

On the Double Descent of Random Features Models Trained with SGD

Fanghui Liu (*EPFL*), **Johan A.K. Suykens** (*KU Leuven*), **Volkan Cevher** (*EPFL*)

Laboratory for Information and Inference Systems (LIONS)
École Polytechnique Fédérale de Lausanne (EPFL)
Switzerland

25th Nov. 2021



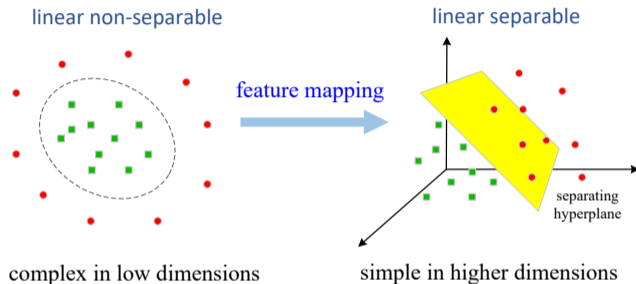
Outline

Research overview

Random features in double descent

Conclusion

Research Overview: Kernel approximation



Scalability of kernel methods: n -by- n kernel matrix.

Solution: approximate the kernel by a low-rank representation

- ▶ Nyström approximation: approximate the kernel matrix
- ▶ Random Fourier features¹: approximate the kernel function

¹Rahimi A, Recht B. Random features for large-scale kernel machines, NeurIPS2007. (the test-of-time award in NeurIPS2017)

Research Overview: Random Fourier features

$$k(\mathbf{x}, \mathbf{x}') = \langle \phi(\mathbf{x}), \phi(\mathbf{x}') \rangle_{\mathcal{H}} \approx \varphi^\top(\mathbf{x})\varphi(\mathbf{x}'),$$

where $\varphi(\mathbf{x}) : \mathbb{R}^d \rightarrow \mathbb{R}^s$ is an **explicit** feature mapping

Bochner's theorem [1]

For a shift-invariant $k(\mathbf{x}, \mathbf{x}') = k(\mathbf{x} - \mathbf{x}')$ and positive definite kernel,

$$\begin{aligned} k(\mathbf{x}, \mathbf{x}') &= \int_{\mathbb{R}^d} p(\omega) \exp\left(i\omega^\top(\mathbf{x} - \mathbf{x}')\right) d\omega \\ &\approx \frac{1}{s} \sum_{j=1}^s \exp(i\omega_j^\top \mathbf{x}) \exp(i\omega_j^\top \mathbf{x}')^* = \varphi(\mathbf{x})^\top \varphi(\mathbf{x}') \end{aligned}$$

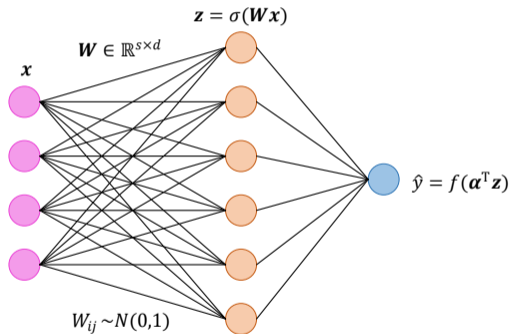
the explicit feature mapping:

$$\varphi(\mathbf{x}) := \frac{1}{\sqrt{s}} \left[\exp(-i\omega_1^\top \mathbf{x}), \dots, \exp(-i\omega_s^\top \mathbf{x}) \right]^\top.$$

Research Overview: Neural network view

RF model: a two-layer, (infinite)-width, fully-connected neural network

$$k(\mathbf{x}, \mathbf{x}') = \mathbb{E}_{\omega \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_d)}[\sigma(\omega^\top \mathbf{x})\sigma(\omega^\top \mathbf{x}')]$$



- ▶ Gaussian kernel: $\sigma(x) = [\cos(x), \sin(x)]^\top$
- ▶ the 1st-order arc-cosine kernel: $\sigma(x) = \max\{0, x\}$
- ▶ soft-max in attention: $\sigma(x) = \exp(x)$

Research Overview: Applied to Linearized Attention in Transformers

self attention

$$\text{Attention}(\mathbf{Q}, \mathbf{K}, \mathbf{V}) = \underbrace{\text{softmax}(\mathbf{Q}\mathbf{K}^\top)}_{:=\mathbf{A}} \mathbf{V} \approx \mathbf{Q}'\mathbf{K}'^\top \mathbf{V},$$

where $\mathbf{A}_{ij} = k(\mathbf{q}_i, \mathbf{k}_j) = \mathbb{E}[\sigma(\mathbf{q}_i)^\top \sigma(\mathbf{k}_j)]$

