Can Uniform Convergence Explain Interpolation Learning?

D.J. Sutherland TTI-Chicago → UBC

based on arXiv:2006.05942 (NeurIPS 2020), with:

Lijia Zhou U Chicago



Nati Srebro TTI-Chicago



Penn State Statistics Seminar, October 8 2020

• 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} :

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

• 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} :

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

• 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} :

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

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

• Standard approaches based on empirical risk minimization:

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

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	

Second Edition

🖉 Springer

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



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
Incention	1 640 402	yes yes	yes no	100.0 100.0	89.05 89.31
inception	1,049,402	no no	yes no	100.0 100.0	86.03 85.75
ang et al., "Rethinking	generalizatio	on", 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"



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	no 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 $L_{\mathcal{D}}(f^*) > 0$)

Zhang et



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

yes	100.0	80.05
•		09.05
no	100.0	89.31
yes	100.0	86.03
no	100.0	85.75
	yes no	yes 100.0 no 100.0

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



, ICIVIE 2010

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

Misha Belkin Simons Institute July 2019

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

A constant factor of 2 invalidates the bound!



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])

 $L_{\mathcal{D}}(\hat{f}) \leq L_{\mathbf{S}}(\hat{f}) + \text{bound}$

WYSIWYG bounds:

Misha Belkin Simons Institute _{Oracle bounds} July 2019

expected loss ≈
optimal loss

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



Generalization theory for interpolation?

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

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



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

*One exception-ish [Negrea/Dziugaite/Roy, ICML 2020]: relates \hat{f} to a surrogate predictor, shows uniform convergence for the surrogate

We're only going to worry about *consistency*:

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

We're only going to worry about *consistency*:

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

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

We're only going to worry about *consistency*:

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

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

Is it possible to show consistency of an interpolator with

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

We're only going to worry about *consistency*:

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

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

Is it possible to show consistency of an interpolator with

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

This requires tight constants!

Our testbed problem



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

Our testbed problem



 $\mathbf{X}\mathbf{w} = \mathbf{v}$

As $d_J \to \infty$, $\hat{\mathbf{w}}_{MN}$ approaches ridge regression on the signal: $\langle \hat{\mathbf{w}}_{MN}, \mathbf{x} \rangle \xrightarrow{d_J \to \infty} \langle \hat{\mathbf{w}}_{\lambda_n}, \mathbf{x}_S \rangle$ for almost all \mathbf{x} $\hat{\mathbf{w}}_{\lambda_n} = \underset{\mathbf{w}_S \in \mathbb{R}^{d_S}}{\operatorname{argmin}} \|\mathbf{X}_S \mathbf{w}_S^2 - \mathbf{y}\|^2 + \lambda_n \|\mathbf{w}_S\|^2$

As $d_J \to \infty$, $\hat{\mathbf{w}}_{MN}$ approaches ridge regression on the signal: $\langle \hat{\mathbf{w}}_{MN}, \mathbf{x} \rangle \xrightarrow{d_J \to \infty} \langle \hat{\mathbf{w}}_{\lambda_n}, \mathbf{x}_S \rangle$ for almost all \mathbf{x} $\hat{\mathbf{w}}_{\lambda_n} = \underset{\mathbf{w}_S \in \mathbb{R}^{d_S}}{\operatorname{argmin}} \|\mathbf{X}_S \mathbf{w}_S^2 - \mathbf{y}\|^2 + \lambda_n \|\mathbf{w}_S\|^2$ If $\lambda_n = o(n)$, $\hat{\mathbf{w}}_{\lambda_n}$ is consistent: $L_{\mathcal{D}}(\hat{\mathbf{w}}_{\lambda_n}) \xrightarrow{n \to \infty} \sigma^2$

