$  minimizes the  $\|\cdot\|_2^2$  norm among the  $\theta$ 's that minimize  $\|y - X\theta\|_2^2$ .

$$\text{So } \|\hat{y} - X\hat{\theta}^{Ls}\|_2^2 \leq \|y - X\theta^*\|_2^2$$

as since by definition  
 $y = X\theta^* + \varepsilon$

$$\|\hat{y} - X\hat{\theta}^{Ls}\|_2^2 \leq \|\varepsilon\|_2^2. \quad \textcircled{1}$$

$$\text{Also } \|\hat{y} - X\hat{\theta}^{Ls}\|_2^2$$

$$\begin{aligned} &= \|\hat{y} - X\hat{\theta}^{Ls} + X\theta^* - X\theta^*\|_2^2 \\ &= \|\hat{y} - X\theta^*\|_2^2 + \|X(\hat{\theta}^{Ls} - \theta^*)\|_2^2 - 2\langle \hat{y} - X\theta^*, X(\hat{\theta}^{Ls} - \theta^*) \rangle \\ &\quad \textcircled{2} \quad = \|\varepsilon\|_2^2 + \|X(\hat{\theta}^{Ls} - \theta^*)\|_2^2 - 2\langle \varepsilon, X(\hat{\theta}^{Ls} - \theta^*) \rangle \end{aligned}$$

Rearranging  $\textcircled{2}$  we get

$$\|X(\hat{\theta}^{Ls} - \theta^*)\|_2^2 \leq \|\hat{y} - X\hat{\theta}^{Ls}\|_2^2 - \|\varepsilon\|_2^2 + 2\langle \varepsilon, X(\hat{\theta}^{Ls} - \theta^*) \rangle$$

$$\textcircled{1} \quad \leq 2 \langle \varepsilon, X(\hat{\theta}^{Ls} - \theta^*) \rangle.$$

Now although we've made some progress

$\hat{\theta}^{Ls}$  depends on  $\varepsilon$ , & the RHS depends on  $\theta^*$ .

One way out is to consider a worst case scenario that is take supremum over everything but  $\varepsilon$

if

$$\|X(\hat{\theta}^{LS} - \theta^*)\|_2 \leq 2 \left\langle \varepsilon, \frac{X(\hat{\theta}^{LS} - \theta^*)}{\|X(\hat{\theta}^{LS} - \theta^*)\|_2} \right\rangle \cdot \|X(\hat{\theta}^{LS} - \theta^*)\|_2$$

$$\Leftarrow \|X(\hat{\theta}^{LS} - \theta)\|_2 \leq 2 \sup \left\{ \langle \varepsilon, v \rangle : v \in \mathbb{R}^n, \|v\|_2 \leq 1 \right\}$$

$$= 2 \|\varepsilon\|_2.$$

Now if  $\varepsilon = (\varepsilon_1, \dots, \varepsilon_n)$  is iid

$$\text{If } \|\varepsilon\|_2^2 = \sum \sigma_i^2 = n\sigma^2 \text{ where } \sigma^2 = \text{Var}(\varepsilon)$$

$$\text{so } \mathbb{E}[\text{MSE}(X\hat{\theta}^{LS})] \leq \frac{2n\sigma^2}{n} = 2\sigma^2.$$

So it seems that taking a large sample actually produces no visible benefit, the MSE  $\not\rightarrow 0$  as  $n \rightarrow \infty$ .

Q: Is this a real phenomenon? or an inefficiency of our argument?

→ If it is an inefficiency, what conditions can get better results.

Example: Gaussian Sequence:

Let  $y_i = \sqrt{n} \theta^*_i + \varepsilon_i, \quad \varepsilon_i \stackrel{\text{ind}}{\sim} N(0, \sigma^2)$

Then  $n=d$   $X = \sqrt{n} I_n$ .

$$X^T = \sqrt{n} I_n$$

$$X^T X = n I_n$$

$$\text{So } \hat{\theta}^{LS} = (X^T X)^{-1} X^T y =$$

$$= \frac{1}{n} \sqrt{n} y = \frac{1}{\sqrt{n}} y$$

$$\begin{aligned} \text{So } E[\text{MSE}(\hat{\theta}^{LS})] &= \frac{1}{n} E\left[\|X(\hat{\theta}^{LS} - \theta^*)\|_2^2\right] \\ &= \frac{1}{n} E\left[\|\sqrt{n} \cdot \left(\frac{y}{\sqrt{n}} - \theta^*\right)\|_2^2\right] \\ &= \frac{1}{n} E\left[\|y - \sqrt{n}\theta^*\|_2^2\right] = \frac{1}{n} E\left[\|\varepsilon\|_2^2\right] \\ &= \sum_{i=1}^n E \frac{\varepsilon_i^2}{n} \\ &= \sigma^2 \end{aligned}$$

So Real phenomenon

Question: Could it be that under additional assumptions we could do better?

Where could our argument be inefficient (under assumptions)

Let's see where our argument may be  
inefficient

When taking sup, we sup'ed over  
all of the  $\ell_2$ -unit ball in  $\mathbb{R}^n$ .

Q: when could this be wasteful?

looking back we replaced

$\frac{X(\hat{\theta}^L - \theta^*)}{\| \cdot \|}$  by the sup over  $v \in \mathbb{R}^n$   
 $\| v \| \leq 1$

But this assumes implicitly that

$$\text{Im}(X) = \mathbb{R}^n,$$

or that  $\text{Rank}(X) = n$ .

What if  $\text{rank}(X) = r < n$ ?

Could we hope for improved bounds?

We need to figure out a  
way to exploit this additional  
information.

①  $\text{Im}(X)$  is  $r$ -dimensional

$\Rightarrow$  let  $\{\mathbf{e}_i : i=1,\dots,n\}$  be the standard basis

and  $\{h_1, h_2, \dots, h_n\}$  be an orthonormal basis

such that  $h_1, \dots, h_r$  form an orthonormal basis of  $\text{Im}(X)$ .

(How to do this? Start with  $h_1, \dots, h_r$  & complete it with the Hilbert-Schmidt procedure).

So any elt of  $\text{Im}(X)$  expressed as a column vector in the  $H$  basis will take the form

$$\begin{pmatrix} h_1 \\ \vdots \\ h_r \\ \hline 0 \end{pmatrix}$$

n-r elems

Next express the  $H$  basis in the  $E$  basis of

$$h_j = \sum_{k=1}^n h_j^k e_k \quad \rightarrow \text{defining row vectors}$$
$$h_j = (h_j^i)_{i=1}^n \text{ & the matrix}$$

Suppose  $v = \sum_{j=1}^r \alpha_j h_j$  so  $[v]_H = \begin{pmatrix} \alpha_1 \\ \vdots \\ \alpha_r \\ 0 \end{pmatrix}$

Now we can write  $v = \sum_{j=1}^r \alpha_j \sum_{k=1}^n h_j^k e_k$

$$= \sum_{k=1}^n e_k \underbrace{\sum_{j=1}^n \alpha_j h_j}_k$$

So

$$\begin{pmatrix} \sum_{j=1}^n h_j^1 \alpha_j \\ \sum_{j=1}^n h_j^2 \alpha_j \\ \vdots \\ \sum_{j=1}^n h_j^n \alpha_j \end{pmatrix} = \begin{bmatrix} h_1^1 h_2^1 h_3^1 \dots h_n^1 \\ h_1^2 h_2^2 h_3^2 \dots h_n^2 \\ \vdots \\ h_1^n h_2^n \dots h_n^n \end{bmatrix} \begin{pmatrix} \alpha_1 \\ \vdots \\ \alpha_n \end{pmatrix}$$

$$= \begin{bmatrix} h_1 \ h_2 \ \dots \ h_n \end{bmatrix} \begin{pmatrix} \alpha_1 \\ \vdots \\ \alpha_n \end{pmatrix}$$

Now this means .

$$[\psi]_E = [h_1 \dots h_n] [\psi]_H$$

also  $\langle h_i, h_j \rangle = \delta_{ij}$  so

$$\begin{bmatrix} h_1 & \dots & h_n \end{bmatrix}^\top =$$

$$\begin{bmatrix} h_1 \\ \vdots \\ h_n \end{bmatrix} [\psi]_E = [\psi]_H$$

let  $O = \begin{bmatrix} h_1 \\ \vdots \\ h_n \end{bmatrix}$ , obviously  $O^T O = O O^T = I_n$ .

50 if  $v \in \text{Im}(x)$  then

$$Ov = \begin{pmatrix} \alpha_1 \\ \vdots \\ \alpha_r \\ 0 \end{pmatrix} \text{ and thus } P_r O v = O v$$


---

Now we have write  $\varphi$  for  $\hat{\theta}^L - \theta^*$   
since  $O$  orthogonal

$$\|x(\hat{\theta}^L - \theta^*)\|_2 \leq 2 \frac{\langle \varepsilon, x\varphi \rangle}{\|x\varphi\|} \leq 2 \frac{\langle O\varepsilon, OX\varphi \rangle}{\|X\varphi\|}$$

$$= 2 \frac{\langle P_r O\varepsilon, P_r OX\varphi \rangle}{\|X\varphi\|} = \textcircled{*}$$

where  $P_r = \begin{pmatrix} I_r & | & 0 \\ 0 & | & 0 \end{pmatrix}$

Since  $OX\varphi$  will have in last  $(n-r)$  entries 0.

$$\textcircled{*} \leq 2 \|P_r O\varepsilon\|_2$$

Now notice first that if

$\varepsilon \sim N(0, \sigma^2 I_n)$  then.

$$\textcircled{1} \quad \tilde{\varepsilon} = O\varepsilon \sim N(0, \sigma^2 I_n)$$

$$\tilde{\varepsilon} := P_r O\varepsilon \sim N(0, \sigma^2 I_r) \quad \&$$

thus finally  $\mathbb{E} \|P_r \Omega \varepsilon\|_2^2 = \mathbb{E} \sum_{i=1}^r \tilde{\varepsilon}_i^2 = r \sigma^2$ .

