Adaptive rates for trend filtering using dual certificates

Sara van de Geer

December, 2020

CIRM 2020

Joint work with **Francesco Ortelli**



< □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □

Contents

- Lasso with structured design
- Part I Non-adaptive bounds
- Best N-term approximations
- Application: 1-hidden-layer neural network
- Entropy of convex hulls based on covering numbers based on best *N*-term approximations
- Examples
- Part II Adaptive bounds for trend filtering

Graphs

(ロ) (同) (三) (三) (三) (三) (○) (○)

- Effective sparsity (dual certificates)
- 1-dimensional case
- Higherdimensional case

This course is based on ideas from

- o Koltchinskii, Lounici and Tsybakov [2011]
- o Dalalyan, Hebiri and Lederer [2016]
- Fuchs [2004]
- o Candès and Fernandez-Granda [2014]

◆□▶ ◆□▶ ▲□▶ ▲□▶ ■ ののの

o . . .

Lasso

- "Structured" design matrix $\Psi \in \mathbb{R}^{n \times p}$
- ▶ Response vector $Y \in \mathbb{R}^n$, $Y \sim \mathcal{N}(f^0, I)$

Goal: estimate the unknown *f*⁰ using the Lasso

$$\hat{f} := \arg\min_{f_b = \Psi b, \ b \in \mathbb{R}^p} \left\{ \|Y - \Psi b\|_2^2/n + 2\lambda \|b\|_1
ight\}.$$

▲□▶ ▲□▶ ▲ 三▶ ▲ 三▶ - 三 - のへぐ

Let $\Psi \subset (L_2(Q), \|\cdot\|_Q)$. Notation

For a linear space $\mathcal{V} \subset L_2(Q)$ and $f \in L_2(Q)$ we let $f_{\mathcal{V}}$ be the projection of f on \mathcal{V} and $f_{\mathcal{V}^{\perp}} := f - f_{\mathcal{V}}$. We let

$$\mathsf{u}(\mathcal{V},\Psi) := \sup_{\psi \in \Psi} \|\psi_{\mathcal{V}^{\perp}}\|_Q.$$

Definition

For $N \in \mathbb{N}$ define the best N-term approximation

$$\gamma(N,\Psi) := \min \left\{ \mathbf{u}(\mathcal{V},\Psi) : \dim(\mathcal{V}) = N \right\}.$$

▲□▶▲□▶▲□▶▲□▶ □ のQ@

Let $\lambda_0(u) := \sqrt{2 \log(2p) + 2u}$. Theorem

Suppose $\gamma(N, \Psi) \leq AN^{-1/W}$. With probability at least $1 - \exp[-u] - \exp[-v]$, and for

$$\delta_n^2 = \frac{(2A\lambda_0(u))^W}{n\lambda^W} + \frac{1+2v}{n}$$

we have

$$\|\hat{f} - f^{0}\|_{2}^{2}/n + \lambda \|\hat{\beta}\|_{1} \leq \underbrace{\left\{ \|f^{*} - f^{0}\|_{2}^{2}/n + 3\lambda \|\beta^{*}\|_{1} \right\}}_{\text{approximation error}} + \delta_{n}^{2}.$$

$$\uparrow_{\text{estimation}}$$

◆□▶ ◆□▶ ▲□▶ ▲□▶ □ のQ@

Application: one-hidden-layer NN

Definition

The hinge, or ReLU, function is

$$z\mapsto z_+: egin{cases} z, & z\geq 0\ 0, & z<0 \end{cases}$$



The Rell activation function

▲□▶ ▲□▶ ▲ 三▶ ▲ 三▶ - 三 - のへぐ

Notation

For $v \in \mathbb{R}^n$ and $c \in \mathbb{R}$

$$(\mathbf{v}-\mathbf{c})_+:=egin{pmatrix} (\mathbf{v}_1-\mathbf{c})_+\ dots\ (\mathbf{v}_n-\mathbf{c})_+ \end{pmatrix}.$$

Construction of the dictionary $\boldsymbol{\Psi}$

• *d* := input dimension

•
$$X = \begin{pmatrix} x_1 \\ \vdots \\ x_n \end{pmatrix} \in \mathbb{R}^{n \times d}$$
 input matrix

•
$$W := (w_1, \dots, w_p) \in \mathbb{R}^{r \times p}$$
 matrix of weights

$$\|w_j\|_2 = 1, j \in \{1, \dots, p\}$$

•
$$\mathbf{c} = (c_1, \dots, c_p) \in \mathbb{R}^p$$
 vector of biases

$$\circ \psi_j := (Xw_j - c_j)_+ \in \mathbb{R}^n, j = 1, \dots, p$$

$$\circ \Psi := \{\psi_1, \dots, \psi_p\} = (XW - \mathbf{c})_+ \in \mathbb{R}^{n imes p}$$

Lemma Suppose $\max_{1 \le i \le n} \|x_i\|_2 \le 1$. Then $\gamma(N, \Psi) \le AN^{-1/d}$.

