Mitigating multiple descents:

A model-agnostic framework for risk monotonization

Pratik Patil

University of California Berkeley

Rutgers University Seminar November 2023

Based on joint work with the following amazing collaborators:

- Arun Kuchibhotla (Carnegie Mellon University)
- Yuting Wei (University of Pennsylvania)
- Alessandro Rinaldo (University of Texas)
- Jin-Hong Du (Carnegie Mellon University)

- Mitigating multiple descents: A model-agnostic framework for risk monotonization (joint with Arun Kuchibhotla, Yuting Wei, Alessandro Rinaldo) [benefits of subsampling]
- 2. Bagging in overparameterized learning: Risk characterization and risk monotonization (joint with Jin-Hong Du, Arun Kuchibhotla) [benefits of ensembling]
- 3. Generalized equivalences between subsampling and ridge regularization (joint with Jin-Hong Du) [connections to ridge]

- Mitigating multiple descents: A model-agnostic framework for risk monotonization (joint with Arun Kuchibhotla, Yuting Wei, Alessandro Rinaldo) [benefits of subsampling]
- 2. Bagging in overparameterized learning: Risk characterization and risk monotonization (joint with Jin-Hong Du, Arun Kuchibhotla) [benefits of ensembling]
- 3. Generalized equivalences between subsampling and ridge regularization (joint with Jin-Hong Du) [connections to ridge]

- 1. Mitigating multiple descents: A model-agnostic framework for risk monotonization (joint with Arun Kuchibhotla, Yuting Wei, Alessandro Rinaldo) [benefits of subsampling]
- 2. Bagging in overparameterized learning: Risk characterization and risk monotonization (joint with Jin-Hong Du, Arun Kuchibhotla) [benefits of ensembling]
- 3. Generalized equivalences between subsampling and ridge regularization (joint with Jin-Hong Du) [connections to ridge]

- 1. Mitigating multiple descents: A model-agnostic framework for risk monotonization (joint with Arun Kuchibhotla, Yuting Wei, Alessandro Rinaldo) [benefits of subsampling]
- 2. Bagging in overparameterized learning: Risk characterization and risk monotonization (joint with Jin-Hong Du, Arun Kuchibhotla) [benefits of ensembling]
- 3. Generalized equivalences between subsampling and ridge regularization (joint with Jin-Hong Du) [connections to ridge]

Outline

Overview of overparameterization

Double descent Current theoretical understanding Case study of linear regression

Risk monotonization

Motivation Zero-step procedure Takeaways and extensions

Bagging analysis

Motivation Risk characterization Optimal subsample size

Connections to ridge regularization

Risk and structural equivalences Implications of equivalences Discussion and extensions

Conclusion

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

- 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

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

- 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

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

- 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

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

- 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

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

- 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

An influential experiment



"Understanding deep learning requires rethinking generalization" Zhang, Bengio, Hardt, Recht, Vinyals, 2017

- CIFAR10 data (60,000 images [32 imes 32]) with artificial label noise
- Three neural network architectures (with number of parameters): Inception (1,649,402), AlexNet (1,387,786), MLP 1x512 (1,209,866

An influential experiment



"Understanding deep learning requires rethinking generalization" Zhang, Bengio, Hardt, Recht, Vinyals, 2017

- CIFAR10 data (60,000 images $[32 \times 32]$) with artificial label noise
- Three neural network architectures (with number of parameters): Inception (1,649,402), AlexNet (1,387,786), MLP 1x512 (1,209,866

An influential experiment



"Understanding deep learning requires rethinking generalization" Zhang, Bengio, Hardt, Recht, Vinyals, 2017

- CIFAR10 data (60,000 images [32 imes 32]) with artificial label noise
- Three neural network architectures (with number of parameters): Inception (1,649,402), AlexNet (1,387,786), MLP 1x512 (1,209,866)

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.

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.

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.

- 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 ...

- 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 ...

- 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 ...

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 ...

- 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 ...

• In nearly all applications, current practice suggests we should design models to be massively overparametrized

- Once trained (typically by SGD), these models interpolate the training data (achieve zero training error)
- Still they are capable of having (often do have) good test error

Current understanding of this? In full theoretical rigor, not great.

- 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"

- In nearly all applications, current practice suggests we should design models to be massively overparametrized
- Once trained (typically by SGD), these models interpolate the training data (achieve zero training error)
- Still they are capable of having (often do have) good test error

Current understanding of this? In full theoretical rigor, not great.

- 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"

- In nearly all applications, current practice suggests we should design models to be massively overparametrized
- Once trained (typically by SGD), these models interpolate the training data (achieve zero training error)
- Still they are capable of having (often do have) good test error

Current understanding of this? In full theoretical rigor, not great.

- 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"

- In nearly all applications, current practice suggests we should design models to be massively overparametrized
- Once trained (typically by SGD), these models interpolate the training data (achieve zero training error)
- Still they are capable of having (often do have) good test error

Current understanding of this? In full theoretical rigor, not great.

- 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"

- In nearly all applications, current practice suggests we should design models to be massively overparametrized
- Once trained (typically by SGD), these models interpolate the training data (achieve zero training error)
- Still they are capable of having (often do have) good test error

Current understanding of this? In full theoretical rigor, not great.

- 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"

- In nearly all applications, current practice suggests we should design models to be massively overparametrized
- Once trained (typically by SGD), these models interpolate the training data (achieve zero training error)
- Still they are capable of having (often do have) good test error

Current understanding of this? In full theoretical rigor, not great.

- 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"

- In nearly all applications, current practice suggests we should design models to be massively overparametrized
- Once trained (typically by SGD), these models interpolate the training data (achieve zero training error)
- Still they are capable of having (often do have) good test error

Current understanding of this? In full theoretical rigor, not great.

- 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"

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

Given i.i.d. training data $(x_i, y_i) \in \mathbb{R}^p \times \mathbb{R}$, i = 1, ..., n, where

$$y_i = \underbrace{x_i^\top \beta}_{f(x_i)} + \epsilon_i, \quad x_i^\top \beta \perp \epsilon_i$$

"Ridgeless" least squares estimator of y on X (which has rows x_i):

$$\widehat{\beta} = (X^{\top}X)^{\dagger}X^{\top}y = \lim_{\lambda \to 0^+} \operatorname*{arg\,min}_{b \in \mathbb{R}^p} \left\{ \frac{1}{n} \|y - Xb\|_2^2 + \lambda \|b\|_2^2 \right\}$$

Let $\sigma^2 = \operatorname{Var}(\epsilon_i)$ [noise energy], $\rho^2 = \mathbb{E}f(x_i)^2$ [signal energy]. Under simplifying assumptions, as $n, p \to \infty$ with $p/n \to \gamma$:

$$\mathbb{E}(f(x_0) - x_0^\top \widehat{\beta})^2 \to \begin{cases} \sigma^2 \frac{\gamma}{1 - \gamma} & \gamma < 1\\ \rho^2 \frac{\gamma - 1}{\gamma} + \sigma^2 \frac{1}{\gamma - 1} & \gamma > 1 \end{cases}$$

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

Given i.i.d. training data $(x_i, y_i) \in \mathbb{R}^p \times \mathbb{R}$, $i = 1, \dots, n$, where

$$y_i = \underbrace{x_i^\top \beta}_{f(x_i)} + \epsilon_i, \quad x_i^\top \beta \perp \epsilon_i$$

"Ridgeless" least squares estimator of y on X (which has rows x_i):

$$\widehat{\beta} = (X^{\top}X)^{\dagger}X^{\top}y = \lim_{\lambda \to 0^+} \operatorname*{arg\,min}_{b \in \mathbb{R}^p} \left\{ \frac{1}{n} \|y - Xb\|_2^2 + \lambda \|b\|_2^2 \right\}$$

Let $\sigma^2 = \operatorname{Var}(\epsilon_i)$ [noise energy], $\rho^2 = \mathbb{E}f(x_i)^2$ [signal energy]. Under simplifying assumptions, as $n, p \to \infty$ with $p/n \to \gamma$:

$$\mathbb{E}(f(x_0) - x_0^\top \widehat{\beta})^2 \to \begin{cases} \sigma^2 \frac{\gamma}{1 - \gamma} & \gamma < 1\\ \rho^2 \frac{\gamma - 1}{\gamma} + \sigma^2 \frac{1}{\gamma - 1} & \gamma > 1 \end{cases}$$

