Learning structured data with random features and kernel methods

Theodor Misiakiewicz

Stanford University

April 28th, 2021

ML lunch

Joint work with Behrooz Ghorbani (Google Research), Song Mei (UC Berkeley) and Andrea Montanari (Stanford)

- 'Mysteries' of Deep Learning:
 - Optimization: training is a highly non-convex problem, but we are still able to find near global optimizers.
 - Generalization: *overparametrized* models. The solution often interpolates the training data while generalizing well on test data.

In this talk, I will focus on the connection between neural networks (NNs) and kernel machines.

- 'Mysteries' of Deep Learning:
 - Optimization: training is a highly non-convex problem, but we are still able to find near global optimizers.
 - Generalization: overparametrized models. The solution often interpolates the training data while generalizing well on test data.

In this talk, I will focus on the connection between neural networks (NNs) and kernel machines.

Neural network: NN_p(x; θ), x ∈ ℝ^d, θ ∈ ℝ^p. e.g., fully-connected neural network:

$$\mathsf{NN}_{\rho}(\boldsymbol{x};\boldsymbol{\theta}) = \langle \boldsymbol{a}, \sigma(\ldots \boldsymbol{W}_{2}\sigma(\boldsymbol{W}_{1}\boldsymbol{x})) \rangle.$$

Linearization around random initialization θ₀:

 $NN_{\rho}(x; \theta) = NN_{\rho}(x; \theta_{0}) + \langle \theta - \theta_{0}, \nabla_{\theta} NN_{\rho}(x; \theta_{0}) \rangle + o(\|\theta - \theta_{0}\|_{2}).$

Neural Tangent (NT) model: [Chizat et al., '18], [Du et al., '18]

 $\mathsf{NT}_{p}(\boldsymbol{x};\boldsymbol{\beta},\boldsymbol{\theta}_{0}) = \langle \boldsymbol{\beta}, \nabla_{\boldsymbol{\theta}} \mathsf{NN}_{p}(\boldsymbol{x};\boldsymbol{\theta}_{0}) \rangle.$

Finite-width approximation of limiting kernel $(p \rightarrow \infty + \text{iid initialization})$:

 $K_{\mathsf{NT}}(x, y) = \mathbb{E}_{\boldsymbol{\theta}_{\mathbf{0}}}[\langle \nabla_{\boldsymbol{\theta}} \mathsf{NN}_{\boldsymbol{\rho}}(x; \boldsymbol{\theta}_{\mathbf{0}}), \nabla_{\boldsymbol{\theta}} \mathsf{NN}_{\boldsymbol{\rho}}(y; \boldsymbol{\theta}_{\mathbf{0}}) \rangle]$

Neural network: NN_p(x; θ), x ∈ ℝ^d, θ ∈ ℝ^p. e.g., fully-connected neural network:

$$\mathsf{NN}_{\rho}(\boldsymbol{x};\boldsymbol{\theta}) = \langle \boldsymbol{a}, \sigma(\ldots \boldsymbol{W}_{2}\sigma(\boldsymbol{W}_{1}\boldsymbol{x})) \rangle.$$

Linearization around random initialization θ₀:

 $\mathsf{NN}_{p}(\boldsymbol{x};\boldsymbol{\theta}) = \mathsf{NN}_{p}(\boldsymbol{x};\boldsymbol{\theta}_{0}) + \underline{\langle \boldsymbol{\theta} - \boldsymbol{\theta}_{0}, \nabla_{\boldsymbol{\theta}} \mathsf{NN}_{p}(\boldsymbol{x};\boldsymbol{\theta}_{0}) \rangle} + o(\|\boldsymbol{\theta} - \boldsymbol{\theta}_{0}\|_{2}).$

Neural Tangent (NT) model: [Chizat et al., '18], [Du et al., '18]

 $\mathsf{NT}_{\rho}(\boldsymbol{x};\boldsymbol{\beta},\boldsymbol{\theta}_0) = \langle \boldsymbol{\beta}, \nabla_{\boldsymbol{\theta}} \mathsf{NN}_{\rho}(\boldsymbol{x};\boldsymbol{\theta}_0) \rangle.$

Finite-width approximation of limiting kernel ($p \rightarrow \infty$ + iid initialization):

 $K_{\mathsf{NT}}(x, y) = \mathbb{E}_{\boldsymbol{\theta}_{0}}[\langle \nabla_{\boldsymbol{\theta}} \mathsf{NN}_{\boldsymbol{\rho}}(x; \boldsymbol{\theta}_{0}), \nabla_{\boldsymbol{\theta}} \mathsf{NN}_{\boldsymbol{\rho}}(y; \boldsymbol{\theta}_{0}) \rangle]$

Neural network: NN_p(x; θ), x ∈ ℝ^d, θ ∈ ℝ^p. e.g., fully-connected neural network:

$$\mathsf{NN}_{\rho}(\boldsymbol{x};\boldsymbol{\theta}) = \langle \boldsymbol{a}, \sigma(\ldots \boldsymbol{W}_{2}\sigma(\boldsymbol{W}_{1}\boldsymbol{x})) \rangle.$$

Linearization around random initialization θ₀:

$$\mathsf{NN}_{\rho}(\boldsymbol{x};\boldsymbol{\theta}) = \mathsf{NN}_{\rho}(\boldsymbol{x};\boldsymbol{\theta}_{0}) + \underline{\langle \boldsymbol{\theta} - \boldsymbol{\theta}_{0}, \nabla_{\boldsymbol{\theta}}\mathsf{NN}_{\rho}(\boldsymbol{x};\boldsymbol{\theta}_{0}) \rangle} + o(\|\boldsymbol{\theta} - \boldsymbol{\theta}_{0}\|_{2}).$$

Neural Tangent (NT) model: [Chizat et al., '18], [Du et al., '18]

 $\mathsf{NT}_{\rho}(\boldsymbol{x};\boldsymbol{\beta},\boldsymbol{\theta}_0) = \langle \boldsymbol{\beta}, \nabla_{\boldsymbol{\theta}} \mathsf{NN}_{\rho}(\boldsymbol{x};\boldsymbol{\theta}_0) \rangle.$

Finite-width approximation of limiting kernel ($p \rightarrow \infty$ + iid initialization):

 $K_{\mathsf{NT}}(x, y) = \mathbb{E}_{\theta_0}[\langle \nabla_{\theta} \mathsf{NN}_p(x; \theta_0), \nabla_{\theta} \mathsf{NN}_p(y; \theta_0) \rangle]$

Neural network: NN_ρ(x; θ), x ∈ ℝ^d, θ ∈ ℝ^p. e.g., fully-connected neural network:

$$\mathsf{NN}_{\rho}(\boldsymbol{x};\boldsymbol{\theta}) = \langle \boldsymbol{a}, \sigma(\ldots \boldsymbol{W}_{2}\sigma(\boldsymbol{W}_{1}\boldsymbol{x})) \rangle.$$

Linearization around random initialization θ₀:

$$\mathsf{NN}_{\rho}(\boldsymbol{x};\boldsymbol{\theta}) = \mathsf{NN}_{\rho}(\boldsymbol{x};\boldsymbol{\theta}_{0}) + \langle \boldsymbol{\theta} - \boldsymbol{\theta}_{0}, \nabla_{\boldsymbol{\theta}}\mathsf{NN}_{\rho}(\boldsymbol{x};\boldsymbol{\theta}_{0}) \rangle + o(\|\boldsymbol{\theta} - \boldsymbol{\theta}_{0}\|_{2}).$$

