

# Facets of regularization in high-dimensional learning:

Cross-validation, risk monotonization, and model complexity

Pratik Patil

Carnegie Mellon University

November 2022

## **Committee:**

Ryan Tibshirani (Chair)

Alessandro Rinaldo

Arun Kumar Kuchibhotla

Yuting Wei (University of Pennsylvania)

Arian Maleki (Columbia University)

# Outline

## Overview

### Cross-validation

- Distribution estimation
- Functional estimation
- Discussion and extensions

### Risk monotonization

- Motivation
- Zero-step procedure
- Discussion and extensions

### Model complexity

- Fixed-X degrees of freedom
- Random-X degrees of freedom
- Discussion and extensions

## Conclusion

# Overparametrization in machine learning

Modern machine learning models typically fit a huge number of parameters. Such overparameterization seems to be useful for<sup>1</sup>:

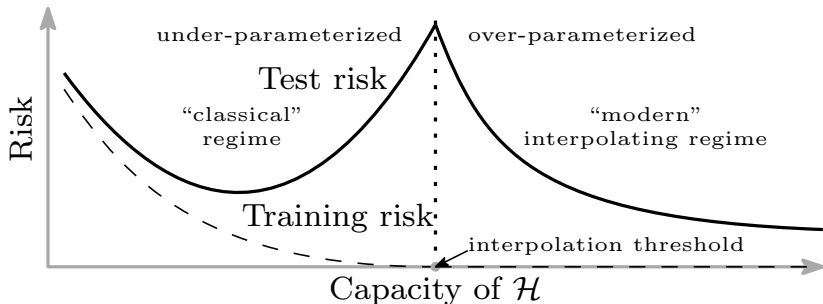
- **Representation**: allows rich, expressive models for diverse real data
- **Optimization**: simple, local optimization methods often find near-optimal solutions to empirical risk minimization problem
- **Generalization**: despite overfitting, models generalize well in practice

This talk is about generalization aspect in overparameterized learning.

---

<sup>1</sup>Credits to Ryan for this nice partition of distinct benefits of overparameterization.

## Peculiar generalization behavior: double descent



Belkin, Hsu, Ma, Mandal, 2018: "Reconciling modern machine learning practice and the bias variance tradeoff"

- The phenomenon is dubbed "double descent" in the risk curve.
- This trend holds for many model classes including linear regression, kernel regression, random forest, boosting, neural networks, etc.

## Recent theoretical developments

Understanding generalization of interpolators in simpler settings:

- Linear regression
  - Hastie, Montanari, Rosset, Tibshirani, 2019
  - Belkin, Hsu, Xu, 2019
  - Muthukumar, Vodrahalli, Sahai, 2019
  - Bartlett, Long, Lugosi, Tsigler, 2019
  - Mei, Montanari, 2019
- Kernel regression
  - Liang, Rakhlin, 2018
  - Liang, Rakhlin, Zhai, 2019
- Local methods
  - Belkin, Hsu, Mitra, 2018
  - Belkin, Rakhlin, Tsybakov, 2018
- and many more ...

Nice survey papers:

- Bartlett, Montanari, and Rakhlin, 2021: “Deep learning: a statistical viewpoint”
- Belkin, 2021: “Fit without fear: remarkable mathematical phenomena of deep learning through the prism of interpolation”

## Motivating questions

We study three operational aspects of overparameterized learning:  
1) cross-validation, 2) risk monotonization, 3) model complexity.

Motivating questions:

1. Does cross-validation still “work” in the overparameterized regime, especially when optimal regularization and train error can be zero?
2. Is it possible to modify any given prediction procedure to mitigate double descent behavior and achieve a monotonic risk behavior?
3. Is there a better and more principled measure of model complexity in general for overparameterized models?

Short answers: YES.

Long answers: Rest of the talk.

# Outline

Overview

## Cross-validation

Distribution estimation

Functional estimation

Discussion and extensions

## Risk monotonization

Motivation

Zero-step procedure

Discussion and extensions

## Model complexity

Fixed-X degrees of freedom

Random-X degrees of freedom

Discussion and extensions

Conclusion

## Motivation and main punchlines

- Given  $\mathcal{D} = \{(x_i, y_i) \in \mathbb{R}^p \times \mathbb{R}, 1 \leq i \leq n\}$ , let  $\hat{\beta}_\lambda$  be **ridge estimator**:

$$\underset{\beta \in \mathbb{R}^p}{\text{minimize}} \sum_{i=1}^n (y_i - x_i^\top \beta)^2 / n + \lambda \|\beta\|_2^2$$

- The **out-of-sample error** of  $\hat{\beta}_\lambda$  is  $y_0 - x_0^\top \hat{\beta}_\lambda$  for a test point  $(x_0, y_0)$

Key question: can we reliably estimate the entire out-of-sample error distribution and its linear and non-linear functionals in high dimensions?

