

Tutorial:

BO, Double Descent and RKHS ridge regression made easy!

Theodor Misiakiewicz (Stanford)

May 26th, 2022

New Interactions Between Statistics and Optimization workshop, BIRS.

- ▶ **New regime for statistics:** overparametrized models, no explicit regularization and capacity control, train until (near-)interpolation even with noisy data.

New phenomenology: benign overfitting, double descent, non-monotonic error curves...

- ▶ Phenomena already present in **linear models** [\[Belkin, Ma, Mandal, '18\]](#).
- ▶ This talk: **kernel ridge regression** (KRR) in the **high dimension regime**.

Goal of this tutorial is to show:

1. How to derive quickly asymptotics for kernel/random features ridge regression.
2. How the above phenomena have very precise explanations in this regime.

- ▶ **New regime for statistics:** overparametrized models, no explicit regularization and capacity control, train until (near-)interpolation even with noisy data.

New phenomenology: benign overfitting, double descent, non-monotonic error curves...

- ▶ Phenomena already present in **linear models** [\[Belkin, Ma, Mandal, '18\]](#).
- ▶ This talk: **kernel ridge regression** (KRR) in the **high dimension regime**.

Goal of this tutorial is to show:

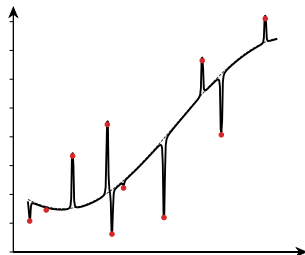
1. How to derive quickly asymptotics for kernel/random features ridge regression.
2. How the above phenomena have very precise explanations in this regime.

Benign overfitting, self-induced regularization and double descent

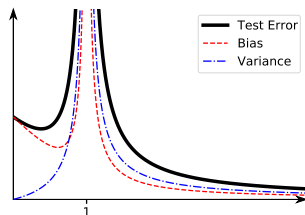
- **Benign overfitting:** interpolator generalizes well.
Idea: $\hat{f} = f_0 + \Delta$ with f_0 smooth solution + spike part with $\|\Delta\|_{L^2} \ll 1$.

In linear models: **self-induced regularization**.
Non-smooth part of the kernel plays the role of an effective ridge regularization.

(This is HD phenomena.)



- **Double descent** and non-monotonic curves.



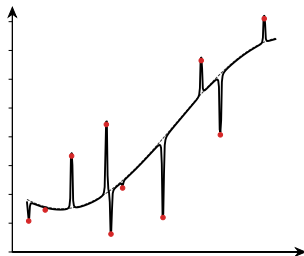
Need exact test error that holds for a given function and is exact up to an additive vanishing constant.

Benign overfitting, self-induced regularization and double descent

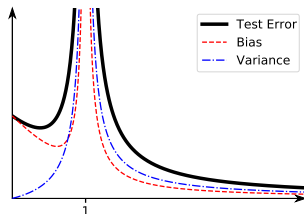
- **Benign overfitting:** interpolator generalizes well.
Idea: $\hat{f} = f_0 + \Delta$ with f_0 smooth solution + spike part with $\|\Delta\|_{L^2} \ll 1$.

In linear models: **self-induced regularization**.
Non-smooth part of the kernel plays the role of an effective ridge regularization.

(This is HD phenomena.)



- **Double descent** and non-monotonic curves.



Need exact test error that holds for a given function and is exact up to an additive vanishing constant.

A subset of references:

- ▶ Benign overfitting: [Liang,Rakhlin,'18], [Ghorbani,Mei,M,Montanari,'19], [Bartlett,Long,Lugosi,Tsigler,'20].
- ▶ Double descent: [Mei,Montanari,'19], [Hastie,Montanari,Rosset,Tibshirani,'20], [Gerace,Loureiro,Krzakala,Mezard,Zdeborova,'20].
- ▶ Linear models: [Tsigler,Bartlett,'20], [Cui,Loureiro,Krzakala,Zdeborova,'21], [Liao,Couillet,Mahoney,'20], [Richards,Mourtada,Rosasco,'21], [Wu,Xu,'20].
- ▶ KRR: [Jacot,Simsek,Spadaro,Hongler,Gabriel,'20], [Canatar,Bordelon,Pehlevan,'21], [Mei,M,Montanari,'21], [Bartlett,Montanari,Rakhlin,'21], [Liu,Liao,Suykens,'21], [Liang,Rakhlin,Zhai,'21], [Hu,Lu,'22].

Quick background on KRR (1)

Covariates: $\mathbf{x} \in (\mathcal{X}, \nu)$.

Kernel function: $K : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$ PSD, associated kernel operator: $\mathbb{K} : L^2(\mathcal{X}) \rightarrow L^2(\mathcal{X})$:

$$\mathbb{K}f(\mathbf{x}) = \int_{\mathcal{X}} K(\mathbf{x}, \mathbf{x}') f(\mathbf{x}') \nu(d\mathbf{x}').$$

Diagonalization: $\{\phi_j\}_{j \geq 1}$ orthonormal basis of $L^2(\mathcal{X})$ and $\{\lambda_j\}_{j \geq 0}$ nonincreasing ($\lambda_j > 0$)

$$\mathbb{K} = \sum_{j \geq 1} \lambda_j \phi_j \phi_j^*, \quad K(\mathbf{x}_1, \mathbf{x}_2) = \sum_{j \geq 1} \lambda_j \phi_j(\mathbf{x}_1) \phi_j(\mathbf{x}_2).$$

Feature map: $\mathbf{x} \mapsto \Phi(\mathbf{x}) = (\sqrt{\lambda_j} \phi_j(\mathbf{x}))_{j \geq 1}$ so that $K(\mathbf{x}_1, \mathbf{x}_2) = \langle \Phi(\mathbf{x}_1), \Phi(\mathbf{x}_2) \rangle_{\ell_2}$