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

Given i.i.d. training data $(x_i, y_i) \in \mathbb{R}^p \times \mathbb{R}$, $i = 1, \dots, n$, where

$$y_i = \underbrace{x_i^\top \beta}_{f(x_i)} + \epsilon_i, \quad x_i^\top \beta \perp \epsilon_i$$

"Ridgeless" least squares estimator of y on X (which has rows x_i):

$$\widehat{\beta} = (X^{\top}X)^{\dagger}X^{\top}y = \lim_{\lambda \to 0^+} \operatorname*{arg\,min}_{b \in \mathbb{R}^p} \left\{ \frac{1}{n} \|y - Xb\|_2^2 + \lambda \|b\|_2^2 \right\}$$

Let $\sigma^2 = \text{Var}(\epsilon_i)$ [noise energy], $\rho^2 = \mathbb{E}f(x_i)^2$ [signal energy]. Under simplifying assumptions, as $n, p \to \infty$ with $p/n \to \gamma$:

$$\mathbb{E}(f(x_0) - x_0^\top \widehat{\beta})^2 \to \begin{cases} \sigma^2 \frac{\gamma}{1-\gamma} & \gamma < 1\\ \rho^2 \frac{\gamma-1}{\gamma} + \sigma^2 \frac{1}{\gamma-1} & \gamma > 1 \end{cases}$$

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

Given i.i.d. training data $(x_i, y_i) \in \mathbb{R}^p \times \mathbb{R}$, $i = 1, \dots, n$, where

$$y_i = \underbrace{x_i^\top \beta}_{f(x_i)} + \epsilon_i, \quad x_i^\top \beta \perp \epsilon_i$$

"Ridgeless" least squares estimator of y on X (which has rows x_i):

$$\widehat{\beta} = (X^{\top}X)^{\dagger}X^{\top}y = \lim_{\lambda \to 0^+} \operatorname*{arg\,min}_{b \in \mathbb{R}^p} \left\{ \frac{1}{n} \|y - Xb\|_2^2 + \lambda \|b\|_2^2 \right\}$$

Let $\sigma^2 = \text{Var}(\epsilon_i)$ [noise energy], $\rho^2 = \mathbb{E}f(x_i)^2$ [signal energy]. Under simplifying assumptions, as $n, p \to \infty$ with $p/n \to \gamma$:

$$\mathbb{E}(f(x_0) - x_0^\top \widehat{\beta})^2 \to \begin{cases} \sigma^2 \frac{\gamma}{1 - \gamma} & \gamma < 1\\ \rho^2 \frac{\gamma - 1}{\gamma} + \sigma^2 \frac{1}{\gamma - 1} & \gamma > 1 \end{cases}$$

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

Given i.i.d. training data $(x_i, y_i) \in \mathbb{R}^p \times \mathbb{R}$, $i = 1, \dots, n$, where

$$y_i = \underbrace{x_i^\top \beta}_{f(x_i)} + \epsilon_i, \quad x_i^\top \beta \perp \epsilon_i$$

"Ridgeless" least squares estimator of y on X (which has rows x_i):

$$\widehat{\beta} = (X^{\top}X)^{\dagger}X^{\top}y = \lim_{\lambda \to 0^+} \operatorname*{arg\,min}_{b \in \mathbb{R}^p} \left\{ \frac{1}{n} \|y - Xb\|_2^2 + \lambda \|b\|_2^2 \right\}$$

Let $\sigma^2 = \text{Var}(\epsilon_i)$ [noise energy], $\rho^2 = \mathbb{E}f(x_i)^2$ [signal energy]. Under simplifying assumptions, as $n, p \to \infty$ with $p/n \to \gamma$:

$$\mathbb{E}(f(x_0) - x_0^\top \widehat{\beta})^2 \to \begin{cases} \sigma^2 \frac{\gamma}{1-\gamma} & \gamma < 1\\ \rho^2 \frac{\gamma-1}{\gamma} + \sigma^2 \frac{1}{\gamma-1} & \gamma > 1 \end{cases}$$

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

Given i.i.d. training data $(x_i, y_i) \in \mathbb{R}^p \times \mathbb{R}$, $i = 1, \dots, n$, where

$$y_i = \underbrace{x_i^\top \beta}_{f(x_i)} + \epsilon_i, \quad x_i^\top \beta \perp \epsilon_i$$

"Ridgeless" least squares estimator of y on X (which has rows x_i):

$$\widehat{\beta} = (X^{\top}X)^{\dagger}X^{\top}y = \lim_{\lambda \to 0^+} \operatorname*{arg\,min}_{b \in \mathbb{R}^p} \left\{ \frac{1}{n} \|y - Xb\|_2^2 + \lambda \|b\|_2^2 \right\}$$

Let $\sigma^2 = \operatorname{Var}(\epsilon_i)$ [noise energy], $\rho^2 = \mathbb{E}f(x_i)^2$ [signal energy]. Under simplifying assumptions, as $n, p \to \infty$ with $p/n \to \gamma$:

$$\mathbb{E}(f(x_0) - x_0^\top \widehat{\beta})^2 \to \begin{cases} \sigma^2 \frac{\gamma}{1-\gamma} & \gamma < 1\\ \rho^2 \frac{\gamma-1}{\gamma} + \sigma^2 \frac{1}{\gamma-1} & \gamma > 1 \end{cases}$$

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

Given i.i.d. training data $(x_i, y_i) \in \mathbb{R}^p \times \mathbb{R}$, $i = 1, \dots, n$, where

$$y_i = \underbrace{x_i^\top \beta}_{f(x_i)} + \epsilon_i, \quad x_i^\top \beta \perp \epsilon_i$$

"Ridgeless" least squares estimator of y on X (which has rows x_i):

$$\widehat{\beta} = (X^{\top}X)^{\dagger}X^{\top}y = \lim_{\lambda \to 0^+} \operatorname*{arg\,min}_{b \in \mathbb{R}^p} \left\{ \frac{1}{n} \|y - Xb\|_2^2 + \lambda \|b\|_2^2 \right\}$$

Let $\sigma^2 = \text{Var}(\epsilon_i)$ [noise energy], $\rho^2 = \mathbb{E}f(x_i)^2$ [signal energy]. Under simplifying assumptions, as $n, p \to \infty$ with $p/n \to \gamma$:

$$\mathbb{E}(f(x_0) - x_0^\top \widehat{\beta})^2 \to \begin{cases} \sigma^2 \frac{\gamma}{1-\gamma} & \gamma < 1\\ \rho^2 \frac{\gamma-1}{\gamma} + \sigma^2 \frac{1}{\gamma-1} & \gamma > 1 \end{cases}$$
Double in linear regression



Here $\sigma^2 = 1$, thus signal-to-noise ratio (SNR) is ρ^2 , and $\gamma = p/n$. Hastie, Montanari, Rosset, Tibshirani, 2019: "Surprises in high-dimensional ridgeless least squares interpolation"

Isotropic features

- The risk first increases as p/n increases up to some threshold and then decreases.
- There are two ways to view this:
 - If p is thought of as fixed (large value), this implies that as sample size increases the risk first decreases and then increases.
 More data hurts.
 - If n is thought of as fixed (large value), this implies that as the number of features/covariates increase the risk first increases and then decreases.

More features do not hurt.

- The risk first increases as p/n increases up to some threshold and then decreases.
- There are two ways to view this:
 - If p is thought of as fixed (large value), this implies that as sample size increases the risk first decreases and then increases.
 More data hurts.
 - If n is thought of as fixed (large value), this implies that as the number of features/covariates increase the risk first increases and then decreases.

More features do not hurt.

- The risk first increases as p/n increases up to some threshold and then decreases.
- There are two ways to view this:
 - If p is thought of as fixed (large value), this implies that as sample size increases the risk first decreases and then increases.
 More data hurts.
 - If n is thought of as fixed (large value), this implies that as the number of features/covariates increase the risk first increases and then decreases.

More features do not hurt.