We show, that under proportional asymptotics, almost surely:

- the empirical distributions of re-weighted in-sample errors from leave-one-out and generalized cross-validation converge weakly to the out-of-sample error distribution, even when  $\lambda = 0$
- the plug-in estimators of these empirical distributions consistent for a broad class of linear and non-linear functionals of error distribution



## Overview of high-dimensional ridge regression

- Let  $X \in \mathbb{R}^{n \times p}$  denote feature matrix,  $y \in \mathbb{R}^n$  denote response vector
- Let  $\hat{\beta}_\lambda \in \mathbb{R}^p$  denote the ridge estimator at regularization level  $\lambda$ :

$$\hat{\beta}_\lambda := \arg \min_{\beta \in \mathbb{R}^p} \|y - X\beta\|_2^2/n + \lambda \|\beta\|_2^2$$

- if  $\lambda > 0$ , the problem is convex in  $\beta$  and has an explicit solution:

$$\hat{\beta}_\lambda = (X^T X/n + \lambda I_p)^{-1} X^T y/n$$

- for any  $\lambda \in \mathbb{R}$ , extend the solution using **Moore-Penrose inverse**:

$$\hat{\beta}_\lambda = (X^T X/n + \lambda I_p)^+ X^T y/n$$

- when  $\lambda = 0$ , this reduces to least squares sol with minimum  $\ell_2$  norm; in particular, when  $\text{rank}(X) = n \leq p$ , the solution interpolates data, i.e.  $X\hat{\beta} = y$ , and has minimum  $\ell_2$  norm among all interpolators

## Out-of-sample error distribution and its functionals

- Let  $P_\lambda$  denote distribution of out-of-sample error of  $\widehat{\beta}_\lambda$ :

$$P_\lambda = \mathcal{L}(y_0 - x_0^\top \widehat{\beta}_\lambda \mid X, y),$$

where  $(x_0, y_0)$  is sampled indep from the same training distribution

- a random distribution (conditional on observed data  $X$  and  $y$ )

- Let  $\psi$  denote a functional such that  $P \mapsto \psi(P) \in \mathbb{R}$ :

- Linear functional:

$$\psi(P_\lambda) = \int t(z) dP_\lambda(z) = \mathbb{E}[t(y_0 - x_0^\top \widehat{\beta}_\lambda) \mid X, y],$$

where  $t: \mathbb{R} \rightarrow \mathbb{R}$  is an error function (e.g., squared or absolute error)

- Nonlinear functional:

$$\psi(P_\lambda) = \text{Quantile}(P_\lambda; \tau) = \inf\{z : F_\lambda(z) \geq \tau\},$$

where  $F_\lambda$  denotes the cumulative distribution function of  $P_\lambda$

We construct estimators of  $P_\lambda$  and  $\psi(P_\lambda)$  by suitably extending [leave-one-out cross-validation](#) and [generalized cross-validation](#) procedures.

## Standard leave-one-out and generalized cross-validation

- Leave-one-out cross-validation (LOOCV):
  - for every  $i$ , train on all data except  $(x_i, y_i)$ , call the estimate  $\hat{\beta}_\lambda^{-i}$
  - compute test error on the  $i^{\text{th}}$  point and take average

$$\begin{aligned} \text{loo}(\lambda) &= \frac{1}{n} \sum_{i=1}^n \left( y_i - x_i^T \hat{\beta}_\lambda^{-i} \right)^2 \\ &\stackrel{\text{(shortcut)}}{=} \frac{1}{n} \sum_{i=1}^n \left( \frac{y_i - x_i^T \hat{\beta}_\lambda}{1 - [L_\lambda]_{ii}} \right)^2 \end{aligned}$$

where  $L_\lambda = X(X^T X/n + \lambda I_p)^+ X^T/n$  is the ridge smoothing matrix

- Generalized cross-validation (GCV)
  - same as leave-one-out shortcut but a single re-weighting

$$\text{gcv}(\lambda) = \frac{1}{n} \sum_{i=1}^n \left( \frac{y_i - x_i^T \hat{\beta}_\lambda}{1 - \text{tr}[L_\lambda]/n} \right)^2$$

- Standard LOOCV and GCV are consistent for the expected squared out-of-sample prediction error

## Proposed estimators

Natural estimators for  $P_\lambda$  and  $\psi(P_\lambda)$  building off from GCV and LOOCV.