L^2 space: for any $f_* \in L^2(\mathcal{X})$,

$$f_*(\mathbf{x}) = \sum_{j \geq 1} \beta_j \phi_j(\mathbf{x}) = \langle \boldsymbol{\theta}_*, \Phi(\mathbf{x}) \rangle_{\ell_2}, \quad \boldsymbol{\theta}_* = (\theta_j)_{j \geq 1}, \quad \theta_j = \beta_j / \sqrt{\lambda_j}.$$

Associated RKHS: $\mathcal{H} = \{f \in L^2(\mathcal{X}) : \|f\|_{\mathcal{H}} < \infty\}$,

$$\|f\|_{\mathcal{H}}^2 = \|\mathbb{K}^{-1/2} f\|_{L^2}^2 = \sum_{j \geq 1} \frac{\beta_j^2}{\lambda_j} = \|\boldsymbol{\theta}_*\|_{\ell_2}^2.$$

Quick background on KRR (2)

Data: $\{(y_i, \mathbf{x}_i)\}_{i \in [n]}$ where $\mathbf{x}_i \sim_{iid} (\mathcal{X}, \nu)$ and $y_i = f_*(\mathbf{x}_i) + \varepsilon_i$, with $f_* \in L^2(\mathcal{X})$ and independent noise ε_i , $\mathbb{E}[\varepsilon_i] = 0$, $\mathbb{E}[\varepsilon_i^2] = \sigma_\varepsilon^2$.

Kernel ridge regression: fit the data with

$$\hat{f}_\lambda = \arg \min_f \left\{ \sum_{i \in [n]} (y_i - f(\mathbf{x}_i))^2 + \lambda \|f\|_{\mathcal{H}}^2 \right\}.$$

Equivalently: $\hat{f}_\lambda = \langle \hat{\boldsymbol{\theta}}_\lambda, \boldsymbol{\Phi}(\cdot) \rangle_{\ell_2}$ where for $\boldsymbol{\Phi} = [\boldsymbol{\Phi}(\mathbf{x}_1), \dots, \boldsymbol{\Phi}(\mathbf{x}_n)]^T \in \mathbb{R}^{n \times \infty}$,

$$\hat{\boldsymbol{\theta}}_\lambda = \arg \min_{\boldsymbol{\theta}} \left\{ \sum_{i \in [n]} (y_i - \langle \boldsymbol{\Phi}(\mathbf{x}_i), \boldsymbol{\theta} \rangle)^2 + \lambda \|\boldsymbol{\theta}\|_{\ell_2}^2 \right\} = \boldsymbol{\Phi}^T (\boldsymbol{\Phi} \boldsymbol{\Phi}^T + \lambda)^{-1} \mathbf{y}.$$

[Representer thm: $\hat{f}_\lambda(\mathbf{x}) = \sum_i \hat{a}_i K(\mathbf{x}, \mathbf{x}_i)$ with $\hat{\mathbf{a}} = (\mathbf{K} + \lambda)^{-1} \mathbf{y}$, $\mathbf{K} = (K(\mathbf{x}_i, \mathbf{x}_j))_{ij \in [n]}$]

Goal: compute the test error $R(f_*, \hat{f}_\lambda) = \mathbb{E}_{\mathbf{x}}[(f_*(\mathbf{x}) - \hat{f}_\lambda(\mathbf{x}))^2]$ in the high dimensional regime $\mathbf{x} \in \mathbb{R}^d$ and $\log(n) \asymp \log(d)$.

Quick background on KRR (2)

Data: $\{(y_i, \mathbf{x}_i)\}_{i \in [n]}$ where $\mathbf{x}_i \sim_{iid} (\mathcal{X}, \nu)$ and $y_i = f_*(\mathbf{x}_i) + \varepsilon_i$, with $f_* \in L^2(\mathcal{X})$ and independent noise ε_i , $\mathbb{E}[\varepsilon_i] = 0$, $\mathbb{E}[\varepsilon_i^2] = \sigma_\varepsilon^2$.

Kernel ridge regression: fit the data with

$$\hat{f}_\lambda = \arg \min_f \left\{ \sum_{i \in [n]} (y_i - f(\mathbf{x}_i))^2 + \lambda \|f\|_{\mathcal{H}}^2 \right\}.$$

Equivalently: $\hat{f}_\lambda = \langle \hat{\boldsymbol{\theta}}_\lambda, \boldsymbol{\Phi}(\cdot) \rangle_{\ell_2}$ where for $\boldsymbol{\Phi} = [\boldsymbol{\Phi}(\mathbf{x}_1), \dots, \boldsymbol{\Phi}(\mathbf{x}_n)]^\top \in \mathbb{R}^{n \times \infty}$,

$$\hat{\boldsymbol{\theta}}_\lambda = \arg \min_{\boldsymbol{\theta}} \left\{ \sum_{i \in [n]} (y_i - \langle \boldsymbol{\Phi}(\mathbf{x}_i), \boldsymbol{\theta} \rangle)^2 + \lambda \|\boldsymbol{\theta}\|_{\ell_2}^2 \right\} = \boldsymbol{\Phi}^\top (\boldsymbol{\Phi} \boldsymbol{\Phi}^\top + \lambda \mathbf{I})^{-1} \mathbf{y}.$$

[Representer thm: $\hat{f}_\lambda(\mathbf{x}) = \sum_i \hat{a}_i K(\mathbf{x}, \mathbf{x}_i)$ with $\hat{\mathbf{a}} = (\mathbf{K} + \lambda \mathbf{I})^{-1} \mathbf{y}$, $\mathbf{K} = (K(\mathbf{x}_i, \mathbf{x}_j))_{ij \in [n]}$.]