Neural Tangent (NT) model: [Chizat et al., '18], [Du et al., '18]

$$\mathsf{NT}_{\rho}(\boldsymbol{x};\boldsymbol{\beta},\boldsymbol{\theta}_{0}) = \langle \boldsymbol{\beta}, \nabla_{\boldsymbol{\theta}} \mathsf{NN}_{\rho}(\boldsymbol{x};\boldsymbol{\theta}_{0}) \rangle.$$

Finite-width approximation of limiting kernel $(p \rightarrow \infty + \text{iid initialization})$:

$$\mathcal{K}_{\mathsf{NT}}(\mathbf{x}, \mathbf{y}) = \mathbb{E}_{\boldsymbol{\theta}_{\mathbf{0}}}[\langle \nabla_{\boldsymbol{\theta}} \mathsf{NN}_{\boldsymbol{\rho}}(\mathbf{x}; \boldsymbol{\theta}_{\mathbf{0}}), \nabla_{\boldsymbol{\theta}} \mathsf{NN}_{\boldsymbol{\rho}}(\mathbf{y}; \boldsymbol{\theta}_{\mathbf{0}}) \rangle]$$

GD trained NNs \approx NT model in some regime

Coupled gradient descent on empirical squared loss:

$$\begin{split} \frac{\mathrm{d}}{\mathrm{d}t} \boldsymbol{\theta}^t &= -\nabla_{\boldsymbol{\theta}} \hat{\mathbb{E}}[(\boldsymbol{y} - \mathsf{NN}_{\boldsymbol{\rho}}(\boldsymbol{x}; \boldsymbol{\theta}^t))^2], \qquad \boldsymbol{\theta}^0 = \boldsymbol{\theta}_0, \\ \frac{\mathrm{d}}{\mathrm{d}t} \boldsymbol{\beta}^t &= -\nabla_{\boldsymbol{\beta}} \hat{\mathbb{E}}[(\boldsymbol{y} - \mathsf{NT}_{\boldsymbol{\rho}}(\boldsymbol{x}; \boldsymbol{\beta}^t, \boldsymbol{\theta}_0))^2], \qquad \boldsymbol{\beta}^0 = 0. \end{split}$$

Theorem (informal)

For any $\varepsilon > 0$, for large enough NNs and proper random initialization θ_0 (Xavier initialization), we have with high probability

$$\sup_{t\geq 0} \sup_{\|x\|_{2}=1} |\mathsf{NN}_{p}(x; \boldsymbol{\theta}^{t}) - \mathsf{NT}_{p}(x; \boldsymbol{\beta}^{t}, \boldsymbol{\theta}_{0})| \leq \varepsilon.$$

 \implies Intuition: weights do not change much and $\|\boldsymbol{\theta}^t - \boldsymbol{\theta}_0\|_2$ remains small.

[Jacot, Gabriel, Hongler,'18] , [Du, Zhai, Poczos, Singh,'18], [Du, Lee, Li, Wang, Zhai,'18], [Allen-Zhu, Li, Song,'18] , [Zou, Cao, Zhou, Gu,'18], [Oymak, Soltanolkotabi,'18], [Chizat, Oyallon, Bach,'19], ...

GD trained NNs \approx NT model in some regime

Coupled gradient descent on empirical squared loss:

$$\begin{split} \frac{\mathrm{d}}{\mathrm{d}t} \boldsymbol{\theta}^t &= -\nabla_{\boldsymbol{\theta}} \hat{\mathbb{E}}[(\boldsymbol{y} - \mathsf{NN}_{\boldsymbol{\rho}}(\boldsymbol{x}; \boldsymbol{\theta}^t))^2], \qquad \boldsymbol{\theta}^0 = \boldsymbol{\theta}_0, \\ \frac{\mathrm{d}}{\mathrm{d}t} \boldsymbol{\beta}^t &= -\nabla_{\boldsymbol{\beta}} \hat{\mathbb{E}}[(\boldsymbol{y} - \mathsf{NT}_{\boldsymbol{\rho}}(\boldsymbol{x}; \boldsymbol{\beta}^t, \boldsymbol{\theta}_0))^2], \qquad \boldsymbol{\beta}^0 = 0. \end{split}$$

Theorem (informal)

For any $\varepsilon > 0$, for large enough NNs and proper random initialization θ_0 (Xavier initialization), we have with high probability

$$\sup_{t\geq 0}\sup_{\|\boldsymbol{x}\|_{2}=1}|\mathsf{NN}_{\rho}(\boldsymbol{x};\boldsymbol{\theta}^{t})-\mathsf{NT}_{\rho}(\boldsymbol{x};\boldsymbol{\beta}^{t},\boldsymbol{\theta}_{0})|\leq\varepsilon.$$

 \implies Intuition: weights do not change much and $\|\boldsymbol{\theta}^t - \boldsymbol{\theta}_0\|_2$ remains small.

[Jacot, Gabriel, Hongler,'18], [Du, Zhai, Poczos, Singh,'18], [Du, Lee, Li, Wang, Zhai,'18], [Allen-Zhu, Li, Song,'18], [Zou, Cao, Zhou, Gu,'18], [Oymak, Soltanolkotabi,'18], [Chizat, Oyallon, Bach,'19], ...

The 'Kernel Regime'

Optimization success:

With sufficient over-parametrization and proper initialization, gradient descent on training loss of NN converges linearly to global minimum.

... but:

- Lazy Training: weights hardly change, there is 'no feature learning'. [Chizat et al.,'18]
- Empirically, NNs perform often better than their linearized counterparts.

However, offer theoretical insights [Bartlett, Montanari, Rakhlin,'21]:

- Tractability via overparametrization.
- Explain some of the trade-off between overparametrization and generalization (benign overfitting, double descent).

The 'Kernel Regime'

Optimization success:

With sufficient over-parametrization and proper initialization, gradient descent on training loss of NN converges linearly to global minimum.

... but:

Lazy Training: weights hardly change, there is 'no feature learning'. [Chizat et al.,'18]

Empirically, NNs perform often better than their linearized counterparts.

However, offer theoretical insights [Bartlett, Montanari, Rakhlin,'21]:

Tractability via overparametrization.

Explain some of the trade-off between overparametrization and generalization (benign overfitting, double descent).

The 'Kernel Regime'

Optimization success:

With sufficient over-parametrization and proper initialization, gradient descent on training loss of NN converges linearly to global minimum.

... but:

- Lazy Training: weights hardly change, there is 'no feature learning'. [Chizat et al.,'18]
- Empirically, NNs perform often better than their linearized counterparts.

However, offer theoretical insights [Bartlett, Montanari, Rakhlin,'21]:

- Tractability via overparametrization.
- Explain some of the trade-off between overparametrization and generalization (benign overfitting, double descent).

New approach for kernel engineering

Progress on CIFAR-10:

Paper	Method	Accuracy
[Coates, Lee, Ng,'11]	feature learning + linear regression	77%
[Arora et al.,'19]	CNTK (data independent)	77%
[Li et al.,'19]	CRFK + data dependent preprocessing	89%
[Shankar et al.,'20]	CRFK (data independent)	90%
-	CNN	> 99%

NTK achieves near state-of-the-art results on UCI dataset [Arora et al.,'19].