In fact the same is true more generally

if  $\varepsilon$  is a sub-Gaussian vector.

why? here we need to be

a little more careful, let  $\varphi = \hat{\theta}^{LS} - \theta^*$

we have

$$\|X\varphi\| \leq 2 \left\langle \varepsilon, \frac{X\varphi}{\|X\varphi\|} \right\rangle$$

let as before  $h_1, \dots, h_r$  be an orthonormal basis

of  $\text{Im}(x)$ , & suppose as before

$h_j = \sum_{k=1}^n h_j^k e_k$ , so the means the vectors

$q_1 = \begin{pmatrix} h_1 \\ h_1 \\ \vdots \\ h_1 \end{pmatrix}, \dots, q_r = \begin{pmatrix} h_r \\ h_r \\ \vdots \\ h_r \end{pmatrix}$  are an o.n. basis for  $\text{Im}(x)$ .

let  $\Phi = [q_1, \dots, q_r] \in \mathbb{R}^{n \times r}$ .

In particular if  $v \in \text{Im}(x)$  then

$$v = \sum_{j=1}^r \alpha_j h_j = \sum_{j=1}^r \alpha_j \sum_{k=1}^n h_j^k e_k =$$

$$= \sum_{k=1}^n e_k \left( \sum_{j=1}^n \alpha_j h_j^{(k)} \right)$$

(So)  $[v]_e = \Phi [v]_h$

Thus for any  $v \in \text{Im}(x)$ ,  $\exists v \in \mathbb{R}^r$  such

$$v = \Phi v_h$$

Also the o.n. property of  $\{h_{ij}\}_{j=1}^n$

$$\Rightarrow (\Phi^\top \Phi)_{ij} = \langle q_i, q_j \rangle = \langle h_{1i}, h_{1j} \rangle = \delta_{ij}$$

(But  $\Phi \Phi^\top$  makes no sense!)

Thus:  $\frac{\langle \varepsilon, X(\theta^* - \theta^*) \rangle}{\|X(\theta^* - \theta^*)\|} = \frac{\langle \varepsilon, \Phi v \rangle}{\|\Phi v\|_2} = \frac{\langle \Phi^\top \varepsilon, v \rangle}{\|v\|_2}$

Since  $\Phi^\top \Phi = I_r$

$$= \frac{\langle \tilde{\varepsilon}, v \rangle}{\|v\|_2}, \text{ where } \tilde{\varepsilon} = \Phi^\top \varepsilon \in \mathbb{R}^r.$$

$$\leq \sup_{u \in B_2^r} \langle \tilde{\varepsilon}, u \rangle.$$

Suppose now that  $\varepsilon \sim \text{sub-G}(\mathbb{R}^n)$ .

Let  $u \in \mathbb{S}^{r-1}$ , then

$$\mathbb{E}[\exp(\gamma \langle u, \tilde{\varepsilon} \rangle)] = \mathbb{E}[\exp(\gamma \langle u, \Phi^\top \varepsilon \rangle)]$$

$$= \mathbb{E} \left\{ \exp(\lambda \langle \Phi u, \varepsilon \rangle) \right\}$$

Notice that if  $u \in \mathbb{S}^{r-1}$ , then  
 $\|\Phi u\|_2^2 = \langle \Phi u, \Phi u \rangle = \langle u, \Phi^\top \Phi u \rangle = \|u\|_2^2$ .

So  $\Phi u$  is a unit vector in  $\mathbb{R}^n$

$$\Rightarrow \mathbb{E} \left[ \exp(\lambda \langle \Phi u, \varepsilon \rangle) \right] = e^{\frac{\lambda^2 \sigma^2}{2}}$$

So  $\tilde{\varepsilon}$  is  $\sigma^2$ -sub-Gaussian ( $\mathbb{R}^r$ ).

We are now interested in

$$\sup_{u \in B_2^r} \langle \tilde{\varepsilon}, u \rangle, \quad \tilde{\varepsilon} \text{ is } \sigma^2\text{-sub-Gaussian}(\mathbb{R}^r).$$

We will make use of the following theorem.

Theorem: let  $X \in \mathbb{R}^d$   $\sigma^2$ -sub-Gaussian vector.

Then for any  $\delta > 0$ , w.p. at least  $1 - \delta$

$$\max_{\theta \in B_2^d} \theta^\top X = \max_{\theta \in B_2^d} \|\theta^\top X\| \leq 4\sigma\sqrt{d} + 2\sigma\sqrt{2\log(1/\delta)}.$$

Proof: the idea is the following.

We will cover  $B_2^d$  with balls of radius  $1/2$ , centred

out the elements of a finite set  $\mathcal{N}$ .

We will control the  $\sup_{B_2^d}$  with the  $\overbrace{\sup \text{ over } \mathcal{N}}$   
 + a correction term  
 $\underbrace{\text{Simple union bound.}}$

Let  $\mathcal{N}$  be a  $1/2$ -net of  $B_2^d$ ,

i.e.  $\forall z \in B_2^d \exists y \in \mathcal{N}$  s.t.  $\|y - z\| < 1/2$ .

Then:

$$\sup_{\theta \in B_2^d} \langle \theta, x \rangle = \sup_{\theta \in \mathcal{N}} \left\{ \langle \theta, x \rangle \mid \theta = z + u, z \in \mathcal{N}, \|u\| \leq 1/2 \right\}$$

$$\stackrel{*}{\leq} \sup_{\theta \in \mathcal{N}} \langle \theta, x \rangle + \sup_{\theta \in \frac{1}{2} B_2^d} \langle \theta, x \rangle$$

Now

$$\theta \in \frac{1}{2} B_2^d \Leftrightarrow \|\theta\|_2 \leq 1/2.$$

$$\text{so } \langle \theta, x \rangle = \frac{1}{2} \langle 2\theta, x \rangle \text{ when } 2\theta \in B_2^d$$

$$\text{so } \sup_{\theta \in \frac{1}{2} B_2^d} \langle \theta, x \rangle = \frac{1}{2} \sup_{\theta \in B_2^d} \langle \theta, x \rangle \text{ so}$$

$$\textcircled{1} \quad \frac{1}{2} \sup_{\theta \in B_2^d} \langle \theta, x \rangle \leq \sup_{\theta \in \mathcal{N}} \langle \theta, x \rangle + \frac{1}{2} \sup_{\theta \in B_2^d} \langle \theta, x \rangle$$

$$\text{so } \sup_{\theta \in B_2^d} \langle \theta, x \rangle \leq 2 \sup_{\theta \in \mathcal{N}} \langle \theta, x \rangle.$$

$$\text{Thus } \mathbb{P} \left[ \max_{\theta \in B_2^d} \langle \theta, x \rangle > t \right] \leq \mathbb{P} \left[ 2 \sup_{\theta \in \mathcal{N}} \langle \theta, x \rangle > t \right]$$

$$= \mathbb{P} \left[ \max_{\theta \in N} \langle \theta, X \rangle > t \right]$$

this is now a union over  $|N|$  r.v's  
so a simple union bound gives Lemma 1.3.2 notes

$$\leq |N| e^{-t^2/8\sigma^2}. \quad \textcircled{X}$$

how big is  $|N|$ ?

Lemma: let  $\varepsilon \in (0,1)$  &  $N$  be an  $\varepsilon$ -net of  $B_2^d$ .

Theorem: Start with  $N = \{0\}$ ,  $X = B_2^d \setminus \left( \bigcup_{z \in N} B_\varepsilon(z) \right)$ .  
 While  $X \neq \emptyset$  choose point in  $X$  & add it to  $N$ .

This will eventually terminate if  $\forall u, y \in N$   
 $\|u-y\| > \varepsilon$ . So if we replace  $\varepsilon$  with  $\varepsilon/2$

then the  $\varepsilon$  balls centred at  $N$  will be disjoint.

$$\text{So } \bigcup_{z \in N} \left\{ z + \frac{\varepsilon}{2} B_2^d \right\} \subset \left( 1 + \frac{\varepsilon}{2} \right) B_2^d$$

& compute their volumes

$$\text{vol} \left( \left( 1 + \frac{\varepsilon}{2} \right) B_2^d \right) \geq \text{vol} \left( \bigcup_{z \in N} \left\{ z + \frac{\varepsilon}{2} B_2^d \right\} \right) = \sum_{z \in N} \text{vol} \left( z + \frac{\varepsilon}{2} B_2^d \right)$$

$$\text{so } \left( 1 + \frac{\varepsilon}{2} \right)^d \geq |N| \left( \frac{\varepsilon}{2} \right)^d$$

$$\Rightarrow |\mathcal{N}| \leq \left( \frac{1 + \frac{\varepsilon}{2}}{\varepsilon/2} \right)^d = \left( \frac{2 + \varepsilon}{\varepsilon} \right)^d \leq \left( \frac{3}{\varepsilon} \right)^d.$$

letting  $\varepsilon = 1/2$  we get that  $|\mathcal{N}| = 6^d$

$$\text{so } \mathbb{P} \left[ \max_{\theta \in \mathcal{B}_2} \langle \theta, X \rangle > t \right] \leq 6^d e^{-\frac{t^2}{8\sigma^2}} \leq 5$$

$$\Rightarrow \exp \left( -\frac{t^2}{8\sigma^2} + d \log 6 \right) \leq 5$$

$$-\frac{t^2}{8\sigma^2} + d \log 6 \leq \log(5)$$

$$\frac{t^2}{8\sigma^2} \geq \log(1/5) + d \log 6.$$

$$t^2 \geq 8\sigma^2 \log(6)d + 8\sigma^2 \log(1/5).$$

$$\text{so take } t = \sqrt{8 \log(6)} \sqrt{d} + 2\sqrt{2 \log(1/5)}.$$

Since  $MSE(\hat{X}_{\theta^{**}}) = \frac{1}{n} \|X_{\theta^{**}} - X_{\theta^*}\|^2$

we conclude that

wp at least  $1-\delta$

$$MSE(\hat{X}_{\theta^{**}}) \leq \frac{\sigma^2 r + \log(1/\delta)}{n}.$$

rank(X)

So we did a lot of work to go from

$$\frac{d}{n} \text{ to } \frac{r}{n}.$$

in fact if  $\text{rank}(X) = d$  then the bounds are the same.