Lasso for 1-hidden-layer NN

Define

$$\mathcal{NN} := \Big\{ f_b = (XW - \mathbf{c})_+ b : b \in \mathbb{R}^p \Big\}.$$

Let

$$\hat{f} = f_{\hat{\beta}} := \arg \min_{f_b \in \mathcal{NN}} \bigg\{ \|\mathbf{Y} - f_b\|_2^2 / n + 2\lambda \|b\|_1 \bigg\}.$$

[Parhi and Nowak , 2020]

▲□▶ ▲□▶ ▲ 三▶ ▲ 三▶ - 三 - のへぐ

Simulation by Peter Hinz d := 100 p := 10000 $\sigma^2 := 1$



◆□▶ ◆□▶ ◆臣▶ ◆臣▶ 臣 のへぐ



◆□▶ ◆□▶ ◆ □▶ ◆ □▶ ● □ ● ● ●

Question:

How good are the bounds using best *N*-term approximations? Note:

∘ { Ψb : $b \in \mathbb{R}^{p}$, $||b||_{1} = 1$ } = conv(± Ψ) is the convex hull of ± Ψ

 $\circ \Psi$ "structured": it has polynomial covering numbers $N(\cdot)$

- [Ball and Pajor, 1990] \rightsquigarrow entropy $H(\cdot)$ of conv(± Ψ)
- Estimation error for LSE over $\operatorname{conv}(\pm \Psi)$ is $\delta_n^2 \simeq H(\delta_n)/n$

Covering and entropy

Let (S, d) be a metric space. Definition For $\delta > 0$, δ -covering number $N(\delta, S, d) :=$ minimum # balls with radius δ necessary to cover S



(ロ) (同) (三) (三) (三) (三) (○) (○)

Definition Entropy $H(\cdot, S) := \log N(\cdot, S)$.

Entropy of convex hulls

Let $\Psi \in (L_2(Q), \|\cdot\|_Q)$. Theorem [Ball & Pajor, 1990] ["Covering Theorem"] Suppose for some positive constants A and W

$$N(\delta, \Psi) \leq A\delta^{-W}, \ \delta > 0.$$

Then for a constant C

$$H(\delta,\operatorname{conv}(\Psi)) \leq C \delta^{-rac{2W}{2+W}}, \; \delta > 0.$$

The entropy bound based on covering numbers can be too loose

Definition The hinge, or ReLU, function is $\begin{bmatrix} 7 & 7 > 0 \end{bmatrix}$

$$z\mapsto z_+:\begin{cases} 2, & 2\leq 0\\ 0, & z<0 \end{cases}$$



Let

$$X := \left(\frac{1}{n}, \frac{2}{n}, \cdots, 1\right)^{\top},$$

$$\psi_1 \equiv 1, \ \psi_2 := X - 1/n, \ \psi_j := \left(X - \frac{j-1}{n}\right)_+, \ j \in \{3, \dots, n\}.$$

Define the space of linear functions

$$\mathcal{N} := \{ \boldsymbol{z} = \psi_1 \boldsymbol{b}_1 + \psi_2 \boldsymbol{b}_2 : (\boldsymbol{b}_1, \boldsymbol{b}_2)^\top \in \mathbb{R}^2 \}$$

Then

$$\left\{z+f: z \in \mathcal{N}, f \in \operatorname{conv}(\{\psi_j\}_{j=3}^n)\right\} = \left\{f \in \mathbb{R}^n: n \|\Delta^2 f\|_1 = 1\right\}$$

with

$$(\Delta^2 f)_j := f_j - 2f_{j-1} + f_{j-2}, \ j \in \{3, \ldots, n\}.$$

Let

$$\|v\|_{Q_n}^2 := \frac{1}{n} \sum_{i=1}^n v_i^2, \ v \in \mathbb{R}^n.$$

By the [Covering Theorem] with W = 1

$$H(\delta,\operatorname{conv}(\{\psi_j\}_{j=3}^n),\|\cdot\|_{Q_n})\leq C\delta^{-\frac{2}{3}},\delta>0.$$

But by [Babenko, 1979]

$$H(\delta,\operatorname{conv}(\{\psi_j\}_{j=3}^n), \|\cdot\|_{Q_n}) \leq C\delta^{-\frac{1}{2}}, \delta > 0!$$

▲□▶ ▲□▶ ▲ 三▶ ▲ 三▶ - 三 - のへぐ

Notation For R > 0, $\mathcal{F}(R) :=$ $\{f \in \overline{\text{conv}}(\Psi) : \|f\|_Q \le R\}$



◆□▶ ◆□▶ ◆臣▶ ◆臣▶ ─臣 ─のへで

Theorem ["Projection Theorem"] Suppose that for some constants positive A and W

$$\gamma(\mathbf{N}, \Psi) \leq \mathbf{A} \mathbf{N}^{-\frac{1}{\mathbf{W}}}$$

For all R > 0 and $\delta > 0$,



◆□▶ ◆□▶ ▲□▶ ▲□▶ ■ ののの

In other words,

$$H(\delta, \mathcal{F}(R)) \lesssim \delta^{-rac{2W}{2+W}} \log(1/\delta)$$

and for the local entropy where $R \simeq \delta$ we get

$$H(\delta, \mathcal{F}(R)) \lesssim \delta^{-rac{2W}{2+W}} \log^{rac{W}{2+W}}(1/\delta).$$

Example: second order discrete derivatives.

$$\mathcal{F} := \left\{ f \in \mathbb{R}^n : n \| \Delta^2 f \|_1 = 1 \right\} = \left\{ \text{linear} + \text{conv}\left(\underbrace{\{\psi_j\}_{j=3}^n}_{\text{ReLU functions}}\right) \right\}$$