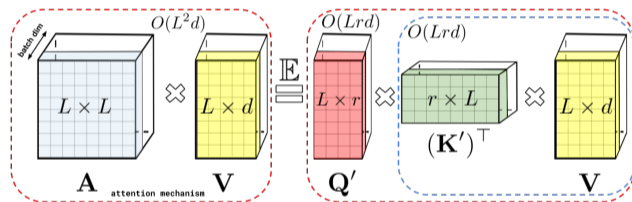


Figure: Approximation of self-attention. source: [2].

- soft-max in attention: $\exp(\mathbf{x}^\top \mathbf{x}') = \mathbb{E}_{\omega \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_d)} \left[\exp\left(\omega^\top \mathbf{x} - \frac{\|\mathbf{x}\|_2^2}{2}\right) \exp\left(\omega^\top \mathbf{x}' - \frac{\|\mathbf{x}'\|_2^2}{2}\right) \right]$

Research Overview: Taxonomy

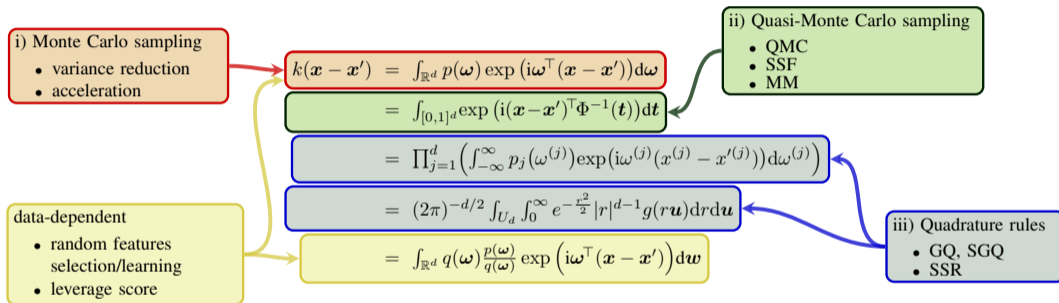


Figure: Taxonomy of random features based algorithms².

²Fanghui Liu, Xiaolin Huang, Yudong Chen, and Johan A.K. Suykens. *Random Features for Kernel Approximation: A Survey on Algorithms, Theory, and Beyond*. TPAMI2021.

Outline

Research overview

Random features in double descent

Conclusion

Background: Double descent

over-parameterized models, e.g., neural networks, random features

- ▶ high dimensions: large n and d
- ▶ abnormal phenomena: training error can be zero but still generalize well

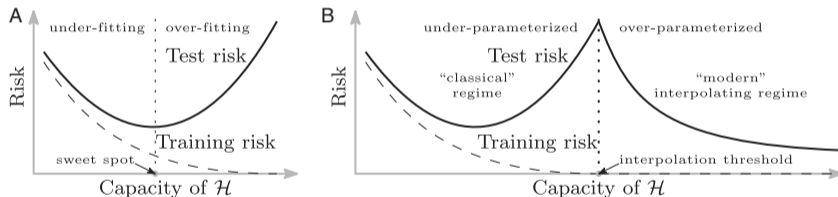


Figure: Bias-variance trade-off [3] (Belkin et al. PNAS2019).

Research Overview: Motivation

- ▶ interplay between optimization and excess risk: trained by SGD
- ▶ bias-variance decomposition for understanding multiple randomness sources

	data assumption	solution	result
(Hastie et al., 2019)	Gaussian	closed-form	variance ↗ ↘
(Ba et al., 2020)	Gaussian	GD	variance ↗ ↘
(Mei & Montanari, 2019)	i.i.d on sphere	closed-form	variance, bias ↗ ↘
(d'Ascoli et al., 2020a)	Gaussian	closed-form	refined ²
(Gerace et al., 2020)	Gaussian	closed-form	↗ ↘
(Adlam & Pennington, 2020)	Gaussian	closed-form	refined
(Dhifallah & Lu, 2020)	Gaussian	closed-form	↗ ↘
(Hu & Lu, 2020)	Gaussian	closed-form	↗ ↘
(Liao et al., 2020)	general	closed-form	↗ ↘
(Lin & Dobriban, 2021)	isotropic features with finite moments	closed form	refined
(Li et al., 2021)	correlated features with polynomial decay on Σ_d	closed form	interpolation learning
Ours	(at least) sub-exponential data	SGD	variance ↗ ↘, bias ↘

¹ A refined decomposition on variance is conducted by sources of randomness on data sampling, initialization, label noise to possess each term (d'Ascoli et al., 2020b) or their full decomposition in (Adlam & Pennington, 2020; Lin & Dobriban, 2021).

