## **Can Uniform Convergence Explain Interpolation Learning?**

#### Danica J. Sutherland (she/her)

University of British Columbia (UBC) / Alberta Machine Intelligence Institute (Amii)

#### based on arXiv:2006.05942 and 2106.09276, with:

UChicago

Lijia Zhou Frederic Koehler Nati Srebro  $MIT \rightarrow Simons Institute$ 

**TTI-Chicago** 







NYLL Center for Data Science – November 10 2021 The HTML version is the "official" version, though this PDF is basically the same.

• Given i.i.d. samples  $\mathbf{S} = \{(\mathbf{x}_i, y_i)\}_{i=1}^n \sim \mathcal{D}^n$ 

• features/covariates  $\mathbf{x}_i \in \mathbb{R}^d$ , labels/targets  $y_i \in \mathbb{R}$ 

• Given i.i.d. samples  $\mathbf{S} = \{(\mathbf{x}_i, y_i)\}_{i=1}^n \sim \mathcal{D}^n$ 

• features/covariates  $\mathbf{x}_i \in \mathbb{R}^d$ , labels/targets  $y_i \in \mathbb{R}$ 

• Want f such that  $f(\mathbf{x}) pprox y$  for new samples from  $\mathcal D$ 

• Given i.i.d. samples  $\mathbf{S} = \{(\mathbf{x}_i, y_i)\}_{i=1}^n \sim \mathcal{D}^n$ 

• features/covariates  $\mathbf{x}_i \in \mathbb{R}^d$ , labels/targets  $y_i \in \mathbb{R}$ 

• Want f such that  $f(\mathbf{x}) pprox y$  for new samples from  $\mathcal{D}$ :

$$f^* = rgmin igg[ L_{\mathcal{D}}(f) := \mathop{\mathbb{E}}_{(\mathbf{x}, \mathrm{y}) \sim \mathcal{D}} L(f(\mathbf{x}), \mathrm{y}) igg]$$

• Given i.i.d. samples  $\mathbf{S} = \{(\mathbf{x}_i, y_i)\}_{i=1}^n \sim \mathcal{D}^n$ 

- features/covariates  $\mathbf{x}_i \in \mathbb{R}^d$ , labels/targets  $y_i \in \mathbb{R}$ 

• Want f such that  $f(\mathbf{x}) pprox y$  for new samples from  $\mathcal{D}$ :

$$f^* = rgmin igg[ L_{\mathcal{D}}(f) := \mathop{\mathbb{E}}_{(\mathbf{x}, \mathrm{y}) \sim \mathcal{D}} L(f(\mathbf{x}), \mathrm{y}) igg]$$

• e.g. squared loss:  $L(\hat{y},y)=(\hat{y}-y)^2$ 

• Given i.i.d. samples  $\mathbf{S} = \{(\mathbf{x}_i, y_i)\}_{i=1}^n \sim \mathcal{D}^n$ 

- features/covariates  $\mathbf{x}_i \in \mathbb{R}^d$ , labels/targets  $y_i \in \mathbb{R}$ 

• Want f such that  $f(\mathbf{x}) pprox y$  for new samples from  $\mathcal{D}$ :

$$f^* = rgmin igg[ L_{\mathcal{D}}(f) := \mathop{\mathbb{E}}_{(\mathbf{x}, \mathrm{y}) \sim \mathcal{D}} L(f(\mathbf{x}), \mathrm{y}) igg]$$

• e.g. squared loss:  $L(\hat{y},y)=(\hat{y}-y)^2$ 

• Standard approaches based on empirical risk minimization:

$$\hat{f} pprox rgmin \left[ L_{\mathbf{S}}(f) := rac{1}{n} \sum_{i=1}^n L(f(\mathbf{x}_i), y_i) 
ight]$$

We have lots of bounds like: with probability  $\geq 1-\delta$ ,

$$\sup_{f\in\mathcal{F}}\left|L_{\mathcal{D}}(f)-L_{\mathbf{S}}(f)
ight|\leq\sqrt{rac{C_{\mathcal{F},\delta}}{n}}$$

We have lots of bounds like: with probability  $\geq 1-\delta$ ,

$$\sup_{f\in\mathcal{F}}\left|L_{\mathcal{D}}(f)-L_{\mathbf{S}}(f)
ight|\leq\sqrt{rac{C_{\mathcal{F},\delta}}{n}}$$

 $C_{\mathcal{F},\delta}$  could be from VC dimension, covering number, RKHS norm, Rademacher complexity, fat-shattering dimension, ...

We have lots of bounds like: with probability  $\geq 1-\delta$ ,

$$\sup_{f\in\mathcal{F}}\left|L_{\mathcal{D}}(f)-L_{\mathbf{S}}(f)
ight|\leq\sqrt{rac{C_{\mathcal{F},\delta}}{n}}$$

We have lots of bounds like: with probability  $\geq 1-\delta$ ,

$$\sup_{f\in\mathcal{F}}\left|L_{\mathcal{D}}(f)-L_{\mathbf{S}}(f)
ight|\leq\sqrt{rac{C_{\mathcal{F},\delta}}{n}}$$

Then for large n,  $L_{\mathcal{D}}(f)pprox L_{\mathbf{S}}(f)$ , so  $\hat{f}pprox f^{*}$ 

We have lots of bounds like: with probability  $\geq 1-\delta$ ,

$$\sup_{f\in\mathcal{F}}\left|L_{\mathcal{D}}(f)-L_{\mathbf{S}}(f)
ight|\leq\sqrt{rac{C_{\mathcal{F},\delta}}{n}}$$

Then for large n,  $L_{\mathcal{D}}(f)pprox L_{\mathbf{S}}(f)$ , so  $\hat{f}\,pprox f^*$ 

$$L_{\mathcal{D}}(\widehat{f}\,) \leq L_{\mathbf{S}}(\widehat{f}\,) + \sup_{f \in \mathcal{F}} \left| L_{\mathcal{D}}(f) - L_{\mathbf{S}}(f) 
ight|$$

Classical wisdom: "a model with zero training error is overfit to the training data and will typically generalize poorly"

Springer Series in Statistics					
Trevor Hastie Robert Tibshirani Jerome Friedman					
The Elements of Statistical Learning					

cond Edition

🖄 Springer

Classical wisdom: "a model with zero training error is overfit to the training data and will typically generalize poorly"

Zhang et



Table 1: The training and test accuracy (in percentage) of various models on the CIFAR10 dataset.

model	# params	random crop	weight decay	train accuracy	test accuracy
Inception 1,649		yes	yes	100.0	89.05
	1 640 402	yes	no	100.0	89.31
	1,049,402	no	yes	100.0	86.03
		no	no	100.0	85.75
l., "Rethinking	generalizatio	n", ICLR 2017	$L_{\mathbf{S}}(\hat{f})$ :	$= 0; L_{\mathcal{D}}($	$(\hat{f})pprox 11$

Classical wisdom: "a model with zero training error is overfit to the training data and will typically generalize poorly"

Zhang et



Table 1: The training and test accuracy (in percentage) of various models on the CIFAR10 dataset.

model	# params	random crop	weight decay	train accuracy	test accuracy
Inception		yes	yes	100.0	89.05
	1 640 402	yes	no	100.0	89.31
	1,049,402	no	yes	100.0	86.03
		no	no	100.0	85.75
l., "Rethinking	generalizatio	n", ICLR 2017	$L_{\mathbf{S}}(\hat{f})$ :	$= 0; L_{\mathcal{D}}($	$(\hat{f})pprox 11$

We'll call a model with  $L_{f S}(f)=0$  an *interpolating* predictor

Classical wisdom: "a model with zero training error is overfit to the training data and will typically generalize poorly" (when  $L_{\mathcal{D}}(f^*) > 0$ )



Table 1: The training and test accuracy (in percentage) of various models on the CIFAR10 dataset.

model	# params	random crop	weight decay	train accuracy	test accuracy
Inception	1,649,402	yes yes no	yes no yes	100.0 100.0 100.0	89.05 89.31 86.03
ang et al "Rethinking	generalizatio	<b>no</b> n″ ICLR 2017	$L_{\mathbf{G}}(\hat{f})$	$= 0 \cdot L_{\mathcal{D}}$	$(\hat{f}) \approx 1^{1}$

We'll call a model with  $L_{f S}(f)=0$  an *interpolating* predictor

## Interpolation does not overfit even for very noisy data

All methods (except Bayes optimal) have zero training square loss.



Belkin/Ma/Mandal, ICML 2018



Belkin/Ma/Mandal, ICML 2018

correct 
$$\sqrt{\frac{C_{\mathcal{F},\delta}}{n}}$$
 nontrivial  $n \to \infty$ 

# There are no bounds like this and no reason they should exist.

A constant factor of 2 invalidates the bound!



Misha Belkin

Simons Institute

July 2019

#### Generalization theory for interpolation?

What theoretical analyses do we have?

- VC-dimension/Rademacher complexity/covering/margin bounds.
  - Cannot deal with interpolated classifiers when Bayes risk is non-zero.
  - > Generalization gap cannot be bound when empirical risk is zero.
- Regularization-type analyses Tikhonov, early stopping/SGD, etc.)
  - Diverge as  $\lambda \to 0$  for fixed n.
- Algorithmic stability.
  - > Does not apply when empirical risk is zero, expected risk nonzero.
- Classical smoothing methods (i.e., Nadaraya-Watson).
  - Most classical analyses do not support interpolation.
  - But 1-NN! (Also Hilbert regression Scheme, [Devroye, et al. 98])

 $\neg \, L_{\mathcal{D}}(\hat{f}) \leq L_{\mathbf{S}}(\hat{f}) \! + \! \mathrm{bound}$ 

WYSIWYG bounds:

expected loss

Misha Belkin Simons Institute <sub>Oracle bounds</sub> July 2019

expected loss  $\approx$ optimal loss

 $L_{\mathcal{D}}(\hat{f}) \leq L_{\mathcal{D}}(f^*) + ext{bound}$ 



What theoretical analyses do we have?

#### Lots of recent theoretical work on interpolation.

[Belkin+ NeurIPS 2018], [Belkin+ AISTATS 2018], [Belkin+ 2019], [Hastie+ 2019],

[Muthukumar+ JSAIT 2020], [Bartlett+ PNAS 2020], [Liang+ COLT 2020], [Montanari+ 2019], many more...

## None\* bound $\sup_{f\in\mathcal{F}}|L_{\mathcal{D}}(f)-L_{\mathbf{S}}(f)|.$

Is it possible to find such a bound?

Can uniform convergence explain interpolation learning?

But 1-NN! (Also Hilbert regression Scheme, [Devroye, et al. 98])

optimal loss

 $L_{\mathcal{D}}(\hat{f}) \leq L_{\mathcal{D}}(f^*) + \text{bound}$ 



What theoretical analyses do we have?