Goal: compute the test error $R(f_*, \hat{f}_\lambda) = \mathbb{E}_{\mathbf{x}}[(f_*(\mathbf{x}) - \hat{f}_\lambda(\mathbf{x}))^2]$ in the high dimensional regime $\mathbf{x} \in \mathbb{R}^d$ and $\log(n) \asymp \log(d)$.

Gaussian equivalent model and universality of feature maps (1)

Ridge regression with features $\phi(\mathbf{x}_i)$: function of random matrix resolvent.

For some “high dimensional” feature map, expect universality to happen: can replace $\phi(\mathbf{x})$ by Gaussian vector \mathbf{z} with matching first two moments.

Covariance matrix: $\Sigma = \mathbb{E}_{\mathbf{x}}[\Phi(\mathbf{x})\Phi(\mathbf{x})^T] = \text{diag}((\lambda_j)_{j \geq 1})$.

| Model | KRR: | Gaussian covariates model: |
|--------------|--|--|
| Distribution | $\Phi(\mathbf{x}) = (\sqrt{\lambda_j} \phi_j(\mathbf{x}))_{j \geq 1}$ with $\mathbf{x} \sim \nu$. | $\mathbf{z} \sim \mathcal{N}(0, \Sigma)$ |
| Data | $(\Phi(\mathbf{x}_i))_{i \in [n]}$ iid, $y_i = \langle \Phi(\mathbf{x}_i), \theta_* \rangle + \varepsilon_i$ | $(\mathbf{z}_i)_{i \in [n]}$ iid, $y_i = \langle \mathbf{z}_i, \theta_* \rangle + \varepsilon_i$ |
| Feature mat. | $\Phi = [\phi(\mathbf{x}_1), \dots, \phi(\mathbf{x}_n)]^T \in \mathbb{R}^{n \times \infty}$ | $\mathbf{Z} = [\mathbf{z}_1, \dots, \mathbf{z}_n]^T \in \mathbb{R}^{n \times \infty}$ |
| Kernel fct | $K(\mathbf{x}_i, \mathbf{x}_j) = \langle \Phi(\mathbf{x}_1), \Phi(\mathbf{x}_2) \rangle_{\ell_2}$ | $K(\mathbf{z}_i, \mathbf{z}_j) = \langle \mathbf{z}_1, \mathbf{z}_2 \rangle_{\ell_2}$ |
| Solution | $\hat{f}_\lambda = \langle \Phi(\mathbf{x}), \hat{\theta}_\lambda \rangle_{\ell_2}$ $\hat{\theta}_\lambda = \Phi^T (\Phi \Phi^T + \lambda)^{-1} \mathbf{y}$ | $\hat{f}_\lambda(\mathbf{z}) = \langle \mathbf{z}, \hat{\theta}_\lambda^G \rangle_{\ell_2}$ $\hat{\theta}_\lambda^G = \mathbf{Z}^T (\mathbf{Z} \mathbf{Z}^T + \lambda)^{-1} \mathbf{y}$ |
| Test error | $R(f_*, \hat{f}_\lambda) = \ \Sigma^{1/2}(\theta_* - \hat{\theta}_\lambda)\ _{\ell_2}^2$ | $R_G(f_*, \hat{f}_\lambda) = \ \Sigma^{1/2}(\theta_* - \hat{\theta}_\lambda^G)\ _{\ell_2}^2$ |

Universality: $R(f_*, \hat{f}_\lambda) - R_G(f_*, \hat{f}_\lambda) \xrightarrow{\mathbb{P}} 0$ (already conjectured previously).

Gaussian equivalent model and universality of feature maps (1)

Ridge regression with features $\phi(\mathbf{x}_i)$: function of random matrix resolvent.

For some “high dimensional” feature map, expect universality to happen: can replace $\phi(\mathbf{x})$ by Gaussian vector \mathbf{z} with matching first two moments.

Covariance matrix: $\Sigma = \mathbb{E}_{\mathbf{x}}[\Phi(\mathbf{x})\Phi(\mathbf{x})^T] = \text{diag}((\lambda_j)_{j \geq 1})$.

| Model | KRR: | Gaussian covariates model: |
|--------------|---|---|
| Distribution | $\Phi(\mathbf{x}) = (\sqrt{\lambda_j}\phi_j(\mathbf{x}))_{j \geq 1}$ with $\mathbf{x} \sim \nu$. | $\mathbf{z} \sim \mathcal{N}(0, \Sigma)$ |
| Data | $(\Phi(\mathbf{x}_i))_{i \in [n]}$ iid, $y_i = \langle \Phi(\mathbf{x}_i), \theta_* \rangle + \varepsilon_i$ | $(\mathbf{z}_i)_{i \in [n]}$ iid, $y_i = \langle \mathbf{z}_i, \theta_* \rangle + \varepsilon_i$ |
| Feature mat. | $\Phi = [\phi(\mathbf{x}_1), \dots, \phi(\mathbf{x}_n)]^T \in \mathbb{R}^{n \times \infty}$ | $\mathbf{Z} = [\mathbf{z}_1, \dots, \mathbf{z}_n]^T \in \mathbb{R}^{n \times \infty}$ |
| Kernel fct | $K(\mathbf{x}_i, \mathbf{x}_j) = \langle \Phi(\mathbf{x}_1), \Phi(\mathbf{x}_2) \rangle_{\ell_2}$ | $K(\mathbf{z}_i, \mathbf{z}_j) = \langle \mathbf{z}_1, \mathbf{z}_2 \rangle_{\ell_2}$ |
| Solution | $\hat{f}_\lambda = \langle \Phi(\mathbf{x}), \hat{\theta}_\lambda \rangle_{\ell_2}$ $\hat{\theta}_\lambda = \Phi^T(\Phi\Phi^T + \lambda)^{-1}\mathbf{y}$ | $\hat{f}_\lambda(\mathbf{z}) = \langle \mathbf{z}, \hat{\theta}_\lambda^G \rangle_{\ell_2}$ $\hat{\theta}_\lambda^G = \mathbf{Z}^T(\mathbf{Z}\mathbf{Z}^T + \lambda)^{-1}\mathbf{y}$ |
| Test error | $R(f_*, \hat{f}_\lambda) = \ \Sigma^{1/2}(\theta_* - \hat{\theta}_\lambda)\ _{\ell_2}^2$ | $R_G(f_*, \hat{f}_\lambda) = \ \Sigma^{1/2}(\theta_* - \hat{\theta}_\lambda^G)\ _{\ell_2}^2$ |

Universality: $R(f_*, \hat{f}_\lambda) - R_G(f_*, \hat{f}_\lambda) \xrightarrow{\mathbb{P}} 0$ (already conjectured previously).