- Empirical distributions of the GCV, LOO re-weighted errors:

$$\hat{P}_\lambda^{\text{gcv}} = \frac{1}{n} \sum_{i=1}^n \delta \left( \frac{y_i - \mathbf{x}_i^\top \hat{\beta}_\lambda}{1 - \text{tr}[L_\lambda]/n} \right) \quad \text{and} \quad \hat{P}_\lambda^{\text{loo}} = \frac{1}{n} \sum_{i=1}^n \delta \left( \frac{y_i - \mathbf{x}_i^\top \hat{\beta}_\lambda}{1 - [L_\lambda]_{ii}} \right)$$

- When  $\hat{\beta}_\lambda$  is an interpolator, i.e.  $L_\lambda = I_n$ , both estimates are “0/0”<sup>2</sup>; we then define the estimates as their respective limits as  $\lambda \rightarrow 0$ :

$$\hat{P}_0^{\text{gcv}} = \frac{1}{n} \sum_{i=1}^n \delta \left( \frac{[(\mathbf{X}\mathbf{X}^\top)^\dagger \mathbf{y}]_i}{\text{tr}[(\mathbf{X}\mathbf{X}^\top)^\dagger]/n} \right) \quad \text{and} \quad \hat{P}_0^{\text{loo}} = \frac{1}{n} \sum_{i=1}^n \delta \left( \frac{[(\mathbf{X}\mathbf{X}^\top)^\dagger \mathbf{y}]_i}{[(\mathbf{X}\mathbf{X}^\top)^\dagger]_{ii}} \right)$$

- Plug-in GCV and LOO estimators:

$$\hat{\psi}_\lambda^{\text{gcv}} = \psi(\hat{P}_\lambda^{\text{gcv}}) \quad \text{and} \quad \hat{\psi}_\lambda^{\text{loo}} = \psi(\hat{P}_\lambda^{\text{loo}})$$

---

<sup>2</sup>The idea of analytic continuation at  $\lambda = 0$  is from Hastie, Montanari, Rosset, Tibshirani, 2019: “Surprises in high-dimensional ridgeless least squares interpolation”

## Distribution estimation

**Theorem.** Under i.i.d. sampling of  $(x_i, y_i)$ ,  $i = 1, \dots, n$  with

1. feature  $x_i$  decomposable into  $x_i = \Sigma^{1/2} z_i$  where  $z_i$  contains i.i.d. entries with mean 0, variance 1, and finite 4+ moment, and spectrum of  $\Sigma$  is uniformly away from  $r_{\min} > 0$  and  $r_{\max} < \infty$ ,
2. response  $y_i$  with bounded 4+ moment,

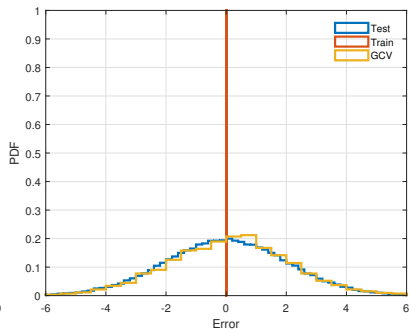
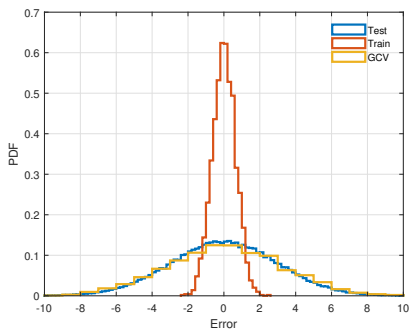
as  $n, p \rightarrow \infty$  such that  $p/n \rightarrow \gamma \in (0, \infty)$ , almost surely, for each  $\lambda > \lambda_{\min} := -(1 - \sqrt{\gamma})^2 r_{\min} \leq 0$ ,

$$\hat{P}_\lambda^{\text{gcv}} \xrightarrow{d} P_\lambda, \quad \text{and} \quad \hat{P}_\lambda^{\text{loo}} \xrightarrow{d} P_\lambda.$$

Remarks:

- Almost sure convergence with respect to the training data
- The regression function does not need to be linear in  $x$
- Amazingly, this results also holds when  $\lambda = 0$  (min-norm estimator)

## Distribution estimation: illustration



- $n = 2500$ ,  $p = 2000$ ,  $p/n = 0.8$
- $\lambda = 0$ , i.e., least squares
- $n = 2500$ ,  $p = 5000$ ,  $p/n = 2$
- $\lambda = 0$ , i.e., min-norm estimator, zero in-sample errors