Then

$$\begin{split} \mathsf{N}(\delta, \{\psi_j\}_{j=3}^n, \mathsf{Q}_n) &\asymp \quad \delta^{-1} : \ \mathsf{W}_{\mathrm{cov}} = 1\\ \overset{[\operatorname{Covering Thm}]}{\Rightarrow} \ \mathsf{H}(\delta, \operatorname{conv}\{\psi_j\}_{j=3}^n) &\lesssim \quad \delta^{-\frac{2\mathsf{W}_{\mathrm{cov}}}{2+\mathsf{W}_{\mathrm{cov}}}} = \delta^{-\frac{2}{3}}.\\ \gamma(\mathsf{N}, \{\psi_j\}_{j=3}^n) &\leq \quad \mathsf{AN}^{-\frac{1}{2/3}} : \ \mathsf{W}_{\mathrm{proj}} = \frac{2}{3}\\ \overset{[\operatorname{Projection Thm}]}{\Rightarrow} \ \mathsf{H}(\delta, \operatorname{conv}\{\psi_j\}_{j=3}^n) &\lesssim \quad \delta^{-\frac{2\mathsf{W}_{\mathrm{proj}}}{2+\mathsf{W}_{\mathrm{proj}}}} \log n = \delta^{-\frac{1}{2}} \log n. \end{split}$$

▲□▶▲□▶▲□▶▲□▶ □ のへで

Example: k^{th} order discrete derivatives

$$\mathcal{F} := \left\{ f \in \mathbb{R}^n : \ n^{k-1} \|\Delta^k f\|_1 = 1 \right\}$$
$$= \left\{ (k-1)^{th} \text{ order polynomial} + \operatorname{conv}(\{\psi_j\}_{j=k+1}^n) \right\},$$

where $\{\psi_j\}_{j=k+1}^n$ is the falling factorial basis. [Wang, Smola, Tibshirani, 2014]

Then

.

$$N(\delta, \{\psi_j\}_{j=k+1}^n) \approx \delta^{-1} : W_{\text{cov}} = 1$$

$$\stackrel{[\text{Covering Thm}]}{\Rightarrow} H(\delta, \{\psi_j\}_{j=k+1}^n) \lesssim \delta^{-\frac{2W_{\text{cov}}}{2+W_{\text{cov}}}} = \delta^{-\frac{2}{3}}.$$

$$\gamma(N, \{\psi_j\}_{j=k+1}^n) \leq AN^{-\frac{1}{2/(2k-1)}} : W_{\text{proj}} = \frac{2}{2k-1}$$

$$\stackrel{[\text{Projection Thm}]}{\Rightarrow} H(\delta, \{\psi_j\}_{j=k+1}^n) \lesssim \delta^{-\frac{2W_{\text{proj}}}{2+W_{\text{proj}}}} \log n = \delta^{-\frac{1}{k}} \log n.$$

Higher-dimensional extension (MARS)

Let $x_i = (\xi_{i,1}, \ldots, \xi_{i,r}) \in \{1/n_0, \ldots, 1\}^d$ and let $n = n_0^r$. Let $\{\psi_j\}_{j=1}^p$ be real-valued functions on $\{1/n_0, \ldots, 1\}$ with $\max_{1 \le j \le p} \|\psi_j\|_{\infty} \le 1$. Let for $\mathbf{j} = (j_1, \ldots, j_r) \in \{1, \ldots, p\}^d$

$$\psi_{\mathbf{j}}(\mathbf{x}_i) = \prod_{t=1}^d \psi_{j_t}(\xi_{i,t})$$

Then functions of the form

$$f = \sum_{\mathbf{j} \in \{1, \dots, p\}^d} \psi_{\mathbf{j}} b_{\mathbf{j}}$$

with ψ_j the falling factorial basis functions corresponds to those used in the context of multiplicative adaptive regression splines (MARS) [Friedman, 1991].

Let

$$\mathcal{F} := \bigg\{ f = \sum_{\mathbf{j} \in \{1, \dots, p\}^d} \psi_{\mathbf{j}} \mathbf{b}_{\mathbf{j}} : \sum_{\mathbf{j} \in \{1, \dots, p\}^d} |\mathbf{b}_{\mathbf{j}}| \leq 1 \bigg\}.$$

.....

. . . .

Then

$$\gamma(N, \{\psi_{\mathbf{j}}\}) \leq AN^{-\frac{2k-1}{2[1+\frac{1}{2}+\cdots+\frac{1}{d}]}}.$$

$$W = W_{\text{proj}} = \frac{2[1+\frac{1}{2}+\cdots+\frac{1}{d}]}{2k-1}$$

$$\left(W_{\text{cov}} = 2d\right)$$
which gives

$$\frac{2W}{2+W} = \frac{2h}{2k+h-1}, h = 1 + \frac{1}{2} + \dots + \frac{1}{d}.$$

We then find by [Projection Theorem]

$$H_2(\delta, \mathcal{F}(R)) \lesssim \delta^{-\frac{2h}{2k+h-1}} \log p$$

For k = 1 still not tight

$$\frac{1}{\delta} \sim H_2(\delta, \mathcal{F}(R)) \lesssim \delta^{-\frac{2h}{h+1}} \log p$$

Part II: Adaptation for the trend filtering problem

Trend filtering problem:

$$\hat{f} := \min_{f \in \mathbb{R}^n} \left\{ \|Y - f\|_n^2 + 2\lambda \|\Delta^k f\|_1 \right\}$$

$$= \min_{f_b = \Psi b} \left\{ \|Y - f_b\|_n^2 + 2\lambda \|b_{-\{1,...,k\}}\|_1 \right\}$$

where Ψ is the falling factorial basis.

[Wang, Smola, Tibshirani, 2014]

Notation $S \subset \{k + 1, \ldots, n\}, s := |S|$

$$\mathcal{N}_{-S} := \{f \in \mathbb{R}^n : (\Delta^k f)_j = 0 \forall j \notin S\}$$

$$r_S := \dim(\mathcal{N}_{-S}) = k + s = \boxed{\text{sparsity}}$$

$$\|f^* - f^0\|_2^2 / n := \text{ approximation error, } f^* \in \mathcal{N}_{-S}$$

$$r_S / n := \text{ estimation error}$$

Goal: Show that $\|\hat{f} - f^0\|_n^2$ trades off approximation error and estimation error r_S/n .