Today, we're mainly going to worry about *consistency*:

$$\mathbb{E}[L_{\mathcal{D}}(\hat{f}) - L_{\mathcal{D}}(f^*)] o 0$$

Today, we're mainly going to worry about *consistency*:

$$\mathbb{E}[L_{\mathcal{D}}(\hat{f})-L_{\mathcal{D}}(f^*)] o 0$$

...in a *noisy* setting:  $L_{\mathcal{D}}(f^*) > 0$ 

Today, we're mainly going to worry about *consistency*:

$$\mathbb{E}[L_{\mathcal{D}}(\hat{f}) - L_{\mathcal{D}}(f^*)] o 0$$

...in a *noisy* setting:  $L_{\mathcal{D}}(f^*) > 0$ 

...for Gaussian linear regression:

$$\mathbf{x} \sim \mathcal{N}(\mathbf{0}, \Sigma) \quad y = \langle \mathbf{x}, w^* 
angle + \mathcal{N}(\mathbf{0}, \sigma^2) \quad L(y, \hat{y}) = (y - \hat{y})^2$$

Today, we're mainly going to worry about *consistency*:

$$\mathbb{E}[L_{\mathcal{D}}(\hat{f}\,)-L_{\mathcal{D}}(f^*)] o 0$$

...in a *noisy* setting:  $L_{\mathcal{D}}(f^*) > 0$ 

... for Gaussian linear regression:

$$\mathbf{x} \sim \mathcal{N}(\mathbf{0}, \Sigma) \quad y = \langle \mathbf{x}, w^* 
angle + \mathcal{N}(\mathbf{0}, \sigma^2) \quad L(y, \hat{y}) = (y - \hat{y})^2$$

Is it possible to show consistency of an interpolator with

$$L_{\mathcal{D}}(\widehat{f}\,) \leq \underbrace{L_{\mathbf{S}}(\widehat{f}\,)}_{0} + \sup_{f \in \mathcal{F}} |L_{\mathcal{D}}(f) - L_{\mathbf{S}}(f)|?$$

Today, we're mainly going to worry about *consistency*:

$$\mathbb{E}[L_{\mathcal{D}}(\hat{f})-L_{\mathcal{D}}(f^*)] o 0$$

...in a *noisy* setting:  $L_{\mathcal{D}}(f^*) > 0$ 

... for Gaussian linear regression:

$$\mathbf{x} \sim \mathcal{N}(\mathbf{0}, \Sigma) \quad y = \langle \mathbf{x}, w^* 
angle + \mathcal{N}(\mathbf{0}, \sigma^2) \quad L(y, \hat{y}) = (y - \hat{y})^2$$

Is it possible to show consistency of an interpolator with

$$L_{\mathcal{D}}(\hat{f}) \leq \underbrace{L_{\mathbf{S}}(\hat{f})}_{f \in \mathcal{F}} + \sup_{f \in \mathcal{F}} |L_{\mathcal{D}}(f) - L_{\mathbf{S}}(f)|?$$
  
This requires tight constants!

## A testbed problem: "junk features"



 $\lambda_n$  controls scale of junk:  $\mathbb{E}\|\mathbf{x}_J\|_2^2 = \lambda_n$ Linear regression:  $L(y, \hat{y}) = (y - \hat{y})^2$ 

## A testbed problem: "junk features"

$$\begin{array}{c|c} \text{"signal", } d_S & \text{"junk", } d_J \to \infty \\ \mathbf{x} & \mathbf{x}_S \sim \mathcal{N}\left(\mathbf{0}_{d_S}, \mathbf{I}_{d_S}\right) & \mathbf{x}_J \sim \mathcal{N}\left(\mathbf{0}_{d_J}, \frac{\lambda_n}{d_J} \mathbf{I}_{d_J}\right) \\ \mathbf{w}^* & \mathbf{w}_S^* & \mathbf{0} \\ & y = \langle \mathbf{x}, \mathbf{w}^* \rangle + \mathcal{N}(\mathbf{0}, \sigma^2) \\ & \swarrow \mathbf{x}_S, \mathbf{w}_S^* \rangle \end{array}$$

 $\lambda_n$  controls scale of junk:  $\mathbb{E} \| \mathbf{x}_J \|_2^2 = \lambda_n$ 

Linear regression: 
$$L(y, \hat{y}) = (y - \hat{y})^2$$

Min-norm interpolator:  $\hat{\mathbf{w}}_{MN} = \operatorname*{arg\,min}_{\mathbf{X}\mathbf{w}=\mathbf{y}} \|\mathbf{w}\|_2 = \mathbf{X}^\dagger \mathbf{y}$ 

## Consistency of $\hat{\mathbf{w}}_{MN}$ $\hat{\mathbf{w}}_{MN} = \underset{\mathbf{X}\mathbf{w}=\mathbf{y}}{\operatorname{arg\,min}} \|\mathbf{w}\|_2 = \mathbf{X}^{\dagger} \mathbf{y}$

As  $d_J o\infty$ ,  $\hat{\mathbf{w}}_{MN}$  approaches ridge regression on the signal  $\hat{\mathbf{w}}_{MN}$  is consistent when  $d_S$  fixed,  $d_J o\infty$ ,  $\lambda_n=o(n)$ 

## Consistency of $\hat{\mathbf{w}}_{MN}$ $\hat{\mathbf{w}}_{MN} = \underset{\mathbf{X}\mathbf{w}=\mathbf{y}}{\operatorname{arg\,min}} \|\mathbf{w}\|_2 = \mathbf{X}^{\dagger} \mathbf{y}$

As  $d_J o\infty$ ,  $\hat{\mathbf{w}}_{MN}$  approaches ridge regression on the signal  $\hat{\mathbf{w}}_{MN}$  is consistent when  $d_S$  fixed,  $d_J o\infty$ ,  $\lambda_n=o(n)$ 

Could we have shown that with uniform convergence?

## Consistency of $\hat{\mathbf{w}}_{MN}$ $\hat{\mathbf{w}}_{MN} = \underset{\mathbf{X}\mathbf{w}=\mathbf{y}}{\operatorname{arg\,min}} \|\mathbf{w}\|_2 = \mathbf{X}^{\dagger} \mathbf{y}$

As  $d_J o \infty$ ,  $\hat{\mathbf{w}}_{MN}$  approaches ridge regression on the signal  $\hat{\mathbf{w}}_{MN}$  is consistent when  $d_S$  fixed,  $d_J o \infty$ ,  $\lambda_n = o(n)$ 

Could we have shown that with uniform convergence?

### No uniform convergence on norm balls

<u>Theorem:</u> In junk features with  $\lambda_n = o(n)$ ,

$$\lim_{n o \infty} \lim_{d_J o \infty} \mathbb{E} \left[ \sup_{\|\mathbf{w}\|_2 \leq \|\hat{\mathbf{w}}_{MN}\|_2} |L_{\mathcal{D}}(\mathbf{w}) - L_{\mathbf{S}}(\mathbf{w})| 
ight] = \infty.$$

### No uniform convergence on norm balls

<u>Theorem:</u> In junk features with  $\lambda_n = o(n)$ ,

$$\lim_{n o \infty} \lim_{d_J o \infty} \mathbb{E} \left[ \sup_{\|\mathbf{w}\|_2 \leq \|\hat{\mathbf{w}}_{MN}\|_2} |L_{\mathcal{D}}(\mathbf{w}) - L_{\mathbf{S}}(\mathbf{w})| 
ight] = \infty.$$

### No uniform convergence on norm balls

<u>Theorem:</u> In junk features with  $\lambda_n = o(n)$ ,

$$\lim_{n o \infty} \lim_{d_J o \infty} \mathbb{E} \left[ \sup_{\|\mathbf{w}\|_2 \leq \|\hat{\mathbf{w}}_{MN}\|_2} |L_\mathcal{D}(\mathbf{w}) - L_\mathbf{S}(\mathbf{w})| 
ight] = \infty.$$

## A more refined uniform convergence analysis?

 $\{\mathbf{w}: \|\mathbf{w}\| \leq B\}$  is no good.

## A more refined uniform convergence analysis?

 $\{\mathbf{w}: \|\mathbf{w}\| \leq B\}$  is no good. Maybe  $\{\mathbf{w}: A \leq \|\mathbf{w}\| \leq B\}$ ?
$\{\mathbf{w}: \|\mathbf{w}\| \leq B\}$  is no good. Maybe  $\{\mathbf{w}: A \leq \|\mathbf{w}\| \leq B\}$ ?

# Uniform convergence may be unable to explain generalization in deep learning

Vaishnavh Nagarajan Department of Computer Science Carnegie Mellon University Pittsburgh, PA vaishnavh@cs.cmu.edu J. Zico Kolter Department of Computer Science Carnegie Mellon University & Bosch Center for Artificial Intelligence Pittsburgh, PA zkolter@cs.cmu.edu

 $\frac{\text{Theorem}}{\text{Theorem}} \text{ (à la [Nagarajan/Kolter, NeurIPS 2019]):} \\ \text{In junk features, for each } \delta \in (0, \frac{1}{2}), \text{ let } \Pr\left(\mathbf{S} \in \mathcal{S}_{n,\delta}\right) \geq 1 - \delta, \\ \end{array}$ 

 $\frac{\text{Theorem}}{\text{Theorem}} \text{ (à la [Nagarajan/Kolter, NeurIPS 2019]):}$   $\text{In junk features, for each } \delta \in (0, \frac{1}{2}), \text{ let } \Pr(\mathbf{S} \in \mathcal{S}_{n,\delta}) \geq 1 - \delta,$   $\hat{\mathbf{w}} \text{ a natural consistent interpolator,}$ 

Natural interpolators:  $\hat{\mathbf{w}}_S$  doesn't change if  $\mathbf{X}_J$  flips to  $-\mathbf{X}_J$ . Examples:  $\hat{\mathbf{w}}_{MN}$ ,  $\underset{\mathbf{w}:\mathbf{X}\mathbf{w}=\mathbf{y}}{\operatorname{arg\,min}} \|\mathbf{w}\|_1$ ,  $\underset{\mathbf{w}:\mathbf{X}\mathbf{w}=\mathbf{y}}{\operatorname{arg\,min}} \|\mathbf{w} - \mathbf{w}^*\|_2$ ,  $\underset{\mathbf{w}:\mathbf{X}\mathbf{w}=\mathbf{y}}{\operatorname{arg\,min}} f_S(\mathbf{w}_S) + f_J(\mathbf{w}_J)$  with each f convex,  $f_J(-\mathbf{w}_J) = f_J(\mathbf{w}_J)$  $\mathbf{w}:\mathbf{X}\mathbf{w}=\mathbf{y}$ 

Theorem (à la [Nagarajan/Kolter, NeurIPS 2019]):

In junk features, for each  $\delta \in (0, rac{1}{2})$ , let  $\Pr{(\mathbf{S} \in \mathcal{S}_{n,\delta})} \geq 1 - \delta$ ,