As $d_J \to \infty$, $\hat{\mathbf{w}}_{MN}$ approaches ridge regression on the signal: $\langle \hat{\mathbf{w}}_{MN}, \mathbf{x} \rangle \xrightarrow{d_J \to \infty} \langle \hat{\mathbf{w}}_{\lambda_n}, \mathbf{x}_S \rangle$ for almost all \mathbf{x} $\hat{\mathbf{w}}_{\lambda_n} = \underset{\mathbf{w}_S \in \mathbb{R}^{d_S}}{\operatorname{argmin}} \|\mathbf{X}_S \mathbf{w}_S^2 - \mathbf{y}\|^2 + \lambda_n \|\mathbf{w}_S\|^2$ If $\lambda_n = o(n)$, $\hat{\mathbf{w}}_{\lambda_n}$ is consistent: $L_{\mathcal{D}}(\hat{\mathbf{w}}_{\lambda_n}) \xrightarrow{n \to \infty} \sigma^2$

 $\hat{f w}_{MN}$ is consistent when d_S fixed, $d_J o\infty$, $\lambda_n=o(n)$. Could we have shown that with uniform convergence?

As $d_J \to \infty$, $\hat{\mathbf{w}}_{MN}$ approaches ridge regression on the signal: $\langle \hat{\mathbf{w}}_{MN}, \mathbf{x} \rangle \xrightarrow{d_J \to \infty} \langle \hat{\mathbf{w}}_{\lambda_n}, \mathbf{x}_S \rangle$ for almost all \mathbf{x} $\hat{\mathbf{w}}_{\lambda_n} = \underset{\mathbf{w}_S \in \mathbb{R}^{d_S}}{\operatorname{argmin}} \|\mathbf{X}_S \mathbf{w}_S^2 - \mathbf{y}\|^2 + \lambda_n \|\mathbf{w}_S\|^2$ If $\lambda_n = o(n)$, $\hat{\mathbf{w}}_{\lambda_n}$ is consistent: $L_{\mathcal{D}}(\hat{\mathbf{w}}_{\lambda_n}) \xrightarrow{n \to \infty} \sigma^2$

 $\hat{\mathbf{w}}_{MN}$ is consistent when d_S fixed, $d_J \to \infty$, $\lambda_n = o(n)$. Could we have shown that with uniform convergence?

Theorem: If
$$\lambda_n = o(n)$$
,

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

Theorem: If
$$\lambda_n = o(n)$$
,

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

Theorem: If
$$\lambda_n = o(n)$$
,

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

Theorem: If
$$\lambda_n = o(n)$$
,

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

Proof idea:

Theorem: If
$$\lambda_n = o(n)$$
,

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

Proof idea:

$$L_{\mathcal{D}}(\mathbf{w}) = (\mathbf{w} - \mathbf{w}^*)^{\mathsf{T}} \mathbf{\Sigma}(\mathbf{w} - \mathbf{w}^*) + L_{\mathcal{D}}(\mathbf{w}^*)$$
Theorem: If
$$\lambda_n = o(n)$$
,

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

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

Theorem: If
$$\lambda_n = o(n)$$
,

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

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

Theorem: If
$$\lambda_n = o(n)$$
,

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

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

Theorem: If
$$\lambda_n = o(n)$$
,

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

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

Theorem: If
$$\lambda_n = o(n)$$
,

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

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

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

 $\{\mathbf{w}: \|\mathbf{w}\| \leq B\}$ 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

<u>Theorem</u> (à la [Nagarajan/Kolter, NeurIPS 2019]): For each $\delta \in (0, \frac{1}{2})$, let $\Pr(\mathbf{S} \in \mathcal{S}_{n,\delta}) \geq 1 - \delta$,

<u>Theorem</u> (à la [Nagarajan/Kolter, NeurIPS 2019]): For each $\delta \in (0, \frac{1}{2})$, let $\Pr(\mathbf{S} \in \mathcal{S}_{n,\delta}) \ge 1 - \delta$, $\hat{\mathbf{w}}$ a *natural* consistent interpolator,