< □ > < 同 > < 三 > < 三 > < 三 > < ○ < ○ </p>

Notation:

Thus

$$Df := \Delta^k f,$$

 $D_S f := \{ (Df)_j : j \in S \}, \ D_{-S} f := \{ (Df)_j : j \notin S \}$

$$\mathcal{N}_{-\mathcal{S}} = \{ f \in \mathbb{R}^n : D_{-\mathcal{S}} f = 0 \}$$

Choice of the tuning parameter

Let $\mathcal{V} \supset \mathcal{N}_{-S}$ be a linear space. Recall

$$\mathbf{u}(\mathcal{V}, \Psi) := \sup_{\psi \in \Psi} \|\psi_{\mathcal{V}^{\perp}}\|_{Q_n}.$$

Definition For $N \ge r_S$ the best N-term approximation is

$$\begin{split} \gamma_{\mathcal{S}}(N,\Psi) &:= \min \left\{ \mathbf{u}(\mathcal{V},\Psi) : \ \mathcal{V} \supset \mathcal{N}_{-\mathcal{S}}, \ \dim(\mathcal{V}) = N \right\} \\ &:= \mathbf{u}(\mathcal{V}_{N},\Psi). \end{split}$$

▲□▶ ▲□▶ ▲□▶ ▲□▶ = 三 のへで

Define for u > 0

$$\lambda_0(u) := \sqrt{\frac{2(\log(2n) + u)}{n}}.$$

Take



▲□▶ ▲□▶ ▲ 三▶ ▲ 三▶ - 三 - のへぐ

Effective sparsity without noise

Definition

Let $q \in \mathbb{R}^{n-k}$ where $q_S \in \{\pm 1\}^s$ (s := |S|) is a fixed sign vector. We call q an interpolating vector

(that interpolates the signs at S).

Definition The noiseless effective sparsity is

$$\Gamma^2(q_S) := \min \bigg\{ n \| D^\top q \|_2^2 : \ q \text{ interpolating } |q_j| \leq 1 \ \forall \ j \notin S \bigg\}.$$



Remark This is a quantified version of the dual certificates used in noiseless compressed sensing.

Definition Define the normalized noise weights

$$\mathbf{w}_{j} := \lambda_{0} \| \psi_{\mathcal{V}_{N}^{\perp}, j} \|_{Q_{n}} / \lambda, \ j \notin S.$$

The effective sparsity is

$$\Gamma^2(q_{\mathcal{S}}, w_{-\mathcal{S}}) := \min\left\{n\|D^\top q\|_2^2: \ q \text{ interpolating } |q_j| \leq 1 - w_j \ \forall \ j \notin \mathcal{S}\right\}.$$



Oracle inequality

Theorem Let $f^* \in \mathbb{R}^n$ be arbitrary and $S := \{j : (Df^*)_j \neq 0\}$. Let $\lambda \ge \lambda_0(u)\mathcal{U}_S(N, \Psi)$. For all u > 0 and v > 0, we have with probability at least $1 - \exp[-u] - \exp[-v]$

$$\|\hat{f} - f^{0}\|_{2}^{2}/n \leq \underbrace{\|f^{*} - f^{0}\|_{2}^{2}/n}_{estimation \ error} + \underbrace{\left(\sqrt{\frac{N}{n} + \sqrt{\frac{2v}{n}} + \lambda\Gamma(q_{S}, w_{-S})}\right)^{2}}_{estimation \ error}$$

where

$$\Gamma^2(\boldsymbol{q}_{\mathcal{S}}, \boldsymbol{w}_{-\mathcal{S}}), \ \boldsymbol{q}_{\mathcal{S}} := \operatorname{sign}(\boldsymbol{D}_{\mathcal{S}}f^*), \ j \in \mathcal{S}$$

is the effective sparsity.

[Dalalyan, Hebiri and Lederer; 2017]

$\begin{array}{c} \text{Example} \\ \text{Total variation (TV) in } \mathbb{R} \end{array}$

$$\hat{f} := \arg\min_{f \in \mathbb{R}^n} \left\{ \|Y - f\|_{Q_n}^2 + 2\lambda \|Df\|_1 \right\}$$

with

$$||Df||_1 := \sum_{j=2}^n |f_j - f_{j-1}| = \mathrm{TV}(f)$$

Thus

$$D = \begin{pmatrix} -1 & +1 & 0 & \cdots & 0 & 0 \\ 0 & -1 & +1 & \cdots & 0 & 0 \\ \vdots & \vdots & \vdots & \ddots & \vdots & \vdots \\ 0 & 0 & 0 & \cdots & -1 & +1 \end{pmatrix} \Rightarrow$$
$$D^{\mathsf{T}} = -\begin{pmatrix} +1 & 0 & \cdots & 0 \\ -1 & +1 & \cdots & 0 \\ 0 & -1 & \cdots & 0 \\ \vdots & \vdots & \vdots & \vdots \\ 0 & 0 & \cdots & +1 \\ 0 & 0 & \cdots & +1 \end{pmatrix} \Rightarrow$$