Gaussian equivalent model and universality of feature maps (2)

[Hu, Lu, '20], [Montanari, Saeed, '22], [Gerace, Loureiro, Krzakala, Mezard, Zdeborova, '20] when $n \asymp d$ and or universality of covariates.

Here universality of the entire feature map and polynomial scaling $\log(n) \asymp \log(d)$.

Such an equivalence is not obvious: coordinates of $\Phi(\mathbf{x})$ are not subgaussian or weakly dependent. Here will present some cases, where it can be shown rigorously.

Test error in the Gaussian covariates model

Test error: $R_G(f_*, \hat{f}_\lambda) = \|\Sigma^{1/2}(\theta_* - \hat{\theta}_\lambda^G)\|_{\ell_2}^2$.

$$\text{Bias} = \|\Sigma^{1/2}\theta_* - \Sigma^{1/2}\mathbf{Z}^\top(\mathbf{Z}\mathbf{Z}^\top + \lambda)^{-1}\mathbf{Z}\theta_*\|_{\ell_2}^2,$$

$$\text{Variance} = \sigma_\varepsilon^2 \text{Tr}[(\mathbf{Z}\mathbf{Z}^\top + \lambda)^{-2}\mathbf{Z}\Sigma\mathbf{Z}^\top].$$

Different than previous Gaussian design work ($n \asymp p$, eigenvalues of same order).

Here for simplicity, we assume: $\exists \delta > 0$ and a sequence $m(n)$ such that $m \leq n^{1-\delta}$ and

$$\lambda_{m+1} \cdot n^{1+\delta} \leq \sum_{j=m+1}^{\infty} \lambda_j.$$

(a ‘spectral gap’, which will happen for models with a lot of symmetries.)

Key quantity: the kernel matrix

Random kernel matrix: $K = (\langle \mathbf{z}_i, \mathbf{z}_j \rangle)_{i,j \in [n]} \in \mathbb{R}^{n \times n}$.

$$K = \sum_{j=1}^{\infty} \lambda_j \mathbf{u}_j \mathbf{u}_j^T = K_{\leq m} + K_{> m}, \quad \mathbf{u}_j = (\mathbf{z}_{ij})_{i \in [n]}.$$

Main intuition:

- **High-frequency part:** $K_{> m} = \mathbf{Z}_{> m} \mathbf{Z}_{> m}^T$ with $\mathbf{Z}_{> m}$ iid Gaussian rows. Denote $\lambda_{> m} = \sum_{j> m} \lambda_j$. Then w.h.p.,

$$\|\mathbf{Z}_{> m} \mathbf{Z}_{> m}^T - \lambda_{> m} \mathbf{I}\|_{\text{op}} \lesssim \lambda_{m+1} + \lambda_{> m}/n \lesssim n^{-1} \cdot \lambda_{> m}.$$

- **Low-frequency part:** $K_{\leq m} = \mathbf{Z}_{\leq m} \mathbf{Z}_{\leq m}^T = \mathbf{G}_m \mathbf{\Sigma}_m \mathbf{G}_m^T$ where $\mathbf{G}_m = \mathbf{Z}_{\leq m} \mathbf{\Sigma}_m^{-1/2}$, $\mathbf{G}_m \in \mathbb{R}^{n \times m}$ iid $N(0, 1)$ with $m \ll n$. Then \mathbf{G}_m almost orthogonal:

$$\|\mathbf{G}_m^T \mathbf{G}_m / n - \mathbf{I}_m\|_{\text{op}} \lesssim \frac{m}{n} = n^{-\delta}.$$

The kernel matrix: $\lambda_{\text{eff}} = \lambda_{> m} + \lambda$,

$$K + \lambda \mathbf{I} = \mathbf{G}_m \mathbf{\Sigma}_m \mathbf{G}_m^T + \lambda_{\text{eff}} \mathbf{I}_n.$$

Asymptotics:

$$\begin{aligned}\text{Bias} &= \|\beta_m - (\Sigma_m + (\lambda_{\text{eff}}/n)I_m)^{-1}\Sigma_m\beta_m\|_2^2 + \|\beta_{>m}\|_{\ell_2}^2 + o_{d,\mathbb{P}}(1), \\ \text{Variance} &= o_{d,\mathbb{P}}(1).\end{aligned}$$

Test error:

$$R_G(f_*, \hat{f}_\lambda) = \|\beta - (\Sigma + (\lambda_{\text{eff}}/n)I)^{-1}\Sigma\beta\|_{\ell_2}^2 + o_{d,\mathbb{P}}(1).$$

For KRR:

$$R(f_*, \hat{f}_\lambda) = \|f_* - \mathbb{S}_\lambda f_*\|_{L^2}^2 + o_{d,\mathbb{P}}(1),$$

with shrinkage operator:

$$\hat{f}_\lambda(\mathbf{x}) \approx \mathbb{S}_\lambda f_*(\mathbf{x}) = \sum_{j \geq 1} \frac{\lambda_j}{\lambda_j + \lambda_{\text{eff}}/n} \langle f_*, \phi_j \rangle_{L^2} \phi_j(\mathbf{x}).$$