## Linear functional estimation (pointwise in $\lambda$ )

- Let  $T_\lambda$  be a linear functional of the out-of-sample error distribution:

$$T_\lambda = \mathbb{E}[t(y_0 - x_0^T \hat{\beta}_\lambda) \mid X, y]$$

- Let  $\hat{T}_\lambda^{\text{gcv}}$  and  $\hat{T}_\lambda^{\text{loo}}$  be plug-in estimators from GCV and LOOCV:

$$\hat{T}_\lambda^{\text{gcv}} = \frac{1}{n} \sum_{i=1}^n t \left( \frac{y_i - x_i^T \hat{\beta}_\lambda}{1 - \text{tr}[L_\lambda]/n} \right) \quad \text{and} \quad \hat{T}_\lambda^{\text{loo}} = \frac{1}{n} \sum_{i=1}^n t \left( \frac{y_i - x_i^T \hat{\beta}_\lambda}{1 - [L_\lambda]_{ii}} \right)$$

**Theorem.** For error functions  $t : \mathbb{R} \rightarrow \mathbb{R}$

- that are continuous
- have quadratic growth, i.e., there exist constants  $a, b, c > 0$  such that  $|t(z)| \leq az^2 + b|z| + c$  for any  $z \in \mathbb{R}$ ,

as  $n, p \rightarrow \infty$  with  $p/n \rightarrow \gamma \in (0, \infty)$ , for  $\lambda > \lambda_{\min} := -(1 - \sqrt{\gamma})^2 r_{\min}$ ,

$$\hat{T}_\lambda^{\text{gcv}} \xrightarrow{\text{a.s.}} T_\lambda, \quad \text{and} \quad \hat{T}_\lambda^{\text{loo}} \xrightarrow{\text{a.s.}} T_\lambda.$$

## Linear functional estimation (uniform in $\lambda$ )

**Theorem.** For error functions  $t : \mathbb{R} \rightarrow \mathbb{R}$

1. that are differentiable
2. have derivative with linear growth rate, i.e., there exist constants  $g, h > 0$  such that  $|t'(z)| \leq g|z| + h$  for any  $z \in \mathbb{R}$

as  $n, p \rightarrow \infty$  with  $p/n \rightarrow \gamma \in (0, \infty)$  for any compact set  $\Lambda \subseteq (\lambda_{\min}, \infty)$ ,

$$\sup_{\lambda \in \Lambda} |\widehat{T}_{\lambda}^{\text{gcv}} - T_{\lambda}| \xrightarrow{\text{a.s.}} 0, \quad \text{and} \quad \sup_{\lambda \in \Lambda} |\widehat{T}_{\lambda}^{\text{loo}} - T_{\lambda}| \xrightarrow{\text{a.s.}} 0.$$

Remarks:

- Special case of  $t(r) = r^2$  exploits bias-variance decomposition
- No bias-variance decomposition for general error functions and result requires a different proof technique via leave-one-out arguments
- Using uniformity arguments, the result can be extended for non-linear variational functionals



## Application: quantile estimation

- Quantile of the out-of-sample error distribution:

$$Q_\lambda(\tau) = \text{Quantile}(y_0 - x_0^T \hat{\beta}_\lambda; \tau) = \arg \min_{u \in \mathbb{R}} \mathbb{E}[t_u(y_0 - x_0^T \hat{\beta}_\lambda; \tau) \mid X, y]$$

where  $t_u(y - x_0^T \hat{\beta}_\lambda; \tau)$  is  $\tau$ -tilted pin-ball loss function with shift  $u$

- Empirical quantiles  $\hat{Q}^{\text{gcv}}$  and  $\hat{Q}^{\text{loo}}$  (of  $\hat{P}_\lambda^{\text{gcv}}$  and  $\hat{P}_\lambda^{\text{loo}}$ )  $\xrightarrow{\text{a.s.}}$   $Q_\lambda$
- Estimated quantiles can be used to construct prediction intervals:

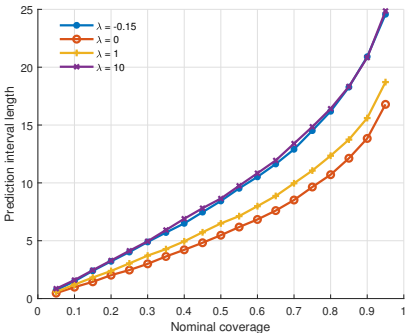
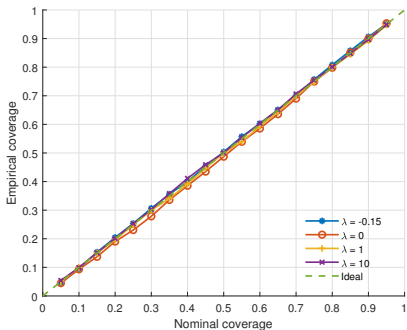
$$\mathcal{I}_\lambda^{\text{gcv}} = [x_0^T \hat{\beta}_\lambda - \hat{Q}_\lambda^{\text{gcv}}(\tau_l), x_0^T \hat{\beta}_\lambda + \hat{Q}_\lambda^{\text{gcv}}(\tau_u)] \quad \text{and} \quad \mathcal{I}_\lambda^{\text{loo}}$$