- The risk first increases as p/n increases up to some threshold and then decreases.
- There are two ways to view this:
 - If p is thought of as fixed (large value), this implies that as sample size increases the risk first decreases and then increases.
 More data hurts.
 - If n is thought of as fixed (large value), this implies that as the number of features/covariates increase the risk first increases and then decreases.

More features do not hurt.

- The risk first increases as p/n increases up to some threshold and then decreases.
- There are two ways to view this:
 - If p is thought of as fixed (large value), this implies that as sample size increases the risk first decreases and then increases.
 More data hurts.
 - If n is thought of as fixed (large value), this implies that as the number of features/covariates increase the risk first increases and then decreases.

More features do not hurt.

- The risk first increases as p/n increases up to some threshold and then decreases.
- There are two ways to view this:
 - If p is thought of as fixed (large value), this implies that as sample size increases the risk first decreases and then increases.
 More data hurts.
 - If n is thought of as fixed (large value), this implies that as the number of features/covariates increase the risk first increases and then decreases.

More features do not hurt.

Outline

Overview of overparameterization

Double descent Current theoretical understanding Case study of linear regression

Risk monotonization

Motivation Zero-step procedure Takeaways and extensions

Bagging analysis

Motivation Risk characterization Optimal subsample size

Connections to ridge regularization

Risk and structural equivalences Implications of equivalences Discussion and extensions

Conclusion

- 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?

- 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?

- 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?

- 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?

- 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?

Method overview and the problem



Isotropic features

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.

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.

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.

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

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

- 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

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

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

- 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

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

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

- 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

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

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

- 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

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

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

- 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

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

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

- 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

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

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

- 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

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

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

- 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

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

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

- 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 $\min_{m \le n} R_m$.



Sample size (n)

Split sample cross-validation

- Given data \mathcal{D}_n of n i.i.d. observations and a prediction procedure \tilde{f} , split \mathcal{D}_n into training data \mathcal{D}_{tr} with $n(1-1/\log n)$ observations and test data \mathcal{D}_{te} with $n/\log n$ observations.
- Note that

$$\lim_n \frac{p}{n} = \lim_n \frac{p}{n(1-1/\log n)}.$$

- For $n^{1/2} \leq k \leq |\mathcal{D}_{tr}|$, obtain a predictor \tilde{f}_k by training \tilde{f} on a subset of \mathcal{D}_{tr} with k observations.
- If p/n converges to γ as $n \to \infty$, then

$$\left\{\frac{p}{n^{1/2}},\frac{p}{n^{1/2}+1},\ldots,\frac{p}{|\mathcal{D}_{\mathrm{tr}}|}\right\} \quad " \to " \quad [\gamma,\infty].$$

The set of aspect ratios for the predictors f_k covers $[\gamma, \infty]$.

Choose one out of *f*_k, n^{1/2} ≤ k ≤ |D_{tr}| using an estimate of out-of-sample risk computed from D_{te} This is split sample cross-validation.

Cross-validation risk estimate

 Traditionally, the risk of a predictor based on a test data is done via average loss. For example, with squared error loss, the traditional estimate of (prediction) risk of a predictor f_k

$$\widehat{R}(\widetilde{f}_k) := rac{1}{|\mathcal{D}_{ ext{te}}|} \sum_{j \in \mathcal{D}_{ ext{te}}} (Y_j - \widetilde{f}_k(X_j))^2.$$

- For a good performance simultaneously over O(n) predictors and also to avoid strong tail assumptions on the loss, we also consider the median-of-means estimator.
- With either the average or median-of-means estimator of risk, we return the predictor $\widehat{f} := \widetilde{f_{\mu}}$ where

$$\widehat{k} := \operatorname*{argmin}_{n^{1/2} \leq k \leq |\mathcal{D}_{\mathrm{tr}}|} \widehat{R}(\widetilde{f}_k).$$

• \hat{k} represents the "best" sample size to use for the given number of features in the dataset and $\tilde{f}_{\hat{k}}$ is what we call a zero-step predictor that achieves risk monotonization.

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 ζ converges to $R^{\text{det}}(\zeta; \hat{f})$, we show:

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

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 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 ζ converges to $R^{\text{det}}(\zeta; \hat{f})$, we show:

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

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)



Takeaways:

- 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.

- 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)

Takeaways:

- 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.

- 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)

Takeaways:

- 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.

- 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)

Takeaways:

- 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.

- 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)
Takeaways and extensions

Takeaways:

- 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)

Takeaways and extensions

Takeaways:

- 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 of overparameterization

Double descent Current theoretical understanding Case study of linear regression

Risk monotonization

Motivation Zero-step procedure Takeaways and extensions

Bagging analysis

Motivation Risk characterization Optimal subsample size

Connections to ridge regularization

Risk and structural equivalences Implications of equivalences Discussion and extensions

Conclusion

Motivation beyond bagging analysis

Key question: How much improvement do we get if we use an ensemble of M > 1 subsampled datasets, rather than just a single subsampled dataset?



We provide precise risk characterization for ridgeless (and ridge) ensemles.

• Let $\mathcal{D}_n = \{(\mathbf{x}_j, y_j) \in \mathbb{R}^p \times \mathbb{R} : j \in [n]\}$ denote a dataset. The ridge estimator fitted on subsampled dataset \mathcal{D}_I with $I \subseteq [n], |I| = k$ is:

$$\widehat{\beta}_k^{\lambda}(\mathcal{D}_l) = \operatorname*{arg\,min}_{\boldsymbol{\beta} \in \mathbb{R}^p} \frac{1}{k} \sum_{j \in I} (y_j - \boldsymbol{x}_j^{\top} \boldsymbol{\beta})^2 + \lambda \|\boldsymbol{\beta}\|_2^2.$$

• For $\lambda \ge 0$ fixed, ensemble ridge estimator is:

$$\widetilde{eta}_{k,M}^{\lambda}(\mathcal{D}_n; \{I_\ell\}_{\ell=1}^M) := rac{1}{M}\sum_{\ell\in [M]}\widehat{eta}_k^{\lambda}(\mathcal{D}_{I_\ell}),$$

with $I_1, \ldots, I_M \sim \mathcal{I}_k := \{\{i_1, \ldots, i_k\} : 1 \leq i_1 < \ldots < i_k \leq n\}$. The *full-ensemble* ridge estimator is defined by letting $M \to \infty$.

• The goal is to quantify and estimate the conditional prediction risk:

$$R_{k,M}^{\lambda} := \mathbb{E}[(y - \boldsymbol{x}^{\top} \widetilde{\beta}_{k,M}^{\lambda})^2 \mid \mathcal{D}_n, \{I_{\ell}\}_{\ell=1}^M]$$

• Let $\mathcal{D}_n = \{(\mathbf{x}_j, y_j) \in \mathbb{R}^p \times \mathbb{R} : j \in [n]\}$ denote a dataset. The ridge estimator fitted on subsampled dataset \mathcal{D}_I with $I \subseteq [n], |I| = k$ is:

$$\widehat{eta}_k^\lambda(\mathcal{D}_I) = \operatorname*{arg\,min}_{oldsymbol{eta}\in\mathbb{R}^p} rac{1}{k} \sum_{j\in I} (y_j - oldsymbol{x}_j^ opoldsymbol{eta})^2 + \lambda \|oldsymbol{eta}\|_2^2.$$

• For $\lambda \ge 0$ fixed, ensemble ridge estimator is:

$$\widetilde{eta}_{k,M}^{\lambda}(\mathcal{D}_n; \{I_\ell\}_{\ell=1}^M) := rac{1}{M} \sum_{\ell \in [M]} \widehat{eta}_k^{\lambda}(\mathcal{D}_{I_\ell}),$$

with $l_1, \ldots, l_M \sim \mathcal{I}_k := \{\{i_1, \ldots, i_k\} : 1 \leq i_1 < \ldots < i_k \leq n\}$. The *full-ensemble* ridge estimator is defined by letting $M \to \infty$.