Observations:

- NNs often outperform kernel methods.
- Kernel methods sometimes achieve comparable accuracy to NNs.
- Kernels induced by NNs perform much better than previous handcrafted kernels.

Neural Networks and Kernel Methods

New approach for kernel engineering

Progress on CIFAR-10:

Paper	Method	Accuracy
[Coates, Lee, Ng,'11]	feature learning + linear regression	77%
[Arora et al.,'19]	CNTK (data independent)	77%
[Li et al.,'19]	CRFK + data dependent preprocessing	89%
[Shankar et al.,'20]	CRFK (data independent)	90%
-	CNN	> 99%

NTK achieves near state-of-the-art results on UCI dataset [Arora et al.,'19].

Observations

- NNs often outperform kernel methods.
- Kernel methods sometimes achieve comparable accuracy to NNs.
- Kernels induced by NNs perform much better than previous handcrafted kernels.

Neural Networks and Kernel Methods

New approach for kernel engineering

Progress on CIFAR-10:

Paper	Method	Accuracy
[Coates, Lee, Ng,'11]	feature learning + linear regression	77%
[Arora et al.,'19]	CNTK (data independent)	77%
[Li et al.,'19]	CRFK + data dependent preprocessing	89%
[Shankar et al.,'20]	CRFK (data independent)	90%
-	CNN	> 99%

▶ NTK achieves near state-of-the-art results on UCI dataset [Arora et al.,'19].

Observations:

- NNs often outperform kernel methods.
- Kernel methods sometimes achieve comparable accuracy to NNs.
- Kernels induced by NNs perform much better than previous handcrafted kernels.

Outline of the talk

Questions

- Why are kernel methods not as good as NNs in general?
- When can we expect kernel methods to perform well, comparable to NNs?
- Can we quantify the benefits of using convolutional kernels against inner-product kernels (e.g., Gaussian kernel)?

Three 'stylized' scenarios:

- A. Isotropic covariates model:
 - large gap between NN and kernel methods.
- B. Spiked covariates model:
 - smaller gap between NN and kernel methods.

C. Invariant function estimation:

quantify the benefits of convolutional kernels over inner-product kernels.

Outline of the talk

Questions

- Why are kernel methods not as good as NNs in general?
- When can we expect kernel methods to perform well, comparable to NNs?
- Can we quantify the benefits of using convolutional kernels against inner-product kernels (e.g., Gaussian kernel)?

Three 'stylized' scenarios:

- A. Isotropic covariates model:
 - large gap between NN and kernel methods.
- B. Spiked covariates model:
 - smaller gap between NN and kernel methods.

C. Invariant function estimation:

quantify the benefits of convolutional kernels over inner-product kernels.

A. Isotropic covariates model

► Given *n* iid samples
$$\{(y_i, \mathbf{x}_i)\}_{i \in [n]}$$
:
 $y_i = f_*(\mathbf{x}_i) + \varepsilon_i, \quad \mathbf{x}_i \sim_{iid} \text{Unif}(\mathbb{S}^{d-1}(\sqrt{d})), \quad \varepsilon_i \sim_{iid} N(0, \tau^2).$

► Two-layers neural networks: $NN_N(x; \theta) = \sum_{i=1}^N a_i \sigma(\langle w_i, x \rangle).$

Random Features model [Rahimi, Recht,'08]:

$$\mathsf{RF}_{N}(x; a, W) = \sum_{i=1}^{N} a_{i}\sigma(\langle w_{i}, x \rangle) = \langle a, \nabla_{a}\mathsf{NN}_{p}(x; \theta_{0}) \rangle.$$

▶ Ridge regression with RF model (with fixed $W = (w_i)_{i \in [N]} \sim_{iid} Unif(S^{d-1}))$

$$\hat{\boldsymbol{a}}(\lambda) = \arg\min_{\boldsymbol{a}} \left\{ \frac{1}{n} \sum_{i=1}^{n} \left(y_i - \sum_{j=1}^{N} a_j \sigma(\langle \boldsymbol{w}_j, \boldsymbol{x}_i \rangle) \right)^2 + \lambda \|\boldsymbol{a}\|_2^2 \right\}.$$

Denote the solution $\widehat{\mathsf{RF}}_{n,N,\lambda} = \mathsf{RF}_N(\cdot; \hat{a}(\lambda), W).$

► Given *n* iid samples $\{(y_i, \mathbf{x}_i)\}_{i \in [n]}$: $y_i = f_*(\mathbf{x}_i) + \varepsilon_i, \quad \mathbf{x}_i \sim_{iid} \text{Unif}(\mathbb{S}^{d-1}(\sqrt{d})), \quad \varepsilon_i \sim_{iid} N(0, \tau^2).$

• Two-layers neural networks: $NN_N(\mathbf{x}; \boldsymbol{\theta}) = \sum_{i=1}^N a_i \sigma(\langle \mathbf{w}_i, \mathbf{x} \rangle).$

Random Features model [Rahimi, Recht,'08]:

$$\mathsf{RF}_{N}(x; a, W) = \sum_{i=1}^{N} a_{i}\sigma(\langle w_{i}, x \rangle) = \langle a, \nabla_{a}\mathsf{NN}_{p}(x; \theta_{0}) \rangle.$$

▶ Ridge regression with RF model (with fixed $W = (w_i)_{i \in [N]} \sim_{iid} \text{Unif}(\mathbb{S}^{d-1})$)

$$\hat{\boldsymbol{a}}(\lambda) = \arg\min_{\boldsymbol{a}} \left\{ \frac{1}{n} \sum_{i=1}^{n} \left(y_i - \sum_{j=1}^{N} a_j \sigma(\langle \boldsymbol{w}_j, \boldsymbol{x}_i \rangle) \right)^2 + \lambda \|\boldsymbol{a}\|_2^2 \right\}.$$

Denote the solution $\widehat{\mathsf{RF}}_{n,N,\lambda} = \mathsf{RF}_N(\cdot; \hat{a}(\lambda), W).$

► Given *n* iid samples $\{(y_i, x_i)\}_{i \in [n]}$: $y_i = f_*(x_i) + \varepsilon_i, \quad x_i \sim_{iid} \text{Unif}(\mathbb{S}^{d-1}(\sqrt{d})), \quad \varepsilon_i \sim_{iid} N(0, \tau^2).$

• Two-layers neural networks: $NN_N(\mathbf{x}; \boldsymbol{\theta}) = \sum_{i=1}^N a_i \sigma(\langle \mathbf{w}_i, \mathbf{x} \rangle).$

Random Features model [Rahimi, Recht,'08]:

$$\mathsf{RF}_N(\mathbf{x}; \mathbf{a}, \mathbf{W}) = \sum_{i=1}^N \mathbf{a}_i \sigma(\langle \mathbf{w}_i, \mathbf{x} \rangle) = \langle \mathbf{a}, \nabla_{\mathbf{a}} \mathsf{NN}_p(\mathbf{x}; \boldsymbol{\theta}_0) \rangle.$$