 $\hat{\mathbf{w}}$  a *natural* consistent interpolator,

and  $\mathcal{W}_{n,\delta} = \{ \hat{\mathbf{w}}(\mathbf{S}) : \mathbf{S} \in \mathcal{S}_{n,\delta} \}.$ 

*Natural* interpolators: 
$$\mathbf{\hat{w}}_S$$
 doesn't change if  $\mathbf{X}_J$  flips to  $-\mathbf{X}_J$ . Examples:  
 $\mathbf{\hat{w}}_{MN}$ ,  $\underset{\mathbf{w}:\mathbf{X}\mathbf{w}=\mathbf{y}}{\operatorname{arg\,min}} \|\mathbf{w}\|_1$ ,  $\underset{\mathbf{w}:\mathbf{X}\mathbf{w}=\mathbf{y}}{\operatorname{arg\,min}} \|\mathbf{w} - \mathbf{w}^*\|_2$ ,  
 $\underset{\mathbf{w}:\mathbf{X}\mathbf{w}=\mathbf{y}}{\operatorname{arg\,min}} f_S(\mathbf{w}_S) + f_J(\mathbf{w}_J)$  with each  $f$  convex,  $f_J(-\mathbf{w}_J) = f_J(\mathbf{w}_J)$   
 $\mathbf{w}:\mathbf{X}\mathbf{w}=\mathbf{y}$ 

Theorem (à la [Nagarajan/Kolter, NeurIPS 2019]):

In junk features, for each  $\delta \in (0, rac{1}{2})$ , let  $\Pr{(\mathbf{S} \in \mathcal{S}_{n,\delta})} \geq 1 - \delta$ ,

 $\hat{\mathbf{w}}$  a *natural* consistent interpolator,

and  $\mathcal{W}_{n,\delta} = \{ \hat{\mathbf{w}}(\mathbf{S}) : \mathbf{S} \in \mathcal{S}_{n,\delta} \}$ . Then, almost surely,

Natural interpolators:  $\hat{\mathbf{w}}_S$  doesn't change if  $\mathbf{X}_J$  flips to  $-\mathbf{X}_J$ . Examples:  $\hat{\mathbf{w}}_{MN}$ ,  $\underset{\mathbf{w}:\mathbf{X}\mathbf{w}=\mathbf{y}}{\operatorname{arg\,min}} \|\mathbf{w}\|_1$ ,  $\underset{\mathbf{w}:\mathbf{X}\mathbf{w}=\mathbf{y}}{\operatorname{arg\,min}} \|\mathbf{w} - \mathbf{w}^*\|_2$ ,  $\underset{\mathbf{w}:\mathbf{X}\mathbf{w}=\mathbf{y}}{\operatorname{arg\,min}} f_S(\mathbf{w}_S) + f_J(\mathbf{w}_J)$  with each f convex,  $f_J(-\mathbf{w}_J) = f_J(\mathbf{w}_J)$  $\mathbf{w}:\mathbf{X}\mathbf{w}=\mathbf{y}$ 

Theorem (à la [Nagarajan/Kolter, NeurIPS 2019]):

In junk features, for each  $\delta \in (0, rac{1}{2})$ , let  $\Pr{(\mathbf{S} \in \mathcal{S}_{n,\delta})} \geq 1 - \delta$ ,

 $\mathbf{\hat{w}}$  a *natural* consistent interpolator,

and 
$$\mathcal{W}_{n,\delta} = \{ \hat{\mathbf{w}}(\mathbf{S}) : \mathbf{S} \in \mathcal{S}_{n,\delta} \}$$
. Then, almost surely,

$$\lim_{n o \infty} \lim_{d_J o \infty} \sup_{\mathbf{S} \in \mathcal{S}_{n,\delta}} \sup_{\mathbf{w} \in \mathcal{W}_{n,\delta}} |L_\mathcal{D}(\mathbf{w}) - L_\mathbf{S}(\mathbf{w})| \geq 3\sigma^2.$$

([Negrea/Dziugaite/Roy, ICML 2020] had a very similar result for  $\hat{\mathbf{w}}_{MN}$ )

*Natural* interpolators:  $\hat{\mathbf{w}}_S$  doesn't change if  $\mathbf{X}_J$  flips to  $-\mathbf{X}_J$ . Examples:  $\hat{\mathbf{w}}_{MN}$ ,  $\underset{\mathbf{w}:\mathbf{X}\mathbf{w}=\mathbf{y}}{\operatorname{arg\,min}} \|\mathbf{w}\|_1$ ,  $\underset{\mathbf{w}:\mathbf{X}\mathbf{w}=\mathbf{y}}{\operatorname{arg\,min}} \|\mathbf{w} - \mathbf{w}^*\|_2$ ,  $\underset{\mathbf{w}:\mathbf{X}\mathbf{w}=\mathbf{y}}{\operatorname{arg\,min}} f_S(\mathbf{w}_S) + f_J(\mathbf{w}_J)$  with each f convex,  $f_J(-\mathbf{w}_J) = f_J(\mathbf{w}_J)$  $\mathbf{w}:\mathbf{X}\mathbf{w}=\mathbf{y}$ 

 $\begin{array}{l} \underline{\text{Theorem}} \text{ (à la [Nagarajan/Kolter, NeurIPS 2019]):} \\ \text{In junk features, for each } \delta \in (0, \frac{1}{2}) \text{, let } \Pr \left( \mathbf{S} \in \mathcal{S}_{n,\delta} \right) \geq 1 - \delta, \\ \widehat{\mathbf{w}} \text{ a natural consistent interpolator,} \\ \text{ and } \mathcal{W}_{n,\delta} = \{ \widehat{\mathbf{w}}(\mathbf{S}) : \mathbf{S} \in \mathcal{S}_{n,\delta} \}. \text{ Then, almost surely,} \\ \lim_{n \to \infty} \lim_{d_J \to \infty} \sup_{\mathbf{S} \in \mathcal{S}_{n,\delta}} \sup_{\mathbf{w} \in \mathcal{W}_{n,\delta}} \left| L_{\mathcal{D}}(\mathbf{w}) - L_{\mathbf{S}}(\mathbf{w}) \right| \geq 3\sigma^2. \end{array}$ 

Proof shows that for most  $\mathbf{S}$ , there's a typical predictor  $\mathbf{w}$  (in  $\mathcal{W}_{n,\delta}$ ) that's good on most inputs ( $L_{\mathcal{D}}(\mathbf{w}) \to \sigma^2$ ), but very bad on *specifically*  $\mathbf{S}$  ( $L_{\mathbf{S}}(\mathbf{w}) \to 4\sigma^2$ )

• Convergence of surrogates [Negrea/Dziugaite/Roy, ICML 2020]?

- Convergence of surrogates [Negrea/Dziugaite/Roy, ICML 2020]?
  - Nice, but not really the same thing...

- Convergence of surrogates [Negrea/Dziugaite/Roy, ICML 2020]?
  - Nice, but not really the same thing...
- Only do analyses based on e.g. exact form of  $\hat{\mathbf{w}}_{MN}$ ?

- Convergence of surrogates [Negrea/Dziugaite/Roy, ICML 2020]?
  - Nice, but not really the same thing...
- Only do analyses based on e.g. exact form of  $\hat{\mathbf{w}}_{MN}$ ?
- We'd like to keep good things about uniform convergence:
  - Apply to more than just one specific predictor
  - Tell us more about "why" things generalize
  - Easier to apply without a nice closed form

- Convergence of surrogates [Negrea/Dziugaite/Roy, ICML 2020]?
  - Nice, but not really the same thing...
- Only do analyses based on e.g. exact form of  $\hat{\mathbf{w}}_{MN}$ ?
- We'd like to keep good things about uniform convergence:
  - Apply to more than just one specific predictor
  - Tell us more about "why" things generalize
  - Easier to apply without a nice closed form
- Or...

#### A broader view of uniform convergence

So far, used  $L_\mathcal{D}(\mathbf{w}) - L_\mathbf{S}(\mathbf{w}) \leq \sup_{\|\mathbf{w}\|_2 \leq B} |L_\mathcal{D}(\mathbf{w}) - L_\mathbf{S}(\mathbf{w})|$ 

But we only care about interpolators. How about

$$\sup_{\|\mathbf{w}\|_2 \leq B, \ L_{\mathbf{S}}(\mathbf{w}) = \mathbf{0}} |L_{\mathcal{D}}(\mathbf{w}) - L_{\mathbf{S}}(\mathbf{w})|?$$

But we only care about interpolators. How about

$$\sup_{\|\mathbf{w}\|_2 \leq B, \ \boldsymbol{L_S}(\mathbf{w}) = \mathbf{0}} |L_\mathcal{D}(\mathbf{w}) - L_\mathbf{S}(\mathbf{w})|?$$

Is this "uniform convergence"?

But we only care about interpolators. How about

$$\sup_{\|\mathbf{w}\|_2 \leq B, \; \boldsymbol{L}_{\mathbf{S}}(\mathbf{w}) = \mathbf{0}} |L_{\mathcal{D}}(\mathbf{w}) - L_{\mathbf{S}}(\mathbf{w})|?$$

Is this "uniform convergence"?

It's the standard notion for noiseless ( $L_{\mathcal{D}}(w^*)=0$ ) analyses...

# A broader view of uniform convergence

Used at least since [Vapnik 1982] and [Valiant 1984]

From [Devroye/Györfi/Lugosi 1996]:

PROOF. For  $n\epsilon \leq 2$ , the inequality is clearly true. So, we assume that  $n\epsilon > 2$ . First observe that since  $\inf_{\phi \in C} L(\phi) = 0$ ,  $\widehat{L}_n(\phi_n^*) = 0$  with probability one. It is easily seen that

$$L(\phi_n^*) \leq \sup_{\phi:\widehat{L}_n(\phi)=0} |L(\phi) - \widehat{L}_n(\phi)|.$$

It's the standard notion for noiseless ( $L_{\mathcal{D}}(w^*)=0$ ) analyses...

# <u>A broader view of uniform convergence</u>



It's the stand

In the example of axis-aligned rectangles that we examined, the hypothesis  $h_S$  returned by the algorithm was always *consistent*, that is, it admitted no error on the training sample S. In this section, we present a general sample complexity bound, or equivalently, a generalization bound, for consistent hypotheses, in the case where the cardinality |H| of the hypothesis set is finite. Since we consider consistent hypotheses, we will assume that the target concept c is in H.