Infact Gaussian sequence model

$$y_i = \sqrt{n} \theta^*_i + \varepsilon_i \quad i=1, \dots, n$$

$$n=d \quad X = \sqrt{n} I_n, \quad \hat{\theta}^{LS} = n^{-1} X^T y = (X^T X)^{-1} X^T y = (n I_n)^{-1} X^T y = \frac{1}{n} \sqrt{n} y = \frac{y}{\sqrt{n}}$$

$$\text{So } \|X(\hat{\theta}^{LS} - \theta^*)\|_2^2 = \sum_{i=1}^n \frac{\varepsilon_i^2}{n} = SIC(1).$$

So it seems

if we are to make any progress we must exploit special structure, e.g. sparsity etc

$$\xrightarrow{\hspace{1cm}} \xrightarrow{\hspace{1cm}} \xrightarrow{\hspace{1cm}} \xrightarrow{\hspace{1cm}}$$

Noiseless Recovery: we will now start considering the effect of sparsity.

To build some intuition about how this can possibly help us we will first consider exact recovery in

the noiseless problem

$$y = X\theta^*.$$

Suppose we know :  $\theta^* \in K \left( \begin{array}{l} k-\text{space} \\ \|\theta^*\|_1 \leq m \\ \text{etc} \end{array} \right)$

We can then consider

$$\hat{\theta} = \underset{\theta \in K}{\operatorname{argmin}} \|y - X\theta\|_2^2.$$

Suppose further that  $\boxed{\theta^* \in K = B_1(1) = \{x \in \mathbb{R}^d \mid \|x\|_1 \leq 1\}}$

Continuing from the previous calculation

$$\begin{aligned} \|X(\theta^* - \hat{\theta})\|_2^2 &\leq 2 \langle \varepsilon, X(\hat{\theta} - \theta^*) \rangle \\ &\leq 2 \sup_{v, w \in K} \langle \varepsilon, X(v - w) \rangle \\ &= 2 \sup_{v' \in B_1(2)} \langle \varepsilon, Xv' \rangle \\ &\leq 4 \sup_{v' \in B_1(1)} \langle \varepsilon, Xv' \rangle. \end{aligned}$$

Now  $\langle \varepsilon, Xv \rangle$  is a linear form in  $v$  &  
the set  $B_1(1)$  is convex, so it can

be written as the convex hull of its extreme points

$$\text{so } \sup_{v \in B_1(1)} \langle \varepsilon, X_v \rangle$$

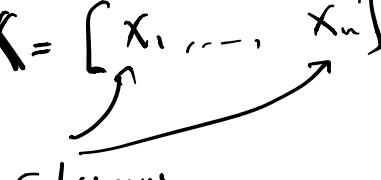
$v \in B_1(1)$

$$= \sup \left\{ \langle \varepsilon, X_v \rangle \mid v \text{ is an extreme point of } B_1(1) \right\}$$

thinking about  $B_1(1)$  we get

$$= \max \left\{ \langle \varepsilon, X_v \rangle : v = \pm e_j, j = 1, \dots, d \right\}$$

if  $X = [x_1, \dots, x_n]$  then



$$X e_j = X \begin{bmatrix} 0 \\ \vdots \\ 1 \\ \vdots \\ 0 \end{bmatrix} = [x_j]$$

thus  $\langle \varepsilon, X e_j \rangle$  is  $\|x_j\|_2 \sigma^2$ -sub Gaussian

so  $\sup_{v \in B_1(1)} \langle \varepsilon, X_v \rangle$  is the max of  $\|x_j\|_2$  for  $j = 1, \dots, d$

$$\left( \max_{j=1, \dots, d} \|x_j\|_2 \right) \sigma^2 \text{-sub Gaussian r.v.}$$

thus

$$\mathbb{E} \left[ \text{MSE}(X \hat{\theta}_k^{LS}) \right] \leq \frac{\|X(\hat{\theta}_k^{LS} - \theta^*)\|_2^2}{n}$$

$$\leq \mathbb{E} \left[ \frac{\max_{j=1 \dots d} \langle \varepsilon, \pm X - e_j \rangle}{n} \right]$$

so

$$\leq \frac{C \sigma \max_j \|X_j\|_2}{n} \sqrt{2 \log(2d)}$$

& if  $t \rightarrow 0$

$$\mathbb{P}[MSE > t] \leq \mathbb{P} \left[ \max_{\substack{0 = \varepsilon + e_j \\ j=1 \dots d}} \varepsilon^T X \geq nt/\sigma \right]$$

$$\leq 2d \exp \left( \frac{-n^2 t^2}{16 \sigma^2 \max_j \|X_j\|_2} \right)$$

If we assume that  $X$  is normalized so that

$$\max_j \|X_j\|_2 \leq \sqrt{n}$$

$$dth \leq 2d \exp \left( \frac{-n t^2}{16 \sigma^2} \right)$$

& concentration occurs at rate  $n$ .

Compare: with previous result in the case  
where  $\text{rank}(X) = n \ll d, m$

So sparsity can get good results even if  $\text{rank}(X) \ll d, m$ .

Similarly if  $k = \text{Bo}(k)$

$\nearrow$   
k-sparse vectors

But to solve over  $B_0(k)$  we need to solve for

$B_0(1), \dots, B_0(k)$  & for each of these you need  
to consider  $\binom{d}{j}$  subsets as possible supports.

So comp. it becomes quickly infeasible.

---

Rather than asking a priori that  $\theta^* \in K = B_p(1)$  s.t.  
we instead penalise solutions with large  $\|\cdot\|_p$  norm.

$$\text{e.g. } \hat{\theta} \in \arg \min \left\{ \frac{1}{2n} \|y - X\theta\|_2^2 + \lambda \|\theta\|_p \right\}$$

user-specified regularization parameter

Whence does this allow us  
to get closer to  $\theta^*$ ?

To understand the effect the regularization has  
w/o the influence of the noise  $\sim \mathcal{N}(0, \sigma^2 I)$  consider noiseless model

e.g.  $Y = X\theta^*$ ,  $\theta^*$  is sparse but we don't know  
a priori how sparse

---

Try to solve

$\min \|x\|_0 \quad \text{s.t. } X\theta = y$

$\theta \in \mathbb{R}^d$

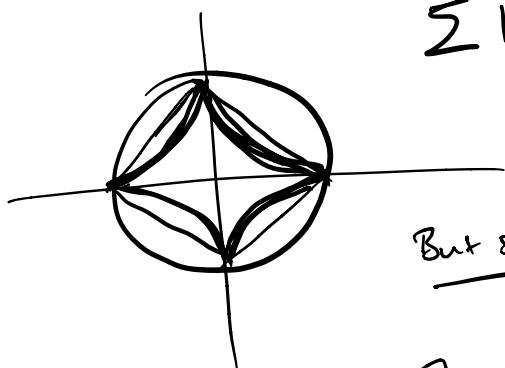
Problem non-convex so again we'd have  
to search  $B_0(u)$  for all  $u$ .

Idea: replace  $\|x\|_0$  by  $\|x\|_1$ .

We will now try to develop some intuition about geometry.

$$\sum \|e_i\| \quad \text{vs} \quad \sum \|e_i\|^2$$

$$\|e\|_2 \leq \|e\|_1 \leq \sqrt{n} \|e\|_2$$



But suppose  $e_i$  for  $i \in I$   
are  $< \delta$  (small)

$$\text{The penalty in } L_2 = \sum_{i \in I} e_i^2$$

which can be much less than the penalty  
in  $L_1$ .

Suppose then

$$y = X\theta^*, \quad \text{and} \quad \theta_j \neq 0 \quad j \in S \subseteq [1:d]$$

Let  $\theta_j = 0$ ,  $j \in S^c \subseteq [1:d]$  support vector  $\theta$ .

Q: When does the solution to

$$\min_{\theta \in \mathbb{R}^d} \|\theta\|_1 \text{ s.t. } X\theta = y \quad (2)$$

give us the solution to

$$\min_{\theta \in \mathbb{R}^d} \|\theta\|_2 \quad X\theta = y \quad (1)$$

First let us think about the space of solutions to  $X\theta = y$ .

$$S(X, y) := \{\theta \in \mathbb{R}^d \mid X\theta = y\}$$

We know  $\theta^*$  is a solution so

$$S(X, y) = \theta^* + \ker(X)$$

where

$$\ker(X) = \{\theta \in \mathbb{R}^d \mid X\theta = 0\}.$$

Now (2)  $\Leftrightarrow \min_{\theta \in S(X, y)} \|\theta\|_1$  so

since  $\theta^* \in S(X, y)$  the solution to (2)

will be  $\theta^*$  iff

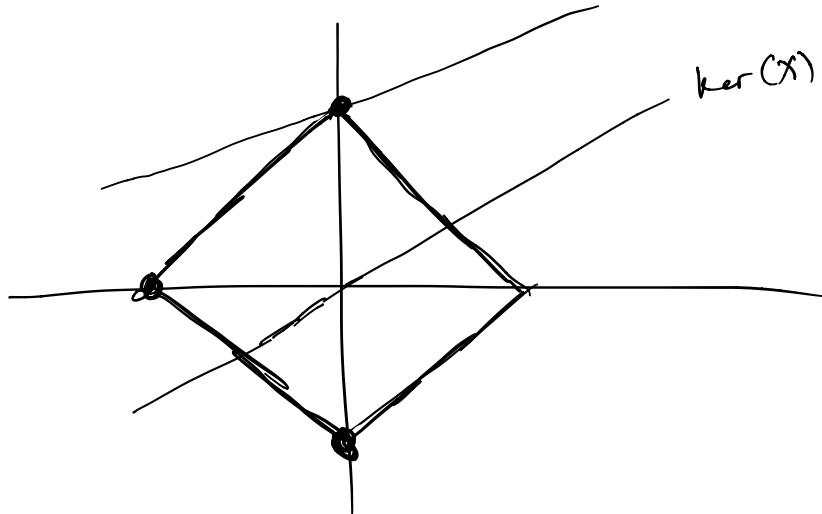
$\theta^*$  has minimal  $\|\cdot\|_1$  norm in  $S(X, y)$

But if  $x \in \ker(X)$  then  $\theta^* + x \in S(X_{\text{cy}})$  also  
 so  $\text{sol}^n$  to (2) =  $\theta^*$   
 iff  $\|\theta^* + \text{null}\|_1 \geq \|\theta^*\|_1$  if  $x \in \ker(X)$

To visualise the situation consider

$$B_1(\|\theta^*\|_1)$$

Note: on  $\theta^{\perp}$  space it will be an extreme point of  
 $B_1(\|\theta^*\|_1)$ , i.e. a corner.