▶ Ridge regression with RF model (with fixed $W = (w_i)_{i \in [N]} \sim_{iid} \text{Unif}(\mathbb{S}^{d-1})$)

$$\hat{\boldsymbol{a}}(\lambda) = \arg\min_{\boldsymbol{a}} \left\{ \frac{1}{n} \sum_{i=1}^{n} \left(y_i - \sum_{j=1}^{N} a_j \sigma(\langle \boldsymbol{w}_j, \boldsymbol{x}_i \rangle) \right)^2 + \lambda \|\boldsymbol{a}\|_2^2 \right\}.$$

Denote the solution $\widehat{\mathsf{RF}}_{n,N,\lambda} = \mathsf{RF}_N(\cdot; \hat{a}(\lambda), W).$

► Given *n* iid samples $\{(y_i, \mathbf{x}_i)\}_{i \in [n]}$: $y_i = f_*(\mathbf{x}_i) + \varepsilon_i, \quad \mathbf{x}_i \sim_{iid} \text{Unif}(\mathbb{S}^{d-1}(\sqrt{d})), \quad \varepsilon_i \sim_{iid} N(0, \tau^2).$

• Two-layers neural networks: $NN_N(\mathbf{x}; \boldsymbol{\theta}) = \sum_{i=1}^N \mathbf{a}_i \sigma(\langle \mathbf{w}_i, \mathbf{x} \rangle).$

Random Features model [Rahimi, Recht,'08]:

$$\mathsf{RF}_{N}(\boldsymbol{x};\boldsymbol{a},\boldsymbol{W}) = \sum_{i=1}^{N} \boldsymbol{a}_{i} \sigma(\langle \boldsymbol{w}_{i},\boldsymbol{x} \rangle) = \langle \boldsymbol{a}, \nabla_{\boldsymbol{a}} \mathsf{NN}_{P}(\boldsymbol{x};\boldsymbol{\theta}_{0}) \rangle.$$

▶ Ridge regression with RF model (with fixed $W = (w_i)_{i \in [N]} \sim_{iid} \text{Unif}(\mathbb{S}^{d-1})$)

$$\hat{\boldsymbol{a}}(\boldsymbol{\lambda}) = \arg\min_{\boldsymbol{a}} \left\{ \frac{1}{n} \sum_{i=1}^{n} \left(y_i - \sum_{j=1}^{N} a_j \sigma(\langle \boldsymbol{w}_j, \boldsymbol{x}_i \rangle) \right)^2 + \lambda \|\boldsymbol{a}\|_2^2 \right\}.$$

Denote the solution $\widehat{\mathsf{RF}}_{n,N,\lambda} = \mathsf{RF}_N(\cdot; \hat{a}(\lambda), \mathcal{W}).$

Theorem (Mei, Misiakiewicz, Montanari, 2021)

Assume $d^{\ell+\delta} \leq \min(n, N) \leq d^{\ell+1-\delta}$ and $|\log(n/N)| \geq \delta \log d$, σ satisfies some generic conditions and λ is small enough. Then,

$$\|f_{\star} - \widehat{\mathsf{RF}}_{n,N,\lambda}\|_{L^2}^2 = \|\mathsf{P}_{>\ell}f_{\star}\|_{L^2}^2 + o_{d,\mathbb{P}}(\cdot).$$

 $P_{>\ell}$: projection orthogonal to the space of degree- ℓ polynomials.

- For d^{ℓ} parameters and samples, RF fits the best degree- ℓ polynomial approximation.
- **n** and **N** play a symmetric role:

Test error $(n, N) = \max \left\{ \text{approx. error } (n = \infty, N), \text{ statistical error } (n, N = \infty) \right\}$

Theorem (Mei, Misiakiewicz, Montanari, 2021)

Assume $d^{\ell+\delta} \leq \min(n, N) \leq d^{\ell+1-\delta}$ and $|\log(n/N)| \geq \delta \log d$, σ satisfies some generic conditions and λ is small enough. Then,

$$\|f_{\star} - \widehat{\mathsf{RF}}_{n,N,\lambda}\|_{L^2}^2 = \|\mathsf{P}_{>\ell}f_{\star}\|_{L^2}^2 + o_{d,\mathbb{P}}(\cdot).$$

 $\mathsf{P}_{>\ell}$: projection orthogonal to the space of degree- ℓ polynomials.

For d^{ℓ} parameters and samples, RF fits the best degree- ℓ polynomial approximation.

n and **N** play a symmetric role:

Test error $(n, N) = \max \left\{ \text{approx. error } (n = \infty, N), \text{ statistical error } (n, N = \infty) \right\}$

Theorem (Mei, Misiakiewicz, Montanari, 2021)

Assume $d^{\ell+\delta} \leq \min(n, N) \leq d^{\ell+1-\delta}$ and $|\log(n/N)| \geq \delta \log d$, σ satisfies some generic conditions and λ is small enough. Then,

$$\|f_{\star} - \widehat{\mathsf{RF}}_{n,N,\lambda}\|_{L^2}^2 = \|\mathsf{P}_{>\ell}f_{\star}\|_{L^2}^2 + o_{d,\mathbb{P}}(\cdot).$$

 $\mathsf{P}_{>\ell}$: projection orthogonal to the space of degree- ℓ polynomials.

- For d^{ℓ} parameters and samples, RF fits the best degree- ℓ polynomial approximation.
- n and N play a symmetric role:

Test error $(n, N) = \max \{ \text{approx. error } (n = \infty, N), \text{ statistical error } (n, N = \infty) \}.$

Theorem (Mei, Misiakiewicz, Montanari, 2021)

Assume $d^{\ell+\delta} \leq \min(n, N) \leq d^{\ell+1-\delta}$ and $|\log(n/N)| \geq \delta \log d$, σ satisfies some generic conditions and λ is small enough. Then,

$$\|f_{\star}-\widehat{\mathsf{RF}}_{n,\boldsymbol{N},\lambda}\|_{L^{2}}^{2}=\|\mathsf{P}_{>\ell}f_{\star}\|_{L^{2}}^{2}+o_{d,\mathbb{P}}(\cdot).$$

 $\mathsf{P}_{>\ell}$: projection orthogonal to the space of degree- ℓ polynomials.

- For d^{ℓ} parameters and samples, RF fits the best degree- ℓ polynomial approximation.
- n and N play a symmetric role:

Test error $(n, N) = \max \{ \text{approx. error } (n = \infty, N), \text{ statistical error } (n, N = \infty) \}.$

The truncated staircase decay



 $f = \mathsf{P}_0 f + \mathsf{P}_1 f + \mathsf{P}_2 f + \mathsf{P}_3 f + \mathsf{P}_4 f.$

Figure: Test error (in red) v.s. $\log(N)/\log(d)$ for $n = d^{2.4}$ and $\lambda = 0^+$. Train error in green.

Theodor Misiakiewicz (Stanford)

Neural Networks and Kernel Methods

Performance gap between NN and kernel methods

We expect NNs to perform way better on 'low complexity' functions E.g., $f_{\star}(\mathbf{x}) = \sigma(\langle \mathbf{w}_{\star}, \mathbf{x} \rangle)$: if $N \propto 1$ and $n \propto d$, then GD solution: $\|f_{\star} - \widehat{\text{NN}}_{N,n}\|_{L^{2}}^{2} = o_{d,\mathbb{P}}(\cdot)$ [Bai et al.,'16]

while for KRR:

$$\|f_{\star} - \widehat{\operatorname{KR}}_{N,n}\|_{L^2}^2 = \|\mathsf{P}_{>1}f_{\star}\|_{L^2}^2 + o_{d,\mathbb{P}}(\cdot).$$

More generally we expect NNs to vastly outperform kernel methods for target functions that only depend on a low dimensional subspace of the data [Bach,'17].