◆□▶ ◆□▶ ▲□▶ ▲□▶ ■ ののの

• $n_1 = t_1 - 1$, $n_2 = t_2 - t_1$, ..., $n_{s+1} = n + 1 - t_s$ be the distances between jumps.
We take

$$N = r_S, \ \mathcal{V}_N = \mathcal{N}_{-S}.$$

Then have noise weights

$$\|\psi_{\mathcal{V}_{N}^{\perp},t_{j-1}+i}\|_{2} = \sqrt{\frac{i(n_{j}-i)}{d_{j}}}, \ i \in [1,n_{j}],$$
$$\gamma_{\mathcal{S}}(N,\Psi) = \sqrt{\frac{n_{\max}}{4n}} \qquad (1/W = 1/2)$$

Choice of tuning parameter:

$$\lambda \geq \lambda_0(u) \sqrt{\frac{n_{\max}}{4n}}.$$

◆□▶ ◆□▶ ◆ □▶ ◆ □▶ ● □ ● ● ● ●



(日)







◆□▶ ◆□▶ ◆臣▶ ◆臣▶ 臣 のへぐ

+1

Łs

Application of the theorem to the total variation penalty:

Take

$$\lambda \geq \underbrace{\frac{\sqrt{2(\log(2n)+u)}}{n}}_{\asymp \sqrt{\log n}/n} \sqrt{\frac{n_{\max}}{4}}.$$

Then with probability at least $1 - \exp[-u] - \exp[-v]$

$$\begin{aligned} \|\hat{f} - f^{0}\|_{2}^{2}/n &\leq \|f^{*} - f^{0}\|_{2}^{2}/n \\ &+ \left(\sqrt{\frac{s+1}{n}} + \sqrt{\frac{2v}{n}} + \lambda \Gamma(q_{s}, w_{-s})\right)^{2} \\ &+ 4\lambda \|D_{-s}f^{*}\|_{1} \end{aligned}$$

where

$$\Gamma^2(q_S, w_{-S}) \leq \frac{n \log n_1}{n_1} + \sum_{j=2}^s \frac{4n \log(n_j/2)}{n_j} + \frac{n \log n_{s+1}}{n_{s+1}}$$

▲□▶ ▲□▶ ▲□▶ ▲□▶ ▲□ ● のへぐ

Special case: equally spaced jumps Suppose $d_j = \frac{n}{s+1}$. We get

$$\Gamma^2(q_S, w_{-S}) \leq 6(s+1)^2 \log n$$

We may take

$$\lambda \asymp \underbrace{\lambda_0(u)}_{\asymp \sqrt{\frac{\log n}{n}}} \underbrace{\sqrt{\frac{1}{s+1}}}_{\asymp \gamma_S(r_S, \Psi)}.$$

So we get

$$\lambda^{2} \Gamma^{2}(q_{S}, w_{-S}) = \mathcal{O}\left(\frac{s+1}{n}\right) \log^{2} n$$
$$= \mathcal{O}\left(\frac{\text{number of parameters}}{\text{number of observations}}\right) \log^{2} n$$

Example total variation of first discrete derivative

$$\hat{f} := \arg\min_{f \in \mathbb{R}^n} \left\{ \|Y - f\|_{Q_n}^2 + 2\lambda \|Df\|_1 \right\}$$

with

$$Df := \Delta^2 f$$

the second order discrete derivative

$$(\Delta^2 f)_j = f_j - 2f_{j-1} + f_{j-2}, \ j \ge 3.$$

[Tibshirani; 2014], [Guntuboyina, Lieu, Chatterjee & Sen; 2020] Then

$$\|Df\|_1 := \|\Delta^2 f\|_1 = \mathrm{TV}(\Delta f).$$

▲□▶▲□▶▲□▶▲□▶ □ のQ@



◆□▶ ◆□▶ ▲□▶ ▲□▶ □ のQ@

Let

$$\mathcal{V}_{N} := \mathcal{N}_{-S} \cup \operatorname{span} \{ \psi_{t_{1}+1}, \dots, \psi_{t_{s}+1} \},$$
$$N = r_{S} + s = 2(s+1)$$

We get

$$\gamma_{\mathcal{S}}(N,\Psi) \leq A n_{\max}^{3/2} / \sqrt{n} \quad (1/W = 3/2).$$

◆□▶ ◆□▶ ◆ □▶ ◆ □▶ ● □ ● ● ● ●



▲□▶ ▲□▶ ▲□▶ ▲□▶ ▲□ ● のへぐ

Noise and tuning: noise weights.



◆□▶ ◆□▶ ◆臣▶ ◆臣▶ 三臣 - のへで

Application of the theorem to the total variation penalty on Δf : Theorem For all u > 0, v > 0 and for

$$\lambda \geq c_2 n \lambda_0(u) \sqrt{n_{ ext{max}}^3/n_{ ext{max}}^3}$$

we have with probability at least $1 - \exp[-u] - \exp[-v]$

$$\|\hat{f} - f^{0}\|_{2}^{2}/n \le \|f^{*} - f^{0}\|_{2}^{2}/n$$
$$+ \left(\sqrt{\frac{2(s+1)}{n}} + \sqrt{\frac{2v}{n}} + \lambda\Gamma(S, w_{-S})\right)^{2}$$

with

$$\underbrace{\Gamma^2(q_S, w_{-S})}_{\text{"effective sparsity"}} \leq C_2^2 \sum_{j=1}^{s+1} \frac{n \log(n_j)}{n_j^3}.$$

◆□▶ ◆□▶ ◆三▶ ◆三▶ 三三 のへで

Example total variation penalty on higher order differences Fix $k \in \mathbb{N}$.