KRR acts as a shrinkage operator

With spectral gap, KRR with finite data :

$$\hat{f}_\lambda = \arg \min_f \left\{ \frac{1}{n} \sum_{i \in [n]} (y_i - f(\mathbf{x}_i))^2 + \frac{\lambda}{n} \|f\|_{\mathcal{H}}^2 \right\},$$

replaced by effective problem with $n = \infty$ and $\lambda_{\text{eff}} = \lambda_{>m} + \lambda$:

$$\hat{f}_\lambda^{\text{eff}} = \arg \min_f \left\{ \mathbb{E}[(f_*(\mathbf{x}) - f(\mathbf{x}))^2] + \frac{\lambda_{\text{eff}}}{n} \|f\|_{\mathcal{H}}^2 \right\}.$$

Components: $\lambda_j \gg \lambda_{\text{eff}}/n$ perfectly fitted, $\lambda_j \ll \lambda_{\text{eff}}/n$ not fitted at all.

Phenomenology:

- 1) $\lambda_{>m}$ self-induced regularization from high-degree part of kernel,
- 2) Interpolator $\lambda = 0$ are optimal,
- 3) KRR learns $P_{\leq m} f_*$ (smooth part) and doesn't learn at all $P_{>m} f_*$ (and $P_{>m} \hat{f}_\lambda$ spiky part for interpolation).

KRR acts as a shrinkage operator

With spectral gap, KRR with finite data :

$$\hat{f}_\lambda = \arg \min_f \left\{ \frac{1}{n} \sum_{i \in [n]} (y_i - f(\mathbf{x}_i))^2 + \frac{\lambda}{n} \|f\|_{\mathcal{H}}^2 \right\},$$

replaced by effective problem with $n = \infty$ and $\lambda_{\text{eff}} = \lambda_{>m} + \lambda$:

$$\hat{f}_\lambda^{\text{eff}} = \arg \min_f \left\{ \mathbb{E}[(f_*(\mathbf{x}) - f(\mathbf{x}))^2] + \frac{\lambda_{\text{eff}}}{n} \|f\|_{\mathcal{H}}^2 \right\}.$$

Components: $\lambda_j \gg \lambda_{\text{eff}}/n$ perfectly fitted, $\lambda_j \ll \lambda_{\text{eff}}/n$ not fitted at all.

Phenomenology:

- 1) $\lambda_{>m}$ self-induced regularization from high-degree part of kernel,
- 2) Interpolator $\lambda = 0$ are optimal,
- 3) KRR learns $P_{\leq m} f_*$ (smooth part) and doesn't learn at all $P_{>m} f_*$ (and $P_{>m} \hat{f}_\lambda$ spiky part for interpolation).

Proving universality

What properties on $\phi(\mathbf{x})$ allow you to show universality with Gaussian model?

[Mei,M.,Montanari,'21]: hypercontractivity of the top eigenfunctions.

Assumptions on $(\phi^{(n)})_{n \geq 1}$:

- **Spectral gap:** exists $m(n)_{n \geq 1}$ and $\delta > 0$ such that $m \leq n^{1-\delta}$ and

$$\lambda_{m+1} \cdot n^{1+\delta} \leq \sum_{j>m} \lambda_j .$$

- **Hypercontractivity:** for any $q \geq 1$, there exists C_q such that

$$\|h\|_{L^{2q}} \leq C_q \|h\|_{L^2} , \quad \forall h \in \text{span}(\psi_s^{(n)} : 1 \leq s \leq m) .$$

E.g., low-degree polynomials for \mathbf{x} Gaussian vector/uniform on hypercube/uniform on hypersphere.

Other abstract assumptions that show universality. However, difficult to check these assumptions in practice.

One example: inner-product kernel on the sphere

- $\mathbf{x}_1, \mathbf{x}_2 \sim \text{Unif}(\mathbb{S}^{d-1}(\sqrt{d}))$ and $h : [-1, +1] \rightarrow \mathbb{R}$ PD, non-polynomial.
Eigendecomposition of the kernel:

$$K(\mathbf{x}_1, \mathbf{x}_2) = h(\langle \mathbf{x}_1, \mathbf{x}_2 \rangle / d) = \sum_{k=0}^{\infty} \xi_k \sum_{s \in [B(d, k)]} Y_{ks}(\mathbf{x}_1) Y_{ks}(\mathbf{x}_2),$$

where Y_{ks} degree- k spherical harmonics and $\xi_k = \Theta_d(B(d, k)^{-1}) = \Theta_d(d^{-k})$.

- For $d^{\ell+\delta} \leq n \leq d^{\ell+1-\delta}$, can take $m = \sum_{k \leq \ell} B(d, k) = \Theta(d^\ell)$:

$$n \cdot \lambda_{m+1} = n \cdot \xi_{\ell+1} \geq d^\delta, \quad \sum_{j > m} \lambda_j = \sum_{k > \ell} \xi_k B(d, k) = \Theta_d(1).$$

- For $j \leq m$, $\lambda_j \gg \lambda^{\text{eff}}/n$ are perfectly learned, $\lambda_j \ll \lambda^{\text{eff}}/n$ are not learned at all, i.e.,

$$\hat{f}_\lambda(\mathbf{x}) = P_{\leq \ell} f_*(\mathbf{x}) + o_{d, \mathbb{P}}(1).$$

What about $n \asymp d^\ell$?

No spectral gap when $n \asymp d^\ell$:

$$K \approx K_{\leq \ell-1} + \mu_\ell \frac{Y_\ell Y_\ell^\top}{B(d, \ell)} + \mu_{>\ell} I_n.$$

- ▶ $K_{\leq \ell-1}$: low-rank spike matrix.
- ▶ $K_{>\ell} \approx \mu_{>\ell} I_n$: self-induced reg from high-degree part.
- ▶ $Y_\ell = (Y_{\ell s}(\mathbf{x}_i))_{i \in [n], s \in [B(d, \ell)]} \in \mathbb{R}^{n \times B(d, \ell)}$, iid rows. Covariance matrix. Spectrum converges to a Marchenko-Pastur law.

