

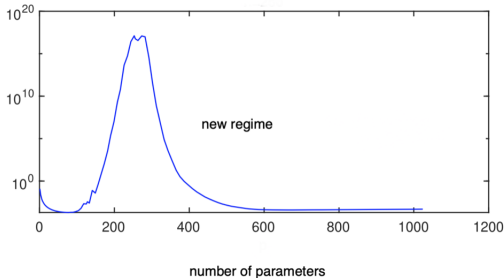
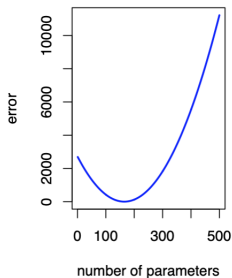
# Overparametrization and the bias-variance dilemma

**Johannes Schmidt-Hieber**

joint work with Alexis Derumigny

<https://arxiv.org/abs/2006.00278.pdf>

## double descent and implicit regularization



overparametrization generalizes well  $\rightsquigarrow$  implicit regularization

# can we defy the bias-variance trade-off?

Geman et al. '92: "the fundamental limitations resulting from the bias-variance dilemma apply to all nonparametric inference methods, including neural networks"

Because of the double descent phenomenon, there is some doubt whether this statement is true. Recent work includes

## Statistics > Machine Learning

*[Submitted on 28 Dec 2018 (v1), last revised 10 Sep 2019 (this version, v2)]*

### **Reconciling modern machine learning practice and the bias-variance trade-off**

Mikhail Belkin, Daniel Hsu, Siyuan Ma, Soumik Mandal

---

## Computer Science > Machine Learning

*[Submitted on 19 Oct 2018 (v1), last revised 18 Dec 2019 (this version, v4)]*

### **A Modern Take on the Bias-Variance Tradeoff in Neural Networks**

Brady Neal, Sarthak Mittal, Aristide Baratin, Vinayak Tantia, Matthew Scicluna, Simon Lacoste-Julien, Ioannis Mitliagkas

## lower bounds on the bias-variance trade-off

Similar to minimax lower bounds we want to establish a general mathematical framework to derive lower bounds on the bias-variance trade-off that hold for all estimators.

**given such bounds we can answer many interesting questions**

- are there methods (e.g. deep learning) that can defy the bias-variance trade-off?
- lower bounds for the  $U$ -shaped curve of the classical bias-variance trade-off

## related literature

- Low '95 provides complete characterization of bias-variance trade-off for functionals in the Gaussian white noise model
- Pfanzagl '99 shows that estimators of functionals satisfying an asymptotic unbiasedness property must have unbounded variance

No general treatment of lower bounds for the bias-variance trade-off yet.

## Cramér-Rao inequality

for parametric problems:

$$V(\theta) \geq \frac{(1 + B'(\theta))^2}{F(\theta)}$$

- $V(\theta)$  the variance
- $B'(\theta)$  the derivative of the bias
- $F(\theta)$  the Fisher information

## change of expectation inequalities

- probability measures  $P_0, \dots, P_M$
- $\chi^2(P_0, \dots, P_M)$  the matrix with entries

$$\chi^2(P_0, \dots, P_M)_{j,k} = \int \frac{dP_j}{dP_0} dP_k - 1$$

- any random variable  $X$
- $\Delta := (E_{P_1}[X] - E_{P_0}[X], \dots, E_{P_M}[X] - E_{P_0}[X])^\top$

then,

$$\Delta^\top \chi^2(P_0, \dots, P_M)^{-1} \Delta \leq \text{Var}_{P_0}(X)$$

## pointwise estimation

**Gaussian white noise model:** We observe  $(Y_x)_x$  with

$$dY_x = f(x) dx + n^{-1/2} dW_x$$

- estimate  $f(x_0)$  for a fixed  $x_0$
- $\mathcal{C}^\beta(R)$  denotes ball of Hölder  $\beta$ -smooth functions
- for any estimator  $\hat{f}(x_0)$ , we obtain the **bias-variance lower bound**

$$\inf_{\hat{f}} \sup_{f \in \mathcal{C}^\beta(R)} |\text{Bias}_f(\hat{f}(x_0))|^{1/\beta} \sup_{f \in \mathcal{C}^\beta(R)} \text{Var}_f(\hat{f}(x_0)) \gtrsim \frac{1}{n}$$