$$\hat{f} = \arg\min_{f \in \mathbb{R}^n} \left\{ \|Y - f\|_2^2 / n + 2\lambda \|Df\|_1 \right\}$$

with

$$\begin{aligned} \|Df\|_1 &:= \|\Delta^k f\|_1 \\ &:= \operatorname{TV}(\Delta^{k-1} f). \end{aligned}$$

くりょう 小田 マイビット 日 うくの

Let - $S := \{t_1, \dots, t_s\}, t_1 < \dots < t_s, t_0 := N, t_{s+1} := n+1,$ - $n_j := t_j - t_{j-1}, j = 1, \dots, s+1,$ - $n_{\max} := \max_{1 < j < s+1} n_j$

$$\mathcal{V}_{N} := \mathcal{N}_{-S} \cup \{\psi_{t_{j}+1}, \dots, \psi_{t_{j}+k-1}\}_{j=1}^{s}$$
$$N = r_{S} + (k-1)s = k(s+1)$$
$$\gamma_{S}(N, \Psi) \le An_{\max}^{\frac{2k-1}{2}}/\sqrt{n} \quad \left(1/W = \frac{2k-1}{2}\right)$$

▲□▶ ▲□▶ ▲ 三▶ ▲ 三▶ - 三■ - のへぐ

Application of the theorem to the total variation penalty on $\Delta^{k-1}f$: Theorem For all u > 0, v > 0 and for

$$\lambda \geq c_k n^{k-1} \lambda_0(u) \sqrt{n_{\max}^{2k-1}/n}$$

we have with probability at least $1 - \exp[-u] - \exp[-v]$

$$\|\hat{f} - f^{0}\|_{2}^{2}/n \le \|f^{*} - f^{0}\|_{2}^{2}/n$$
$$+ \left(\sqrt{\frac{k(s+1)}{n}} + \sqrt{\frac{2v}{n}} + \lambda\Gamma(q_{S}, w_{-S})\right)^{2}$$

with



◆□▶ ◆□▶ ◆□▶ ◆□▶ ● ● ● ●

Interpolating vector for k = 4



・ロト・四ト・モート ヨー うへの

Corollary If the jumps of $\Delta^{k-1} f^0$ are "roughly" equidistant the rate of convergence is

$$\|\hat{f} - f^0\|_2^2/n = \mathcal{O}_{\mathbb{P}}\left(\frac{(\mathbf{s}_0 + 1)\log^2 n}{n}\right)$$

provided

$$\lambda \asymp n^{k-1} \sqrt{\frac{\log n}{n}} \left(\frac{1}{\frac{s_0}{s_0}+1}\right)^{\frac{2k-1}{2}}.$$

▲□▶ ▲□▶ ▲ 三▶ ▲ 三▶ - 三 - のへぐ

Example Total variation in higher dimensions [Fang, Guntuboyina & Sen; 2019]

Let

 $f \in \mathbb{R}^{n_1 \times n_2}$ (a matrix)

and

$$\|D_1 f D_2^{\top}\|_1 := \sum_{j=2}^{n_1} \sum_{k=2}^{n_2} |f_{j,k} - f_{j-1,k} - f_{j,k-1} + f_{j-1,k-1}|$$

Hardy-Krause variation

◆□ ▶ ◆□ ▶ ◆ □ ▶ ◆ □ ▶ ● ● ● ● ●

Useful for pictures ····









◆□▶ ◆□▶ ◆臣▶ ◆臣▶ ─臣 ─のへ⊙



Effective sparsity

Let $\{t_1, \ldots, t_s\} \subset \mathbb{R}^2$ be the locations of the jumps. Construct a tiling of the rectangle $[1 : n_1] \times [1 : n_2]$ into *s* sub-rectangles, such that each sub-rectangle contains one jump.



・ロットロット キョン キョン・ヨー のへで

Interpolating vector q







◆□▶ ◆□▶ ◆□▶ ◆□▶ ●□ ● ● ●

Let $(n_j^{--}, n_j^{+-}, n_j^{-+}, n_j^{++})$ contain the areas between t_j and the corners of its sub-rectangle, j = 1, ..., s. Then

$$\Gamma^2(q_{\mathcal{S}}, w_{-\mathcal{S}}) \leq n \|D_1^{\top} q D_2\|_2^2$$

$$\leq \sum_{j=1}^{s} \left(\frac{n}{n_{j}^{--}} + \frac{n}{n_{j}^{+-}} + \frac{n}{n_{j}^{-+}} + \frac{n}{n_{j}^{++}} \right) \left(\log(1+n_{1}) + \log(1+n_{2}) \right)$$

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

The areas
$$(n_j^{--}, n_j^{+-}, n_j^{-+}, n_j^{++}), j = 1, \dots, s.$$





(日)

Application of the theorem to the 2dim total variation penalty Theorem For all u > 0, v > 0 and for

$$\lambda \ge 2\sqrt{\frac{2(\log(2n)+2u)}{n}}\sqrt{n_{1,\max}/n_1+n_{2,\max}/n_2}$$

we have with probability at least $1 - \exp[-u] - \exp[-v]$

$$\|\hat{f} - f^{0}\|_{2}^{2}/n \le \|f^{*} - f^{0}\|_{2}^{2}/n$$
$$+ \left(\sqrt{\frac{(s+1)}{n}} + \sqrt{\frac{2v}{n}} + \lambda\Gamma(q_{S}, w_{-S})\right)^{2}$$

with



$$\leq \sum_{j=1}^{s} \left(\frac{n}{n_{j}^{--}} + \frac{n}{n_{j}^{+-}} + \frac{n}{n_{j}^{-+}} + \frac{n}{n_{j}^{++}} \right) \left(\log(1+n_{1}) + \log(1+n_{2}) \right)$$