(Generalization of [El Karoui, '10] to the polynomial scaling.)

Fits completely degree- $(\ell - 1)$ polynomial approximation, none of the degree $> \ell$ components, and partially degree- ℓ components.

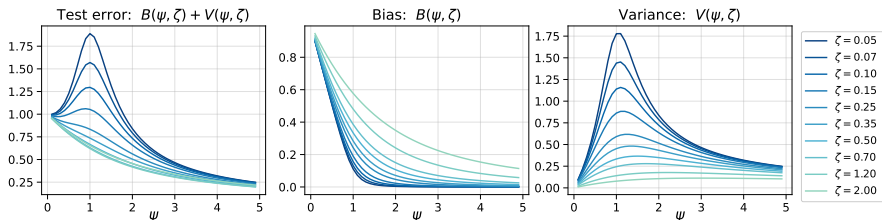
Test error for $n \asymp d^\ell$

Test error = $\|P_{>\ell} f_*\|_{L^2}^2$ + the test error of ridge regression model with $\mathbf{x}_i \sim N(0, I_B)$ and $y_i = \langle \mathbf{x}, \boldsymbol{\beta} \rangle + \varepsilon_i$ with $\|\boldsymbol{\beta}\|_2 = \|P_\ell f_*\|_{L^2}$ and noise $\mathbb{E}[\varepsilon_i^2] = \|P_{>\ell} f_*\|_{L^2}^2 + \sigma_\varepsilon^2$ and regularization $\xi_\ell = (\mu_{>\ell} + \lambda)/\mu_\ell$:

$$\min_{\boldsymbol{\beta}} \left\{ \|\mathbf{y} - \mathbf{X}\boldsymbol{\beta}\|_2^2 + \xi_\ell \|\boldsymbol{\beta}\|_2^2 \right\}.$$

As $n/B(d, \ell) \rightarrow \psi$:

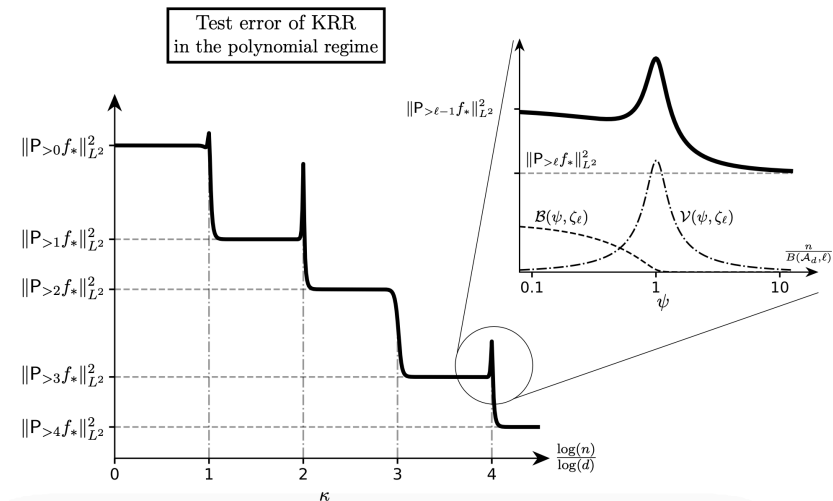
$$R(f_*; \hat{f}_\lambda) = \|P_\ell f_*\|_{L^2}^2 \cdot \mathcal{B}(\psi, \zeta_\ell) + (\|P_{>\ell} f_*\|_{L^2}^2 + \sigma_\varepsilon^2) \cdot \mathcal{V}(\psi, \zeta_\ell) + \|P_{>\ell} f_*\|_{L^2}^2 + o_{d, \mathbb{P}}(1).$$



Asymptotics of KRR on the sphere in polynomial scaling

$\mathbf{x} \sim \text{Unif}(\mathbb{S}^{d-1}(\sqrt{d}))$, $K(\mathbf{x}; \mathbf{z}) = h(\langle \mathbf{x}, \mathbf{z} \rangle / d)$.

Asymptotics in polynomial scaling $n/d^\kappa \rightarrow \psi$ for any $\kappa, \psi > 0$:



$K(\mathbf{x}_1, \mathbf{x}_2) = \mathbb{E}_{\mathbf{w} \sim \nu} [\sigma(\mathbf{x}_1; \mathbf{w}) \sigma(\mathbf{x}_2; \mathbf{w})]$ with $\sigma \in L^2(\mathcal{X} \times \mathcal{V})$.

Random feature approx: sample $(\mathbf{w}_s) \sim_{iid} \nu$ and replace K by

$$K_N(\mathbf{x}_1, \mathbf{x}_2) = \frac{1}{N} \sum_{s \in [N]} \sigma(\mathbf{x}_1; \mathbf{w}_s) \sigma(\mathbf{x}_2; \mathbf{w}_s).$$

Random Features Ridge Regression (RFRR): fit a model $\hat{f}_{N,\lambda}(\mathbf{x}) = \frac{1}{N} \sum_{s \in [N]} \hat{a}_s \sigma(\mathbf{x}; \mathbf{w}_s)$ with

$$\hat{\mathbf{a}} = \arg \min_{\mathbf{a} \in \mathbb{R}^N} \left\{ \sum_{i \in [n]} (y_i - f_N(\mathbf{x}_i; \mathbf{a}))^2 + \frac{\lambda}{N} \|\mathbf{a}\|_2^2 \right\}.$$

When $N \rightarrow \infty$, $\hat{f}_{N,\lambda} \rightarrow \hat{f}_\lambda$.