Intuition: fit $\varphi(x_1)$ using features $\sigma(\langle w_i, x \rangle)$

For RF, w_i 's are random, $\operatorname{corr}(x_1e_1, \langle w_i, x \rangle)^2 = \langle w_i, e_1 \rangle^2 / ||w_i||_2^2$ is small in HD.

For NN, w_i 's can be chosen to have a large correlation with e_1 .

NN can 'adaptively learn' w;'s while RF cannot

Performance gap between NN and kernel methods

We expect NNs to perform way better on 'low complexity' functions E.g., $f_{\star}(\mathbf{x}) = \sigma(\langle \mathbf{w}_{\star}, \mathbf{x} \rangle)$: if $N \propto 1$ and $n \propto d$, then GD solution: $\|f_{\star} - \widehat{NN}_{N,n}\|_{L^2}^2 = o_{d,\mathbb{P}}(\cdot)$ [Bai et al.,'16]

while for KRR:

$$\|f_{\star} - \widehat{\operatorname{KR}}_{N,n}\|_{L^2}^2 = \|\mathsf{P}_{>1}f_{\star}\|_{L^2}^2 + o_{d,\mathbb{P}}(\cdot).$$

More generally we expect NNs to vastly outperform kernel methods for target functions that only depend on a low dimensional subspace of the data [Bach,'17].

ntuition: fit $\varphi(x_1)$ using features $\sigma(\langle w_i, x \rangle)$

For RF, w_i 's are random, $\operatorname{corr}(x_1e_1, \langle w_i, x \rangle)^2 = \langle w_i, e_1 \rangle^2 / ||w_i||_2^2$ is small in HD.

For NN, w_i 's can be chosen to have a large correlation with e_1 .

NN can 'adaptively learn' w_i's while RF cannot.

Performance gap between NN and kernel methods

We expect NNs to perform way better on 'low complexity' functions E.g., f_{*}(x) = σ(⟨w_{*}, x⟩): if N ∝ 1 and n ∝ d, then GD solution: ||f_{*} − NN_{N,n}||²_{L²} = o_{d,P}(·) [Bai et al.,'16]

while for KRR:

$$\|f_{\star} - \widehat{\operatorname{KR}}_{N,n}\|_{L^2}^2 = \|\mathsf{P}_{>1}f_{\star}\|_{L^2}^2 + o_{d,\mathbb{P}}(\cdot).$$

More generally we expect NNs to vastly outperform kernel methods for target functions that only depend on a low dimensional subspace of the data [Bach,'17].

Intuition: fit $\varphi(x_1)$ using features $\sigma(\langle \mathbf{w}_i, \mathbf{x} \rangle)$

- For RF, \boldsymbol{w}_i 's are random, $\operatorname{corr}(x_1\boldsymbol{e}_1, \langle \boldsymbol{w}_i, \boldsymbol{x} \rangle)^2 = \langle \boldsymbol{w}_i, \boldsymbol{e}_1 \rangle^2 / \|\boldsymbol{w}_i\|_2^2$ is small in HD.
- For NN, w_i 's can be chosen to have a large correlation with e_1 .

NN can 'adaptively learn' w_i 's while RF cannot.

Theodor Misiakiewicz (Stanford)

Neural Networks and Kernel Methods

In the isotropic covariates model:

- Kernel methods suffer from the **curse of dimensionality**.
- For target functions that depend on a low dimensional projection of the data, NNs can adaptively learn a good representation of the data and vastly outperform kernel methods.

B. Spiked covariates model

Spiked covariates model

Covariate vector: orthogonal matrix $[\boldsymbol{U}, \boldsymbol{U}^{\perp}]$

$$\mathbf{x} = \mathbf{U}\mathbf{z}_1 + \mathbf{U}^{\perp}\mathbf{z}_2, \qquad \mathbf{z}_1 \in \mathbb{R}^{d_s}, \mathbf{z}_2 \in \mathbb{R}^{d-d_s}.$$

Signal part: $z_1 \sim \text{Unif}\left(\mathbb{S}^{d_s-1}(\sqrt{\operatorname{snr}_c \cdot d_s})\right)$. Noise part: $z_2 \sim \text{Unif}\left(\mathbb{S}^{d-d_s-1}(\sqrt{d-d_s})\right)$

 $d_s = \text{signal dimension}.$

 $snr_c = covariate SNR.$

Target function: $f_{\star}(\mathbf{x}) = \varphi(\mathbf{z}_1)$.



Figure: Isotropic covariates: $snr_c = 1.$



Figure: Anisotropic covariates: $snr_c > 1$.

Test error of KRR in the spiked covariates model

- Given iid samples $(\{x_i, y_i\})_{i \in [n]}$ from the spiked covariates model.
- Effective dimension: $d_{eff} = d_s \vee (d/\operatorname{snr}_c)$.

Theorem (Ghorbani, Mei, **Misiakiewicz**, Montanari, 2020)

Assume $d_{\text{eff}}^{\ell+\delta} \leq \min(n, N) \leq d_{\text{eff}}^{\ell+1-\delta}$ and $|\log(n/N)| \geq \delta \log d_{\text{eff}}$, σ satisfies some generic conditions and λ is small enough. Then,

$$\|f_{\star}-\widehat{\mathsf{RF}}_{n,\boldsymbol{N},\lambda}\|_{L^2}^2=\|\mathsf{P}_{>\ell}f_{\star}\|_{L^2}^2+o_{d,\mathbb{P}}(\cdot).$$

d_{eff} capture the 'effective low-dimensionality' of the data.

- d_{eff} decreases with snr_c (data more anisotropic) and kernel methods will perform better.
- On the contrary, we expect NNs to learn features w_i aligned with z₁ and to depend mildly on snr_c.

Test error of KRR in the spiked covariates model

- Given iid samples $(\{x_i, y_i\})_{i \in [n]}$ from the spiked covariates model.
- Effective dimension: $d_{eff} = d_s \lor (d/\operatorname{snr}_c)$.

Theorem (Ghorbani, Mei, **Misiakiewicz**, Montanari, 2020)

Assume $d_{\text{eff}}^{\ell+\delta} \leq \min(n, N) \leq d_{\text{eff}}^{\ell+1-\delta}$ and $|\log(n/N)| \geq \delta \log d_{\text{eff}}$, σ satisfies some generic conditions and λ is small enough. Then,

$$\|f_{\star}-\widehat{\mathsf{RF}}_{n,N,\lambda}\|_{L^2}^2=\|\mathsf{P}_{>\ell}f_{\star}\|_{L^2}^2+o_{d,\mathbb{P}}(\cdot).$$

d_{eff} capture the 'effective low-dimensionality' of the data.

- d_{eff} decreases with snr_c (data more anisotropic) and kernel methods will perform better.
- On the contrary, we expect NNs to learn features w_i aligned with z₁ and to depend mildly on snr_c.