Problem settings: Random features regression model

data: $y = f_\rho(\mathbf{x}) + \varepsilon$

- ▶ training data: $\{(\mathbf{x}_i, y_i)\}_{i=1}^n \sim \rho$
Assumption: sub-exponential data and $\|\mathbf{x}\|_2^2 \sim \mathcal{O}(d)$
- ▶ target function: $f_\rho(\mathbf{x}) = \int_Y y \, d\rho(y | \mathbf{x})$
- ▶ noise: $\mathbb{E}(\varepsilon) = 0$ and $\mathbb{E}(\varepsilon^2) = \tau^2$

function space

define the random features mapping $\varphi(\mathbf{x}) := \frac{1}{\sqrt{m}} \sigma(\mathbf{W}\mathbf{x} / \sqrt{d})$,

$$\mathcal{H} := \left\{ f \in L^2_{\rho_X} \mid f(\mathbf{x}) = \langle \theta, \varphi(\mathbf{x}) \rangle \right\}, \quad \mathbf{W}_{ij} \sim \mathcal{N}(0, 1)$$

covariance operator: $\Sigma_m := \int_X [\varphi(\mathbf{x}) \otimes \varphi(\mathbf{x})] d\rho_X(\mathbf{x})$

expected covariance operator: $\tilde{\Sigma}_m := \mathbb{E}_{\mathbf{x}, \mathbf{W}} [\varphi(\mathbf{x}) \otimes \varphi(\mathbf{x})]$

Problem settings: averaged SGD under adaptive step-size setting

$$\theta_t = \theta_{t-1} + \gamma_t [y_t - \langle \theta_{t-1}, \varphi(\mathbf{x}_t) \rangle] \varphi(\mathbf{x}_t), \quad t = 1, 2, \dots, n,$$

- ▶ averaged output: $\bar{\theta}_n := \frac{1}{n} \sum_{t=0}^{n-1} \theta_t \implies \bar{f}_n = \langle \varphi(\cdot), \bar{\theta}_n \rangle$
- ▶ adaptive step-size: $\gamma_t := \gamma_0 t^{-\zeta}, \zeta \in [0, 1)$
- ▶ optimal solution: $f^* = \arg \min_{f \in \mathcal{H}} \|f - f_\rho\|_{L_{\rho, X}}^2$
- ▶ averaged excess risk: $\mathbb{E} \|\bar{f}_n - f^*\|_{L_{\rho, X}}^2 = \mathbb{E}_{\mathbf{X}, \mathbf{W}, \varepsilon} \langle \bar{f}_n - f^*, \Sigma_m(\bar{f}_n - f^*) \rangle$

Assumptions

- ▶ **boundedness of f^* :** $\|f^*\|_{\mathcal{H}} < \infty$
- ▶ **high dimension:** $c \leq \{d/n, m/n\} \leq C$, $\|\mathbf{x}\|_2^2 \sim \mathcal{O}(d)$, $\Sigma_d := \mathbb{E}_{\mathbf{x}}[\mathbf{x} \otimes \mathbf{x}]$ with $\|\Sigma_d\|_2 < \infty$
- ▶ **activation function:** $\sigma(\cdot)$: Lipschitz continuous
- ▶ **noise condition:** $\Xi := \mathbb{E}_{\mathbf{x}}[\varepsilon^2 \varphi(\mathbf{x}) \otimes \varphi(\mathbf{x})] \leq \tau^2 \Sigma_m$.
uniformly bounded noise, sub-Gaussian noise
- ▶ **fourth moment condition:**
for any PSD operator A , we have $\mathbb{E}_{\mathbf{W}}[\Sigma_m A \Sigma_m] \leq r' \mathbb{E}_{\mathbf{W}}[\text{Tr}(\Sigma_m A) \Sigma_m] \leq r \text{Tr}(\tilde{\Sigma}_m A) \tilde{\Sigma}_m$.
 - 1) The special case $A := I$ can be proved.
 - 2) holds for sub-Gaussian/exponential data.