Such intervals have correct coverage conditional on the training data:

**Corollary.** Under proportional asymptotics, almost surely

$$\mathbb{P}(y_0 \in \mathcal{I}_\lambda^{\text{gcv}} \mid X, y) \xrightarrow{\text{a.s.}} 1 - \alpha, \quad \text{and} \quad \mathbb{P}(y_0 \in \mathcal{I}_\lambda^{\text{loo}} \mid X, y) \xrightarrow{\text{a.s.}} 1 - \alpha.$$

## Prediction intervals: illustration (coverage and length)



- $n = 2500, p = 5000$
- Features: autoregressive feature covariance structure
- Signal: latent signal aligned with the principal eigenvector
- Coverage nearly exact, even for  $\lambda = 0$ !
- The case of  $\lambda = 0$  provides the minimum interval length!

## Discussion and extensions

Take-away from this work: empirical distributions of GCV and LOOCV track out-of-sample error distribution and a wide class of its functionals for ridge regression under proportional asymptotics framework

Key relation that we exploit:

$$y_i - x_i^\top \hat{\beta}_{-i,\lambda} = \frac{y_i - x_i^\top \hat{\beta}_\lambda}{1 - [L_\lambda]_{ii}} \approx \frac{y_i - x_i^\top \hat{\beta}_\lambda}{1 - \text{tr}[L_\lambda]/n}$$
$$y_i - x_i^\top \hat{\beta}_{-i,0} = \frac{[(XX^\top)^\dagger y]_i}{[(XX^\top)^\dagger]_{ii}} \approx \frac{[(XX^\top)^\dagger y]_i}{\text{tr}[(XX^\top)^\dagger]/n}$$

Extensions:

- Generalized ridge/less regression through structural equivalences
- Kernel ridge/less regression through risk equivalences

# Outline

Overview

Cross-validation

- Distribution estimation

- Functional estimation

- Discussion and extensions

Risk monotonization

- Motivation

- Zero-step procedure

- Discussion and extensions

Model complexity

- Fixed-X degrees of freedom

- Random-X degrees of freedom

- Discussion and extensions

Conclusion

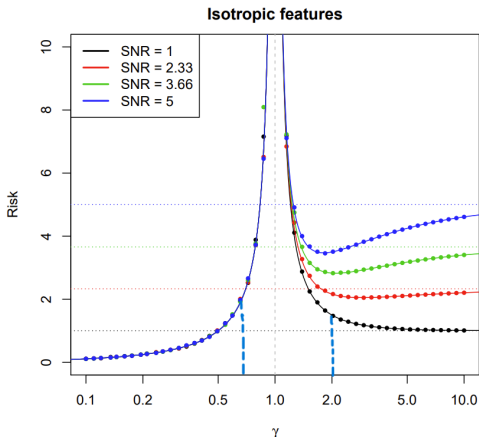
## Motivation and main punchlines

- When the data comprises of i.i.d. observations, we expect that more data will help in prediction or estimation.
- Double or multiple descent behaviour implies that for fixed feature size  $p$  (large value), as sample size increases the risk first decreases and then increases. **More data can hurt!**
- A procedure leading to worse risk as the number of observations increases is not using the data properly.

Key question: Can we modify any prediction procedure to mitigate the double or multiple descent behavior and achieve a monotonic risk behavior?

We propose two methods, dubbed zero-step and one-step, that take an input an arbitrary procedure and return a modified procedure that has a monotonic risk behavior. The main idea is that of subsampling.

# Motivation and the problem



**Figure:** Risk of the minimum  $\ell_2$ -norm least squares as a function of  $p/n \approx \gamma$ .

Hastie, Montanari, Rosset, Tibshirani, 2019: "Surprises in high-dimensional ridgeless least squares interpolation"

## The problem

- Given a number of observations ( $n$ ) and a number of features ( $p$ ), how do we know if a lesser number of observations would actually yield a better risk?
- What is the best sample size to reduce the dataset in order to attain the best possible risk?