Theorem 2.1 Learning bounds — finite H, consistent case Let H be a finite set of functions mapping from  $\mathcal{X}$  to  $\mathcal{Y}$ . Let  $\mathcal{A}$  be an algorithm that for any target concept  $c \in H$  and i.i.d. sample S returns a consistent hypothesis  $h_S$ :  $\widehat{R}(h_S) = 0$ . Then, for any  $\epsilon, \delta > 0$ , the inequality  $\Pr_{S \sim D^m}[R(h_S) \leq \epsilon] \geq 1 - \delta$  holds if

$$m \ge \frac{1}{\epsilon} \Big( \log |H| + \log \frac{1}{\delta} \Big). \tag{2.8}$$

This sample complexity result admits the following equivalent statement as a generalization bound: for any  $\epsilon, \delta > 0$ , with probability at least  $1 - \delta$ ,

$$R(h_S) \le \frac{1}{m} \Big( \log |H| + \log \frac{1}{\delta} \Big).$$
(2.9)

**Proof** Fix  $\epsilon > 0$ . We do not know which consistent hypothesis  $h_S \in H$  is selected by the algorithm  $\mathcal{A}$ . This hypothesis further depends on the training sample S. Therefore, we need to give a *uniform convergence bound*, that is, a bound that holds for the set of all consistent hypotheses, which a fortiori includes  $h_S$ . Thus,

But we only care about interpolators. How about

$$\sup_{\|\mathbf{w}\|_2 \leq B, \; \boldsymbol{L}_{\mathbf{S}}(\mathbf{w}) = \mathbf{0}} |L_{\mathcal{D}}(\mathbf{w}) - L_{\mathbf{S}}(\mathbf{w})|?$$

Is this "uniform convergence"?

It's the standard notion for noiseless ( $L_{\mathcal{D}}(w^*)=0$ ) analyses...

But we only care about interpolators. How about

$$\sup_{\|\mathbf{w}\|_2 \leq B, \; \boldsymbol{L}_{\mathbf{S}}(\mathbf{w}) = \mathbf{0}} |L_{\mathcal{D}}(\mathbf{w}) - L_{\mathbf{S}}(\mathbf{w})|?$$

Is this "uniform convergence"?

It's the standard notion for noiseless ( $L_{\mathcal{D}}(w^*)=0$ ) analyses...

# The interpolator ball in linear regression

What does  $\{\mathbf{w}: \|\mathbf{w}\|_2 \leq B, L_{\mathbf{S}}(\mathbf{w}) = 0\}$  look like?

# The interpolator ball in linear regression

What does  $\{\mathbf{w}: \|\mathbf{w}\|_2 \leq B, L_{\mathbf{S}}(\mathbf{w}) = 0\}$  look like?



Intersection of *d*-ball

# The interpolator ball in linear regression What does $\{\mathbf{w} : \|\mathbf{w}\|_2 \leq B, L_{\mathbf{S}}(\mathbf{w}) = 0\}$ look like? $\{\mathbf{w} : L_{\mathbf{S}}(\mathbf{w}) = \frac{1}{n} \|\mathbf{X}\mathbf{w} - \mathbf{y}\|_2^2 = 0\}$ is the plane $\mathbf{X}\mathbf{w} = \mathbf{y}$



Intersection of d-ball with (d - n)-hyperplane:

# The interpolator ball in linear regression What does $\{\mathbf{w} : \|\mathbf{w}\|_2 \leq B, L_{\mathbf{S}}(\mathbf{w}) = 0\}$ look like? $\{\mathbf{w} : L_{\mathbf{S}}(\mathbf{w}) = \frac{1}{n} \|\mathbf{X}\mathbf{w} - \mathbf{y}\|_2^2 = 0\}$ is the plane $\mathbf{X}\mathbf{w} = \mathbf{y}$



Intersection of d-ball with (d - n)-hyperplane: (d - n)-ball

# The interpolator ball in linear regression What does $\{\mathbf{w} : \|\mathbf{w}\|_2 \leq B, L_{\mathbf{S}}(\mathbf{w}) = 0\}$ look like? $\{\mathbf{w} : L_{\mathbf{S}}(\mathbf{w}) = \frac{1}{n} \|\mathbf{X}\mathbf{w} - \mathbf{y}\|_2^2 = 0\}$ is the plane $\mathbf{X}\mathbf{w} = \mathbf{y}$



Intersection of d-ball with (d - n)-hyperplane: (d - n)-ball centered at  $\hat{\mathbf{w}}_{MN}$ 

Applying [Srebro/Sridharan/Tewari 2010]: for all  $\|\mathbf{w}\|_2 \leq B$ ,  $\psi_n$ : high-prob bound on  $\max_{i=1,\dots,n} \|\mathbf{x}_i\|_2^2$  $L_{\mathcal{D}}(\mathbf{w}) - L_{\mathbf{S}}(\mathbf{w}) \leq \tilde{\mathcal{O}}_P\left(\frac{B^2\psi_n}{n} + \sqrt{L_{\mathbf{S}}(\mathbf{w})\frac{B^2\psi_n}{n}}\right)$ 

Applying [Srebro/Sridharan/Tewari 2010]: for all  $\|\mathbf{w}\|_2 \leq B$ ,  $\psi_n$ : high-prob bound on  $\max_{i=1,\dots,n} \|\mathbf{x}_i\|_2^2$  $L_{\mathcal{D}}(\mathbf{w}) - L_{\mathbf{S}}(\mathbf{w}) \leq \tilde{\mathcal{O}}_P\left(\frac{B^2\psi_n}{n} + \sqrt{L_{\mathbf{S}}(\mathbf{w})\frac{B^2\psi_n}{n}}\right)$ 

$$\sup_{\|\mathbf{w}\|_2 \leq B, \ L_{\mathbf{S}}(\mathbf{w}) = \mathbf{0}} L_{\mathcal{D}}(\mathbf{w}) \leq rac{B^2 \psi_n}{n} + o_P(1)$$

Applying [Srebro/Sridharan/Tewari 2010]: for all  $\|\mathbf{w}\|_2 \leq B$ ,  $\psi_n$ : high-prob bound on  $\max_{i=1,\ldots,n} \|\mathbf{x}_i\|_2^2$  $L_{\mathcal{D}}(\mathbf{w}) - L_{\mathbf{S}}(\mathbf{w}) \leq \tilde{\mathcal{O}}_P\left(\frac{B^2\psi_n}{n} + \sqrt{L_{\mathbf{S}}(\mathbf{w})\frac{B^2\psi_n}{n}}\right)$ 

$$egin{aligned} \sup_{\|\mathbf{w}\|_2 \leq B, \ L_{\mathbf{S}}(\mathbf{w}) = \mathbf{0}} L_{\mathcal{D}}(\mathbf{w}) \leq & oldsymbol{c} rac{B^2 \psi_n}{n} + o_P(1) \ & ext{if } 1 \ll \lambda_n \ll n, \ B = \|\hat{\mathbf{w}}_{MN}\|_{2}, \ & o c \ L_{\mathcal{D}}(\mathbf{w}^*) \end{aligned}$$

Applying [Srebro/Sridharan/Tewari 2010]: for all  $\|\mathbf{w}\|_2 \leq B$ ,  $\psi_n$ : high-prob bound on  $\max_{i=1,\ldots,n} \|\mathbf{x}_i\|_2^2$ 

$$L_{\mathcal{D}}(\mathbf{w}) - L_{\mathbf{S}}(\mathbf{w}) \leq \tilde{\mathcal{O}}_{P}\left(\frac{B^{2}\psi_{n}}{n} + \sqrt{L_{\mathbf{S}}(\mathbf{w})\frac{B^{2}\psi_{n}}{n}}
ight)$$

$$egin{aligned} \sup_{\|\mathbf{w}\|_2 \leq B, \ L_{\mathbf{S}}(\mathbf{w}) = \mathbf{0}} L_{\mathcal{D}}(\mathbf{w}) \leq & rac{B^2 \psi_n}{n} + o_P(1) \ & ext{if } 1 \ll \lambda_n \ll n, \ B = \|\hat{\mathbf{w}}_{MN}\|_2, \ o & c \, L_{\mathcal{D}}(\mathbf{w}^*) \end{aligned}$$

If this holds with c = 1 (and maybe  $\psi_n = \mathbb{E} \|\mathbf{x}\|_2^2$ ), would explain consistency on junk features, and predict that  $B = \alpha \|\hat{\mathbf{w}}_{MN}\|_2$  gives  $\alpha^2 L_{\mathcal{D}}(\mathbf{w}^*)$ 

 $\begin{array}{l} \text{Applying [Srebro/Sridharan/Tewari 2010]: for all } \|\mathbf{w}\|_2 \leq B, \\ \psi_n: \text{high-prob bound on } \max_{i=1,\ldots,n} \|\mathbf{x}_i\|_2^2 \\ L_{\mathcal{D}}(\mathbf{w}) - L_{\mathbf{S}}(\mathbf{w}) \leq \tilde{\mathcal{O}}_P \left(\frac{B^2 \psi_n}{n} + \sqrt{L_{\mathbf{S}}(\mathbf{w}) \frac{B^2 \psi_n}{n}}\right) \\ c \leq 200,000 \log^3(n) \\ \sup_{\|\mathbf{w}\|_2 \leq B, \ L_{\mathbf{S}}(\mathbf{w})=0} L_{\mathcal{D}}(\mathbf{w}) \leq \frac{B^2 \psi_n}{n} + o_P(1) \\ \text{if } 1 \ll \lambda_n \ll n, B = \|\hat{\mathbf{w}}_{MN}\|_2, \quad \rightarrow c L_{\mathcal{D}}(\mathbf{w}^*) \end{array}$ 

If this holds with c = 1 (and maybe  $\psi_n = \mathbb{E} \|\mathbf{x}\|_2^2$ ), would explain consistency on junk features, and predict that  $B = \alpha \|\hat{\mathbf{w}}_{MN}\|_2$  gives  $\alpha^2 L_{\mathcal{D}}(\mathbf{w}^*)$ 

## **Conjecture holds (for Gaussian linear regression)**

Specifically, our more general bound implies that w.h.p.

$$\sup_{\|\mathbf{w}\|_2 \leq B, \ L_{\mathbf{S}}(\mathbf{w}) = 0} L_{\mathcal{D}}(\mathbf{w}) \leq (1 + o(1)) rac{B^2 \operatorname{Tr}(\Sigma_2)}{n}$$

 $\Sigma = \Sigma_1 \oplus \Sigma_2$  splits up covariance eigenvectors;  $\mathrm{Tr}(\Sigma_2) \leq \mathrm{Tr}(\Sigma) = \mathbb{E} \|\mathbf{x}\|^2$ 

## **Conjecture holds (for Gaussian linear regression)**

Specifically, our more general bound implies that w.h.p.

$$\sup_{\|\mathbf{w}\|_2 \leq B, \ L_{\mathbf{S}}(\mathbf{w}) = 0} L_{\mathcal{D}}(\mathbf{w}) \leq (1 + o(1)) rac{B^2 \operatorname{Tr}(\Sigma_2)}{n}$$

 $\Sigma = \Sigma_1 \oplus \Sigma_2$  splits up covariance eigenvectors;  $\mathrm{Tr}(\Sigma_2) \leq \mathrm{Tr}(\Sigma) = \mathbb{E} \|\mathbf{x}\|^2$ 

For this to mean anything, need  $B \geq \hat{\mathbf{w}}_{MN}$ 

## **Conjecture holds (for Gaussian linear regression)**

Specifically, our more general bound implies that w.h.p.