Properties of covariance operators

$\sigma(\cdot) : \mathbb{R} \mapsto \mathbb{R}$ Lipschitz continuous

covariance operator $\Sigma_m := \mathbb{E}_{\mathbf{x}}[\varphi(\mathbf{x}) \otimes \varphi(\mathbf{x})]$

expected covariance operator $\tilde{\Sigma}_m := \mathbb{E}_{\mathbf{x}, \mathbf{W}}[\varphi(\mathbf{x}) \otimes \varphi(\mathbf{x})]$

eigenvalue of $\tilde{\Sigma}_m$

the same diagonal/non-diagonal elements: $\mathcal{O}(1/m)$

two distinct eigenvalues: $\tilde{\lambda}_1 \sim \mathcal{O}(1)$, $\tilde{\lambda}_2 \sim \mathcal{O}(1/m)$

sub-exponential random variables

$\|\Sigma_m\|_2$, $\|\Sigma_m - \tilde{\Sigma}_m\|_2$, $\text{Tr}(\Sigma_m)$, and $\left\| \tilde{\Sigma}_m^{-1} \mathbb{E}_{\mathbf{W}}(\Sigma_m^2) \right\|_2$ with $\mathcal{O}(1)$ sub-exponential norm order

Bias-variance decomposition

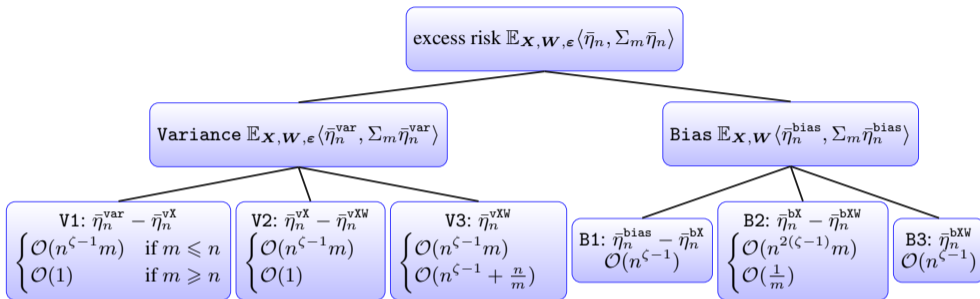
Define $\eta_t := f_t - f^*$, we have

$$\begin{aligned}\eta_t &= [I - \gamma_t \varphi(\mathbf{x}_t) \otimes \varphi(\mathbf{x}_t)](f_{t-1} - f^*) + \gamma_t \varepsilon_t \varphi(\mathbf{x}_t), \\ \eta_t^{\text{bias}} &= [I - \gamma_t \varphi(\mathbf{x}_t) \otimes \varphi(\mathbf{x}_t)] \eta_{t-1}^{\text{bias}}, \quad \eta_0^{\text{bias}} = f^*, \\ \eta_t^{\text{var}} &= [I - \gamma_t \varphi(\mathbf{x}_t) \otimes \varphi(\mathbf{x}_t)] \eta_{t-1}^{\text{var}} + \gamma_t \varepsilon_t \varphi(\mathbf{x}_t), \quad \eta_0^{\text{var}} = 0.\end{aligned}$$

Bias-variance decomposition

$$\mathbb{E} \|\bar{f}_n - f^*\|_{L^2_{\rho_X}}^2 = \underbrace{\mathbb{E}_{\mathbf{X}, \mathbf{W}} \langle \bar{\eta}_n^{\text{bias}}, \Sigma_m \bar{\eta}_n^{\text{bias}} \rangle}_{:= \text{Bias}} + \underbrace{\mathbb{E}_{\mathbf{X}, \mathbf{W}, \varepsilon} \langle \bar{\eta}_n^{\text{var}}, \Sigma_m \bar{\eta}_n^{\text{var}} \rangle}_{:= \text{Variance}}.$$

Proof framework

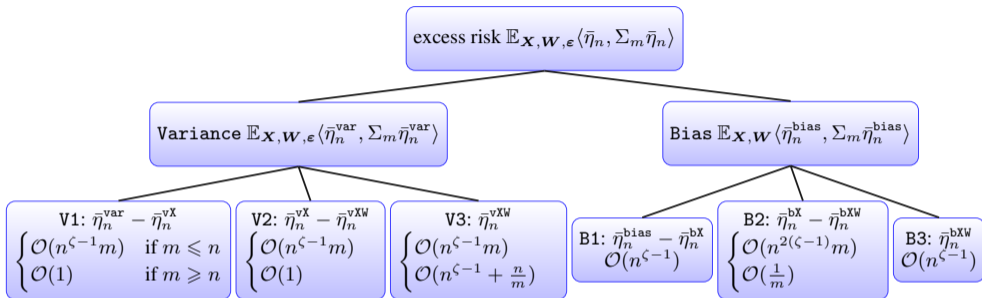


$$\text{Bias : } \eta_t^{\text{bias}} = [I - \gamma_t \varphi(\mathbf{x}_t) \otimes \varphi(\mathbf{x}_t)] \eta_{t-1}^{\text{bias}}$$