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

<u>Theorem</u> (à la [Nagarajan/Kolter, NeurIPS 2019]): For each $\delta \in (0, \frac{1}{2})$, 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: $\hat{\mathbf{w}}_S$ doesn't change if \mathbf{X}_J flips to $-\mathbf{X}_J$. Examples: $\hat{\mathbf{w}}_{MN}$, $\underset{\mathbf{w}:\mathbf{X}\mathbf{w}=\mathbf{y}}{\operatorname{argmin}} \|\mathbf{w}\|_1$, $\underset{\mathbf{w}:\mathbf{X}\mathbf{w}=\mathbf{y}}{\operatorname{argmin}} \|\mathbf{w} - \mathbf{w}^*\|_2$, $\underset{\mathbf{w}:\mathbf{X}\mathbf{w}=\mathbf{y}}{\operatorname{argmin}} f_S(\mathbf{w}_S) + f_J(\mathbf{w}_J)$ with each f convex, $f_J(-\mathbf{w}_J) = f_J(\mathbf{w}_J)$ $\mathbf{w}:\mathbf{X}\mathbf{w}=\mathbf{y}$

<u>Theorem</u> (à la [Nagarajan/Kolter, NeurIPS 2019]): For each $\delta \in (0, \frac{1}{2})$, 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,

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

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

([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: $\hat{\mathbf{w}}_{MN}$, $\underset{\mathbf{w}:\mathbf{X}\mathbf{w}=\mathbf{y}}{\operatorname{argmin}} \|\mathbf{w}\|_1$, $\underset{\mathbf{w}:\mathbf{X}\mathbf{w}=\mathbf{y}}{\operatorname{argmin}} \|\mathbf{w} - \mathbf{w}^*\|_2$, $\underset{\mathbf{w}:\mathbf{X}\mathbf{w}=\mathbf{y}}{\operatorname{argmin}} f_S(\mathbf{w}_S) + f_J(\mathbf{w}_J)$ with each f convex, $f_J(-\mathbf{w}_J) = f_J(\mathbf{w}_J)$ $\mathbf{w}:\mathbf{X}\mathbf{w}=\mathbf{y}$

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

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...
- Give up?

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

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 ∞ for norm balls in our testbed
 - $\lambda_{\max}(\mathbf{\Sigma}-\hat{\mathbf{\Sigma}})$ instead of $\|\mathbf{\Sigma}-\hat{\mathbf{\Sigma}}\|$

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 ∞ for norm balls in our testbed

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

• 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 $\mathcal{F} = \{f : L_{\mathcal{D}}(f) \le L_{\mathcal{D}}(f^*) + \epsilon_{n,\delta}\}$

A broader view of uniform convergence

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

But we only care about interpolators. How about

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

But we only care about interpolators. How about

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

But we only care about interpolators. How about

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

Is this "uniform convergence"?

But we only care about interpolators. How about

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

Is this "uniform convergence"?

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

A broader view of uniform convergence



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

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

Is this "uniform convergence"?

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

But we only care about interpolators. How about

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

Is this "uniform convergence"?

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

Are there analyses like this for $L_{\mathcal{D}}(w^*) > 0$?



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

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

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

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



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

But if we suppose c = 1, would get a novel prediction:

$$\sup_{\|\mathbf{w}\|\leq lpha\|\hat{\mathbf{w}}_{MN}\|,\ L_{\mathbf{S}}(\mathbf{w})=0} L_{\mathcal{D}}(\mathbf{w})\leq lpha^2\left[\sigma^2+o_P(1)
ight]$$
Main result

$$rac{ ext{Theorem:}}{\lim_{n o \infty} \lim_{d_J o \infty} \mathbb{E} \left[egin{array}{c} \sup_{\|\mathbf{w}\| \leq lpha \| \hat{\mathbf{w}}_{MN} \| \ L_{\mathbf{S}}(\mathbf{w}) = 0 \end{array} | L_{\mathcal{D}}(\mathbf{w}) - L_{\mathbf{S}}(\mathbf{w}) |
ight] = lpha^2 \ L_{\mathcal{D}}(\mathbf{w}^*)$$