Test error of KRR in the spiked covariates model

- Given iid samples $(\{x_i, y_i\})_{i \in [n]}$ from the spiked covariates model.
- Effective dimension: $d_{eff} = d_s \lor (d/\operatorname{snr}_c)$.

Theorem (Ghorbani, Mei, Misiakiewicz, Montanari, 2020)

Assume $d_{\text{eff}}^{\ell+\delta} \leq \min(n, N) \leq d_{\text{eff}}^{\ell+1-\delta}$ and $|\log(n/N)| \geq \delta \log d_{\text{eff}}$, σ satisfies some generic conditions and λ is small enough. Then,

$$\|f_{\star} - \widehat{\mathsf{RF}}_{n,\mathbf{N},\lambda}\|_{L^2}^2 = \|\mathsf{P}_{>\ell}f_{\star}\|_{L^2}^2 + o_{d,\mathbb{P}}(\cdot).$$

d_{eff} capture the 'effective low-dimensionality' of the data.

- d_{eff} decreases with snr_c (data more anisotropic) and kernel methods will perform better.
- On the contrary, we expect NNs to learn features w_i aligned with z₁ and to depend mildly on snr_c.
Test error of KRR in the spiked covariates model

- Given iid samples $(\{x_i, y_i\})_{i \in [n]}$ from the spiked covariates model.
- Effective dimension: $d_{eff} = d_s \lor (d/\operatorname{snr}_c)$.

Theorem (Ghorbani, Mei, Misiakiewicz, Montanari, 2020)

Assume $d_{\text{eff}}^{\ell+\delta} \leq \min(n, N) \leq d_{\text{eff}}^{\ell+1-\delta}$ and $|\log(n/N)| \geq \delta \log d_{\text{eff}}$, σ satisfies some generic conditions and λ is small enough. Then,

$$\|f_{\star}-\widehat{\mathsf{RF}}_{n,\boldsymbol{N},\lambda}\|_{L^{2}}^{2}=\|\mathsf{P}_{>\ell}f_{\star}\|_{L^{2}}^{2}+o_{d,\mathbb{P}}(\cdot).$$

d_{eff} capture the 'effective low-dimensionality' of the data.

- ▶ d_{eff} decreases with snr_c (data more anisotropic) and kernel methods will perform better.
- On the contrary, we expect NNs to learn features w_i aligned with z₁ and to depend mildly on snr_c.

Test error of KRR in the spiked covariates model

- Given iid samples $(\{x_i, y_i\})_{i \in [n]}$ from the spiked covariates model.
- Effective dimension: $d_{eff} = d_s \lor (d/\operatorname{snr}_c)$.

Theorem (Ghorbani, Mei, Misiakiewicz, Montanari, 2020)

Assume $d_{\text{eff}}^{\ell+\delta} \leq \min(n, N) \leq d_{\text{eff}}^{\ell+1-\delta}$ and $|\log(n/N)| \geq \delta \log d_{\text{eff}}$, σ satisfies some generic conditions and λ is small enough. Then,

$$\|f_{\star}-\widehat{\mathsf{RF}}_{n,\mathbf{N},\lambda}\|_{L^2}^2=\|\mathsf{P}_{>\ell}f_{\star}\|_{L^2}^2+o_{d,\mathbb{P}}(\cdot).$$

d_{eff} capture the 'effective low-dimensionality' of the data.

- d_{eff} decreases with snr_c (data more anisotropic) and kernel methods will perform better.
- On the contrary, we expect NNs to learn features w_i aligned with z₁ and to depend mildly on snr_c.

Test error of KRR in the spiked covariates model

- Given iid samples $(\{x_i, y_i\})_{i \in [n]}$ from the spiked covariates model.
- Effective dimension: $d_{eff} = d_s \lor (d/\operatorname{snr}_c)$.

Theorem (Ghorbani, Mei, Misiakiewicz, Montanari, 2020)

Assume $d_{\text{eff}}^{\ell+\delta} \leq \min(n, N) \leq d_{\text{eff}}^{\ell+1-\delta}$ and $|\log(n/N)| \geq \delta \log d_{\text{eff}}$, σ satisfies some generic conditions and λ is small enough. Then,

$$\|f_{\star}-\widehat{\mathsf{RF}}_{n,\boldsymbol{N},\lambda}\|_{L^2}^2=\|\mathsf{P}_{>\ell}f_{\star}\|_{L^2}^2+o_{d,\mathbb{P}}(\cdot).$$

d_{eff} capture the 'effective low-dimensionality' of the data.

- d_{eff} decreases with snr_c (data more anisotropic) and kernel methods will perform better.
- On the contrary, we expect NNs to learn features w_i aligned with z₁ and to depend mildly on snr_c.

NNs vs kernel methods in the spiked covariates model

Effective dimension: $d_{eff} = d_s \vee (d/\operatorname{snr}_c)$.

In the spiked covariates model, we expect the generalization error of:

- Kernel methods to depend on d_{eff};
- NNs to depend on d_s.

• $d_s \leq d_{eff} \leq d$:

Test error: $NN \leq Kernel methods$

 \blacktriangleright decreases with snr_c:

lsotropic features: $d_{eff} = d$ (low covariate SNR)

Test error: $NN \ll Kernel methods.$

• Highly anisotropic features: $d_{eff} = d_s$ (high covariate SNR)

Fest error: $\mathsf{NN}\sim\mathsf{Kernel}$ methods.

NNs vs kernel methods in the spiked covariates model

Effective dimension: $d_{eff} = d_s \vee (d/\operatorname{snr}_c)$.

In the spiked covariates model, we expect the generalization error of:

- Kernel methods to depend on d_{eff};
- NNs to depend on d_s.

 $d_s \leq d_{\rm eff} \leq d:$

Test error: $NN \leq Kernel methods$

 \blacktriangleright d_{eff} decreases with snr_c:

lsotropic features: $d_{eff} = d$ (low covariate SNR)

Test error: $NN \ll Kernel methods$.

• Highly anisotropic features: $d_{eff} = d_s$ (high covariate SNR)

Fest error: $\mathsf{NN}\sim\mathsf{Kernel}$ methods.

NNs vs kernel methods in the spiked covariates model

Effective dimension: $d_{eff} = d_s \vee (d/\operatorname{snr}_c)$.

In the spiked covariates model, we expect the generalization error of:

- Kernel methods to depend on d_{eff};
- NNs to depend on d_s.

► $d_s \leq d_{eff} \leq d$: Test error: NN \leq Kernel methods

d_{eff} decreases with snr_c:

lsotropic features: $d_{eff} = d$ (low covariate SNR)

Test error: $NN \ll Kernel methods$.

• Highly anisotropic features: $d_{eff} = d_s$ (high covariate SNR)

Test error: $NN \sim Kernel methods.$

Goal: fitting $\varphi(x_1)$ using features $\sigma(\langle \mathbf{w}_i, \mathbf{x} \rangle)$.

Consider $\mathbb{E}[xx^{\mathsf{T}}] = \operatorname{diag}(\operatorname{snr}_{c}, \operatorname{Id}_{d-1}).$

▶ To fit $\varphi(x_1)$, we need to use features $\sigma(\langle w_i, x \rangle)$ that are correlated to x_1 , i.e.,

$$\operatorname{corr}(x_1 e_1, \langle w_i, x \rangle) = \operatorname{snr}_c \cdot \frac{\langle w_i, e_1 \rangle^2}{\|\Sigma^{1/2} w_i\|_2^2} = \operatorname{snr}_c \cdot O_{d, \mathbb{P}}(d^{-1}) = O_{d, \mathbb{P}}(d_{\operatorname{eff}}^{-1}).$$

For RF, when snr_c is large, all random w_i 's are good.