$$\sup_{\|\mathbf{w}\|_2 \leq B, \ L_{\mathbf{S}}(\mathbf{w}) = 0} L_{\mathcal{D}}(\mathbf{w}) \leq (1 + o(1)) rac{B^2 \operatorname{Tr}(\Sigma_2)}{n}$$

 $\Sigma = \Sigma_1 \oplus \Sigma_2$  splits up covariance eigenvectors;  $\mathrm{Tr}(\Sigma_2) \leq \mathrm{Tr}(\Sigma) = \mathbb{E} \|\mathbf{x}\|^2$ 

For this to mean anything, need  $B \geq \hat{\mathbf{w}}_{MN}$ 

Combine with a new analysis on  $\|\hat{\mathbf{w}}_{MN}\|$ : whp,

$$\|\hat{\mathbf{w}}_{MN}\|_2 \leq \|\mathbf{w}^*\|_2 + (1+o(1)) \; \sqrt{rac{\sigma^2 n}{\operatorname{Tr}(\Sigma_2)}}$$

# Benign overfitting of $\hat{\mathbf{w}}_{MN}$

Plugging the two bounds together:

$$L_{\mathcal{D}}(\hat{\mathbf{w}}) \leq (1+o(1)) \left(\sigma + \|\mathbf{w}^*\| \sqrt{rac{\mathrm{Tr}(\Sigma_2)}{n}}
ight)^2$$

# Benign overfitting of $\hat{\mathbf{w}}_{MN}$

Plugging the two bounds together:

$$L_{\mathcal{D}}(\hat{\mathbf{w}}) \leq (1+o(1)) \left(\sigma + \|\mathbf{w}^*\| \sqrt{rac{\mathrm{Tr}(\Sigma_2)}{n}}
ight)^2$$

Including all the fiddly conditions I didn't mention, we recover the consistency conditions of the landmark paper [Bartlett/Long/Lugosi/Tsigler PNAS 2020]
### Benign overfitting of $\hat{\mathbf{w}}_{MN}$

Plugging the two bounds together:

$$L_{\mathcal{D}}(\hat{\mathbf{w}}) \leq (1+o(1)) \left(\sigma + \|\mathbf{w}^*\| \sqrt{rac{\mathrm{Tr}(\Sigma_2)}{n}}
ight)^2$$

Including all the fiddly conditions I didn't mention, we recover the consistency conditions of the landmark paper [Bartlett/Long/Lugosi/Tsigler PNAS 2020]

Additionally tells us about nearly-minimal-norm interpolators

Generalization error in compact sets <u>Theorem</u>. If  $\Sigma = \Sigma_1 \oplus \Sigma_2$  with  $\mathrm{rank}(\Sigma_1) = o(n)$ , w.h.p.

$$\sup_{\mathbf{w}\in\mathcal{K},\ L_{\mathbf{S}}(\mathbf{w})=0}L_{\mathcal{D}}(\mathbf{w})\leq \left(1+o(1)
ight)rac{W(\Sigma_{2}^{1/2}\mathcal{K})^{2}}{n}$$

where  $W(\mathcal{K}) := \mathbb{E}_{H \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_d)}[\sup_{\mathbf{w} \in \mathcal{K}} |\langle H, \mathbf{w} \rangle|]$ is the Gaussian width (a standard tool)

this is an informal statement, but gets the gist

#### Norm needed to interpolate for general norms

Theorem. Let 
$$\|\cdot\|_*$$
 be the dual norm of  $\|\cdot\|$ .  
Call  $\hat{\mathbf{w}} = \arg\min_{\mathbf{w}: L_{\mathbf{S}}(\mathbf{w})=0} \|\mathbf{w}\|$ .  
Under some conditions, w.h.p.

$$egin{aligned} \|\hat{\mathbf{w}}\| &\leq \|\mathbf{w}^*\| + (1+o(1))rac{\sigma\sqrt{n}}{\mathbb{E}_{H\sim\mathcal{N}(\mathbf{0},\mathbf{I}_d)}\|\Sigma_2^{1/2}H\|_*}. \end{aligned}$$

#### Norm needed to interpolate for general norms

Theorem. Let 
$$\|\cdot\|_*$$
 be the dual norm of  $\|\cdot\|$ .  
Call  $\hat{\mathbf{w}} = \arg\min_{\mathbf{w}: L_{\mathbf{S}}(\mathbf{w})=0} \|\mathbf{w}\|$ .  
Under some conditions, w.h.p.

$$\|\hat{\mathbf{w}}\| \leq \|\mathbf{w}^*\| + (1+o(1))rac{\sigma\sqrt{n}}{\mathbb{E}_{H\sim\mathcal{N}(\mathbf{0},\mathbf{I}_d)}\|\Sigma_2^{1/2}H\|_*}.$$

Plugging them together, get consistency conditions analogous to the [BLLT] ones for minimal-norm interpolators for any norm.

## New application: minimum $\|\mathbf{w}\|_1$

LASSO, Adaboost, compressed sensing, basis pursuit, ...

Much harder to analyze directly, because no closed form! Some analysis in isotropic case; didn't show consistency [Ju/Lin/Liu NeurIPS 2020] [Chinot/Löffler/van de Geer 2021]

## New application: minimum $\|\mathbf{w}\|_1$

LASSO, Adaboost, compressed sensing, basis pursuit, ...

Much harder to analyze directly, because no closed form! Some analysis in isotropic case; didn't show consistency [Ju/Lin/Liu NeurIPS 2020] [Chinot/Löffler/van de Geer 2021]

Our conditions hold in a junk features setting, if  $d=e^{\omega(n)}$ 

## New application: minimum $\|\mathbf{w}\|_1$

LASSO, Adaboost, compressed sensing, basis pursuit, ...

Much harder to analyze directly, because no closed form! Some analysis in isotropic case; didn't show consistency [Ju/Lin/Liu NeurIPS 2020] [Chinot/Löffler/van de Geer 2021]

Our conditions hold in a junk features setting, if  $d=e^{\omega(n)}$ 

Very limited setting, but (as far as we know) first consistency result for  $\sigma>0$ ,  $w^*
eq 0$ 

**On Uniform Convergence and Low-Norm Interpolation Learning** Zhou, Sutherland, Srebro [NeurIPS 2020] [arXiv:2006.05942]

Uniform Convergence of Interpolators: Gaussian Width, Norm Bounds and Benign Overfitting Koehler\*, Zhou\*, Sutherland, Srebro [NeurIPS 2021] [arXiv:2106.09276]

- Junk features example:
  - $\hat{\mathbf{w}}_{MN}$  is consistent; usual uniform convergence can't show that
  - Uniform convergence over norm ball can't show any learning
- Uniform convergence of interpolators does work
  - Matches previously known (nearly necessary) sufficient conditions
  - Applies to general norm balls (though can be hard to evaluate)
  - Our analysis is very specific to Gaussian data
- Coming soon: extension to near-interpolators via optimistic rates

#### **Backup slides**

<u>Theorem:</u> In junk features with  $\lambda_n = o(n)$ ,

$$\lim_{n o \infty} \lim_{d_J o \infty} \mathbb{E} \left[ \sup_{\|\mathbf{w}\|_2 \leq \|\hat{\mathbf{w}}_{MN}\|_2} |L_{\mathcal{D}}(\mathbf{w}) - L_{\mathbf{S}}(\mathbf{w})| 
ight] = \infty.$$

<u>Theorem:</u> In junk features with  $\lambda_n = o(n)$ ,

$$\lim_{n o \infty} \lim_{d_J o \infty} \mathbb{E} \left[ \sup_{\|\mathbf{w}\|_2 \leq \|\hat{\mathbf{w}}_{MN}\|_2} |L_\mathcal{D}(\mathbf{w}) - L_\mathbf{S}(\mathbf{w})| 
ight] = \infty.$$

<u>Theorem:</u> In junk features with  $\lambda_n = o(n)$ ,

$$\lim_{n o \infty} \lim_{d_J o \infty} \mathbb{E} \left[ \sup_{\|\mathbf{w}\|_2 \leq \|\hat{\mathbf{w}}_{MN}\|_2} |L_\mathcal{D}(\mathbf{w}) - L_\mathbf{S}(\mathbf{w})| 
ight] = \infty.$$

<u>Theorem:</u> In junk features with  $\lambda_n = o(n)$ ,

$$\lim_{n o \infty} \lim_{d_J o \infty} \mathbb{E} \left[ \sup_{\|\mathbf{w}\|_2 \leq \|\hat{\mathbf{w}}_{MN}\|_2} |L_\mathcal{D}(\mathbf{w}) - L_\mathbf{S}(\mathbf{w})| 
ight] = \infty.$$

<u>Theorem:</u> In junk features with  $\lambda_n = o(n)$ ,

$$\lim_{n o \infty} \lim_{d_J o \infty} \mathbb{E} \left[ \sup_{\|\mathbf{w}\|_2 \leq \|\hat{\mathbf{w}}_{MN}\|_2} |L_\mathcal{D}(\mathbf{w}) - L_\mathbf{S}(\mathbf{w})| 
ight] = \infty.$$

$$L_{\mathcal{D}}(\mathbf{w}) = (\mathbf{w} - \mathbf{w}^*)^{\mathsf{T}} \mathbf{\Sigma}(\mathbf{w} - \mathbf{w}^*) + L_{\mathcal{D}}(\mathbf{w}^*)$$

<u>Theorem:</u> In junk features with  $\lambda_n = o(n)$ ,

$$\lim_{n o \infty} \lim_{d_J o \infty} \mathbb{E} igg[ \sup_{\|\mathbf{w}\|_2 \leq \|\hat{\mathbf{w}}_{MN}\|_2} |L_\mathcal{D}(\mathbf{w}) - L_\mathbf{S}(\mathbf{w})| igg] = \infty.$$

$$egin{aligned} & L_{\mathcal{D}}(\mathbf{w}) = (\mathbf{w} - \mathbf{w}^*)^\mathsf{T} \mathbf{\Sigma} (\mathbf{w} - \mathbf{w}^*) + L_{\mathcal{D}}(\mathbf{w}^*) \ & L_{\mathcal{D}}(\mathbf{w}) - L_{\mathbf{S}}(\mathbf{w}) = (\mathbf{w} - \mathbf{w}^*) (\mathbf{\Sigma} - \hat{\mathbf{\Sigma}}) (\mathbf{w} - \mathbf{w}^*) \ & + (L_{\mathcal{D}}(\mathbf{w}^*) - L_{\mathbf{S}}(\mathbf{w}^*)) - ext{cross term} \end{aligned}$$