So (2) will recover the sol<sup>n</sup> to (1)

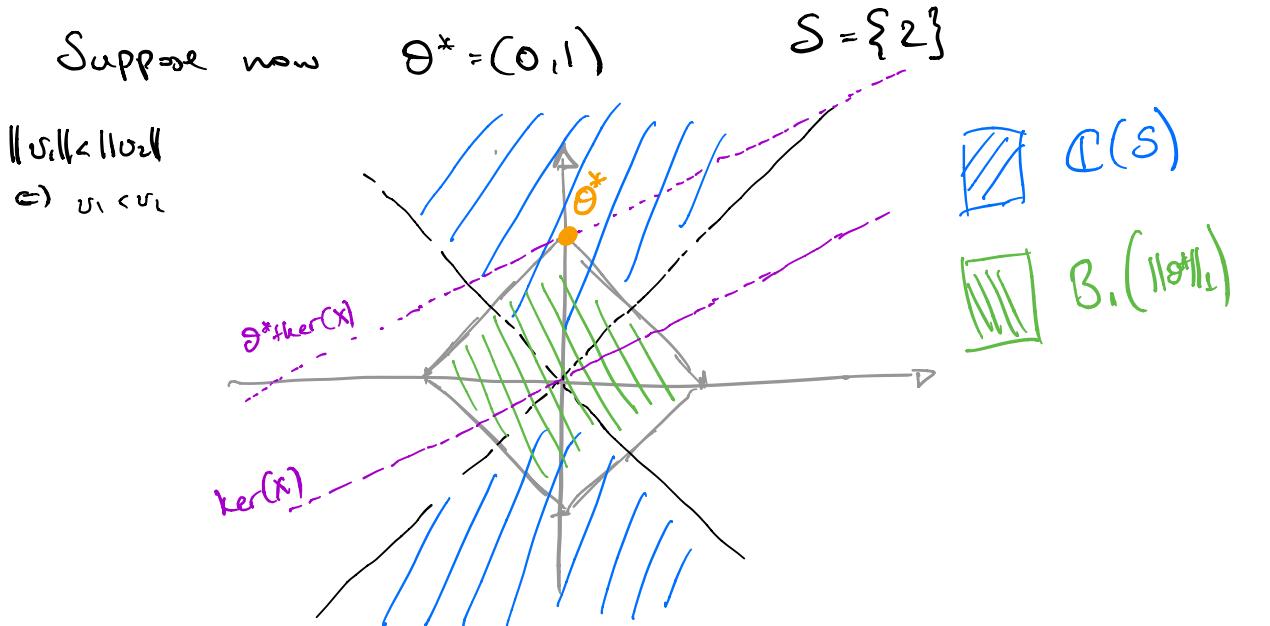
if  $B_1(\|\theta^*\|_1)$  intersects  $\theta^{\perp} \ker(X)$  only at  $\theta^*$ .

Definition set (for  $S \subset [1:d]$ ) CONE

$$\mathbb{C}(S) := \left\{ v \in \mathbb{R}^d \mid \|v_S\|_1 \leq \|v_S\|_1 \right\}$$

where  $(v_S)_i = v_i$  if  $i \in S$  & 0 o.w.

i.e.  $(v_S)_i = v_i \mathbb{1}_{\{i \in S\}}$ .



So  $(\theta^* + \ker X) \cap B_1(\|S\|_1) = \{\theta^*\}$

iff  $\ker(X) \cap C(S) = \{0\}$

Definition (Restricted Nullspace Property)

We say  $X$  satisfies the (RNP) with respect to  $S \subseteq [1:d]$

if  $C(S) \cap \ker(X) = \{0\}$

②  $\min_{\theta} \{ \|\theta\|_1 \mid X\theta = y \}$

Theorem: TFAE

(a)  $\forall \theta^*$  supported on  $S \subseteq [1:d]$ , (2) applied with  $y = X\theta^*$  having solution  $\hat{\theta} = \theta^*$ .

(b)  $X$  satisfies RNP wrt  $S$ .

Pf: (b)  $\Rightarrow$  (a)

We know  $y = X\theta^*$  so

Suffices to show  $\theta'$  solution  $\Rightarrow \|\theta'\|_1 \geq \|\theta^*\|_1$ .

Let  $\theta' = \theta^* + v$ ,  $v \in \ker(x) \neq \emptyset$

$$(b) \Rightarrow \|v_{S^c}\|_1 > \|v_S\|_1$$

$\& \text{supp}(\theta^*) = S$ , so

$$\|\theta'\|_1 = \|\theta^* + v\|_1 = \|\theta^*_{S^c} + v_{S^c}\|_1 + \|\theta^*_S + v_S\|_1.$$

$$= \|v_{S^c}\|_1 + \|\theta^*_S + v_S\|_1 \geq \|\theta^*_S\|_1 - \|v_S\|_1 + \|v_{S^c}\|_1$$

$$\stackrel{\Delta}{\begin{aligned} &\text{ineq} \\ &= \|\theta^*_S\|_1 + \underbrace{(\|v_{S^c}\|_1 - \|v_S\|_1)}_{\geq 0 \text{ by (b)}} \\ &\geq \|\theta^*_S\|_1. \end{aligned}}$$

(a)  $\Rightarrow$  (b) Let  $\theta^* \in \ker(x)$ ,  $\theta^* \neq 0$

Then  $X\theta^* = 0$  & thus. remove coords if needed

$$X \begin{bmatrix} \theta_S^* \\ 0 \end{bmatrix} + X \begin{bmatrix} 0 \\ \theta_{S^c}^* \end{bmatrix} = 0$$

$$\Rightarrow X \begin{bmatrix} \theta_S^*, 0 \end{bmatrix}^\top = X \begin{bmatrix} 0, -\theta_{S^c}^* \end{bmatrix}^\top \Leftrightarrow \tilde{y}$$

$$\Rightarrow \begin{bmatrix} 0 \\ -\theta_{S^c}^* \end{bmatrix} \text{ is a soln to } X\theta = \left\{ X \begin{bmatrix} \theta_S^* \\ 0 \end{bmatrix} \right\}$$

i.e. for problem

$$X\theta = \tilde{g} \quad (= X[\theta_s^*])$$

has obviously the solution  $\begin{bmatrix} \theta_s^* \\ 0 \end{bmatrix}$

but also the solution  $\begin{bmatrix} 0 \\ -\theta_{s^c}^* \end{bmatrix}$ .

by (a) since (2) has a unique soln, so it must be that  
 $\theta_s^*$  uniquely

it must be that

$$\|\theta_{s^c}^*\| > \|\theta_s^*\|_1 \Rightarrow \ker(X) \cap C(S) = \{0\}.$$

---

### Verifiable Condition for RNP

If  $v \in \ker(X) \setminus \{0\}$ , then

$$Xv=0 \Rightarrow v_1 X_1 + \dots + v_d X_d = 0$$

$\nearrow$  columns  $\searrow$   
 $\theta_s^* X$

$\Rightarrow \{x_1, \dots, x_d\}$  not linearly independent!

So one obvious way to ensure  $X$  satisfies RNP( $S$ ) for any  $S \subseteq [1:d]$

is to assume that  $\ker(X) = \{0\}$

$\Rightarrow \{x_1, \dots, x_d\}$  lin. independent.

This is already quite restrictive if  $d \leq n$   
and impossible for  $d > n$ .

Note: that it indeed the columns are lin. ind.

We can transform the model so make them  
 orthonormal i.e.  $\frac{1}{n} \langle \tilde{X}_i, \tilde{X}_j \rangle = \delta_{ij}$

Or

$$\frac{\tilde{X}^T \tilde{X}}{n} = \mathbb{I}_d \quad (\text{ORT})$$

Idea: if we allow ORT to fail in a "controlled"  
 quantifiable way, how we make any progress?

### DEF'n (PAIRWISE INCOHERENCE PARAMETER)

For  $n \times d$  matrix  $X$  define

$$\delta_{\text{PW}}(X) := \max_{j,k=1,\dots,d} \left| \frac{1}{n} \langle X_i, X_j \rangle - \delta_{jk} \right|$$

This measures precisely the degree to which  
ORT fails.

let  $S \subseteq [1:d]$   $|S|=s$  & suppose

Define  $X_S = (X_{i,j})_{i,j \in S}$ .

Suppose that

$\lambda$  is an eigenvalue of  $\frac{X_S^T X_S}{n}$

Then for some  $w \in \mathbb{R}^S$

$$\left( \frac{1}{n} X_S^T X_S - \mathbb{1}_S \mathbb{1}_S^T \right) w = (\lambda - 1) w$$

$$|\lambda - 1| \|w\|_2 \leq \underbrace{\left\| \frac{1}{n} X_S^T X_S - \mathbb{1}_S \mathbb{1}_S^T \right\|_2}_{\text{operator norm}} \|w\|_2$$

$$\leq \left\| \frac{1}{n} X_S^T X_S - \mathbb{1}_S \mathbb{1}_S^T \right\|_F \|w\|_2$$

\* \*

$\|A\|_2 \leq \|A\|_F$   
(otherwise)  
 Frobenius norm

$$\|A\|_F^2 = \sum_{i,j} |A_{i,j}|^2$$

$$\textcircled{*} |\lambda - 1|^2 \leq \sum_{j,k \in S} \left( \frac{1}{n} \langle X_j, X_k \rangle - \delta_{jk} \right)^2$$

$$\leq |S|^2 \sigma_{pw}(x)^2 \leq |S|^2 \frac{\delta^2}{|S|^2}$$

Now if  $\gamma \in (0, 1) \Rightarrow 1-\gamma \leq \gamma$  Alternatively

$$\Rightarrow 2 \geq 1-\gamma > 0.$$

Thus for any  $\theta \in \mathbb{R}^d$

$$\theta_s^\top \frac{X_s^\top X_s}{n} \theta_s \geq (1-\gamma) \|\theta_s\|_2^2$$

Since  $X_s \theta_s = X \theta_s$

$$\Leftrightarrow \|\theta_s\|_2^2 \leq \frac{1}{1-\gamma} \theta_s^\top \frac{X^\top X}{n} \theta_s.$$

Let now  $\theta \in \ker(X)$  ( $\theta = \theta_s + \theta_{s^c}$ )

so  $X \theta_s = -X \theta_{s^c}$

$$\|\theta_s\|_2^2 \leq \frac{1}{1-\gamma} \theta_s^\top \frac{X^\top X}{n} \theta_s$$

$$= \frac{1}{1-\gamma} \theta_s^\top \frac{X^\top X}{n} \theta_{s^c}$$

$$= \frac{1}{1-\gamma} \theta_s^\top \left( I - \frac{X^\top X}{n} \right) \theta_{s^c}$$