Universality and Gaussian equivalence model

Again, σ can be seen as a compact operator and is diagonalizable:

$$\sigma(\mathbf{x}; \mathbf{w}) = \sum_{j \geq 1} \sqrt{\lambda_j} \phi_j(\mathbf{x}) \psi_j(\mathbf{w}) = \langle \Phi(\mathbf{x}), \Psi(\mathbf{w}) \rangle_{\ell_2},$$

$$\Phi(\mathbf{x}) = (\lambda_j^{1/4} \phi_j(\mathbf{x}))_{j \geq 1}, \quad \Psi(\mathbf{w}) = (\lambda_j^{1/4} \psi_j(\mathbf{w}))_{j \geq 1},$$

with $\{\phi_j\}$ orthonormal basis of $L^2(\mathcal{X})$ and $\{\psi_j\}$ orthonormal basis of $L^2(\mathcal{V})$.

Gaussian equivalent model: $\Phi(\mathbf{x}_i) \leftrightarrow \mathbf{z}_i$, $\Psi(\mathbf{w}_j) \leftrightarrow \mathbf{g}_j$ and $\sigma(\mathbf{x}_i; \mathbf{w}_j) \leftrightarrow \langle \mathbf{z}_i, \mathbf{g}_j \rangle_{\ell_2}$.

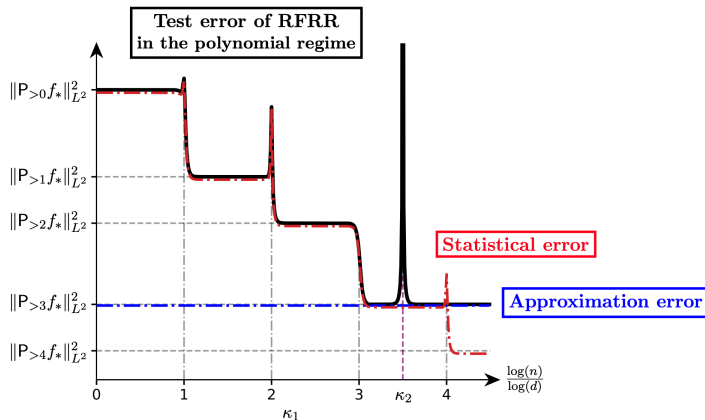
$$\mathbf{F} = (\sigma(\mathbf{x}_i; \mathbf{w}_j))_{i \in [n], j \in [M]}, \quad \mathbf{Z} = [\mathbf{z}_1, \dots, \mathbf{z}_n]^T \in \mathbb{R}^{n \times \infty}, \quad \mathbf{G} = [\mathbf{g}_1, \dots, \mathbf{g}_n]^T \in \mathbb{R}^{N \times \infty}.$$

$$\hat{f}_{N,\lambda}^G(\mathbf{z}) = \mathbf{z}^T \mathbf{G}^T \mathbf{G} \mathbf{Z}^T (\mathbf{Z} \mathbf{G}^T \mathbf{G} \mathbf{Z} + \lambda \mathbf{I}/N)^{-1} \mathbf{y}/N.$$

Asymptotics of RFRR on the sphere in polynomial scaling

$\mathbf{x} \sim \text{Unif}(\mathbb{S}^{d-1}(\sqrt{d}))$, $\mathbf{w} \sim \text{Unif}(\mathbb{S}^{d-1}(\sqrt{d}))$, $\hat{f}_{RF}(\mathbf{x}; \mathbf{a}) = \sum_{i \in [N]} a_i \sigma(\langle \mathbf{x}, \mathbf{w}_i \rangle)$.

Asymptotics in polynomial scaling $n/d^{\kappa_1} \rightarrow \psi_1$, $N/d^{\kappa_2} \rightarrow \psi_2$ for any $\kappa_1, \kappa_2, \psi_1, \psi_2 > 0$:



Test error $\approx \max(\text{approximation error}, \text{statistical error})$.

Application I: learning with group-invariance (1)

- ▶ Data invariant by group action \mathcal{G}_d (subgroup of $\mathcal{O}(d)$):
i.e., $f_*(g \cdot \mathbf{x}) = f_*(\mathbf{x})$ for all $g \in \mathcal{G}_d, \mathbf{x} \in \mathbb{S}^{d-1}(\sqrt{d})$.
- ▶ Comparison between learning with:
 - standard kernel $K(\mathbf{x}_1, \mathbf{x}_2) = h(\langle \mathbf{x}_1, \mathbf{x}_2 \rangle / d)$;
 - \mathcal{G}_d -invariant kernel $K(\mathbf{x}_1, \mathbf{x}_2) = \int_{\mathcal{G}_d} h(\langle \mathbf{x}_1, g \cdot \mathbf{x}_2 \rangle / d) \pi_d(dg)$.

- ▶ **Group \mathcal{G}_d of degeneracy α :** if for any $k \geq \alpha$,

$$\frac{\dim(V_{d,k})}{\dim(V_{d,k}(\mathcal{G}_d))} \asymp d^\alpha,$$

$V_{d,k}$: space of degree- k polynomials; $V_{d,k}(\mathcal{G}_d)$: degree- k \mathcal{G}_d -invariant polynomials.

Cyclic group: $\alpha = 1$ ($g_r \cdot \mathbf{x} = (x_{1+r}, x_{2+r}, \dots, x_d, x_1, \dots, x_r)$).

- ▶ To learn a degree- ℓ polynomial approximation needs $d^{\ell-\alpha}$ samples.
(Gain of a factor d^α in sample size and number of random features.)