<u>Theorem:</u> In junk features with  $\lambda_n = o(n)$ ,

$$\lim_{n o \infty} \lim_{d_J o \infty} \mathbb{E} igg[ \sup_{\|\mathbf{w}\|_2 \leq \|\hat{\mathbf{w}}_{MN}\|_2} |L_{\mathcal{D}}(\mathbf{w}) - L_{\mathbf{S}}(\mathbf{w})| igg] = \infty.$$

$$egin{aligned} &L_{\mathcal{D}}(\mathbf{w}) = (\mathbf{w} - \mathbf{w}^*)^\mathsf{T} \mathbf{\Sigma} (\mathbf{w} - \mathbf{w}^*) + L_{\mathcal{D}}(\mathbf{w}^*) \ &L_{\mathcal{D}}(\mathbf{w}) - L_{\mathbf{S}}(\mathbf{w}) = (\mathbf{w} - \mathbf{w}^*) (\mathbf{\Sigma} - \hat{\mathbf{\Sigma}}) (\mathbf{w} - \mathbf{w}^*) \ &+ (L_{\mathcal{D}}(\mathbf{w}^*) - L_{\mathbf{S}}(\mathbf{w}^*)) - ext{cross term} \ &\sup[\ldots] \geq \|\mathbf{\Sigma} - \hat{\mathbf{\Sigma}}\|_{op} \cdot (\|\hat{\mathbf{w}}_{MN}\|_2 - \|\mathbf{w}^*\|_2)^2 + o(1) \end{aligned}$$

<u>Theorem:</u> In junk features with  $\lambda_n = o(n)$ ,

$$\lim_{n o \infty} \lim_{d_J o \infty} \mathbb{E} \left[ \sup_{\|\mathbf{w}\|_2 \leq \|\hat{\mathbf{w}}_{MN}\|_2} |L_\mathcal{D}(\mathbf{w}) - L_\mathbf{S}(\mathbf{w})| 
ight] = \infty.$$

$$egin{aligned} &L_{\mathcal{D}}(\mathbf{w}) = (\mathbf{w} - \mathbf{w}^*)^\mathsf{T} \mathbf{\Sigma} (\mathbf{w} - \mathbf{w}^*) + L_{\mathcal{D}}(\mathbf{w}^*) \ &L_{\mathcal{D}}(\mathbf{w}) - L_\mathbf{S}(\mathbf{w}) = (\mathbf{w} - \mathbf{w}^*) (\mathbf{\Sigma} - \hat{\mathbf{\Sigma}}) (\mathbf{w} - \mathbf{w}^*) \ &+ (L_{\mathcal{D}}(\mathbf{w}^*) - L_\mathbf{S}(\mathbf{w}^*)) - ext{cross term} \ &\sup[\ldots] \geq \|\mathbf{\Sigma} - \hat{\mathbf{\Sigma}}\|_{op} \cdot (\|\hat{\mathbf{w}}_{MN}\|_2 - \|\mathbf{w}^*\|_2)^2 + o(1) \ &\Theta\left(rac{n}{\lambda_n}
ight) \end{aligned}$$

<u>Theorem:</u> In junk features with  $\lambda_n = o(n)$ ,

$$\lim_{n o \infty} \lim_{d_J o \infty} \mathbb{E} igg[ \sup_{\|\mathbf{w}\|_2 \leq \|\hat{\mathbf{w}}_{MN}\|_2} |L_{\mathcal{D}}(\mathbf{w}) - L_{\mathbf{S}}(\mathbf{w})| igg] = \infty.$$

$$L_{\mathcal{D}}(\mathbf{w}) = (\mathbf{w} - \mathbf{w}^*)^\mathsf{T} \mathbf{\Sigma} (\mathbf{w} - \mathbf{w}^*) + L_{\mathcal{D}}(\mathbf{w}^*)$$

$$L_{\mathcal{D}}(\mathbf{w}) - L_{\mathbf{S}}(\mathbf{w}) = (\mathbf{w} - \mathbf{w}^*)(\mathbf{\Sigma} - \hat{\mathbf{\Sigma}})(\mathbf{w} - \mathbf{w}^*)$$

$$+ (L_{\mathcal{D}}(\mathbf{w}^*) - L_{\mathbf{S}}(\mathbf{w}^*)) - \text{cross term}$$

$$\sup[\ldots] \ge \|\mathbf{\Sigma} - \hat{\mathbf{\Sigma}}\|_{op} \cdot (\|\hat{\mathbf{w}}_{MN}\|_2 - \|\mathbf{w}^*\|_2)^2 + o(1)$$

$$\Theta\left(\sqrt{\frac{\lambda_n}{n}}\right) \qquad \Theta\left(\frac{n}{\lambda_n}\right)$$
Koltchinskii/Lounici, Bernoulli 2017

<u>Theorem:</u> In junk features with  $\lambda_n = o(n)$ ,

$$\lim_{n o \infty} \lim_{d_J o \infty} \mathbb{E} igg[ \sup_{\|\mathbf{w}\|_2 \leq \|\hat{\mathbf{w}}_{MN}\|_2} |L_{\mathcal{D}}(\mathbf{w}) - L_{\mathbf{S}}(\mathbf{w})| igg] = \infty.$$

$$L_{\mathcal{D}}(\mathbf{w}) = (\mathbf{w} - \mathbf{w}^*)^{\mathsf{T}} \mathbf{\Sigma} (\mathbf{w} - \mathbf{w}^*) + L_{\mathcal{D}}(\mathbf{w}^*)$$

$$L_{\mathcal{D}}(\mathbf{w}) - L_{\mathbf{S}}(\mathbf{w}) = (\mathbf{w} - \mathbf{w}^*)(\mathbf{\Sigma} - \hat{\mathbf{\Sigma}})(\mathbf{w} - \mathbf{w}^*)$$

$$+ (L_{\mathcal{D}}(\mathbf{w}^*) - L_{\mathbf{S}}(\mathbf{w}^*)) - \text{cross term}$$

$$\sup[\ldots] \ge \|\mathbf{\Sigma} - \hat{\mathbf{\Sigma}}\|_{op} \cdot (\|\hat{\mathbf{w}}_{MN}\|_2 - \|\mathbf{w}^*\|_2)^2 + o(1) \to \infty$$

$$\Theta\left(\sqrt{\frac{\lambda_n}{n}}\right) \qquad \Theta\left(\frac{n}{\lambda_n}\right)$$
Koltchinskii/Lounici, Bernoulli 2017

We don't really care about small  $L_D$ , big  $L_S$ .... Could we bound  $\sup L_D - L_S$  instead of  $\sup |L_D - L_S|$ ?

We don't really care about small  $L_D$ , big  $L_S$ .... Could we bound  $\sup L_D - L_S$  instead of  $\sup |L_D - L_S|$ ?

- Existing uniform convergence proofs are "really" about  $|L_{\mathcal{D}} - L_{\mathbf{S}}|$  [Nagarajan/Kolter, NeurIPS 2019]

We don't really care about small  $L_D$ , big  $L_S$ .... Could we bound  $\sup L_D - L_S$  instead of  $\sup |L_D - L_S|$ ?

- Existing uniform convergence proofs are "really" about  $|L_{\mathcal{D}} L_{\mathbf{S}}|$  [Nagarajan/Kolter, NeurIPS 2019]
- Strongly expect still  $\infty$  for norm balls in our testbed
  - $\lambda_{\max}(\mathbf{\Sigma} \hat{\mathbf{\Sigma}})$  instead of  $\|\mathbf{\Sigma} \hat{\mathbf{\Sigma}}\|_{op}$

We don't really care about small  $L_D$ , big  $L_S$ .... Could we bound  $\sup L_D - L_S$  instead of  $\sup |L_D - L_S|$ ?

- Existing uniform convergence proofs are "really" about  $|L_{\mathcal{D}} L_{\mathbf{S}}|$  [Nagarajan/Kolter, NeurIPS 2019]
- Strongly expect still  $\infty$  for norm balls in our testbed

• 
$$\lambda_{\max}(\mathbf{\Sigma} - \hat{\mathbf{\Sigma}})$$
 instead of  $\|\mathbf{\Sigma} - \hat{\mathbf{\Sigma}}\|_{op}$ 

• Not possible to show  $\sup_{f \in \mathcal{F}} L_{\mathcal{D}} - L_{\mathbf{S}}$  is big for all  $\mathcal{F}$ • If  $\hat{f}$  consistent and  $\inf_f L_{\mathbf{S}}(f) \ge 0$ , use  $\mathcal{F} = \{f : L_{\mathcal{D}}(f) \le L_{\mathcal{D}}(f^*) + \epsilon_{n,\delta}\}$ 

# 

- Confirms speculation based on c=1 assumption
- Shows consistency with uniform convergence (of interpolators)
- New result for error of not-quite-minimal-norm interpolators
  - Norm  $\|\hat{\mathbf{w}}_{MN}\| + \text{const}$  is asympttically consistent
  - Norm  $1.1 \| \hat{\mathbf{w}}_{MN} \|$  is at worst  $1.21 \, L_{\mathcal{D}}(\mathbf{w}^*)$

#### What does $\{\mathbf{w}: \|\mathbf{w}\| \leq B, \, L_{\mathbf{S}}(\mathbf{w}) = 0\}$ look like?

#### What does $\{\mathbf{w}: \|\mathbf{w}\| \leq B, L_{\mathbf{S}}(\mathbf{w}) = 0\}$ look like?



#### Intersection of *d*-ball

## What does $\{\mathbf{w} : \|\mathbf{w}\| \le B, L_{\mathbf{S}}(\mathbf{w}) = 0\}$ look like? $\{\mathbf{w} : L_{\mathbf{S}}(\mathbf{w}) = \frac{1}{n} \|\mathbf{X}\mathbf{w} - \mathbf{y}\|^2 = 0\}$ is the plane $\mathbf{X}\mathbf{w} = \mathbf{y}$



Intersection of *d*-ball with (d - n)-hyperplane:

## What does $\{\mathbf{w} : \|\mathbf{w}\| \le B, L_{\mathbf{S}}(\mathbf{w}) = 0\}$ look like? $\{\mathbf{w} : L_{\mathbf{S}}(\mathbf{w}) = \frac{1}{n} \|\mathbf{X}\mathbf{w} - \mathbf{y}\|^2 = 0\}$ is the plane $\mathbf{X}\mathbf{w} = \mathbf{y}$



Intersection of d-ball with (d - n)-hyperplane: (d - n)-ball

What does  $\{\mathbf{w} : \|\mathbf{w}\| \le B, L_{\mathbf{S}}(\mathbf{w}) = 0\}$  look like?  $\{\mathbf{w} : L_{\mathbf{S}}(\mathbf{w}) = \frac{1}{n} \|\mathbf{X}\mathbf{w} - \mathbf{y}\|^2 = 0\}$  is the plane  $\mathbf{X}\mathbf{w} = \mathbf{y}$ 



Intersection of d-ball with (d - n)-hyperplane: (d - n)-ball centered at  $\hat{\mathbf{w}}_{MN}$  What does  $\{\mathbf{w}: \|\mathbf{w}\| \leq B, L_{\mathbf{S}}(\mathbf{w}) = 0\}$  look like?  $\{\mathbf{w}: L_{\mathbf{S}}(\mathbf{w}) = \frac{1}{n} \|\mathbf{X}\mathbf{w} - \mathbf{y}\|^2 = 0\}$  is the plane  $\mathbf{X}\mathbf{w} = \mathbf{y}$ 