Corollary Suppose the jumps are roughly on a regular grid. Then the rate of convergence is

$$\|\hat{f} - f^0\|_2^2 / n = \mathcal{O}_{\mathbb{P}}\left(\frac{(s_0 + 1)\log^2 n}{n}\right) \sqrt{s_0 + 1}$$

provided

$$\lambda \asymp \sqrt{\frac{\log n}{n}} \left(\frac{1}{s_0 + 1}\right)^{\frac{1}{4}}$$

Remark The extra factor is

$$\frac{\text{distance in } \mathbb{R}^d}{\text{area in } \mathbb{R}^d} = \frac{\left(\frac{1}{s_0+1}\right)^{\frac{1}{d}}}{\frac{1}{s_0+1}}$$
$$= \begin{cases} 1 & \text{for } d = 1\\ \sqrt{s_0+1} & \text{for } d = 2 \end{cases}$$

▲□▶ ▲□▶ ▲ 三▶ ▲ 三▶ - 三 - のへぐ

Example: Total variation in dimension d



Application of the theorem ...

Corollary Suppose the jumps are roughly on a regular grid. Then the rate of convergence is

$$\|\hat{f} - f^0\|_2^2/n = \mathcal{O}_{\mathbb{P}}\left(\frac{(s_0+1)\log^2 n}{n}\right)(s_0+1)^{1-\frac{1}{d}}$$

provided

$$\lambda \asymp \sqrt{\frac{\log n}{n}} \left(\frac{1}{s_0 + 1}\right)^{\frac{1}{2d}}$$

◆□▶ ◆□▶ ◆□▶ ◆□▶ ● ● ● ●

remark For mesh grid

$$\|\hat{f} - f^0\|_2^2 / n = \mathcal{O}_{\mathbb{P}}\left(\frac{(s_0 + 1)\log^2 n}{n}\right) (s_0 + 1)^{1 - \frac{1}{1 + \frac{1}{2} + \dots + \frac{1}{d}}}$$

Example: MARS in dimension *d* Application of the theorem ...

Example Total variation on graphs

(ロ) (同) (三) (三) (三) (○) (○)

- $\mathcal{G}:=(\textbf{V},\mathcal{E})$ be a connected graph,
- V := set of vertices
- $n := |\mathbf{V}| \#$ vertices
- ${\mathcal E}$ edges of of ${\mathcal G}.$

Signal

$$f^{\mathsf{0}} = \{f_{\mathsf{v}}^{\mathsf{0}}\}_{\mathsf{v}\in\mathsf{V}} \in \mathbb{R}^{n}.$$

Observations

$$Y_{\boldsymbol{\nu}} = f_{\boldsymbol{\nu}}^{0} + \epsilon_{\boldsymbol{\nu}}, \ \boldsymbol{\nu} \in \boldsymbol{\mathsf{V}}, \ \{\epsilon_{\boldsymbol{\nu}}\}_{\boldsymbol{\nu} \in \boldsymbol{\mathsf{V}}} \text{ i.i.d. } \mathcal{N}(0,1).$$

Total variation

$$TV(f) := \sum_{v \sim v'} |f_v - f_{v'}|, \ f = \{f_v\}_{v \in V}$$

Analysis estimator

$$\hat{f} := \arg\min_{f} \left\{ \sum_{\nu \in \mathbf{V}} (Y_{\nu} - f_{\nu})^2 / n + 2\lambda \mathrm{TV}(f) \right\}.$$

▲□▶ ▲□▶ ▲□▶ ▲□▶ ▲□ ● のへぐ

We have considered the path graph

$\bullet \bullet \bullet \bullet \bullet \bullet \bullet$

◆□▶ ◆□▶ ◆ □▶ ◆ □▶ ● □ ● ● ● ●

At each node v is attached a value f_v (v = 1, ..., n)



▲□▶▲圖▶▲≣▶▲≣▶ ■ のへで

We now consider more general graphs



For example cycles



(口)

By cutting it we get back the path graph!
Effective sparsity without noise for graphs

The interpolating vector *q* for the cycle



ヘロト 人間 とくほとくほとう

э

The interpolating vector is as for the path graph

For general graphs, we first construct a generating tree and then split up the tree into path graphs. Both steps depend on where the jumps are.



A 目 > A 目 > A 目 > A

Tuning and noise for graphs

<□ > < @ > < E > < E > E のQ @

We now have $f \in \mathcal{N}_{-S}$ $\Leftrightarrow f$ is constant on certain connected sub-graphs $(f_{\mathcal{N}_{-S}})_{\nu} = \text{local average}$ within the sub-graph , $\nu \in \text{sub-graph}$

 $\sim \rightarrow$

to find the dictionary $\{\psi_e\}_{e\notin S}$ we do a local version of centring *f*.