For NN, w_i 's can be chosen to have a large correlation with e_1 .

Goal: fitting $\varphi(x_1)$ using features $\sigma(\langle \mathbf{w}_i, \mathbf{x} \rangle)$.

Consider $\mathbb{E}[xx^{\mathsf{T}}] = \operatorname{diag}(\operatorname{snr}_{c}, \operatorname{Id}_{d-1}).$

▶ To fit $\varphi(x_1)$, we need to use features $\sigma(\langle \mathbf{w}_i, \mathbf{x} \rangle)$ that are correlated to x_1 , i.e.,

$$\operatorname{corr}(x_1 \boldsymbol{e}_1, \langle \boldsymbol{w}_i, \boldsymbol{x} \rangle) = \operatorname{snr}_c \cdot \frac{\langle \boldsymbol{w}_i, \boldsymbol{e}_1 \rangle^2}{\|\boldsymbol{\Sigma}^{1/2} \boldsymbol{w}_i\|_2^2} = \operatorname{snr}_c \cdot O_{d,\mathbb{P}}(d^{-1}) = O_{d,\mathbb{P}}(d_{\operatorname{eff}}^{-1}).$$

For RF, when snr_c is large, all random w_i 's are good.

For NN, w_i 's can be chosen to have a large correlation with e_1 .

Goal: fitting $\varphi(x_1)$ using features $\sigma(\langle \mathbf{w}_i, \mathbf{x} \rangle)$.

Consider $\mathbb{E}[xx^{\mathsf{T}}] = \operatorname{diag}(\operatorname{snr}_{c}, \operatorname{Id}_{d-1}).$

▶ To fit $\varphi(x_1)$, we need to use features $\sigma(\langle \mathbf{w}_i, \mathbf{x} \rangle)$ that are correlated to x_1 , i.e.,

$$\operatorname{corr}(x_1 \boldsymbol{e}_1, \langle \boldsymbol{w}_i, \boldsymbol{x} \rangle) = \operatorname{snr}_c \cdot \frac{\langle \boldsymbol{w}_i, \boldsymbol{e}_1 \rangle^2}{\|\boldsymbol{\Sigma}^{1/2} \boldsymbol{w}_i\|_2^2} = \operatorname{snr}_c \cdot O_{d,\mathbb{P}}(d^{-1}) = O_{d,\mathbb{P}}(d_{\operatorname{eff}}^{-1}).$$

For RF, when snr_c is large, all random w_i 's are good.

For NN, w_i 's can be chosen to have a large correlation with e_1 .

Goal: fitting $\varphi(x_1)$ using features $\sigma(\langle \mathbf{w}_i, \mathbf{x} \rangle)$.

Consider $\mathbb{E}[xx^{\mathsf{T}}] = \operatorname{diag}(\operatorname{snr}_{c}, \operatorname{Id}_{d-1}).$

▶ To fit $\varphi(x_1)$, we need to use features $\sigma(\langle \mathbf{w}_i, \mathbf{x} \rangle)$ that are correlated to x_1 , i.e.,

$$\operatorname{corr}(x_1 \boldsymbol{e}_1, \langle \boldsymbol{w}_i, \boldsymbol{x} \rangle) = \operatorname{snr}_c \cdot \frac{\langle \boldsymbol{w}_i, \boldsymbol{e}_1 \rangle^2}{\|\boldsymbol{\Sigma}^{1/2} \boldsymbol{w}_i\|_2^2} = \operatorname{snr}_c \cdot O_{d,\mathbb{P}}(d^{-1}) = O_{d,\mathbb{P}}(d_{\operatorname{eff}}^{-1}).$$

For RF, when snr_c is large, all random w_i 's are good.

For NN, w_i's can be chosen to have a large correlation with e₁.

Goal: fitting $\varphi(x_1)$ using features $\sigma(\langle \mathbf{w}_i, \mathbf{x} \rangle)$.

Consider $\mathbb{E}[xx^{\mathsf{T}}] = \operatorname{diag}(\operatorname{snr}_{c}, \operatorname{Id}_{d-1}).$

▶ To fit $\varphi(x_1)$, we need to use features $\sigma(\langle \mathbf{w}_i, \mathbf{x} \rangle)$ that are correlated to x_1 , i.e.,

$$\operatorname{corr}(x_1 \boldsymbol{e}_1, \langle \boldsymbol{w}_i, \boldsymbol{x} \rangle) = \operatorname{snr}_c \cdot \frac{\langle \boldsymbol{w}_i, \boldsymbol{e}_1 \rangle^2}{\|\boldsymbol{\Sigma}^{1/2} \boldsymbol{w}_i\|_2^2} = \operatorname{snr}_c \cdot O_{d,\mathbb{P}}(d^{-1}) = O_{d,\mathbb{P}}(d_{\operatorname{eff}}^{-1}).$$

For RF, when snr_c is large, all random w_i 's are good.

▶ For NN, *w_i*'s can be chosen to have a large correlation with *e*₁.

Insight

Insight

Lower covariate SNR (data more isotropic) should lead to larger generalization gap between NNs and kernel methods.

How to test this insight?

In *image classification*, we expect

- Images to have most of their spectrum concentrated on low-frequencies (z1 part);
- The labels to depend predominantly on the low-frequencies $(y = f_*(x) = \varphi(z_1))$.

Experiment: add noise to the high frequency components (z_2 part) of the images (i.e., decreasing the snr_c).

Adding noise in covariates should increase the generalization gap between NN and kernel methods.

Insight

Insight

Lower covariate SNR (data more isotropic) should lead to larger generalization gap between NNs and kernel methods.

How to test this insight?

In image classification, we expect

- Images to have most of their spectrum concentrated on low-frequencies (z₁ part);
- ► The labels to depend predominantly on the low-frequencies $(y = f_*(x) = \varphi(z_1))$.

Experiment: add noise to the high frequency components (z_2 part) of the images (i.e., decreasing the snr_c).

Adding noise in covariates should increase the generalization gap between NN and kernel methods.

Insight

Insight

Lower covariate SNR (data more isotropic) should lead to larger generalization gap between NNs and kernel methods.

How to test this insight?

In *image classification*, we expect

- Images to have most of their spectrum concentrated on low-frequencies (z₁ part);
- ► The labels to depend predominantly on the low-frequencies $(y = f_*(x) = \varphi(z_1))$.

Experiment: add noise to the high frequency components (z_2 part) of the images (i.e., decreasing the snr_c).

Adding noise in covariates should increase the generalization gap between NN and kernel methods.

Numerical simulations



Figure: Test accuracy on FMNIST: adding noise to the high frequency components.

Spiked covariates model:

A controlling parameter of the performance gap between NN and kernel methods is

 $\operatorname{snr}_{c} = \operatorname{Covariate} \operatorname{SNR} = \frac{\operatorname{Signal covariates variance}}{\operatorname{Noise covariates variance}}.$

- Small snr_c: large separation.
- Large snr_c: kernel methods perform closer to NN.