Intersection of d-ball with (d - n)-hyperplane: (d - n)-ball centered at  $\hat{\mathbf{w}}_{MN}$ 

Can write as  $\{\hat{\mathbf{w}} + \mathbf{F}\mathbf{z} : \mathbf{z} \in \mathbb{R}^{d-n}, \|\hat{\mathbf{w}} + \mathbf{F}\mathbf{z}\| \leq B\}$ where  $\hat{\mathbf{w}}$  is *any* interpolator, **F** is basis for  $\text{ker}(\mathbf{X})$ 

Can change variables in  $\sup_{\mathbf{w}:\|\mathbf{w}\|\leq B, \ L_{\mathbf{S}}(\mathbf{w})=0} L_{\mathcal{D}}(\mathbf{w})$  to

$$L_{\mathcal{D}}(\mathbf{w}^*) + \sup_{\mathbf{z}: \|\hat{\mathbf{w}} + \mathbf{F}\mathbf{z}\|^2 \leq B^2} (\hat{\mathbf{w}} + \mathbf{F}\mathbf{z} - w^*)^\mathsf{T} \mathbf{\Sigma} (\hat{\mathbf{w}} + \mathbf{F}\mathbf{z} - w^*)$$

Can change variables in  $\sup_{\mathbf{w}:\|\mathbf{w}\|\leq B, \ L_{\mathbf{S}}(\mathbf{w})=0} L_{\mathcal{D}}(\mathbf{w})$  to

$$L_{\mathcal{D}}(\mathbf{w}^*) + \sup_{\mathbf{z}: \|\hat{\mathbf{w}} + \mathbf{F}\mathbf{z}\|^2 \leq B^2} (\hat{\mathbf{w}} + \mathbf{F}\mathbf{z} - w^*)^\mathsf{T} \mathbf{\Sigma} (\hat{\mathbf{w}} + \mathbf{F}\mathbf{z} - w^*)$$

Quadratic program, one quadratic constraint: strong duality

Can change variables in  $\sup_{\mathbf{w}:\|\mathbf{w}\|\leq B,\ L_{\mathbf{S}}(\mathbf{w})=0}L_{\mathcal{D}}(\mathbf{w})$  to

$$L_{\mathcal{D}}(\mathbf{w}^*) + \sup_{\mathbf{z}: \|\hat{\mathbf{w}} + \mathbf{F}\mathbf{z}\|^2 \leq B^2} (\hat{\mathbf{w}} + \mathbf{F}\mathbf{z} - w^*)^\mathsf{T} \mathbf{\Sigma} (\hat{\mathbf{w}} + \mathbf{F}\mathbf{z} - w^*)$$

Quadratic program, one quadratic constraint: strong duality

Exactly equivalent to problem in *one* scalar variable:

$$L_{\mathcal{D}}(\hat{\mathbf{w}}) + \inf_{\mu > \|\mathbf{F}^{\mathsf{T}}\boldsymbol{\Sigma}\mathbf{F}\|} \left\| \mathbf{F}^{\mathsf{T}}[\mu \hat{\mathbf{w}} - \Sigma(\hat{\mathbf{w}} - \mathbf{w}^{*})] \right\|_{(\mu \mathbf{I}_{p-n} - \mathbf{F}^{\mathsf{T}}\boldsymbol{\Sigma}\mathbf{F})^{-1}} + \mu (B^{2} - \|\hat{\mathbf{w}}\|^{2})$$

Can change variables in  $\sup_{\mathbf{w}:\|\mathbf{w}\|\leq B,\ L_{\mathbf{S}}(\mathbf{w})=0}L_{\mathcal{D}}(\mathbf{w})$  to

$$L_{\mathcal{D}}(\mathbf{w}^*) + \sup_{\mathbf{z}: \|\hat{\mathbf{w}} + \mathbf{F}\mathbf{z}\|^2 \leq B^2} (\hat{\mathbf{w}} + \mathbf{F}\mathbf{z} - w^*)^\mathsf{T} \mathbf{\Sigma} (\hat{\mathbf{w}} + \mathbf{F}\mathbf{z} - w^*)$$

Quadratic program, one quadratic constraint: strong duality

Exactly equivalent to problem in *one* scalar variable:

$$L_{\mathcal{D}}(\hat{\mathbf{w}}) + \inf_{\mu > \|\mathbf{F}^{\mathsf{T}} \boldsymbol{\Sigma} \mathbf{F}\|} \left\| \mathbf{F}^{\mathsf{T}} [\mu \hat{\mathbf{w}} - \boldsymbol{\Sigma} (\hat{\mathbf{w}} - \mathbf{w}^{*})] \right\|_{(\mu \mathbf{I}_{p-n} - \mathbf{F}^{\mathsf{T}} \boldsymbol{\Sigma} \mathbf{F})^{-1}} + \mu (B^{2} - \| \hat{\mathbf{w}} \|^{2})$$

Can analyze this for different choices of  $\hat{\mathbf{w}}$ ...

#### The minimal-risk interpolator

$$\hat{\mathbf{w}}_{MR} = rgmin_{\mathcal{W}} L_{\mathcal{D}}(\mathbf{w}) \ \mathbf{w}: \mathbf{X} \mathbf{w} = \mathbf{y}$$
#### The minimal-risk interpolator

$$egin{aligned} \hat{\mathbf{w}}_{MR} &= rg\min_{\mathbf{w}:\mathbf{X}\mathbf{w}=\mathbf{y}} L_{\mathcal{D}}(\mathbf{w}) \ &= \mathbf{w}^* + \Sigma^{-1}\mathbf{X}^\mathsf{T}(\mathbf{X}\Sigma^{-1}\mathbf{X}^\mathsf{T})^{-1}(Y-X\mathbf{w}^*) \end{aligned}$$

#### The minimal-risk interpolator

$$egin{aligned} \hat{\mathbf{w}}_{MR} &= rg\min_{\mathbf{w}:\mathbf{X}\mathbf{w}=\mathbf{y}} L_{\mathcal{D}}(\mathbf{w}) \ &= \mathbf{w}^* + \Sigma^{-1}\mathbf{X}^\mathsf{T}(\mathbf{X}\Sigma^{-1}\mathbf{X}^\mathsf{T})^{-1}(Y-X\mathbf{w}^*) \end{aligned}$$

In Gaussian least squares generally, have that

$$\mathbb{E} \, L_{\mathcal{D}}(\hat{\mathbf{w}}_{MR}) = rac{d-1}{d-1-n} \, L_{\mathcal{D}}(\mathbf{w}^*)$$

so  $\hat{\mathbf{w}}_{MR}$  is consistent iff n=o(d).

#### The minimal-risk interpolator

$$egin{aligned} \hat{\mathbf{w}}_{MR} &= rg\min_{\mathbf{w}:\mathbf{X}\mathbf{w}=\mathbf{y}} L_{\mathcal{D}}(\mathbf{w}) \ &= \mathbf{w}^* + \Sigma^{-1}\mathbf{X}^\mathsf{T}(\mathbf{X}\Sigma^{-1}\mathbf{X}^\mathsf{T})^{-1}(Y-X\mathbf{w}^*) \end{aligned}$$

In Gaussian least squares generally, have that

$$\mathbb{E} \, L_{\mathcal{D}}(\hat{\mathbf{w}}_{MR}) = rac{d-1}{d-1-n} \, L_{\mathcal{D}}(\mathbf{w}^*)$$

so  $\hat{\mathbf{w}}_{MR}$  is consistent iff n=o(d).

Very useful for lower bounds! [Muthukumar+ JSAIT 2020]

#### **Restricted eigenvalue under interpolation**

$$\kappa_{\mathbf{X}}(\mathbf{\Sigma}) = \sup_{\|\mathbf{w}\|=1, \ \mathbf{X}\mathbf{w}=\mathbf{0}} \mathbf{w}^{\mathsf{T}} \mathbf{\Sigma} \mathbf{w}$$

Roughly, "how much" of  ${f \Sigma}$  is "missed" by  ${f X}$ 

Analyzing dual with  $\hat{\mathbf{w}}_{MR}$ , get without **any** distributional assumptions that  $\sup_{\substack{1 \le \beta \le 4 \\ L_{\mathcal{D}}(\mathbf{w}) = L_{\mathcal{D}}(\hat{\mathbf{w}}_{MR}) + \beta \kappa_{X}(\Sigma) \left[ \|\hat{\mathbf{w}}_{MR}\|^{2} - \|\hat{\mathbf{w}}_{MN}\|^{2} \right]$ 

Analyzing dual with  $\hat{\mathbf{w}}_{MR}$ , get without **any** distributional assumptions that  $\sup_{\substack{1 \leq \beta \leq 4 \\ L_{\mathcal{D}}(\mathbf{w}) = L_{\mathcal{D}}(\hat{\mathbf{w}}_{MR}) + \beta \kappa_{X}(\Sigma) \left[ \| \hat{\mathbf{w}}_{MR} \|^{2} - \| \hat{\mathbf{w}}_{MN} \|^{2} \right]$ (amount of missed energy) · (available norm)

Analyzing dual with  $\hat{\mathbf{w}}_{MR}$ , get without **any** distributional assumptions that  $\sup_{\substack{1 \leq \beta \leq 4 \\ \mathcal{L}_{\mathcal{D}}(\mathbf{w}) = L_{\mathcal{D}}(\hat{\mathbf{w}}_{MR}) + \beta \kappa_{X}(\Sigma) \left[ \| \hat{\mathbf{w}}_{MR} \|^{2} - \| \hat{\mathbf{w}}_{MN} \|^{2} \right]$  $\| \mathbf{w} \| \leq \| \hat{\mathbf{w}}_{MR} \|$ (amount of missed energy) · (available norm) If  $\hat{\mathbf{w}}_{MR}$  consistent, everything smaller-norm also consistent iff  $\beta$  term  $\rightarrow 0$ 

In the generic results, 
$$L_{\mathcal{D}}$$
 means  $L_{\mathcal{D}}(\mathbf{w}) = L_{\mathcal{D}}(\mathbf{w}^*) + (\mathbf{w} - \mathbf{w}^*)^\mathsf{T}\Sigma(\mathbf{w} - \mathbf{w}^*)$  for some  $\mathbf{w}^*$ 