Since  $\theta_s^\top \theta_{s^c} = 0$

$$\leq \frac{1}{1-\gamma} \|I - \frac{X^\top X}{n}\|_\infty \|\theta_{s^c}\|_1 \|\theta_s\|_1. \quad \textcircled{*}$$

$$\leq \frac{1}{1-\gamma} \delta_{PW}(x) \|\theta_S\|_1 \|\theta_{S^c}\|_1.$$

Using  $\ell_1 - \ell_2$  inequality we get

$$\|\theta_S\|_1^2 \leq |S| \|\theta_S\|_2^2$$

$$\textcircled{*} + \quad \Rightarrow \quad \|\theta_S\|_1^2 \leq S \|\theta_S\|_2^2$$

$$\leq \frac{S}{1-\gamma} \delta_{PW}(x) \|\theta_S\|_1 \|\theta_{S^c}\|_1$$

$$\Rightarrow \|\theta_S\|_1 \leq \frac{\gamma}{1-\gamma} \frac{\gamma}{\gamma} \|\theta_{S^c}\|_1$$

So if  $\frac{\gamma}{1-\gamma} < 1$  we have for RNP wrt  $S$  with  $|S|=s$ .

$$\Leftrightarrow \gamma < 1-\gamma \Leftrightarrow$$

$$\gamma < 1-\gamma$$

$$2\gamma < 1 \quad \gamma < \frac{1}{2} \Leftrightarrow \|\theta_S\|_1 < \|\theta_{S^c}\|_1.$$

for any  $\theta \in \ker(x) \setminus \{\theta\}$ .

This leads to

Prop: If for some  $s \in \mathbb{N}$

$$\delta_{PW}(x) \leq \frac{1}{2s} \quad \text{then } \text{RNP}$$

RNP( $S$ ) holds  $\forall S \subseteq [1:d]$  with  $|S| \leq s$ .

## The LASSO :

We now have some idea about how  $X \theta^*$  may be allowed to intersect if we may hope that (2) will actually recover the solution in the noiseless setting.

$$\hookrightarrow \text{Bachto} \quad Y = X\theta^* + \varepsilon$$

$$\text{Lasso estimator} \quad \hat{\theta}_L \in \underset{\theta}{\operatorname{argmin}} \left\{ \frac{1}{2n} \|y - X\theta\|_2^2 + \lambda \|\theta\|_1 \right\}$$

Tibshirani 1996

Structure : first slow rates  $\leftarrow$  minimal assumptions

then we will see an example where we can obviously do better than these rates

$\longrightarrow$  FAST RATE

& then we will consider some assumptions that can guarantee us the fast rates.

Attempt 1: first we will assume that  $X$  is normalized so that

$$\max_j \|X_j\|_2^2 \leq n.$$

$$\|Y - X\hat{\theta}^L\|_2^2 = \|X\theta^* + \varepsilon - X\hat{\theta}^L\|_2^2 =$$

$$= \|X(\theta^* - \hat{\theta}^L)\|_2^2 + \boxed{2 \langle \varepsilon, X(\theta^* - \hat{\theta}^L) \rangle} + \|\varepsilon\|_2^2$$

therefore  
removing

$$\|X(\theta^* - \hat{\theta}^L)\|_2^2 = \|Y - X\hat{\theta}^L\|_2^2 - 2 \langle \varepsilon, X(\theta^* - \hat{\theta}^L) \rangle - \|\varepsilon\|_2^2$$

(1)

From  $\hat{\theta}^L \in \arg\min \left\{ \frac{1}{n} \|Y - X\theta\|_2^2 + C\|\theta\|_1 \right\}$

we have  $\frac{1}{n} \|Y - X\theta\|_2^2 + 2C\|\hat{\theta}^L\|_1 \leq \frac{1}{n} \|Y - X\theta^*\|_2^2 + 2C\|\theta^*\|_1$

(2)

so (1) & (2) combined give

$$\begin{aligned} \|X(\theta^* - \hat{\theta}^L)\|_2^2 &\leq \boxed{2 \langle \varepsilon, X(\hat{\theta}^L - \theta^*) \rangle} + \boxed{2nC\|\theta^*\|_1} \\ &\quad - \boxed{2nC\|\hat{\theta}^L\|_1} + \boxed{\|\varepsilon\|_2^2} - \boxed{\|\varepsilon\|_2^2} \\ &= 2 \langle X^\top \varepsilon, (\hat{\theta}^L - \theta^*) \rangle + 2nC(\|\theta^*\|_1 - \|\hat{\theta}^L\|_1) \\ &= 2 \langle X^\top \varepsilon, \hat{\theta}^L \rangle - 2nC\|\hat{\theta}^L\|_1 + 2 \left( \langle X^\top \varepsilon, \theta^* \rangle + nC\|\theta^*\|_1 \right) \\ &\leq 2 \left( \|X^\top \varepsilon\|_\infty - nC \right) \|\hat{\theta}^L\|_1 + 2 \left( \|X^\top \varepsilon\|_\infty + nC \right) \|\theta^*\|_1. \end{aligned}$$

Again union bound gives us

$$\begin{aligned} \mathbb{P} \left[ \|X^\top \varepsilon\|_\infty \geq t \right] &= \mathbb{P} \left[ \max_{j=1, \dots, d} |\langle \varepsilon, X_j \rangle| \geq t \right] \\ &\leq \sum_{j=1}^d \mathbb{P} \left[ |\langle \varepsilon, X_j \rangle| \geq t \right] \leq \sum_{j=1}^d 2 \exp \left\{ - \frac{t^2}{2\sigma^2 \|X_j\|_2^2} \right\} \\ &\leq \sum_{j=1}^d 2 \exp \left\{ - \frac{t^2}{2\sigma^2 \max_j \|X_j\|_2^2} \right\} \leq 2d e^{-\frac{t^2}{2n\sigma^2}} \\ &\text{since } \boxed{\max_j \|X_j\|_2^2 \leq n}. \end{aligned}$$

Setting  $\mathbb{P} [ \|X^T \varepsilon\|_\infty \geq t ] \leq \delta$

which can be guaranteed with

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

Choose  $\varepsilon = t/n$  so that:

$$\|X(\hat{\theta}^L - \theta^*)\|_2^2 \leq 2 \left( \|X^T \varepsilon\|_\infty - n\varepsilon \right) \|\hat{\theta}^L\|_1 + 2 (\|X^T \varepsilon\|_\infty + n\varepsilon) \|\theta^*\|_1$$

if  $t = n\varepsilon$

$$\text{then } \|X^T \varepsilon\|_\infty < t \Rightarrow$$

$$\begin{aligned} \|X(\hat{\theta}^L - \theta^*)\|_2^2 &< 2 (\|X^T \varepsilon\|_\infty + n\varepsilon) \|\theta^*\|_1 \\ &\leq 4n\varepsilon \|\theta^*\|_1 \end{aligned}$$

so

Since  $\|X^T \varepsilon\|_2 < t$   $\omega_P > t + \delta$  we're done.  $\square$

Theorem: Suppose  $y = X\theta^* + \varepsilon$  & let  $\hat{\theta}^L$  solve

$$\arg \min \left\{ \frac{1}{2n} \|y - X\theta\|_2^2 + \varepsilon \|\theta\|_1 \right\}$$

with  $\varepsilon = \sigma \sqrt{\frac{2}{n}} \left( \sqrt{\log(2d)} + \sqrt{\log(1/\delta)} \right)$ . Then

$$\omega_P > 1 - \delta$$

$$\text{MSE}(X\hat{\theta}^L) \leq \frac{4\|\theta^*\|_1}{\sqrt{n}} \left( \sqrt{2 \log(d)} + \sqrt{2 \log(1/\delta)} \right).$$

This was obtained from assuming simply that

$$\max_j \|X_j\|_2^2 \leq n.$$

$\frac{1}{\sqrt{n}}$  (slow rate)

## Gaussian Sequence Revisited

$$y = X\theta + \varepsilon.$$

$$n=d, \quad X = \sqrt{n} I, \quad \text{so} \quad y_i = \sqrt{n} \theta^*_i + \varepsilon_i$$

$$\arg \min \left[ \frac{1}{2n} \|y - X\theta\|_2^2 + C\|\theta\|_1 \right] \quad \textcircled{1}$$

We can transform the original model to obtain (multiply by  $\frac{X^T}{n}$ )

$$\begin{aligned} y' &:= \frac{X^T}{n} y = \frac{X^T X}{n} \theta^* + \zeta \\ &= \theta^* + \zeta \end{aligned}$$

$\zeta_i$  independent  $\frac{\sigma^2}{n}$ -sub-Gaussian.