• The goal is to quantify and estimate the conditional prediction risk:

$$R_{k,M}^{\lambda} := \mathbb{E}[(y - \boldsymbol{x}^{\top} \widetilde{\beta}_{k,M}^{\lambda})^2 \mid \mathcal{D}_n, \{I_{\ell}\}_{\ell=1}^M]$$

• Let $\mathcal{D}_n = \{(\mathbf{x}_j, y_j) \in \mathbb{R}^p \times \mathbb{R} : j \in [n]\}$ denote a dataset. The ridge estimator fitted on subsampled dataset \mathcal{D}_I with $I \subseteq [n], |I| = k$ is:

$$\widehat{eta}_k^\lambda(\mathcal{D}_I) = rgmin_{oldsymbol{eta}\in\mathbb{R}^p}rac{1}{k}\sum_{j\in I}(y_j-oldsymbol{x}_j^ opoldsymbol{eta})^2+\lambda\|oldsymbol{eta}\|_2^2.$$

• For $\lambda \ge 0$ fixed, ensemble ridge estimator is:

$$\widetilde{eta}_{k,M}^{\lambda}(\mathcal{D}_{\mathsf{n}}; \{I_{\ell}\}_{\ell=1}^{M}) := rac{1}{M} \sum_{\ell \in [M]} \widehat{eta}_{k}^{\lambda}(\mathcal{D}_{I_{\ell}}),$$

with $I_1, \ldots, I_M \sim \mathcal{I}_k := \{\{i_1, \ldots, i_k\} : 1 \le i_1 < \ldots < i_k \le n\}$. The *full-ensemble* ridge estimator is defined by letting $M \to \infty$.

• The goal is to quantify and estimate the conditional prediction risk:

$$R_{k,M}^{\lambda} := \mathbb{E}[(y - \mathbf{x}^{\top} \widetilde{\beta}_{k,M}^{\lambda})^2 \mid \mathcal{D}_n, \{I_{\ell}\}_{\ell=1}^M]$$

Let D_n = {(x_j, y_j) ∈ ℝ^p × ℝ : j ∈ [n]} denote a dataset. The ridge estimator fitted on subsampled dataset D_I with I ⊆ [n], |I| = k is:

$$\widehat{eta}_k^\lambda(\mathcal{D}_I) = rgmin_{oldsymbol{eta}\in\mathbb{R}^p}rac{1}{k}\sum_{j\in I}(y_j-oldsymbol{x}_j^ opoldsymbol{eta})^2+\lambda\|oldsymbol{eta}\|_2^2.$$

• For $\lambda \ge 0$ fixed, ensemble ridge estimator is:

$$\widetilde{eta}_{k,M}^{\lambda}(\mathcal{D}_{\mathsf{n}}; \{I_{\ell}\}_{\ell=1}^{M}) := rac{1}{M} \sum_{\ell \in [M]} \widehat{eta}_{k}^{\lambda}(\mathcal{D}_{I_{\ell}}),$$

with $I_1, \ldots, I_M \sim \mathcal{I}_k := \{\{i_1, \ldots, i_k\} : 1 \le i_1 < \ldots < i_k \le n\}$. The *full-ensemble* ridge estimator is defined by letting $M \to \infty$.

• The goal is to quantify and estimate the conditional prediction risk:

$$R_{k,M}^{\lambda} := \mathbb{E}[(y - \boldsymbol{x}^{\top} \widetilde{\beta}_{k,M}^{\lambda})^2 \mid \mathcal{D}_n, \{I_{\ell}\}_{\ell=1}^M]$$

- 1. Feature model:
 - Feature structure: x_i = Σ^{1/2}z_i, z_i ∈ ℝ^p is a random vector containing i.i.d. entries with mean 0, variance 1, and bounded moment of order 4 + δ for some δ > 0.
 - Covariance norm: There exist r_{\min}, r_{\max} independent of p with $0 \le r_{\min} \le r_{\max} \le \infty$ such that $r_{\min} t \ge \sum_{i=1}^{\infty} \frac{1}{r_{\min}} t_{i}$
- 2. Response model:
 - Response structure: $y_i = \mathbf{x}_i^\top \boldsymbol{\beta}_0 + \epsilon_i$.
 - Noise structure: ϵ_i is an unobserved error that is assumed to be independent of x_i with mean 0, variance σ^2 , and bounded moment of order $4 + \delta$ for some $\delta > 0$.
 - Signal norm: $\|\beta_0\|_2$ uniformly bounded in p and $\lim_p \|\beta_0\|_2^2 = \rho^2$.
- 3. Convergence of covariance and signal-weighted spectrums:
 - Covariance spectrum: $\boldsymbol{\Sigma} = \boldsymbol{W} \boldsymbol{R} \boldsymbol{W}^{ op}$ is the eigenvalue decomposition.
 - Empirical spectrums: Assume there exist fixed distributions H and G such that the empirical spectral distributions satisfy

$$H_p(r) := \frac{1}{p} \sum_{i=1}^p \mathbb{1}_{\{r_i \le r\}} \xrightarrow{\mathrm{d}} H,$$

$$G_p(r):=rac{1}{\|oldsymbol{eta}_0\|_2^2}\sum_{i=1}^{\cdot}(oldsymbol{eta}_0^{ op}oldsymbol{w}_i)^2~\mathbbm{1}_{\{r_i\leq r\}} extsf{d} \in G.$$

- 1. Feature model:
 - Feature structure: x_i = Σ^{1/2}z_i, z_i ∈ ℝ^p is a random vector containing i.i.d. entries with mean 0, variance 1, and bounded moment of order 4 + δ for some δ > 0.
 - Covariance norm: There exist r_{\min}, r_{\max} independent of p with $0 < r_{\min} \le r_{\max} < \infty$ such that $r_{\min} I_p \le \Sigma \le r_{\max} I_p$.
- 2. Response model:
 - Response structure: $y_i = x_i^{\top} \beta_0 + \epsilon_i$.
 - Noise structure: ϵ_i is an unobserved error that is assumed to be independent of x_i with mean 0, variance σ^2 , and bounded moment of order $4 + \delta$ for some $\delta > 0$.
 - Signal norm: $\|\beta_0\|_2$ uniformly bounded in p and $\lim_p \|\beta_0\|_2^2 = \rho^2$.
- 3. Convergence of covariance and signal-weighted spectrums:
 - Covariance spectrum: $oldsymbol{\Sigma} = oldsymbol{W} oldsymbol{W}^ op$ is the eigenvalue decomposition.
 - Empirical spectrums: Assume there exist fixed distributions H and G such that the empirical spectral distributions satisfy

$$H_p(r) := rac{1}{p} \sum_{i=1}^p \mathbb{1}_{\{r_i \leq r\}} \stackrel{\mathrm{d}}{ o} H,$$

$$G_{p}(r):=rac{1}{\|oldsymbol{eta}_{0}\|_{2}^{2}}\sum_{i=1}^{r}(oldsymbol{eta}_{0}^{ op}oldsymbol{w}_{i})^{2}\ \mathbb{1}_{\{r_{i}\leq r\}}\stackrel{\mathrm{d}}{
ightarrow} G.$$

- 1. Feature model:
 - Feature structure: x_i = Σ^{1/2}z_i, z_i ∈ ℝ^p is a random vector containing i.i.d. entries with mean 0, variance 1, and bounded moment of order 4 + δ for some δ > 0.
 - Covariance norm: There exist r_{\min}, r_{\max} independent of p with $0 < r_{\min} \le r_{\max} < \infty$ such that $r_{\min}I_p \le \Sigma \le r_{\max}I_p$.
- 2. Response model:
 - Response structure: $y_i = \mathbf{x}_i^\top \boldsymbol{\beta}_0 + \epsilon_i$.
 - Noise structure: ϵ_i is an unobserved error that is assumed to be independent of x_i with mean 0, variance σ^2 , and bounded moment of order $4 + \delta$ for some $\delta > 0$.
 - Signal norm: $\|\beta_0\|_2$ uniformly bounded in p and $\lim_p \|\beta_0\|_2^2 = \rho^2$.
- 3. Convergence of covariance and signal-weighted spectrums:
 - Covariance spectrum: $oldsymbol{\Sigma} = W R W^+$ is the eigenvalue decomposition.
 - Empirical spectrums: Assume there exist fixed distributions H and G such that the empirical spectral distributions satisfy

$$H_p(r) := rac{1}{p} \sum_{i=1}^p \mathbb{1}_{\{r_i \leq r\}} \stackrel{\mathrm{d}}{ o} H,$$

$$G_p(r) := rac{1}{\|oldsymbol{eta}_0\|_2^2} \sum_{i=1}^p (oldsymbol{eta}_0^{ op} oldsymbol{w}_i)^2 \ \mathbbm{1}_{\{r_i \leq r\}} \stackrel{\mathrm{d}}{
ightarrow} G.$$