Application I: learning with group-invariance (2)

Cyclic invariant MNIST:

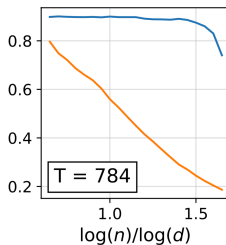
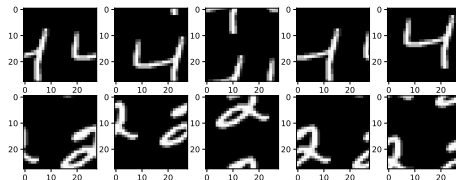


Figure: Test accuracy against number of samples (orange: cyclic kernel, blue: standard kernel).

Application II: learning with convolutional kernels

Covariates $\mathbf{x} = (x_1, \dots, x_d) \in \mathbb{R}^d$, patches of size q : $\mathbf{x}_{(k)} = (x_k, \dots, x_{k+q-1})$.

NTK of 1-layer convolutional kernel followed by local average pooling:

$$H_{q,\omega}(\mathbf{x}, \mathbf{z}) = \frac{1}{d\omega} \sum_{k \in [d]} \sum_{s, s' \in [\omega]} h(\langle \mathbf{x}_{k+s}, \mathbf{z}_{k+s'} \rangle / q).$$

For $\mathbf{x} \sim \text{Unif}(\{+1, -1\}^d)$, $H_{q,\omega}$ can be diagonalized and we can compute sharp asymptotics for the test error.

E.g., target function: $f_*(\mathbf{x}) = \frac{1}{d} \sum_{k \in [d]} P_\ell(\mathbf{x}_{(k)})$.

| To fit f_* | H^{FC} | H_{GP}^{FC} | H^{CK} | H_ω^{CK} | H_{GP}^{CK} |
|-------------------|----------|---------------|---------------|----------------------|---------------|
| Sample complexity | d^ℓ | $d^{\ell-1}$ | $dq^{\ell-1}$ | $dq^{\ell-1}/\omega$ | $q^{\ell-1}$ |

$$\begin{array}{lll} H^{FC}: q = d, \omega = 1; & H_{GP}^{FC}: q = d, \omega = d; & \\ H^{CK}: q, \omega = 1; & H_\omega^{CK}: q, \omega; & H_{GP}^{CK}: q, \omega = d. \end{array}$$

Application III: learning with anisotropic data (1)

Spiked covariates model: orthogonal matrix $[U, U^\perp]$

$$x = Uz_1 + U^\perp z_2, \quad z_1 \in \mathbb{R}^{d_s}, z_2 \in \mathbb{R}^{d-d_s}.$$

Signal part: $z_1 \sim \text{Unif}\left(\mathbb{S}^{d_s-1}(\sqrt{\text{snr}_c \cdot d_s})\right)$.

Noise part: $z_2 \sim \text{Unif}\left(\mathbb{S}^{d-d_s-1}(\sqrt{d-d_s})\right)$

d_s = signal dimension.

snr_c = covariate SNR.

Target function: $f_*(x) = \varphi(z_1)$.

Define effective dimension: $d_{\text{eff}} = d_s \vee (d/\text{snr}_c)$

$$\text{for } d_{\text{eff}}^{\ell+\delta} \leq n \leq d_{\text{eff}}^{\ell+1-\delta}, \quad R(f_*, \hat{f}) = \|P_{>\ell} f_*\|_{L^2}^2 + o_{d,\mathbb{P}}(1).$$

- ▶ Approx. isotropic data: $\text{snr}_c \approx 1$, $d_{\text{eff}} \approx d$.
- ▶ Very anisotropic data: $\text{snr}_c \gg 1$, $d_{\text{eff}} \approx d_s \ll d$. KRR much more efficient.

Application III: learning with anisotropic data (2)

- For images: (1) Spectrum concentrates on low-frequencies;
(2) Labels depend predominantly on low-frequencies.

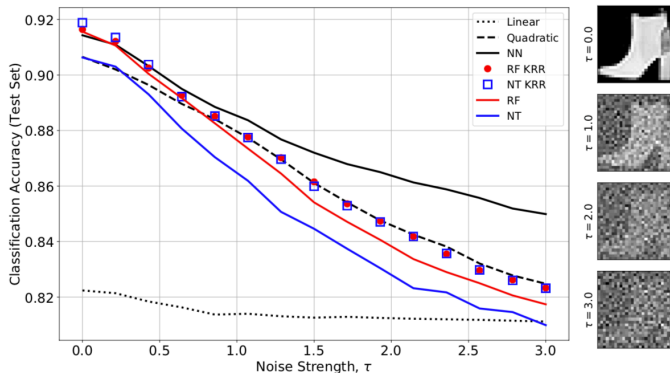


Figure: Test accuracy on FMNIST: adding noise to the high frequency components (decreasing snr_c , increasing $d_{\text{eff}} = d/snr_c$).

- ▶ Gaussian equivalent model for "high-dimensional enough models".
- ▶ Can give very precise results which give clear conceptual picture.
- ▶ **More general type of universality:** entire feature maps + polynomial scaling $\log(n) \asymp \log(d)$.
- ▶ **Limitations:** hard to apply it to specific setting (most of the time, no explicit diagonalization).
- ▶ More general directions (ERM universality).

Thank you!

1. *Linearized two-layers neural networks in high dimension*. Ghorbani, Mei, **M.**, Montanari (2019).
2. *When do neural networks outperform kernel methods?* Ghorbani, Mei, **M.**, Montanari (2020).
3. *Generalization error of random feature and kernel methods: Hypercontractivity and kernel matrix concentration*. Mei, **M.**, Montanari (2021).
4. *Learning with invariances in random features and kernel models*. Mei, **M.**, Montanari (2021).
5. *Learning with convolution and pooling operations in kernel methods*. **M.**, Mei (2021).
6. *Spectrum of inner-product kernel matrices in the polynomial scaling and multiple descent phenomenon in kernel ridge regression*. **M.** (2022).

+ some ongoing work.