Notice that  $\frac{1}{n} \|y - X\theta\|_2^2 = \left\| \frac{1}{n} X^T (y - X\theta) \right\|_2^2 = \|y' - \theta\|_2^2$

$$\begin{aligned} \Rightarrow \textcircled{1}(\leq) \quad &\arg \min \left[ \|y' - \theta\|_2^2 + 2C\|\theta\|_1 \right] \\ &= \arg \min_{\theta \in \mathbb{R}^d} \sum_{i=1}^d (y'_i - \theta_i)^2 + 2C\|\theta\|_1 \end{aligned}$$

$$\begin{aligned} \theta_i > 0 \quad &\Rightarrow \theta_i^2 = 2y'_i \theta_i + y_i'^2 + 2C\theta_i = \theta_i^2 + 2(z + y'_i)\theta_i + y_i'^2 \\ &\sim \text{minimized at } \hat{\theta}_i = y'_i - z \end{aligned}$$

$$\begin{aligned} \theta_i < 0 \quad &\Rightarrow \theta_i^2 - 2(z + y'_i)\theta_i + y_i'^2 \\ &\sim \text{minimized at } \hat{\theta}_i = y'_i + z \end{aligned}$$

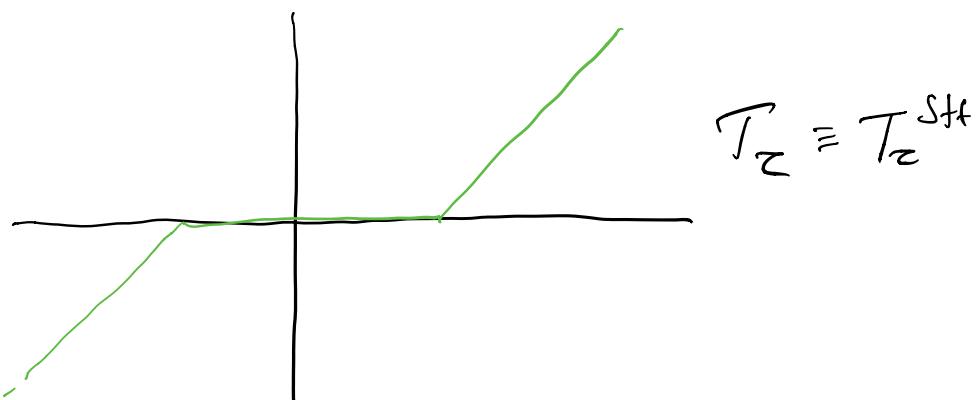
$$\begin{aligned} \text{Thus if } \quad &y'_i > z \quad \hat{\theta}_i = y'_i - z \\ &y'_i < z \quad \hat{\theta}_i = y'_i + z. \end{aligned}$$

else if  $|y'_i| \leq z$  no solution so minimum occurs at

boundary  $\Rightarrow \partial z = 0$

This can be summarized using the  
soft thresholding function

$$T_z^{\text{soft}}(x) := \begin{cases} \text{sign}(x) & |x| \geq z \\ 0 & \text{o.w.} \end{cases}$$



Then  $\hat{\beta}_i = T_z(y_i)$   $i = 1, \dots, d$

$$\text{let } z := z_n := \sigma \sqrt{\frac{2 \log(2d/\delta)}{n}}$$

$$\mathcal{A} := \left\{ \max_i |\hat{\beta}_i| \leq z/2 \right\}$$

Since  $\hat{\beta}_i$  is  $\frac{\sigma^2}{n}$ -Sub Gaussian we get

$$\mathbb{P}[A^c] \leq 2d e^{-\frac{n z^2}{8\sigma^2}} \leq \delta$$

On the event  $\mathcal{A}$  we can estimate

$$\begin{aligned}
\|\hat{\theta} - \theta^*\|_2^2 &= \sum_{j=1}^d \left| T_\varepsilon(y_j') - \theta_j^* \right|^2 \\
&= \sum_{j=1}^d \left[ \mathbb{1}\{y_j' \geq \varepsilon\} (\theta_j^* + \zeta_j - \varepsilon - \theta_j^*) \right. \\
&\quad \left. + \mathbb{1}\{y_j' \leq -\varepsilon\} (\theta_j^* + \zeta_j + \varepsilon - \theta_j^*) \right. \\
&\quad \left. - \mathbb{1}\{|y_j'| \leq \varepsilon\} |\theta_j^*| \right]^2 \\
&= \sum_{j=1}^d \left[ \mathbb{1}\{y_j' \geq \varepsilon\} (\zeta_j - \varepsilon) \right. \\
&\quad \left. + \mathbb{1}\{y_j' \leq -\varepsilon\} (\zeta_j + \varepsilon) \right. \\
&\quad \left. - \mathbb{1}\{|y_j'| \leq \varepsilon\} |\theta_j^*| \right]^2 \\
&\stackrel{\text{on } A}{\leq} \sum_{j=1}^d \left\{ 2c \mathbb{1}\{|y_j'| > \varepsilon\} + \mathbb{1}\{|y_j'| \leq \varepsilon\} |\theta_j^*| \right\}^2 \\
&\quad |y_j'| \leq \frac{3\varepsilon}{2} \\
&\quad \Rightarrow |\zeta_j| \leq \frac{\varepsilon}{2} \\
&\quad \Rightarrow |\theta_j^*| \leq \frac{3\varepsilon}{2} \\
\text{so } \|\hat{\theta} - \theta^*\|_2^2 &\leq \sum_{j=1}^d \left[ 2c \mathbb{1}\{|\theta_j^*| > \varepsilon\} + \mathbb{1}\{|\theta_j^*| \leq \frac{3\varepsilon}{2}\} |\theta_j^*| \right]^2 \\
&\leq \sum_{j=1}^d \left[ 4 \min\{|\theta_j^*|, \frac{3\varepsilon}{2}\} \right]^2 \\
&\leq \sum_{j=1}^d 16 \min\left\{ |\theta_j^*|, \frac{\varepsilon^2}{4} \right\} \leq \\
&\leq \sum_{j \in S} 16 \frac{\varepsilon^2}{4} \leq 4 \|\theta_s^*\|_0 \varepsilon^2
\end{aligned}$$

$$\leq \frac{32 \|x^*\|_0 \sigma^2 \log(2d/\delta)}{n}$$

from def'n of  $\mathcal{L} = \mathcal{L}_n$

So for the Gaussian seg<sup>n</sup> model we get the  $\frac{1}{n}$  rate.

In fact all we used was that we can multiply by

$$\frac{X^T}{n} \rightsquigarrow y' = \frac{X^T y}{n} = \theta^* + \{\}, \quad \{\} = \frac{X^T \varepsilon}{n}$$

So the same rates hold under (ORT)

since we can transform

$$\begin{aligned} y &= X\theta^* + \varepsilon \\ \rightsquigarrow y' &= \frac{1}{n} X^T X \theta^* + \{\} \quad \text{if } \frac{X^T X}{n} = I_n \\ &= \theta^* + \{\}. \end{aligned}$$

Why does ORT help so much?

it transforms the problem  $y = X\theta^* + \varepsilon$

where  $\theta^*$  is only observed after

it has been corrupted by the action of the operator  $X$

to a direct problem  $y' = \theta^* + \{\}$

where there's no corruption

## FAST RATES & RESTRICTED EIGENVALUE CONDITION

So far we have seen that we can obtain fast rates ( $\sqrt{n}$ ) under the **ORT** condition, at least for the prediction error.

Under the much weaker normalization condition  $\max_j \|X_{j\cdot}\|_2 \leq \sqrt{n}$  we were only able to obtain

the slower  $\frac{d}{\sqrt{n}}$  rate

We will now attempt to obtain the  $\sqrt{n}$  rate w/o **ORT**.

Here the noiseless scenario will serve as a guide in the following sense:

In the noiseless scenario the RNP was crucial for recovery.

In this setting we will now use a similar albeit stronger assumption

Let  $\alpha > 1$ ,  $S \subseteq [1:d]$

$$\mathcal{C}_\alpha(S) := \left\{ \theta \in \mathbb{R}^d \mid \|\theta_S\|_1 \leq \alpha \|\theta_S\|_2 \right\}$$

Definition (Restricted Eigenvalue Condition)

$X$  satisfies  $(\kappa, \alpha)$ -REC over  $S \subseteq [1:d]$   
if for all  $v \in \mathcal{C}_\alpha(S)$ , for some  $\kappa > 0$   
 $\kappa \|v\|_2^2 \leq \frac{1}{n} \|Xv\|_2^2$ .

Essentially this attempts to control the minimum eigenvalue of the matrix  $\frac{1}{n} X^T X$

But of course when  $\text{ker}(X) \neq \emptyset$  this is obviously 0,  
since  $\exists v \in \text{ker}(X), Xv=0$ .

Instead it tries to control the "minimum eigenvalue"  
when restricted to the subset  $C_\alpha(S)$ .

We will see later that the  $(\kappa, \alpha)$ -REC condition  
is also useful, and in a certain sense necessary,  
for recovery, i.e. when one wants  $\theta^*$  rather than  
simply prediction.

By that though, let's first see what result we  
can obtain using the  $(\kappa, \alpha)$ -REC.

Then

Suppose  $y = X\theta^* + \varepsilon$ , with  $\text{supp}(\theta^*) = S \subset [1:d]$ ,

$\varepsilon = (\varepsilon_i)_{i=1}^d$   $\sigma^2$ -SubGaussian. Suppose also that

$\max_j \|X_j\|_2^2 \leq n$  & that  $X$  satisfies the

$(\kappa, \beta)$ -REC w.r.t  $S$ . Let  $\hat{\theta}^*$  solve

$$\underset{\theta}{\operatorname{argmin}} \left\{ \frac{1}{2n} \|y - X\theta\|_2^2 + \zeta \|\theta\|_1 \right\}$$

$$\text{with } \zeta := \zeta_n := \sqrt{\frac{8\sigma^2 \log(2d/\delta)}{n}}.$$

Then  $\omega_p > 1 - \delta$

$$MSE(\hat{\theta}^L) \leq \frac{24 \|\theta^*\|_0 \sigma^2 \log(2d/\delta)}{kn}.$$

Pf. Let  $A := \left\{ \max_{j \in n} \frac{1}{n} \langle X_j, \varepsilon \rangle \leq z/2 \right\}$

$$\text{Then } \mathbb{P}(A^c) \leq 2d \exp\left(-\frac{n^2 z^2}{8\sigma^2 \max_j \|X_j\|_2}\right) \leq 2d \exp\left(-\frac{n z^2}{8\sigma^2}\right) \leq \delta$$

since  $\max_j \|X_j\|_2 \leq n$ .

$$\text{Let } L(\theta; z) := L_n(\theta; z) = \frac{1}{2n} \|y - X\theta\|_2^2 + z\|\theta\|_1$$

$$\text{let } \Delta := \hat{\theta}^L - \theta^*.$$

$$\text{Since } y - X\hat{\theta}^L = -X\Delta + \varepsilon$$

$$L(\hat{\theta}^L; z) = \frac{1}{2n} \|X\Delta\|_2^2 - \frac{1}{n} \langle X\Delta, \varepsilon \rangle + \frac{\|\varepsilon\|_2^2}{2n} + z\|\hat{\theta}^L\|_1$$

By definition of  $\hat{\theta}^L$  (argmin)

$$L(\hat{\theta}^L; z) \leq L(\theta^*; z) = \frac{1}{2n} \|\varepsilon\|_2^2 + z\|\theta^*\|_1$$

$$\begin{aligned} \text{so } \frac{1}{2n} \|X\Delta\|_2^2 - \frac{1}{n} \langle X\Delta, \varepsilon \rangle &\stackrel{+ \frac{\|\varepsilon\|_2^2}{2n} + z\|\hat{\theta}^L\|_1}{\cancel{\leq}} \\ &\leq \frac{1}{2n} \|\varepsilon\|_2^2 + z\|\theta^*\|_1 \end{aligned}$$