- 1. Feature model:
 - Feature structure: x_i = Σ^{1/2}z_i, z_i ∈ ℝ^p is a random vector containing i.i.d. entries with mean 0, variance 1, and bounded moment of order 4 + δ for some δ > 0.
 - Covariance norm: There exist r_{\min}, r_{\max} independent of p with $0 < r_{\min} \le r_{\max} < \infty$ such that $r_{\min}I_p \le \Sigma \le r_{\max}I_p$.
- 2. Response model:
 - Response structure: $y_i = \mathbf{x}_i^\top \boldsymbol{\beta}_0 + \epsilon_i$.
 - Noise structure: ϵ_i is an unobserved error that is assumed to be independent of x_i with mean 0, variance σ^2 , and bounded moment of order $4 + \delta$ for some $\delta > 0$.
 - Signal norm: $\|\beta_0\|_2$ uniformly bounded in p and $\lim_p \|\beta_0\|_2^2 = \rho^2$.
- 3. Convergence of covariance and signal-weighted spectrums:
 - Covariance spectrum: $\boldsymbol{\Sigma} = \boldsymbol{W} \boldsymbol{R} \boldsymbol{W}^{ op}$ is the eigenvalue decomposition.
 - Empirical spectrums: Assume there exist fixed distributions H and G such that the empirical spectral distributions satisfy

$$egin{aligned} & H_p(r) := rac{1}{p} \sum_{i=1}^p \mathbbm{1}_{\{r_i \leq r\}} \stackrel{\mathrm{d}}{ o} H, \ & G_p(r) := rac{1}{\|eta_0\|_2^2} \sum_{i=1}^p (eta_0^ o oldsymbol{w}_i)^2 \ \mathbbm{1}_{\{r_i \leq r\}} \stackrel{\mathrm{d}}{ o} \mathcal{G}. \end{aligned}$$

Risk characterization of bagged ridge predictors

Theorem. Under aforementioned assumptions, as $k, n, p \to \infty$ such that $p/n \to \phi \in (0, \infty)$ and $p/k \to \phi_s \in [\phi, \infty]$, the asymptotic risk $\mathscr{R}^{\text{sub}}_{\lambda, M}(\phi, \phi_s)$ is:

data aspect ratio

subsample aspect ratio

$$\mathscr{R}^{\mathrm{sub}}_{\lambda,\mathcal{M}}(\phi,\phi_{s})=\sigma^{2}+\mathscr{R}^{\mathrm{sub}}_{\lambda,\mathcal{M}}(\phi,\phi_{s})+\mathscr{V}^{\mathrm{sub}}_{\lambda,\mathcal{M}}(\phi,\phi_{s}),$$

where the bias and variance terms are given by

$$\begin{split} \mathscr{B}_{\lambda,M}^{\mathrm{sub}}(\phi,\phi_s) &= M^{-1}B_{\lambda}(\phi_s,\phi_s) + (1-M^{-1})B_{\lambda}(\phi,\phi_s), \\ \mathscr{V}_{\lambda,M}^{\mathrm{sub}}(\phi,\phi_s) &= M^{-1}V_{\lambda}(\phi_s,\phi_s) + (1-M^{-1})V_{\lambda}(\phi,\phi_s), \end{split}$$

and the functions $B_\lambda(\cdot,\cdot)$ and $V_\lambda(\cdot,\cdot)$ are defined as

$$B_{\lambda}(\vartheta,\theta) = \rho^{2}(1+\widetilde{v}(-\lambda;\vartheta,\theta))\widetilde{c}(-\lambda;\theta), \qquad V_{\lambda}(\vartheta,\theta) = \sigma^{2}\widetilde{v}(-\lambda;\vartheta,\theta).$$

Here the non-negative constants $\tilde{v}(-\lambda; \vartheta, \theta)$ and $\tilde{c}(-\lambda; \theta)$ are defined as:

$$\begin{split} \widetilde{v}(-\lambda;\vartheta,\theta) &= \frac{\vartheta \int r^2 (1+v(-\lambda;\theta)r)^{-2} \,\mathrm{d}H(r)}{v(-\lambda;\theta)^{-2} - \vartheta \int r^2 (1+v(-\lambda;\theta)r)^{-2} \,\mathrm{d}H(r)},\\ \widetilde{c}(-\lambda;\theta) &= \int \frac{r}{(1+v(-\lambda;\theta)r)^2} \,\mathrm{d}G(r). \end{split}$$

Finally, $v(-\lambda; \theta)$ is the unique nonnegative solution to the fixed-point equation:

$$\frac{1}{\nu(-\lambda;\theta)} = \lambda + \theta \int \frac{r}{1 + \nu(-\lambda;\theta)r} \,\mathrm{d}H(r).$$

Bagged ridge risk characterization (illustration)



Figure: Asymptotic prediction risk curves for bagged ridgeless predictors $(\lambda = 0)$, under AR1 model when $\rho_{ar1} = 0.25$ and $\sigma^2 = 1$, for varying subsample sizes $k = \lfloor p/\phi_s \rfloor$ and numbers of bags M. The null risk is marked as a dotted line. For each value of M, the points denote finite-sample risks averaged over 100 dataset repetitions, with $n = \lfloor p\phi \rfloor$ and p = 500. The left and the right panels correspond to the cases when p < n ($\phi = 0.1$) and p > n ($\phi = 1.1$), respectively.

Optimal bagged ridgeless predictor

Theorem. For any $\phi \ge 0$, the global minimum of $\phi_s \mapsto \mathscr{R}^{\text{sub}}_{0,\infty}(\phi, \phi_s)$ is obtained in $\phi_s^* \in (1,\infty)$. That is





Subagged ridgeless *interpolators* always outperform subagged least squares, even when the full data has more observations than the number of features.

Optimal bagged ridgeless predictor

Theorem. For any $\phi \ge 0$, the global minimum of $\phi_s \mapsto \mathscr{R}^{\text{sub}}_{0,\infty}(\phi, \phi_s)$ is obtained in $\phi_s^* \in (1,\infty)$. That is





Subagged ridgeless *interpolators* always outperform subagged least squares, even when the full data has more observations than the number of features.

Back to risk monotonization

- Risk characterization \rightarrow risk monotonization.
- Data splitting and cross-validation over subsample size.



Figure: Asymptotic excess risk curves for cross-validated bagged ridgeless predictors ($\lambda = 0$), under the isotopic model when $\rho^2 = 1$ for varying SNR, subsample sizes $k = \lfloor p/\phi_s \rfloor$, and numbers of bags M with replacement. For each value of M, the points denote finite-sample risks and the shaded regions denote the values within one standard deviation, with n = 1000, $n_{\rm te} = 63$, and $p = \lfloor n\phi \rfloor$.

Comparison with optimal ridge regularization



Recall here $\gamma = p/n$ is the aspect ratio. The base predictor is ridgeless.

Key question: Is the connection to ridge regularization just coincedental?

Comparison with optimal ridge regularization



Recall here $\gamma = p/n$ is the aspect ratio. The base predictor is ridgeless.

Key question: Is the connection to ridge regularization just coincedental?

Outline

Overview of overparameterization

Double descent Current theoretical understanding Case study of linear regression

Risk monotonization

Motivation Zero-step procedure Takeaways and extensions

Bagging analysis

Motivation Risk characterization Optimal subsample size

Connections to ridge regularization

Risk and structural equivalences Implications of equivalences Discussion and extensions

Conclusion

As p/n → φ and p/k → φ_s, the prediction risk in the full ensemble (M = ∞) converges:

$$R_{k,\infty}^{\lambda} \xrightarrow{\text{a.s.}} \mathscr{R}_{\infty}^{\lambda}(\phi,\phi_{s}).$$

• For $\phi = 0.1$, the risk profile as a function of (λ, ϕ_s) is shown in the figure in the log-log scale.