▲□▶ ▲□▶ ▲□▶ ▲□▶ ▲□ ● ● ●

Global dictionary

Notation: $v \rightsquigarrow v' \subset \mathcal{E}$ path from v to v'

$$n|\psi_{\mathbf{v},\mathbf{e}}^{\text{global}}| := \left| \left\{ \mathbf{v}' \in \mathcal{V} : \mathbf{e} \in \mathbf{v} \rightsquigarrow \mathbf{v}' \right\} \right|$$
$$= \# \text{ paths starting at } \mathbf{v}$$
that contain the edge \mathbf{e}

$$v \in \mathbf{V}, \ \boldsymbol{e} \in \mathcal{E}.$$

Let

$$\overline{f} := \frac{\sum_{v' \in \mathcal{V}} f_{v'}}{n}$$

Then

$$f_{V}-\bar{f}=\frac{\sum_{v'\in\mathbf{V}}(f_{V}-f_{V'})}{n}$$

Define now for all $\textbf{\textit{e}}:=\textbf{\textit{u}}\sim \textbf{\textit{u}}'\in \mathcal{E}$,

$$b_{\boldsymbol{\theta}} := f_{\boldsymbol{U}} - f_{\boldsymbol{U}'}.$$

Then

$$f_{\mathcal{V}}-f_{\mathcal{V}'}=\sum_{e\in\mathcal{V}\rightsquigarrow\mathcal{V}'}b_e.$$

and so

$$f_{\mathbf{V}} - \overline{f} = \sum_{\mathbf{e} \in \mathcal{E}} \psi_{\mathbf{V},\mathbf{e}}^{\text{global}} \mathbf{b}_{\mathbf{e}}.$$

◆□▶ ◆□▶ ◆ □▶ ◆ □▶ ● □ ● ● ●

 \rightsquigarrow global noise weights

$$\|\psi^{\mathrm{global}}_{\boldsymbol{e}}\|_{2}^{2} = \sum_{\boldsymbol{v}\in\boldsymbol{V}} (\psi^{\mathrm{global}}_{\boldsymbol{v},\boldsymbol{e}})^{2}, \ \boldsymbol{e}\in\mathcal{E}.$$

▲□▶ ▲□▶ ▲□▶ ▲□▶ ▲□ ● のへぐ

Now to choose the paths in such a way that this is small is "small".

Example Total variation in 2d bis

[Sadhalana, Wang & Tibshirani; 2016], [Sadhalana, Wang, Sharpnack & Tibshirani; 2017]

< □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □

$$\|Df\|_{1} := \sum_{j=1}^{n_{1}} \sum_{k=1}^{n_{2}} \left\{ |f_{j,k} - f_{j-1,k}| + |f_{j,k} - f_{j,k-1}| \right\}$$



"dimension" (
$$f^0$$
)=2 $\stackrel{?}{\Rightarrow}$ $\|\hat{f} - f^0\|_2^2/n \asymp \frac{2}{n}$?
jumps = \sqrt{n} $\stackrel{?}{\Rightarrow}$ $\|\hat{f} - f^0\|_2^2/n \asymp \frac{\sqrt{n}}{n}$?



◆□▶ ◆□▶ ◆□▶ ◆□▶ ●□ ● ● ●

As paths from v to v' we take grid lines

$$x \mapsto y = \lfloor ax + c \rfloor$$

For example, for v = (j, k) and v' = (j', k') with $j' \le j$, $k' \le k$ we let take the path following the grid line $x \mapsto y$ with

$$k-y = \left\lfloor \left(\frac{k-k'}{j-j'+c}(j-x+c) \right) \right\rfloor.$$



Count the number of paths that contain the horizontal/vertical edge at the blue point

Lemma Let
$$\mathbf{v} = (j, k)$$
, $\mathbf{e} = (j'', k'') \sim (j'' - 1, k'')$. Then
$$n\psi_{\mathbf{v}, \mathbf{e}}^{\text{global}} \leq (j + \mathbf{c}) \left(\frac{j''}{j - j''} \wedge \frac{k''}{k - k''}\right)$$

$$n^{2} \|\psi_{v,e}^{\text{global}}\|^{2} \leq 2j'' k'' (n_{1} - j'' + 1)(n_{1} - j'' + 2c) \\ + 4j''^{3} k'' \left(\log\left(1 + \frac{n_{1} - j''}{c}\right) + \frac{1}{c} \right)$$

◆□ ▶ ◆□ ▶ ◆三 ▶ ◆□ ▶ ◆□ ▶

Corollary Taking c = 1 gives

$$\max_{e \notin S} \|\psi_e\|_2^2 \asymp \log n.$$

 $\sim \rightarrow$

tuning parameter

 $\log n/n \lesssim \lambda$

effective sparsity

$$\Gamma^2(q_{\mathcal{S}},1) symp n\sqrt{n}$$

and the (minimax [Hütter & Rigollet, 2016]) rate

$$\|\hat{f}-f^0\|_2^2/n=\mathcal{O}_{\mathbb{P}}\bigg(\frac{\log^2 n}{\sqrt{n}}\bigg).$$

▲□▶ ▲□▶ ▲ 三▶ ▲ 三▶ - 三 - のへぐ

Corollary Taking $c = \sqrt{n}$ gives

$$\max_{e \notin S} \|\psi_e\|_2^2 \asymp n^{1/4} \log n$$

and similarly for the vertical edges.

 \rightsquigarrow

tuning parameter

$$n^{-rac{7}{8}}\log n\lesssim\lambda$$

effective sparsity

$$\Gamma^2(q_S, w_{-S}) symp n \log n$$

and the (adaptive) rate

$$\|\hat{f} - f^0\|_2^2/n = \mathcal{O}_{\mathbb{P}}\left(\frac{\log^3 n}{n^{3/4}}\right).$$

as in [Chatterjee and Goswami, 2019]

< □ > < 同 > < 三 > < 三 > < 三 > < ○ < ○ </p>

Conclusion

- \circ "Lasso" can also work for highly correlated design
- Overparametrization vs entropy ...
- o interpolating vectors useful for structured problems
 → improvement of restricted eigenvalue conditions

◆□▶ ◆□▶ ▲□▶ ▲□▶ □ のQ@

- TV-type penalties are adaptive
- Remark ∃ extensions to other loss functions (e.g. logistic loss)



▲□▶ ▲□▶ ▲□▶ ▲□▶ = 三 のへで