Main result

$$rac{ ext{Theorem:}}{\sum_{n o \infty} \lim_{d_J o \infty} \mathbb{E} \left[egin{array}{c} \sup_{\|\mathbf{w}\| \leq lpha \| \hat{\mathbf{w}}_{MN} \| \ L_{\mathbf{S}}(\mathbf{w}) = 0 \end{array} | L_{\mathcal{D}}(\mathbf{w}) - L_{\mathbf{S}}(\mathbf{w}) |
ight] = lpha^2 \ L_{\mathcal{D}}(\mathbf{w}^*)$$

- 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



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



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



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



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 $\hat{\mathbf{w}}$...

$$\hat{\mathbf{w}}_{MR} = \operatorname*{argmin}_{\mathbf{w}:\mathbf{X}\mathbf{w}=\mathbf{y}} L_{\mathcal{D}}(\mathbf{w})$$

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

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

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).

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

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 $\sup_{\substack{1 \le \beta \le 4 \\ \mathbf{w}_{MR} \parallel \\ L_{\mathcal{D}}(\mathbf{w}) = L_{\mathcal{D}}(\hat{\mathbf{w}}_{MR}) + \beta \kappa_{X}(\Sigma) \left[\|\hat{\mathbf{w}}_{MR}\|^{2} - \|\hat{\mathbf{w}}_{MN}\|^{2} \right]$ (amount of missed energy) · (available norm) If $\hat{\mathbf{w}}_{MR}$ consistent, everything smaller-norm also consistent iff β term $\rightarrow 0$

Analyzing dual with $\hat{\mathbf{w}}_{MR}$, get without **any** distributional assumptions that $\sup_{\substack{1 \le \beta \le 4 \\ \mathbf{w}_{MR} \parallel \leq \|\hat{\mathbf{w}}_{MR}\| \\ L_{\mathcal{D}}(\mathbf{w}) = L_{\mathcal{D}}(\hat{\mathbf{w}}_{MR}) + \beta \kappa_{X}(\Sigma) \left[\|\hat{\mathbf{w}}_{MR}\|^{2} - \|\hat{\mathbf{w}}_{MN}\|^{2} \right]$ (amount of missed energy) \cdot (available norm) If $\hat{\mathbf{w}}_{MR}$ consistent, everything smaller-norm also consistent iff β 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 $1 \le eta \le 4$ $\sup_{\boldsymbol{\omega}} \quad L_{\mathcal{D}}(\mathbf{w}) = L_{\mathcal{D}}(\hat{\mathbf{w}}_{MR}) + ar{eta} ar{\kappa}_X(\Sigma) \left[\| \hat{\mathbf{w}}_{MR} \|^2 - \| \hat{\mathbf{w}}_{MN} \|^2
ight]$ $\|\mathbf{w}\|{\leq}\|\hat{\mathbf{w}}_{MR}\|$ (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, $L_{\mathcal{D}}(\hat{\mathbf{w}}_{MR})
ightarrow L_{\mathcal{D}}(\mathbf{w}^*)$ $\kappa_X(\mathbf{\Sigma}) pprox rac{\lambda_n}{n} \quad \mathbb{E}ig[\|\hat{\mathbf{w}}_{MR}\|^2 - \|\hat{\mathbf{w}}_{MN}\|^2ig] = rac{\sigma^2 d_S}{\lambda_n} + o\left(1
ight)$

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:

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

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

- "Regular" uniform convergence can't explain consistency of $\hat{\mathbf{w}}_{MN}$
 - Uniform convergence over norm ball can't show any learning
- An "interpolating" uniform convergence bound does
 - Shows low norm is sufficient for interpolation learning here
 - Predicts exact worst-case error as norm grows
- Optimistic/interpolating rates might be able to explain interpolation learning more broadly
 - Need to get the constants on leading terms *exactly* right!
- Analyzing generalization gap via duality may be broadly applicable
 - Can always get upper bounds via weak duality