• Risk equivalence (Theorem 2.3):

$$\underbrace{\min_{\substack{\phi_s \geq \phi}\\ \phi_s \geq \phi}}_{\substack{\text{opt. ridgeless}\\ \text{ensemble}}} \mathcal{R}_{\lambda,\infty}^{\text{sub}}(\phi,\phi) = \underbrace{\min_{\substack{\lambda \geq 0\\\\ \lambda \geq 0}}}_{\substack{\lambda \geq 0}} \mathcal{R}_{\lambda,\infty}^{\text{sub}}(\phi,\phi_s) \cdot \underbrace{\sum_{\substack{\phi_s \geq \phi,\\ \lambda \geq 0\\\\ \text{opt. ridge}\\ \text{ensemble}}}_{\substack{\phi_s \geq \phi,\\ \lambda \geq 0\\\\ \text{opt. ridge}\\ \text{ensemble}}} \mathcal{R}_{\lambda,\infty}^{\text{sub}}(\phi,\phi_s) \cdot \underbrace{\sum_{\substack{\phi_s \geq \phi,\\ \lambda \geq 0\\\\ \text{opt. ridge}\\ \text{ensemble}}}_{\substack{\phi_s \geq \phi,\\ \lambda \geq 0\\\\ \text{opt. ridge}\\ \text{ensemble}}} \mathcal{R}_{\lambda,\infty}^{\text{sub}}(\phi,\phi_s) \cdot \underbrace{\sum_{\substack{\phi_s \geq \phi,\\ \lambda \geq 0\\\\ \text{opt. ridge}\\ \text{ensemble}}}_{\substack{\phi_s \geq \phi,\\ \lambda \geq 0\\\\ \text{opt. ridge}\\ \text{ensemble}}} \mathcal{R}_{\lambda,\infty}^{\text{sub}}(\phi,\phi_s) \cdot \underbrace{\sum_{\substack{\phi_s \geq \phi,\\ \lambda \geq 0\\\\ \text{opt. ridge}\\ \text{ensemble}}}_{\substack{\phi_s \geq \phi,\\ \lambda \geq 0\\\\ \text{opt. ridge}\\ \text{ensemble}}} \mathcal{R}_{\lambda,\infty}^{\text{sub}}(\phi,\phi_s) \cdot \underbrace{\sum_{\substack{\phi_s \geq \phi,\\ \lambda \geq 0\\\\ \text{opt. ridge}\\ \text{ensemble}}}_{\substack{\phi_s \geq \phi,\\ \lambda \geq 0\\\\ \text{opt. ridge}\\ \text{ensemble}}} \mathcal{R}_{\lambda,\infty}^{\text{sub}}(\phi,\phi_s) \cdot \underbrace{\sum_{\substack{\phi_s \geq \phi,\\ \lambda \geq 0\\\\ \text{opt. ridge}\\ \text{ensemble}}}_{\substack{\phi_s \geq \phi,\\ \lambda \geq 0\\\\ \text{opt. ridge}\\ \text{ensemble}}} \mathcal{R}_{\lambda,\infty}^{\text{sub}}(\phi,\phi_s) \cdot \underbrace{\sum_{\substack{\phi_s \geq \phi,\\ \lambda \geq 0\\\\ \text{opt. ridge}}}}_{\substack{\phi_s \geq \phi,\\ \lambda \geq 0\\\\ \text{opt. ridge}}} \mathcal{R}_{\lambda,\infty}^{\text{sub}}(\phi,\phi_s) \cdot \underbrace{\sum_{\substack{\phi_s \geq \phi,\\ \lambda \geq 0\\\\ \text{opt. ridge}}}}_{\substack{\phi_s \geq \phi,\\ \lambda \geq 0\\\\ \text{opt. ridge}}} \mathcal{R}_{\lambda,\infty}^{\text{sub}}(\phi,\phi_s) \cdot \underbrace{\sum_{\substack{\phi_s \geq \phi,\\ \lambda \geq 0\\\\ \text{opt. ridge}}}}_{\substack{\phi_s \geq \phi,\\ \lambda \geq 0\\\\ \text{opt. ridge}}} \mathcal{R}_{\lambda,\infty}^{\text{sub}}(\phi,\phi_s) \cdot \underbrace{\sum_{\substack{\phi_s \geq \phi,\\ \lambda \geq 0\\\\ \text{opt. ridge}}}}_{\substack{\phi_s \geq \phi,\\ \lambda \geq 0\\\\ \text{opt. ridge}}} \mathcal{R}_{\lambda,\infty}^{\text{sub}}(\phi,\phi_s) \cdot \underbrace{\sum_{\substack{\phi_s \geq \phi,\\ \lambda \geq 0\\\\ \text{opt. ridge}}}}_{\substack{\phi_s \geq 0\\\\ \text{opt. ridge}}} \mathcal{R}_{\lambda,\infty}^{\text{sub}}(\phi,\phi_s) \cdot \underbrace{\sum_{\substack{\phi_s \geq \phi,\\ \lambda \geq 0\\\\ \text{opt. ridge}}}}_{\substack{\phi_s \geq 0\\\\ \text{opt. ridge}}}}_{\substack{\phi_s \geq 0\\\\ \text{opt. ridge}}} \mathcal{R}_{\lambda,\infty}^{\text{sub}}(\phi,\phi_s) \cdot \underbrace{\sum_{\substack{\phi_s \geq 0\\\\ \text{opt. ridge}}}_{\substack{\phi_s \geq 0\\\\ \text{opt. ridge}}}}_{\substack{\phi_s \geq 0\\\\ \text{opt. ridge}}}_{\substack{\phi_s \geq 0\\\\ \text{opt. ridge}}} \mathcal{R}_{\lambda,\infty}^{\text{sub}}(\phi,\phi_s) \cdot \underbrace{\sum_{\substack{\phi_s \geq 0\\\\ \text{opt. ridge}}}}_{\substack{\phi_s \geq 0\\\\ \text{opt. ridge}}}_{\substack{\phi_s \geq 0\\\\ \text{opt. ridge}}}}_{\substack{\phi_s \geq 0\\\\ \text{opt. ridge}}}_{\substack{\phi_s \geq 0\\\\ \text{opt. ridge}}} \mathcal{R}_{\lambda,\infty}^{\text{sub}}(\phi,\phi_s) \cdot \underbrace{\sum_{\substack{\phi_s \geq 0\\\\ \text{opt. ridge}}}}_{\substack{\phi_s \geq 0\\\\ \text{opt. ridge}}}_{\substack{\phi_s \geq 0\\\\ \text{opt. ridge}}}_{\substack{\phi_s \geq$$

As p/n → φ and p/k → φ_s, the prediction risk in the full ensemble (M = ∞) converges:

$$R_{k,\infty}^{\lambda} \xrightarrow{\text{a.s.}} \mathscr{R}_{\infty}^{\lambda}(\phi,\phi_{s}).$$

• For $\phi = 0.1$, the risk profile as a function of (λ, ϕ_s) is shown in the figure in the log-log scale.



$$\underbrace{\min_{\substack{\phi_s \geq \phi}} \mathscr{R}^{\mathrm{sub}}_{0,\infty}(\phi,\phi_s)}_{\text{opt. ridgeless}} = \underbrace{\min_{\lambda \geq 0} \mathscr{R}^{\mathrm{sub}}_{\lambda,\infty}(\phi,\phi)}_{\text{opt. ridge}} = \underbrace{\min_{\substack{\phi_s \geq \phi, \\ \lambda \geq 0}} \mathscr{R}^{\mathrm{sub}}_{\lambda,\infty}(\phi,\phi_s)}_{\text{opt. ridge}}.$$



As p/n → φ and p/k → φ_s, the prediction risk in the full ensemble (M = ∞) converges:

$$R_{k,\infty}^{\lambda} \xrightarrow{\text{a.s.}} \mathscr{R}_{\infty}^{\lambda}(\phi,\phi_{s}).$$

• For $\phi = 0.1$, the risk profile as a function of (λ, ϕ_s) is shown in the figure in the log-log scale.