Define "semi-stochastic" version: $\eta_t^{\text{bX}} = (I - \gamma_t \Sigma_m) \eta_{t-1}^{\text{bX}}$, $\eta_t^{\text{bXW}} = (I - \gamma_t \tilde{\Sigma}_m) \eta_{t-1}^{\text{bXW}}$,

- ▶ $B1 := \mathbb{E}_{\mathbf{X}, \mathbf{W}} [\langle \bar{\eta}_n^{\text{bias}} - \bar{\eta}_n^{\text{bX}}, \Sigma_m (\bar{\eta}_n^{\text{bias}} - \bar{\eta}_n^{\text{bX}}) \rangle]$
- ▶ $B2 := \mathbb{E}_{\mathbf{W}} [\langle \bar{\eta}_n^{\text{bX}} - \bar{\eta}_n^{\text{bXW}}, \Sigma_m (\bar{\eta}_n^{\text{bX}} - \bar{\eta}_n^{\text{bXW}}) \rangle]$
- ▶ $B3 := \langle \bar{\eta}_n^{\text{bXW}}, \tilde{\Sigma}_m \bar{\eta}_n^{\text{bXW}} \rangle$

Proof framework



$$\text{Variance: } \eta_t^{\text{var}} = [I - \gamma_t \varphi(\mathbf{x}_t) \otimes \varphi(\mathbf{x}_t)] \eta_{t-1}^{\text{var}} + \gamma_t \epsilon_t \varphi(\mathbf{x}_t)$$

Define "semi-stochastic" version: $\eta_t^{\text{vX}} := (I - \gamma_t \Sigma_m) \eta_{t-1}^{\text{vX}} + \gamma_t \epsilon_t \varphi(\mathbf{x}_t)$, $\eta_t^{\text{vXW}} := (I - \gamma_t \tilde{\Sigma}_m) \eta_{t-1}^{\text{vXW}} + \gamma_t \epsilon_t \varphi(\mathbf{x}_t)$

- ▶ V1 := $\mathbb{E}_{\mathbf{X}, \mathbf{W}, \epsilon} [\langle \bar{\eta}_n^{\text{var}} - \bar{\eta}_n^{\text{vX}}, \Sigma_m (\bar{\eta}_n^{\text{var}} - \bar{\eta}_n^{\text{vX}}) \rangle]$
- ▶ V2 := $\mathbb{E}_{\mathbf{X}, \mathbf{W}, \epsilon} [\langle \bar{\eta}_n^{\text{vX}} - \bar{\eta}_n^{\text{vXW}}, \Sigma_m (\bar{\eta}_n^{\text{vX}} - \bar{\eta}_n^{\text{vXW}}) \rangle]$
- ▶ V3 := $\mathbb{E}_{\mathbf{X}, \mathbf{W}, \epsilon} \langle \bar{\eta}_n^{\text{vXW}}, \Sigma_m \bar{\eta}_n^{\text{vXW}} \rangle$

Results: error bounds

Theorem

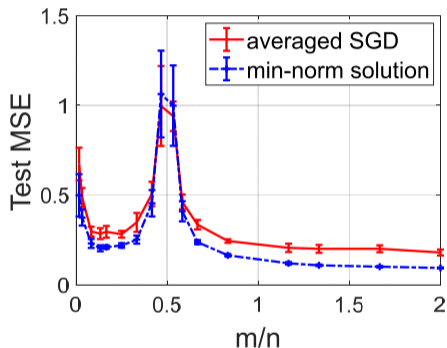
Under the above-mentioned assumptions, if the step-size $\gamma_t := \gamma_0 t^{-\zeta}$ with $\zeta \in [0, 1)$ satisfies $\gamma_0 < C$, we have

$$\text{Bias} \lesssim \frac{\gamma_0 r' n^{\zeta-1}}{\sqrt{\mathbb{E}[1 - \gamma_0 r' \text{Tr}(\Sigma_m)]^4}} \|f^*\|^2 \sim \mathcal{O}(n^{\zeta-1}).$$