Analyzing dual with  $\hat{\mathbf{w}}_{MR}$ , get without **any** distributional assumptions that  $\sup_{\substack{1 \le \beta \le 4 \\ \mathcal{L}_{\mathcal{D}}(\mathbf{w}) = L_{\mathcal{D}}(\hat{\mathbf{w}}_{MR}) + \beta \kappa_{X}(\Sigma) \left[ \|\hat{\mathbf{w}}_{MR}\|^{2} - \|\hat{\mathbf{w}}_{MN}\|^{2} \right]$  $\|\mathbf{w}\| \le \|\hat{\mathbf{w}}_{MR}\|$ (amount of missed energy) · (available norm)  $L_{\mathbf{S}}(\mathbf{w}) = 0$ If  $\hat{\mathbf{w}}_{MR}$  consistent, everything smaller-norm also consistent iff  $\beta$  term  $\rightarrow 0$ 

In our setting:

 $\hat{\mathbf{w}}_{MR}$  is consistent,  $L_{\mathcal{D}}(\hat{\mathbf{w}}_{MR}) 
ightarrow L_{\mathcal{D}}(\mathbf{w}^*)$ 

In the generic results,  $L_{\mathcal{D}}$  means  $L_{\mathcal{D}}(\mathbf{w}) = L_{\mathcal{D}}(\mathbf{w}^*) + (\mathbf{w} - \mathbf{w}^*)^\mathsf{T} \Sigma(\mathbf{w} - \mathbf{w}^*)$  for some  $\mathbf{w}^*$ 

Analyzing dual with  $\hat{\mathbf{w}}_{MR}$ , get without **any** distributional assumptions that  $1 \le eta \le 4$  $\sup \quad L_{\mathcal{D}}(\mathbf{w}) = L_{\mathcal{D}}(\hat{\mathbf{w}}_{MR}) + eta \overline{\kappa}_X(\Sigma) \left[ \| \hat{\mathbf{w}}_{MR} \|^2 - \| \hat{\mathbf{w}}_{MN} \|^2 
ight]$  $\|\mathbf{w}\|{\leq}\|\hat{\mathbf{w}}_{MR}\|$ (amount of missed energy) · (available norm)  $L_{\mathbf{S}}(\mathbf{w})=0$ If  $\hat{\mathbf{w}}_{MR}$  consistent, everything smaller-norm also consistent iff eta term ightarrow 0In our setting:  $\hat{\mathbf{w}}_{MR}$  is consistent,  $L_{\mathcal{D}}(\hat{\mathbf{w}}_{MR}) 
ightarrow L_{\mathcal{D}}(\mathbf{w}^*)$  $\kappa_X(\mathbf{\Sigma}) pprox rac{\lambda_n}{n} \quad \mathbb{E}ig[\|\hat{\mathbf{w}}_{MR}\|^2 - \|\hat{\mathbf{w}}_{MN}\|^2ig] = rac{\sigma^2 d_S}{\lambda_n} + o\left(1
ight)$ 

In the generic results,  $L_{\mathcal{D}}$  means  $L_{\mathcal{D}}(\mathbf{w}) = L_{\mathcal{D}}(\mathbf{w}^*) + (\mathbf{w} - \mathbf{w}^*)^{\mathsf{T}} \Sigma(\mathbf{w} - \mathbf{w}^*)$  for some  $\mathbf{w}^*$ 

Analyzing dual with  $\hat{\mathbf{w}}_{MR}$ , get without **any** distributional assumptions that  $1 \leq \beta \leq 4$  $\sup \quad \ \ L_{\mathcal{D}}(\mathbf{w}) = L_{\mathcal{D}}(\hat{\mathbf{w}}_{MR}) + ar{eta} ar{\kappa}_X(\Sigma) \left[ \| \hat{\mathbf{w}}_{MR} \|^2 - \| \hat{\mathbf{w}}_{MN} \|^2 
ight]$  $\|\mathbf{w}\| \leq \|\hat{\mathbf{w}}_{MR}\|$ (amount of missed energy) · (available norm)  $L_{\mathbf{S}}(\mathbf{w})=0$ If  $\hat{\mathbf{w}}_{MR}$  consistent, everything smaller-norm also consistent iff eta term ightarrow 0In our setting:  $\hat{\mathbf{w}}_{MR}$  is consistent,  $L_{\mathcal{D}}(\hat{\mathbf{w}}_{MR}) 
ightarrow L_{\mathcal{D}}(\mathbf{w}^*)$  $\kappa_X(\mathbf{\Sigma}) pprox rac{\lambda_n}{n} \quad \mathbb{E}ig[\|\hat{\mathbf{w}}_{MR}\|^2 - \|\hat{\mathbf{w}}_{MN}\|^2ig] = rac{\sigma^2 d_S}{\lambda_{ au}} + o\left(1
ight)$  $\mathbb{E} \sup_{\|\mathbf{w}\| \leq \|\hat{\mathbf{w}}_{MR}\|, \ L_{\mathbf{S}}(\mathbf{w}) = 0} L_{\mathcal{D}}(\mathbf{w}) 
ightarrow L_{\mathcal{D}}(\mathbf{w}^{*})$ Plugging in:

In the generic results,  $L_{\mathcal{D}}$  means  $L_{\mathcal{D}}(\mathbf{w}) = L_{\mathcal{D}}(\mathbf{w}^*) + (\mathbf{w} - \mathbf{w}^*)^\mathsf{T} \Sigma(\mathbf{w} - \mathbf{w}^*)$  for some  $\mathbf{w}^*$ 

Analyzing dual with  $\hat{\mathbf{w}}_{MN}$  for  $\hat{\mathbf{w}}, lpha \geq 1$ , get in general:

 $\sup_{\substack{\|\mathbf{w}\| \leq \alpha \|\hat{\mathbf{w}}_{MN}\|\\L_{\mathbf{S}}(\mathbf{w}) = 0}} L_{\mathcal{D}}(\mathbf{w}) = L_{\mathcal{D}}(\hat{\mathbf{w}}_{MN}) + (\alpha^2 - 1) \kappa_X(\Sigma) \|\hat{\mathbf{w}}_{MN}\|^2 + R_n$  $R_n \to 0 \text{ if } \hat{\mathbf{w}}_{MN} \text{ is consistent}$ 

Analyzing dual with  $\hat{\mathbf{w}}_{MN}$  for  $\hat{\mathbf{w}}$ ,  $lpha \geq 1$ , get in general:

 $\sup_{\substack{\|\mathbf{w}\| \leq \alpha \|\hat{\mathbf{w}}_{MN}\|\\L_{\mathbf{S}}(\mathbf{w}) = 0}} L_{\mathcal{D}}(\mathbf{w}) = L_{\mathcal{D}}(\hat{\mathbf{w}}_{MN}) + (\alpha^2 - 1) \kappa_X(\Sigma) \|\hat{\mathbf{w}}_{MN}\|^2 + R_n$  $R_n \to 0 \text{ if } \hat{\mathbf{w}}_{MN} \text{ is consistent}$ 

In our setting:

 $\hat{\mathbf{w}}_{MN}$  is consistent, because  $\|\hat{\mathbf{w}}_{MN}\| \leq \|\hat{\mathbf{w}}_{MR}\|$ 

Analyzing dual with  $\hat{\mathbf{w}}_{MN}$  for  $\hat{\mathbf{w}}$ ,  $lpha \geq 1$ , get in general:

 $\sup_{\substack{\|\mathbf{w}\| \leq \alpha \|\hat{\mathbf{w}}_{MN}\|\\ L_{\mathbf{S}}(\mathbf{w}) = 0}} L_{\mathcal{D}}(\mathbf{w}) = L_{\mathcal{D}}(\hat{\mathbf{w}}_{MN}) + (\alpha^2 - 1) \kappa_X(\Sigma) \|\hat{\mathbf{w}}_{MN}\|^2 + R_n$ 

In our setting:

 $\hat{\mathbf{w}}_{MN}$  is consistent, because  $\|\hat{\mathbf{w}}_{MN}\| \leq \|\hat{\mathbf{w}}_{MR}\|$  $\mathbb{E} \kappa_{\mathbf{X}}(\mathbf{\Sigma}) \|\hat{\mathbf{w}}_{MN}\|^2 o \sigma^2 = L_{\mathcal{D}}(\mathbf{w}^*)$ 

Analyzing dual with  $\hat{\mathbf{w}}_{MN}$  for  $\hat{\mathbf{w}}$ ,  $lpha \geq 1$ , get in general:

 $\sup_{\substack{\|\mathbf{w}\| \leq \alpha \|\hat{\mathbf{w}}_{MN}\|\\L_{\mathbf{S}}(\mathbf{w}) = 0}} L_{\mathcal{D}}(\mathbf{w}) = L_{\mathcal{D}}(\hat{\mathbf{w}}_{MN}) + (\alpha^2 - 1) \kappa_X(\Sigma) \|\hat{\mathbf{w}}_{MN}\|^2 + R_n$ 

In our setting:

 $\hat{\mathbf{w}}_{MN}$  is consistent, because  $\|\hat{\mathbf{w}}_{MN}\| \le \|\hat{\mathbf{w}}_{MR}\|$  $\mathbb{E} \kappa_{\mathbf{X}}(\mathbf{\Sigma}) \|\hat{\mathbf{w}}_{MN}\|^2 \to \sigma^2 = L_{\mathcal{D}}(\mathbf{w}^*)$ Plugging in:  $\mathbb{E} \sup_{\|\mathbf{w}\| \le \alpha \|\hat{\mathbf{w}}_{MN}\|, L_{\mathbf{S}}(\mathbf{w})=0} L_{\mathcal{D}}(\mathbf{w}) \to \alpha^2 L_{\mathcal{D}}(\mathbf{w}^*)$ 

Analyzing dual with  $\hat{\mathbf{w}}_{MN}$  for  $\hat{\mathbf{w}}$ ,  $lpha \geq 1$ , get in general:

 $\sup_{\substack{\|\mathbf{w}\| \leq \alpha \|\hat{\mathbf{w}}_{MN}\|\\L_{\mathbf{S}}(\mathbf{w}) = 0}} L_{\mathcal{D}}(\mathbf{w}) = L_{\mathcal{D}}(\hat{\mathbf{w}}_{MN}) + (\alpha^2 - 1) \kappa_X(\Sigma) \|\hat{\mathbf{w}}_{MN}\|^2 + R_n$ 

In our setting:

$$\begin{split} \hat{\mathbf{w}}_{MN} \text{ is consistent, because } \| \hat{\mathbf{w}}_{MN} \| &\leq \| \hat{\mathbf{w}}_{MR} \| \\ & \mathbb{E} \, \kappa_{\mathbf{X}}(\mathbf{\Sigma}) \, \| \hat{\mathbf{w}}_{MN} \|^2 \to \sigma^2 = L_{\mathcal{D}}(\mathbf{w}^*) \\ \end{split}$$
Plugging in:  $\mathbb{E} \, \sup_{\|\mathbf{w}\| \leq \alpha \| \hat{\mathbf{w}}_{MN} \|, \, L_{\mathbf{S}}(\mathbf{w}) = 0} \, L_{\mathcal{D}}(\mathbf{w}) \to \alpha^2 \, L_{\mathcal{D}}(\mathbf{w}^*) \\ \dots \text{ and we're done!} \end{split}$