$$\underbrace{\min_{\substack{\phi_s \ge \phi}\\ \phi_s \ge \phi}}_{\substack{\text{opt. ridgeless}\\ \text{ensemble}}} \mathcal{R}_{\lambda \ge 0}^{\text{sub}}(\phi, \phi) = \underbrace{\min_{\substack{\lambda \ge 0\\ \lambda \ge 0}}}_{\substack{\lambda \ge 0}} \mathcal{R}_{\lambda, \infty}^{\text{sub}}(\phi, \phi_s) = \underbrace{\min_{\substack{\phi_s \ge \phi,\\ \lambda \ge 0}}}_{\substack{\phi_s \ge \phi,\\ \lambda \ge 0}} \mathcal{R}_{\lambda, \infty}^{\text{sub}}(\phi, \phi_s)$$



As p/n → φ and p/k → φ_s, the prediction risk in the full ensemble (M = ∞) converges:

$$R_{k,\infty}^{\lambda} \xrightarrow{\text{a.s.}} \mathscr{R}_{\infty}^{\lambda}(\phi,\phi_{s}).$$

• For $\phi = 0.1$, the risk profile as a function of (λ, ϕ_s) is shown in the figure in the log-log scale.



$$\underbrace{\min_{\substack{\phi_s \ge \phi}} \mathscr{R}^{\mathrm{sub}}_{0,\infty}(\phi,\phi_s)}_{\text{opt. ridgeless}} = \underbrace{\min_{\lambda \ge 0} \mathscr{R}^{\mathrm{sub}}_{\lambda,\infty}(\phi,\phi)}_{\text{opt. ridge}} = \underbrace{\min_{\substack{\phi_s \ge \phi, \\ \lambda \ge 0}} \mathscr{R}^{\mathrm{sub}}_{\lambda,\infty}(\phi,\phi_s)}_{\text{opt. ridge}}.$$



Generalized risk

- Let $\beta_0 = \mathbb{E}[\mathbf{x}\mathbf{x}^{\top}]^{-1}\mathbb{E}[\mathbf{x}y]$ be the best linear projection of y onto \mathbf{x}
- For a linear functional $L(\beta) = A\beta + b$, we study generalized risks:

$$R(\widehat{\beta}; \boldsymbol{A}, \boldsymbol{b}, \beta_0) = \frac{1}{\operatorname{nrow}(\boldsymbol{A})} \|L(\widehat{\beta} - \beta_0)\|_2^2,$$
(1)

Statistical learning problem	$L(\widehat{eta}-oldsymbol{eta}_0)$	A	b	$\operatorname{nrow}(\boldsymbol{A})$
vector coefficient estimation	$\widehat{eta} - oldsymbol{eta}_{0}$	I_p		р
projected coefficient estimation	$oldsymbol{a}^ op(\widehateta-oldsymbol{eta}_0)$	а		1
training error estimation	$oldsymbol{X}\widehat{eta}-oldsymbol{y}$	X	$-f_{\rm NL}$	п
in-sample prediction	$oldsymbol{X}(\widehat{eta}-oldsymbol{eta}_0)$	X		п
out-of-sample prediction	$\mathbf{x}_{0}^{\top}\widehat{\mathbf{\beta}} - \mathbf{y}_{0}$	\boldsymbol{x}_0	$-\epsilon_0$	1

Generalized risk

- Let $\beta_0 = \mathbb{E}[\mathbf{x}\mathbf{x}^{ op}]^{-1}\mathbb{E}[\mathbf{x}y]$ be the best linear projection of y onto \mathbf{x}
- For a linear functional $L(\beta) = A\beta + b$, we study generalized risks:

$$R(\widehat{\beta}; \boldsymbol{A}, \boldsymbol{b}, \boldsymbol{\beta}_0) = \frac{1}{\operatorname{nrow}(\boldsymbol{A})} \|L(\widehat{\beta} - \boldsymbol{\beta}_0)\|_2^2, \quad (1)$$

Statistical learning problem	$L(\widehat{eta}-oldsymbol{eta}_{0})$	Α	Ь	$\operatorname{nrow}(\boldsymbol{A})$
vector coefficient estimation	$\widehat{eta} - oldsymbol{eta}_0$	I_p	0	p
projected coefficient estimation	$oldsymbol{a}^ op(\widehateta-oldsymbol{eta}_0)$	$a^{ op}$	0	1
training error estimation	$oldsymbol{X}\widehat{eta}-oldsymbol{y}$	X	$-f_{\rm NL}$	п
in-sample prediction	$oldsymbol{X}(\widehat{eta}-oldsymbol{eta}_{0})$	X	0	п
out-of-sample prediction	$\mathbf{x}_{0}^{ op}\widehat{eta} - \mathbf{y}_{0}$	$\mathbf{x}_0^ op$	$-\epsilon_0$	1

Asymptotic equivalence and relaxed assumptions

Asymptotic equivalence:

- Let A_p and B_p be sequences of (additively) conformable matrices of arbitrary dimensions (including vectors and scalars).
- We say that \mathbf{A}_{p} and \mathbf{B}_{p} are asymptotically equivalent, denoted as $\mathbf{A}_{p} \simeq \mathbf{B}_{p}$, if $\lim_{p \to \infty} |\operatorname{tr}[\mathbf{C}_{p}(\mathbf{A}_{p} \mathbf{B}_{p})]| = 0$ almost surely for any sequence of random matrices \mathbf{C}_{p} with bounded trace norm that are (multiplicatively) conformable and independent of \mathbf{A}_{p} and \mathbf{B}_{p} .
- Note that for sequences of scalar random variables, the definition simply reduces to the typical almost sure convergence of sequences of random variables involved.

- Feature distribution: Each feature vector x_i for i ∈ [n] can be decomposed as x_i = Σ^{1/2}z_i, where z_i ∈ ℝ^p contains i.i.d. entries z_{ij} for j ∈ [p] with mean 0, variance 1, and bounded 4 + μ moments for some μ > 0.
- Response distribution: Each response variable y_i for i ∈ [n] has mean 0, and bounded 4 + μ moments.

Asymptotic equivalence and relaxed assumptions

Asymptotic equivalence:

- Let A_p and B_p be sequences of (additively) conformable matrices of arbitrary dimensions (including vectors and scalars).
- We say that A_p and B_p are asymptotically equivalent, denoted as $A_p \simeq B_p$, if $\lim_{p\to\infty} |\operatorname{tr}[C_p(A_p B_p)]| = 0$ almost surely for any sequence of random matrices C_p with bounded trace norm that are (multiplicatively) conformable and independent of A_p and B_p .
- Note that for sequences of scalar random variables, the definition simply reduces to the typical almost sure convergence of sequences of random variables involved.

- Feature distribution: Each feature vector *x_i* for *i* ∈ [*n*] can be decomposed as *x_i* = Σ^{1/2}*z_i*, where *z_i* ∈ ℝ^p contains i.i.d. entries *z_{ij}* for *j* ∈ [*p*] with mean 0, variance 1, and bounded 4 + μ moments for some μ > 0.
- Response distribution: Each response variable y_i for i ∈ [n] has mean 0, and bounded 4 + μ moments.