**Solution:** cross-validation.

## Basic idea of zero-step procedure

Given any arbitrary prediction procedure at a given aspect ratio  $\gamma = p/n$ :

1. Risk estimation: construct a (dense grid of) aspect ratios  $\geq \gamma$  by using datasets of sizes smaller than  $n$ , and estimate risks on test set
2. Model selection: select aspect ratio that delivers the smallest estimated risk and return the corresponding predictor
3. Risk monotonization: show that the risk profile of the resulting procedure is asymptotically monotone in the aspect ratio

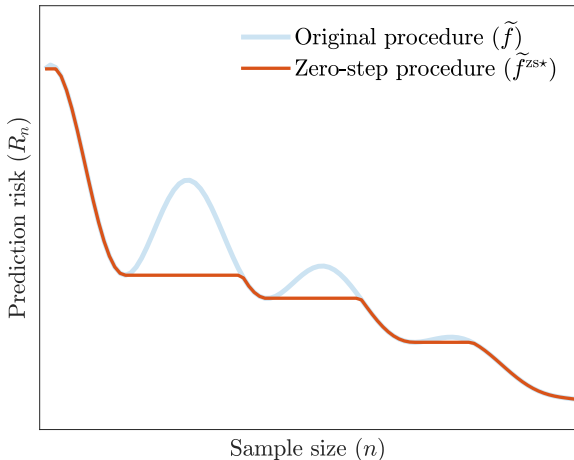
Method highlights:

- applicable to generic (e.g. **black-box**) prediction methods and common classification and regression loss functions
- **model agnostic** and requires **minimal distributional assumptions**
- works for procedures with **diverging risks** at some aspect ratios



## Risk monotonization illustration

If  $R_n$  represents the “risk” of a procedure at sample size  $n$ , then by risk monotonization we mean a procedure with risk  $\min_{m \leq n} R_m$ .



## Risk monotonization guarantee

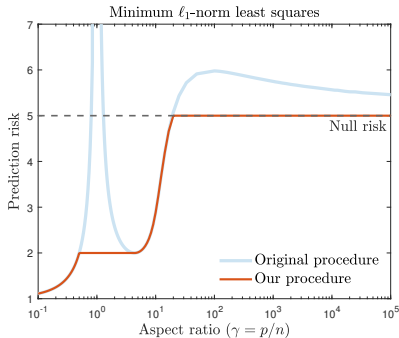
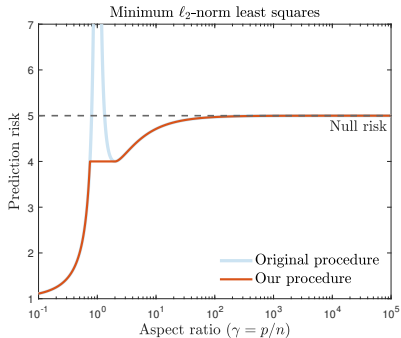
**Theorem.** Under the proportional asymptotics regime ( $p/n \rightarrow \gamma$ ), and a mild assumption on the convergence of the prediction risk of  $\hat{f}$  trained on datasets with a limiting aspect ratio  $\zeta$  converges to  $R^{\text{det}}(\zeta; \hat{f})$ , we show:

$$R(\hat{f}^{\text{cv}}) = \inf_{\zeta \in [\gamma, \infty]} R^{\text{det}}(\zeta; \hat{f}) \times (1 + o_p(1)).$$

This shows that the zero-step predictor has a **monotone risk** in terms of the sample size and hence with respect to the limiting aspect ratio.

This is a **model-free result** in that no parametric model is assumed for the data. This is unlike most results in overparametrized learning which require stringent assumptions.

## Risk monotonization (illustration)



- minimum  $\ell_2$ -norm least squares (ridgeless regression)

- minimum  $\ell_1$ -norm least squares (lassoless regression)

## Discussion and extensions

Take-aways:

- We have introduced the **zero-step prediction procedure** that provably monotonizes the risk of a given predictor.
- The main idea is **cross-validation** based on test data, but with splitting done so as to maintain the limiting aspect ratio.

Extensions:

- We also introduce a **one-step prediction procedure** inspired by classical one-step estimator that improves on zero-step procedure (similar to boosting)
- Both zero-step and one-step procedures can be further improved by **multiple subsamplings and averaging** (similar to bagging)

# Outline

Overview

Cross-validation

- Distribution estimation

- Functional estimation

- Discussion and extensions

Risk monotonization

- Motivation

- Zero-step procedure

- Discussion and extensions

**Model complexity**

- Fixed-X degrees of freedom

- Random-X degrees of freedom

- Discussion and extensions

Conclusion

## Motivation and main punchlines

Key question: is there a principled measure of model complexity in general for overparameterized models?

- Propose measures of model complexity that are:
  - **algorithm-specific** and applies for any prediction algorithm
  - produce a number between **0 and  $n$**  (the number of observations)
- Two variants of model complexities are:
  - **emergent** model complexity that depends on the prediction algorithm as well as underlying the regression function
  - **intrinsic** model complexity that depends on the prediction algorithm only and its adaptability to pure noise

Based on ideas from **optimism theory** and **degrees of freedom**.