- bound is attained by most estimators
- generates  $U$ -shaped curve



## high-dimensional models

### Gaussian sequence model:

- observe independent  $X_i \sim \mathcal{N}(\theta_i, 1)$ ,  $i = 1, \dots, n$
- $\Theta(s)$  the space of  $s$ -sparse vectors (here:  $s \leq \sqrt{n}/2$ )
- bias-variance decomposition

$$E_{\theta}[\|\hat{\theta} - \theta\|^2] = \underbrace{\|E_{\theta}[\hat{\theta}] - \theta\|^2}_{B^2(\theta)} + \sum_{i=1}^n \text{Var}_{\theta}(\hat{\theta}_i)$$

- **bias-variance lower bound:** if  $B^2(\theta) \leq \gamma s \log(n/s^2)$ , then,

$$\sum_{i=1}^n \text{Var}_0(\hat{\theta}_i) \gtrsim n \left(\frac{s^2}{n}\right)^{4\gamma}$$

- bound is matched (up to a factor in the exponent) by soft thresholding
- bias-variance trade-off more extreme than  $U$ -shape
- results also extend to high-dimensional linear regression

**Gaussian white noise model:** We observe  $(Y_x)_x$  with

$$dY_x = f(x) dx + n^{-1/2} dW_x$$

- bias-variance decomposition

$$\begin{aligned} \text{MISE}_f(\hat{f}) &:= E_f [\|\hat{f} - f\|_{L^2[0,1]}^2] \\ &= \int_0^1 \text{Bias}_f^2(\hat{f}(x)) dx + \int_0^1 \text{Var}_f(\hat{f}(x)) dx \\ &=: \text{IBias}_f^2(\hat{f}) + \text{IVar}_f(\hat{f}). \end{aligned}$$

- is there a bias-variance trade-off between  $\text{IBias}_f^2(\hat{f})$  and  $\text{IVar}_f(\hat{f})$ ?
- turns out to be a very hard problem

## $L^2$ -loss (ctd.)

- we propose a two-fold reduction scheme
  - reduction to a simpler model
  - reduction to a smaller class of estimators
- $S^\beta(R)$  Sobolev space of  $\beta$ -smooth functions

**Bias-variance lower bound:** For any estimator  $\hat{f}$ ,

$$\inf_{\hat{f}} \sup_{f \in S^\beta(R)} |\text{Bias}_f(\hat{f})|^{1/\beta} \sup_{f \in S^\beta(R)} \text{IVar}_f(\hat{f}) \geq \frac{1}{8n},$$

- many estimators  $\hat{f}$  can be found with upper bound  $\lesssim 1/n$

## mean absolute deviation

- several extensions of the bias-variance trade-off have been proposed in the literature, e.g. for classification
- the mean absolute deviation (MAD) of an estimator  $\hat{\theta}$  is

$$E_{\theta}[|\hat{\theta} - m|]$$

with  $m$  either the mean or the median of  $\hat{\theta}$

can the general framework be extended to lower bounds on the trade-off between bias and MAD?

- derived change of expectation inequality
- this can be used to obtain a partial answer for pointwise estimation in the Gaussian white noise model

## Summary

- general framework to derive bias-variance lower bounds
- leads to matching bias-variance lower bounds for standard models in nonparametric and high-dimensional statistics
- different types of the bias-variance trade-off occur
- can machine learning methods defy the bias-variance trade-off? **No, there are universal lower bounds that no method can avoid**

for details and more results consult the preprint

<https://arxiv.org/abs/2006.00278.pdf>