Generalized risk equivalences

Theorem. For any $\bar{\psi} \in [\phi, +\infty]$, let $\bar{\lambda}$ be as defined in (4). Then, for any pair of (λ_1, ψ_1) and (λ_2, ψ_2) on the path $\mathcal{P}(\bar{\lambda}; \phi, \bar{\psi})$ as defined in (5), the generalized risk functionals (1) of the full-ensemble estimator are asymptotically equivalent:

$$R(\widehat{\beta}_{\lfloor \boldsymbol{p}/\psi_1 \rfloor,\infty}^{\lambda_1}; \boldsymbol{A}, \boldsymbol{b}, \beta_0) \simeq R(\widehat{\beta}_{\lfloor \boldsymbol{p}/\psi_2 \rfloor,\infty}^{\lambda_2}; \boldsymbol{A}, \boldsymbol{b}, \beta_0).$$
(2)



Structural equivalences

Theorem. For any $\bar{\psi} \in [\phi, +\infty]$, let $\bar{\lambda}$ be as in (4). Then, for any $M \in \mathbb{N} \cup \{\infty\}$ and any pair of (λ_1, ψ_1) and (λ_2, ψ_2) on the path (5), the *M*-ensemble estimators are asymptotically equivalent:

$$\widehat{\beta}_{\lfloor p/\psi_{1}\rfloor,M}^{\lambda_{1}} \simeq \widehat{\beta}_{\lfloor p/\psi_{2}\rfloor,M}^{\lambda_{2}}, \quad \forall (\lambda_{1},\psi_{1}), (\lambda_{2},\psi_{2}) \in \mathcal{P}(\bar{\lambda};\phi,\bar{\psi}).$$
(3)

Equivalence paths

- Given φ ∈ (0,∞) and ψ
 ∈ [φ,∞], our statement of equivalences between different ensemble estimators is defined through certain paths characterized by two endpoints (0, ψ
) and (λ
 , φ).
- Let H_p be the empirical spectral distribution of Σ : $H_p(r) = p^{-1} \sum_{i=1}^p \mathbb{1}_{\{r_i \leq r\}}$, where r_i 's are the eigenvalues of Σ . Consider the following system of equations in $\overline{\lambda}$ and v:

$$\frac{1}{v} = \bar{\lambda} + \phi \int \frac{r}{1 + vr} \mathrm{d}H_{\rho}(r), \quad \text{and} \quad \frac{1}{v} = \bar{\psi} \int \frac{r}{1 + vr} \mathrm{d}H_{\rho}(r). \tag{4}$$

• Now, define a path $\mathcal{P}(\bar{\lambda}; \phi, \bar{\psi})$ that passes through the endpoints $(0, \bar{\psi})$ and $(\bar{\lambda}, \phi)$:

$$\mathcal{P}(\overline{\lambda};\phi,\overline{\psi}) = \left\{ (1-\theta) \cdot (\overline{\lambda},\phi) + \theta \cdot (0,\overline{\psi}) \mid \theta \in [0,1] \right\}.$$
(5)

For any M ∈ N ∪ {∞}, let λ
_n be the value that satisfies the following equation in ensemble ridgeless and ridge gram matrices:

$$\frac{1}{M}\sum_{\ell=1}^{M}\frac{1}{k}\operatorname{tr}\left[\left(\frac{1}{k}\boldsymbol{L}_{l_{\ell}}\boldsymbol{X}\boldsymbol{X}^{\top}\boldsymbol{L}_{l_{\ell}}\right)^{+}\right] = \frac{1}{n}\operatorname{tr}\left[\left(\frac{1}{n}\boldsymbol{X}\boldsymbol{X}^{\top} + \bar{\lambda}_{n}\boldsymbol{I}_{n}\right)^{-1}\right].$$
 (6)

Define the data-dependent path $\mathcal{P}_n = \mathcal{P}(\overline{\lambda}_n; \phi_n, \overline{\psi}_n).$

Implications: Monotonicity of optimal ridge

- An open problem raised by Nakkiran et al. (2021) asks whether the prediction risk of ridge regression with optimal ridge penalty λ* is monotonically increasing in the data aspect ratio φ = p/n.
- Our equivalences imply that the prediction risk of an optimally-tuned ridge estimator is monotonically increasing in the data aspect ratio under mild regularity conditions.
- Under proportional asymptotics, our result settles a recent open question raised by Conjecture 1 of Nakkiran et al. (2021) concerning the monotonicity of optimal ridge regression under anisotropic features and general data models while maintaining a regularity condition that preserves the linearized signal-to-noise ratios across regression problems.

Implications of equivalences: illustration

Theorem. Let $k, n, p \to \infty$ such that $p/n \to \phi \in (0, \infty)$ and $p/k \to \psi \in [\phi, \infty]$. Then, for $\mathbf{A} = \Sigma^{1/2}$ and $\mathbf{b} = \mathbf{0}$, the optimal risk of the ridgeless ensemble, $\min_{\psi \ge \phi} \mathscr{R}(0; \phi, \psi)$, is monotonically increasing in ϕ . Consequently, the optimal risk of the ridge predictor, $\min_{\ge 0} \mathscr{R}(;\phi,\phi)$, is also monotonically increasing in ϕ .



Extension 1: Equivalences for random features

Conjecture. Define $\phi_n = p/n$. Let $k \leq n$ be the subsample size and denote by $\bar{\psi}_n = p/k$. Suppose φ satisfies certain regularity conditions. For any $M \in \mathbb{N} \cup \{\infty\}$, let $\bar{\lambda}_n$ be the value that satisfies

$$\frac{1}{M}\sum_{\ell=1}^{M}\frac{1}{k}\operatorname{tr}\left[\left(\frac{1}{k}\varphi(\boldsymbol{L}_{l_{\ell}}\boldsymbol{X}\boldsymbol{F}^{\top})\varphi(\boldsymbol{L}_{l_{\ell}}\boldsymbol{X}\boldsymbol{F}^{\top})^{\top}\right)^{+}\right] = \frac{1}{n}\operatorname{tr}\left[\left(\frac{1}{n}\varphi(\boldsymbol{X}\boldsymbol{F}^{\top})\varphi(\boldsymbol{X}\boldsymbol{F}^{\top})^{\top} + \bar{\lambda}_{n}\boldsymbol{I}_{n}\right)^{-1}\right]$$

Define the data-dependent path $\mathcal{P}_n = \mathcal{P}(\bar{\lambda}_n; \phi_n, \bar{\psi}_n)$. Then similar equivalences continue to hold along \mathcal{P}_n .



Extension 2: Equivalences for kernel features

Conjecture. Define $\phi_n = p/n$. Suppose the kernel K satisfies certain regularity conditions. Let $k \leq n$ be the subsample size and denote by $\bar{\psi}_n = p/k$. For any $M \in \mathbb{N} \cup \{\infty\}$, let $\bar{\lambda}_n$ be a solution to

$$\frac{1}{M}\sum_{\ell=1}^{M} \operatorname{tr}\left[\boldsymbol{K}_{l_{\ell}}^{+}\right] = \operatorname{tr}\left[\left(\boldsymbol{K}_{[n]} + \frac{n}{p}\bar{\lambda}_{n}\boldsymbol{I}_{n}\right)^{-1}\right]$$

Define the data-dependent path $\mathcal{P}_n = \mathcal{P}(\bar{\lambda}_n; \phi_n, \bar{\psi}_n)$. Then similar equivalences continue to hold along \mathcal{P}_n .



Outline

Overview of overparameterization

Double descent Current theoretical understanding Case study of linear regression

Risk monotonization

Motivation Zero-step procedure Takeaways and extensions

Bagging analysis

Motivation Risk characterization Optimal subsample size

Connections to ridge regularization

Risk and structural equivalences Implications of equivalences Discussion and extensions

Conclusion
- 1. 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.
- 2. Ensembling helps significantly near the interpolator threshold. Subagged ridgeless interpolators always outperform subagged least squares, even when the full data has more observations than the number of features.
- 3. There are connections between the implicit regularization induced by subsampling and explicit ridge regularization for subsampled ridge ensembles.

- 1. 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.
- 2. Ensembling helps significantly near the interpolator threshold. Subagged ridgeless interpolators always outperform subagged least squares, even when the full data has more observations than the number of features.
- 3. There are connections between the implicit regularization induced by subsampling and explicit ridge regularization for subsampled ridge ensembles.

- 1. 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.
- 2. Ensembling helps significantly near the interpolator threshold. Subagged ridgeless interpolators always outperform subagged least squares, even when the full data has more observations than the number of features.
- 3. There are connections between the implicit regularization induced by subsampling and explicit ridge regularization for subsampled ridge ensembles.

- 1. 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.
- 2. Ensembling helps significantly near the interpolator threshold. Subagged ridgeless interpolators always outperform subagged least squares, even when the full data has more observations than the number of features.
- 3. There are connections between the implicit regularization induced by subsampling and explicit ridge regularization for subsampled ridge ensembles.

Thanks for listening!

Questions/comments/thoughts?

What about lasso?



"Mitigating multiple descents: A model-agnostic framework for risk monotonization" P., Kuchibhotla, Wei, Rinaldo, 2021

What about lasso?



More empirical evidence for lasso