## Fixed-X degrees of freedom

Consider data  $(x_i, y_i) \in \mathbb{R}^p \times \mathbb{R}$ ,  $i = 1, \dots, n$  such that  $y_i = f(x_i) + \varepsilon_i$  where  $f : \mathbb{R}^p \rightarrow \mathbb{R}$  is regression function,  $\varepsilon_i$  has mean 0 and variance  $\sigma^2$ .

Let  $\mathcal{A}$  be any fitting algorithm that maps  $\{(x_i, y_i)\}_{i=1}^n \xrightarrow{\mathcal{A}} \hat{f}$ .

The **degrees of freedom** of predictor  $\hat{f}$  is defined as

$$\text{DofF}(\hat{f}) = \sum_{i=1}^n \text{Cov}(y_i, \hat{f}(x_i)) / \sigma^2 = \text{tr} [\text{Cov}(y, \hat{f}(X))] / \sigma^2,$$

where  $y$ : response vector,  $X$ : feature matrix,  $\hat{f}(X)$ : predicted response

Where does squared error loss come into play?

$$\underbrace{\mathbb{E} \left[ \sum_{i=1}^n (\tilde{y}_i - \hat{f}(x_i))^2 \right]}_{\text{fixed-X prediction error} =: \text{ErrF}(\hat{f})} - \underbrace{\mathbb{E} \left[ \sum_{i=1}^n (y_i - \hat{f}(x_i))^2 \right]}_{\text{expected training error} =: \text{ErrT}(\hat{f})} = 2\sigma^2 \text{DofF}(\hat{f})$$

## Fixed-X degrees of freedom in linear regression

- Suppose  $p \leq n$  and  $X$  has full (column) rank, and we take  $\hat{f}$  to be **ordinary least squares** predictor  $\hat{f}(X) = X\hat{\beta}$ , where

$$\hat{\beta} = \arg \min_{\beta \in \mathbb{R}^p} \|y - X\beta\|_2^2 = (X^T X)^{-1} X^T y.$$

$$\text{DofF}(\hat{f}) = \text{tr}[\text{Cov}(y, X\hat{\beta})]/\sigma^2 = \text{tr}[\sigma^2 X(X^T X)^{-1} X^T]/\sigma^2 = p.$$

- Suppose  $p \geq n$  and  $X$  has full (row) rank, and we take  $\hat{f}$  to be **min  $\ell_2$ -norm least squares** predictor  $\hat{f}(X) = X\hat{\beta}$ , where

$$\hat{\beta} = \arg \min_{\beta \in \mathbb{R}^p} \{\|\beta\|_2 : X\beta = y\} = (X^T X)^\dagger X^T y.$$

$$\text{DofF}(\hat{f}) = \text{tr}[\text{Cov}(y, X\hat{\beta})]/\sigma^2 = \text{tr}[\sigma^2 X^T (X X^T)^{-1} X]/\sigma^2 = n.$$

Thus,  $\text{DofF}(\hat{f})$  is  $p$  for  $p \leq n$ , but is always  $n$  for  $p \geq n$  (not meaningful).

This is **fixed-X** degrees of freedom. How to extend for **random-X setting**?



## Re-interpreting fixed-X degrees of freedom

Fixed-X degrees of freedom is a standard algorithm specific measure of complexity, but no notion of random-X degrees of freedom we know of.

We cast fixed-X degrees of freedom from a **different perspective**.

- Define fixed-X optimism of  $\hat{f}$  by  $\text{OptF}(\hat{f}) = \text{ErrF}(\hat{f}) - \text{ErrT}(\hat{f})$ .
- Consider the following family of “reference” models:
  - $\mathcal{A}^{\text{ref}}$  is the **least squares** reference algorithm,
  - $(U_k, \nu)$  is random design with  $k$  features, and noise with level  $\sigma^2$ .
- Recall that  $\text{DofF}(\mathcal{A}^{\text{ref}}(U_k, \nu)) = k$  so long as  $\text{rank}(U_k) = k$ .
- Thus, for a fitting procedure  $\hat{f} = \mathcal{A}(X, y)$ ,  $\text{DofF}(\hat{f})$  is also equal to the value of  $k$  that satisfy the following relation:

$$\text{OptF}(\mathcal{A}(X, y)) = \text{OptF}(\mathcal{A}^{\text{ref}}(U_k, \nu)) \quad (\text{dfF})$$

## Emergent random-X degrees of freedom

“Matching optimism” interpretation can be extended to random-X setting and leads to the definition of random-X degrees of freedom.