④

$$\text{so } 0 \leq \frac{1}{2n} \|X\Delta\|_2^2 \stackrel{\cancel{\leq}}{\leq} \frac{1}{n} \langle X\Delta, \varepsilon \rangle + z(\|\theta^*\|_1 - \|\hat{\theta}^L\|_1)$$

$\text{supp}(\theta^*) = S$  by assumption so

$$\begin{aligned} \|\theta^*\|_1 - \|\hat{\theta}^L\|_1 &= \|\theta_S^*\|_1 - \|\hat{\theta}_S^L\|_1 - \|\hat{\theta}_{S^c}^L\|_1 & (\hat{\theta}^L)_{S^c} &= (\theta^* + \Delta)_{S^c} \\ &= \|\theta_S^*\|_1 - \|\theta_S^* + \Delta_S\|_1 - \|\Delta_{S^c}\|_1 \\ &\leq \|\theta_S^*\|_1 - (\|\theta_S^*\|_1 - \|\Delta_S\|_1) - \|\Delta_{S^c}\|_1 = \|\Delta_S\|_1 - \|\Delta_{S^c}\|_1 \end{aligned}$$

$$\textcircled{S} \Rightarrow O \leq \frac{1}{n} \|X\Delta\|_2^2 \leq \frac{2}{n} \langle X^T \varepsilon, \Delta \rangle + 2c (\|\Delta_s\|_1 - \|\Delta_{s^c}\|_1) \\ \leq 2 \|\Delta\|_1 \max_j \frac{\langle x_j, \varepsilon \rangle}{n} + 2c (\|\Delta_s\|_1 - \|\Delta_{s^c}\|_1)$$

and on the event

$$\mathcal{A} \leq 2 \|\Delta\|_1 \cdot \frac{\varepsilon}{\sqrt{n}} + 2c (\|\Delta_s\|_1 - \|\Delta_{s^c}\|_1) \\ \leq c (3 \|\Delta_s\|_1 - \|\Delta_{s^c}\|_1)$$

Recapping: on the event  $\mathcal{A}$

$$O \leq \frac{1}{n} \|X\Delta\|_2^2 \leq c (3 \|\Delta_s\|_1 - \|\Delta_{s^c}\|_1) \\ \text{& thus } \Delta \in C_S(S). \quad \leftarrow \begin{array}{l} c \|\Delta\|_2 \leq \frac{1}{n} \|X\Delta\|_2^2 \\ \text{L}_1\text{-triangle inequality} \end{array}$$

Continuing from above

$$\frac{1}{n} \|X\Delta\|_2^2 \leq 3c \|\Delta_s\|_1 \leq 3c \sqrt{s} \|\Delta_s\|_2 \\ \leq 3c \sqrt{s} \|\Delta\|_2 \\ \leq 3c \sqrt{s} \frac{\|X\Delta\|_2}{\sqrt{kn}}$$

$$\Rightarrow \frac{1}{\sqrt{n}} \|X\Delta\|_2 \leq 3c \sqrt{\frac{s}{k}}$$

$$\text{Recall that } c = c_n = \sqrt{\frac{8\sigma^2 \log(2d/\delta)}{kn}}$$

$$\|X\Delta\|_2^2 \leq 3 \sqrt{\frac{s}{k}} \cdot \frac{8\sigma^2 \log(2d/\delta)}{\sqrt{kn}}$$

$$\text{& } \text{MSE}(X\hat{\theta}^*) = \frac{1}{n} \|X\Delta\|_2^2 \leq \frac{24 \|\theta^*\|_0 \sigma^2 \log(2d/\delta)}{kn}$$

□

So with the  $(n, \alpha)$ -REC we recover the  $1/n$  error rate.

## {BOUNDS ON L<sub>2</sub>-error}

It is in this situation where the (n, ε)-REC condition really comes into its own and can be fully appreciated.

To simplify things for now suppose we are looking at

$$\hat{\theta} = \underset{\theta \in K}{\operatorname{argmin}} \|y - X\theta\|_2^2 = \underset{\theta \in K = B_2(R)}{\operatorname{argmin}} L(\theta)$$

$$R = \|\theta^*\|_1$$

↑  
true feasible

$$\text{where } L(\theta) := \frac{1}{n} \|y - X\theta\|_2^2$$

$\hat{\theta}$  will minimize the empirical loss above, whereas

$$\begin{aligned} L(\theta) &= E[\|y - X\theta\|_2^2] = E[\|X(\theta^* - \theta) + \varepsilon\|_2^2] \\ &= E[\|\varepsilon\|_2^2] + E[\underbrace{\langle \varepsilon, X(\theta^* - \theta) \rangle}_{\parallel} + \|X(\theta^* - \theta)\|_2^2] \end{aligned}$$

this is clearly minimized at  $\theta = \theta^*$ .

**So** we are essentially minimizing the empirical loss

worrying that the solution (shrunken sample is large enough)

will be close to the minimizer  $\theta^*$  of  $L(\theta)$ , the true loss.

**BUT:** is this really what is happening?

when trying to minimize  $L_n(\theta)$ , the  $\boxed{n \gg 1}$

should guarantee that  $L_n(\hat{\theta}) \approx L(\hat{\theta})$

if so that if  $L_n(\hat{\theta})$  is small  $L(\hat{\theta})$  will also be small.

Think about it differently, if  $L_n \approx L$  and  $L(\theta^*)$  small

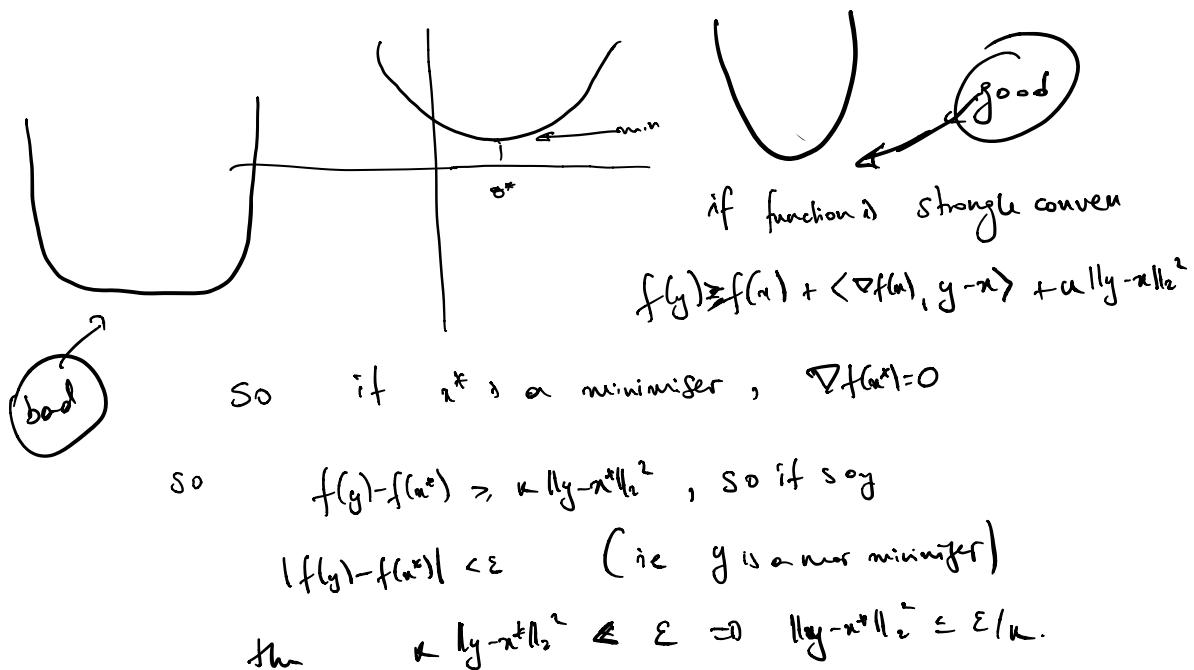
then  $L_n(\theta^*)$  should also be close to its minimum

1.2

$$|\ln(\theta^*) - \ln(\hat{\theta})| \approx \text{small}$$

When does this guarantee that  
 $|\theta^* - \hat{\theta}|$  is also small?

In the 1-d case the picture is quite simple



$$\text{so } f(y) - f(x^*) \geq \kappa \|y - x^*\|^2, \text{ so if say}$$

$$|f(y) - f(x^*)| < \varepsilon \quad (\text{ie } y \text{ is a near minimizer})$$

$$\text{then } \kappa \|y - x^*\|^2 \leq \varepsilon \Rightarrow \|y - x^*\|_2 \leq \varepsilon / \sqrt{\kappa}.$$

Now we know that  $\hat{\theta}$  is a minimizer of  $\ln(\theta)$  so  $\theta^*$  is a near minimizer

so is our function  $\ln(\theta)$  strongly convex?

despite "quadratic" it's not

When does  $\frac{1}{2n} \|y - X\theta\|^2$  if  $X$  has an-trivial kernel

so along the kernel of  $X$ .