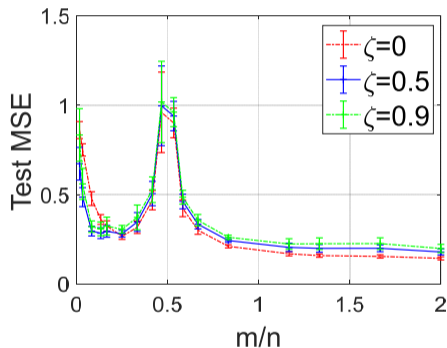
$$\begin{aligned} \text{Variance} &\lesssim \frac{\gamma_0 r' \tau^2}{\sqrt{\mathbb{E}[1 - \gamma_0 r' \text{Tr}(\Sigma_m)]^2}} \begin{cases} mn^{\zeta-1}, & \text{if } m \leq n \\ \gamma_0 \tau^2, & \text{if } m > n \end{cases} \\ &\sim \begin{cases} \mathcal{O}(mn^{\zeta-1}), & \text{if } m \leq n \\ \mathcal{O}(1), & \text{if } m > n. \end{cases} \end{aligned}$$

Experiments on MNIST

Gaussian kernel $k(\mathbf{x}, \mathbf{x}') = \exp\left(-\frac{\|\mathbf{x} - \mathbf{x}'\|_2^2}{2d}\right)$



(a) SGD vs. min-norm solution

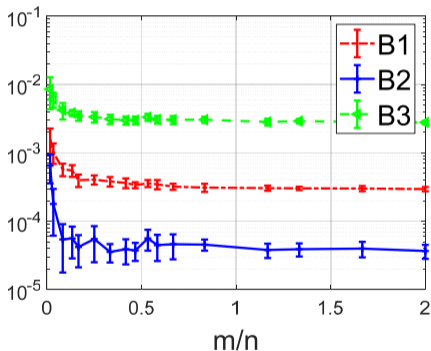


(b) step-size

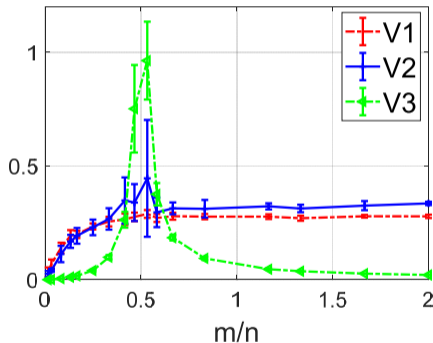
Figure: Test MSE (mean \pm std.) of RF regression as a function of the ratio m/n on MNIST data set (digit 3 vs. 7) for $d = 784$ and $n = 600$.

Validation for bias and variance

- ▶ noise: $\varepsilon \sim \mathcal{N}(0, 1)$
- ▶ $\Sigma_m, \tilde{\Sigma}_m$: sample covariance matrices with Monte Carlo sampling



(a) Bias



(b) Variance

Outline

Research overview

Random features in double descent

Conclusion

Take-away message

high dimensional random features model trained by SGD

findings {
 expected covariance operator $\tilde{\Sigma}_m$ has only two distinct eigenvalues
 bias-variance decomposition: multiple randomness sources
 monotonic decreasing bias and unimodal variance
 optimization effect on excess risk: constant step-size SGD vs. min-norm solution

Future works:

- ▶ SGD: implicit bias/regularization
- ▶ function space, high dimensions

Thanks for your attention!

Q & A

my homepage <http://lfhsgre.org> for more information!



NEW: ERC Advanced Grant E-DUALITY
Exploring duality for future data-driven modelling



References I

- [1] Salomon Bochner.
Harmonic Analysis and the Theory of Probability.
Courier Corporation, 2005.
- [2] Krzysztof Choromanski, Valerii Likhoshesterov, David Dohan, Xingyou Song, Andreea Gane, Tamas Sarlos, Peter Hawkins, Jared Davis, Afroz Mohiuddin, Lukasz Kaiser, and Weller Adrian.
Rethinking attention with performers.
In *International Conference on Learning Representations*, 2021.
- [3] Mikhail Belkin, Daniel Hsu, Siyuan Ma, and Soumik Mandal.
Reconciling modern machine-learning practice and the classical bias–variance trade-off.
the National Academy of Sciences, 116(32):15849–15854, 2019.