- Define **random-X optimism** of  $\hat{f}$  by  $\text{OptR}(\hat{f}) = \text{ErrR}(\hat{f}) - \text{ErrT}(\hat{f})$ , where  $\text{ErrR}(\hat{f}) = \mathbb{E}[(y_0 - \hat{f}(x_0))^2]$  is the random-X prediction error.
- We thus define the random-X degrees of freedom,  $\text{DofR}(\hat{f})$ , of any predictor  $\hat{f} = \mathcal{A}(X, y)$ , as the value of  $k$  for which the following relation holds:

$$\text{OptR}(\mathcal{A}(X, y)) = \text{OptR}(\mathcal{A}^{\text{ref}}(U_k, \nu)) \quad (\text{dfR, emergent})$$

Recall here:

- $\mathcal{A}^{\text{ref}}$  is the **least squares** reference algorithm,
- $(U_k, \nu)$  is random design with  $k$  features, and noise with level  $\sigma^2$ .

We call the measure **emergent** random-X degrees of freedom.

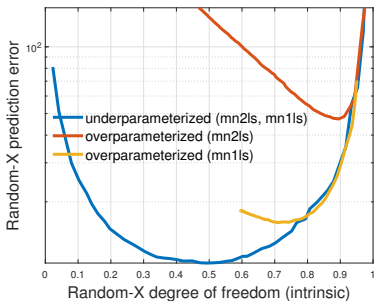
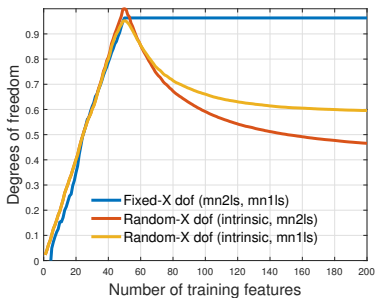
## Intrinsic random-X degrees of freedom

- The emergent random-X degrees of freedom,  $\text{DofR}(\hat{f})$ , depends of both the the predictor  $\hat{f}$  and the underlying regression function  $f$ .
- When matching optimisms, the observed random-X optimism of  $\hat{f}$  consists of bias, which may inflate the degrees of freedom.
- We thus also define **intrinsic** random-X degrees of freedom, denoted by  $\text{DofR}^i$ , as the  $k$  for which the following relation holds:

$$\text{OptR}(\mathcal{A}(X, \nu)) = \text{OptR}(\mathcal{A}^{\text{ref}}(U_k, \nu)) \quad (\text{dfR, intrinsic})$$

The intrinsic random-X degrees of freedom measures the **inherent** complexity of the predictor  $\hat{f}$  in terms of overfitting to “pure noise”.

## Random-X degrees of freedom illustration



- Fixed data with  $n = 50$  and response non-linear in  $p = 200$  features
- Model class: estimators fitted on nested subsets of 1 to 200 features
- Fixed-X: increase then constant; random-X: increase then decrease
- Underparameterized: *U*-curve; overparameterized: also *U*-curve!
- Punchline: reparameterize overparameterized to underparameterized

## Discussion and future directions

A high-level view of the work:

- Suppose we are given a family of models for which we want a complexity measure under a **specific error metric**.
- Construct a family of **“reference” models** spanning same optimisms.
- Find the model in the reference family that is closest to the **observed optimism**. Declare complexity as complexity of that reference model.

Key relation:

$$\text{OptR}(\hat{f}) = \text{OptR}(\hat{f}^{\text{ref}})$$

Future directions:

- Attribute total complexity to various components: bias, variance, covariate shift, etc.
- Other error metrics beyond squared error

# Outline

Overview

Cross-validation

- Distribution estimation

- Functional estimation

- Discussion and extensions

Risk monotonization

- Motivation

- Zero-step procedure

- Discussion and extensions

Model complexity

- Fixed-X degrees of freedom

- Random-X degrees of freedom

- Discussion and extensions

Conclusion

## Motivating thesis questions with take-aways

We studied three operational aspects of overparameterized learning:  
1) cross-validation, 2) risk monotonization, 3) model complexity.

1. Cross-validation still works in the overparameterized regime, especially when optimal regularization and train error can be zero for ridge regression through analytic continuation.
2. It is possible to modify any given prediction procedure to mitigate double descent behavior and achieve a monotonic risk behavior through subsampling and cross-validation.
3. There is a principled measure of model complexity in general for overparameterized models in the form of random-X degrees of freedom.

Thanks for listening!

Questions/comments/thoughts?



# BIG THANKS!

- Ryan
- Committee: Ale, Arun, Yuting, Arian
- Collaborators
- Faculty
- Staff
- Students
- Funding agency