The only hope is if we can control the convexity

on directions  $\perp \ker(X)$  &  $\theta^* \not\perp \ker(X)$

Ex:  $d=2, n=1 \quad X = [1, 0]$

$$y = X[\theta^*, \theta_2^*] + \varepsilon = \theta_1^* + \varepsilon$$

If  $\theta^* = [1, 0]$  & suppose  $X$   $(\kappa, \beta)$ -REC w.r.t  $\delta_{\varepsilon}$   
 $\varepsilon = \frac{\beta}{2}$

If  $v \in C_S(\mathcal{A})$ , i.e.  $\|Xv\|_2 \geq \|v\|_1$ , then

$$\|Xv\|_2^2 = v_1^2 = (1-\varepsilon)v_1^2 + \varepsilon v_1^2 \geq (1-\varepsilon)v_1^2 + \frac{\varepsilon}{\beta} v_2^2 \geq \kappa \|v\|_1^2$$

$$\kappa := \min\left\{1-\varepsilon, \frac{\varepsilon}{\beta}\right\}$$

If  $\hat{\theta}$  is the minimizer

$$(y - X\hat{\theta})^2 + 2\|\hat{\theta}\|_1 = (1+\varepsilon-\hat{\theta}_1)^2 + 2\|\hat{\theta}\|_1 + 2\|\hat{\theta}\|_1$$

$$\text{minimised at } \hat{\theta} = \left[1 + \varepsilon - \frac{\lambda}{2}, 0\right]$$

& prediction & L2 error  $(\varepsilon-\lambda)^2$

so for  $\varepsilon \ll \lambda$   $\lambda \ll 1$  small error is small

if now  $\theta^* = [0, 1]$  so  $X$  no longer satisfies  $(\kappa, \beta)$ -REC w.r.t  $\{\varepsilon\}$

for any  $\lambda > 0$ .

$$\text{Then } y = \varepsilon, \quad (y - X\hat{\theta})^2 + 2\|\hat{\theta}\|_1 = (\varepsilon - \hat{\theta}_1)^2 + 2\|\hat{\theta}\|_1 + 2\|\hat{\theta}\|_1$$

for small  $\lambda$  this is min'd at

$$\hat{\theta} = [0, \varepsilon - \lambda/2]$$

then prediction error  $\|X(\theta^* - \hat{\theta})\|_2^2 = 0$

but L2 error  $\|\theta^* - \hat{\theta}\|_2^2 = (1 - \varepsilon + \frac{\lambda}{2})^2$  which is large  
 $\varepsilon \ll \lambda$  small

Thm Suppose  $y = X\theta^* + \varepsilon$ , with  $\text{supp}(\theta^*) = S \subset [1:d]$ ,

$\varepsilon = (\varepsilon_i)_{i=1}^d$   $\sigma^2$ -SubGaussian. Suppose also that

$\max_j \|X_j\|_2 \leq n$  & that  $X$  satisfies the

$(\kappa, \beta)$ -REC wrt  $S$ . Let  $\hat{\theta}^*$  solve

$$\arg \min_{\theta} \left\{ \frac{1}{2n} \|y - X\theta\|_2^2 + z\|\theta\|_1 \right\}$$

$$\text{with } z := z_n := \sqrt{\frac{8\sigma^2 \log(2d/\delta)}{n}}$$

Then  $\omega_P > 1 - \delta$ .

$$\|\hat{\theta}^* - \theta^*\|_2 \leq \frac{3}{\kappa} \sqrt{\frac{8\|S\|\sigma^2 \log(2d/\delta)}{n}}$$

Pf: essentially the last proof, on the event  $A$  we showed  $\Delta \in C_S(S)$ , ( $\Delta = \hat{\theta}^* - \theta^*$ )

From the  $(\kappa, \beta)$ -condition

$$\sqrt{\kappa} \|\Delta\|_2 \leq \frac{1}{\sqrt{n}} \|X\Delta\|_2 \leq 3 \sqrt{\frac{|S|}{\kappa}}$$

& done



## RANDOM DESIGN

The idea here is to essentially try and

prove that the conditions on  $X$  we have been using so far hold for a random matrix  $\mathbf{X}$  under some assumptions

Thm:  $X \in \mathbb{R}^{n \times d}$ , iid  $N(0,1)$  entries. Then  $\exists$  universal constants  $c_1 < 1 < c_2$  such that

$$\textcircled{**} \quad \frac{\|X\theta\|_2^2}{n} \geq c_1 \|\theta\|_2^2 - c_2 \frac{\log d}{n} \|\theta\|_1^2, \quad \theta \in \mathbb{R}^d$$

$$\text{w.p.} > 1 - \frac{e^{-n/32}}{1 - e^{-n/32}}.$$


---

Pf: Sufficient to look at  $\theta \in \mathbb{S}^{d-1}$

Let  $g(t) := 2\sqrt{\frac{\log(d)}{n}}t$  & define the event

$$\mathcal{E} := \left\{ X \in \mathbb{R}^{n \times d} \mid \inf_{\theta \in \mathbb{S}^{d-1}} \frac{\|X\theta\|_2}{\sqrt{n}} \leq \frac{1}{4} - 2g(\|\theta\|_1) \right\}$$

exercise: prove that on  $\mathcal{E}^c$  the desired bound  $\textcircled{**}$  holds.

We need to bound  $\mathbb{P}[\mathcal{E}]$ .

Although we can now control  $\|\theta\|_2^2$ , since dimension is high

$\|\theta\|_1$  can vary by a factor of  $\sqrt{n}$

→ Split  $\mathcal{E}$  into smaller pieces according to size of  $\|\theta\|_1$ .

For  $0 < r < s$  let

$$K(r,s) := \left\{ \theta \in \mathbb{S}^{d-1} \mid g(\|\theta\|_1) \in [r, s] \right\}$$

and consider  $\mathcal{A}(r,s) := \left\{ X \mid \inf_{\theta \in K(r,s)} \frac{\|X\theta\|_2}{\sqrt{n}} \leq \frac{1}{2} - 2s \right\}$ .

Claim:  $\mathcal{E} \subset \mathcal{A}(0, 1/\alpha) \cup \bigcup_{l=1}^{\infty} \mathcal{A}\left(\frac{2^{l-1}}{\alpha}, \frac{2^l}{\alpha}\right)$

Pf of Claim: Let  $\vartheta$  attain the infimum in  $\mathcal{E}$  ( $\frac{\|X\vartheta\|_2}{\sqrt{n}} = \inf_{\theta \in \mathcal{E}} \frac{\|X\theta\|_2}{\sqrt{n}}$ )

$\vartheta$  must belong to either  $K(0, 1/\alpha)$  or one of the  $K\left(\frac{2^{l-1}}{\alpha}, \frac{2^l}{\alpha}\right)$

If  $\vartheta \in K(0, 1/\alpha)$   $g(\|\vartheta\|_1) \leq 1/\alpha$

$$\text{for } X \in \mathcal{E} \quad \frac{\|X\vartheta\|_2}{\sqrt{n}} \leq \frac{1}{4} - 2g(\|\vartheta\|_1) \leq \frac{1}{4} = \frac{1}{2} - \frac{1}{4}$$

$$\Rightarrow \vartheta \in \mathcal{A}(0, 1/\alpha)$$

Else If  $\vartheta \in K\left(\frac{2^{l-1}}{\alpha}, \frac{2^l}{\alpha}\right)$ , for  $l > 1$

$$\frac{\|X\vartheta\|_2}{\sqrt{n}} \leq \frac{1}{4} - 2g(\|\vartheta\|_1) \leq \frac{1}{8} - 2 \times \frac{2^{l-1}}{4} = \frac{1}{2} - \frac{2^l}{4}$$

$$\Rightarrow \vartheta \in \mathcal{A}\left(\frac{2^{l-1}}{\alpha}, \frac{2^l}{\alpha}\right).$$

Union bound  $\Rightarrow \mathbb{P}[\mathcal{E}] \leq \mathbb{P}[\mathcal{A}(0, 1/\alpha)] + \sum_{l=1}^{\infty} \mathbb{P}\left[\mathcal{A}\left(\frac{2^{l-1}}{\alpha}, \frac{2^l}{\alpha}\right)\right]$

Claim:  $\mathbb{P}(\mathcal{A}(r,s)) \leq e^{-n/s^2} e^{-ns^2/2}$

Notice we can aim for a lower bound on

$$T(r,s) = -\inf_{\theta \in K(r,s)} \frac{\|X\theta\|_2}{\sqrt{n}}$$

$$T(\gamma) = -\inf_{\theta \in K(\gamma)} \sup_{u \in \mathbb{S}^{n-1}} \frac{\langle u, \gamma \rangle}{\sqrt{n}}$$

$$= \sup_{\theta \in K(\gamma)} \inf_{u \in \mathbb{S}^{n-1}} \frac{\langle u, \gamma \rangle}{\sqrt{n}}$$

Write  $X_{u,\theta} := \langle u, \gamma \rangle$

$X_{u,\theta}$  is a Gaussian process on  $\mathbb{S}^{n-1} \times \mathbb{S}^{d-1}$

$$X_{u,\theta} \sim N(0, 1_n)$$

We'll use

Then (Gordon '85) Let  $\{X_{ij} : i \in [n], j \in [m]\}, \{Y_{ij} : i \in [n], j \in [m]\}$   
two Gaussian arrays s.t.

$$\mathbb{E}[(X_{ij} - X_{ik})^2] \leq \mathbb{E}[(Y_{ij} - Y_{ik})^2] \quad \forall i, j, k.$$

$$\mathbb{E}[(X_{ij} - X_{il})^2] \geq \mathbb{E}[(Y_{ij} - Y_{il})^2] \quad i \neq l$$

Then  $\mathbb{E}\left[\min_i \max_j X_{ij}\right] \leq \mathbb{E}\left[\min_i \max_j Y_{ij}\right]$

Our  $X_{u,\theta}$  will be compared with

$$Y_{u,\theta} := \frac{\langle u, \gamma \rangle}{\sqrt{n}} + \frac{\langle \theta, \gamma \rangle}{\sqrt{n}}, \quad \begin{cases} \gamma \sim N(0, 1_n) \\ \theta \sim N(0, 1_d) \end{cases}$$

$$\mathbb{E}[(X_{u,\theta} - Y_{u,\theta})^2] = \mathbb{E}[\langle u, \gamma - (\theta - \alpha) \rangle]^2 = \sum_{i,j} u_i^2 (\theta_j - \alpha_j)^2 = \|\theta - \alpha\|^2$$

$$\mathbb{E}[(Y_{u,\theta} - Y_{u,\alpha})^2] = \mathbb{E}[\langle \theta, \theta - \alpha \rangle^2] = \|\theta - \alpha\|^2 = 1.$$

$$\text{if } u \neq w \quad \mathbb{E}[(y_{u,s} - y_{w,s})^2] = \mathbb{E}[(\langle u, \varphi_s \rangle + \langle s, \varphi_u \rangle)^2]$$

$$= \|u - w\|^2 + \|s - q\|^2$$

$$\mathbb{E}[(x_{u,s} - x_{w,s})^2] = \mathbb{E}[(\langle u, x_s \rangle - \langle w, x_s \rangle)^2]$$

$$= \mathbb{E}\left[\left(\sum_{i,j} x_{ij} (u_i \varphi_j - w_i \varphi_j)\right)^2\right]$$

$$= \sum_{i,j} (u_i \varphi_j - w_i \varphi_j)^2 = \|u \varphi^\top - w \varphi^\top\|_F^2$$

↗  
Frobenius