Intuition:

- ▶ When snr_c is small, RF fail to find the signal covariates.
- ▶ When snr_c is big, RF automatically find the signal covariates.
- NNs always look for the signal covariates.

C. Invariant function estimation on the sphere

In many learning tasks, the target function present some natural symmetries.
 E.g., image classification: labels invariant under translation of the images.

Effectiveness of NNs architectures is often loosely attributed to their ability to encode various symmetries present in the data.

E.g., CNNs with translation invariance.

CNNs perform better than fully-connected networks.

Convolutional kernels perform better than inner-product kernels.

Can we quantify this gain?

- In many learning tasks, the target function present some natural symmetries. E.g., image classification: labels invariant under translation of the images.
- Effectiveness of NNs architectures is often loosely attributed to their ability to encode various symmetries present in the data.

E.g., CNNs with translation invariance.

- CNNs perform better than fully-connected networks.
 Convolutional kernels perform better than inner-product kerne
- Can we quantify this gain?

- In many learning tasks, the target function present some natural symmetries. E.g., image classification: labels invariant under translation of the images.
- Effectiveness of NNs architectures is often loosely attributed to their ability to encode various symmetries present in the data.

E.g., CNNs with translation invariance.

CNNs perform better than fully-connected networks.

Convolutional kernels perform better than inner-product kernels.

Can we quantify this gain?

- In many learning tasks, the target function present some natural symmetries. E.g., image classification: labels invariant under translation of the images.
- Effectiveness of NNs architectures is often loosely attributed to their ability to encode various symmetries present in the data.

E.g., CNNs with translation invariance.

CNNs perform better than fully-connected networks.

Convolutional kernels perform better than inner-product kernels.

Can we quantify this gain?

• Data
$$x \sim \text{Unif}(\mathbb{S}^{d-1}(\sqrt{d})).$$

• Consider $\mathcal{G}_d = \{g_0, g_1, \dots, g_{d-1}\}$ the cyclic group:

$$g_i \cdot \mathbf{x} = (x_{d-i+1}, x_{d-i+2}, \dots, x_d, x_1, x_2, \dots, x_{d-i}).$$

Goal: learn a cyclic invaeriant function f_{*}:
 i.e., f_{*}(g ⋅ x) = f_{*}(x) for all g ∈ G
 E.g., f(x) = ∑^d_{i=1} x_ix_{i+1}.

• Data
$$x \sim \text{Unif}(\mathbb{S}^{d-1}(\sqrt{d})).$$

• Consider
$$\mathcal{G}_d = \{g_0, g_1, \dots, g_{d-1}\}$$
 the cyclic group:

$$g_i \cdot \mathbf{x} = (x_{d-i+1}, x_{d-i+2}, \dots, x_d, x_1, x_2, \dots, x_{d-i}).$$

Goal: learn a cyclic invaeriant function f_{*}:
 i.e., f_{*}(g ⋅ x) = f_{*}(x) for all g ∈ g
 E.g., f(x) = ∑^d_{i=1} x_ix_{i+1}.

• Data
$$x \sim \text{Unif}(\mathbb{S}^{d-1}(\sqrt{d})).$$

• Consider
$$\mathcal{G}_d = \{g_0, g_1, \dots, g_{d-1}\}$$
 the cyclic group:

$$g_i \cdot \mathbf{x} = (x_{d-i+1}, x_{d-i+2}, \dots, x_d, x_1, x_2, \dots, x_{d-i}).$$

Goal: learn a cyclic invaeriant function f_{*}:
 i.e., f_{*}(g ⋅ x) = f_{*}(x) for all g ∈ G_d.
 E.g., f(x) = ∑^d_{i=1} x_ix_{i+1}.

• Data
$$x \sim \text{Unif}(\mathbb{S}^{d-1}(\sqrt{d})).$$

• Consider
$$\mathcal{G}_d = \{g_0, g_1, \dots, g_{d-1}\}$$
 the cyclic group:

$$g_i \cdot \mathbf{x} = (x_{d-i+1}, x_{d-i+2}, \dots, x_d, x_1, x_2, \dots, x_{d-i}).$$

Goal: learn a cyclic invaeriant function f_{*}:
 i.e., f_{*}(g ⋅ x) = f_{*}(x) for all g ∈ G_d.
 E.g., f(x) = ∑^d_{i-1} x_ix_{i+1}.

Invariant random features and kernels

Invariant RF model:

$$\mathsf{RF}_N^{\mathrm{inv}}(\mathbf{x};\mathbf{a},\mathbf{W}) = \sum_{i=1}^N a_i \int_{g \in \mathcal{G}_d} \sigma(\langle \mathbf{w}_i, g \cdot \mathbf{x} \rangle) \pi(\mathrm{d}g) = \frac{1}{d} \sum_{i=1}^N a_i \sum_{k=1}^d \sigma(\langle \mathbf{w}_i, g_k \cdot \mathbf{x} \rangle).$$

Neural tangent model of a CNN with N filters $w_i \in \mathbb{R}^d$. Associated invariant kernel:

$$\mathcal{K}_{\mathrm{inv}}(\mathbf{x},\mathbf{y}) = \int_{g\in\mathcal{G}_d} h(\langle g\cdot\mathbf{x},\mathbf{y}
angle/d)\pi(\mathrm{d}g).$$

Compared to standard RF model:

$$\mathsf{RF}_N(\mathbf{x}; \mathbf{a}, \mathbf{W}) = \sum_{i=1}^N a_i \sigma(\langle \mathbf{x}, \mathbf{w}_i \rangle),$$

and associated 'standard' (inner-product) kernel:

$$K(\mathbf{x},\mathbf{y})=h(\langle \mathbf{x},\mathbf{y}\rangle/d).$$

Invariant random features and kernels

Invariant RF model:

$$\mathsf{RF}_N^{\mathrm{inv}}(\mathbf{x};\mathbf{a},\mathbf{W}) = \sum_{i=1}^N a_i \int_{g \in \mathcal{G}_d} \sigma(\langle \mathbf{w}_i, g \cdot \mathbf{x} \rangle) \pi(\mathrm{d}g) = \frac{1}{d} \sum_{i=1}^N a_i \sum_{k=1}^d \sigma(\langle \mathbf{w}_i, g_k \cdot \mathbf{x} \rangle).$$

Neural tangent model of a CNN with N filters $w_i \in \mathbb{R}^d$. Associated invariant kernel:

$$\mathcal{K}_{\mathrm{inv}}(\boldsymbol{x}, \boldsymbol{y}) = \int_{g \in \mathcal{G}_d} h(\langle g \cdot \boldsymbol{x}, \boldsymbol{y} \rangle / d) \pi(\mathrm{d}g).$$

Compared to standard RF model:

$$\mathsf{RF}_{N}(\boldsymbol{x};\boldsymbol{a},\boldsymbol{W}) = \sum_{i=1}^{N} \boldsymbol{a}_{i} \sigma(\langle \boldsymbol{x},\boldsymbol{w}_{i} \rangle),$$

and associated 'standard' (inner-product) kernel:

$$K(\mathbf{x},\mathbf{y}) = h(\langle \mathbf{x},\mathbf{y}\rangle/d).$$

Test error of KRR with invariant kernel

• Cyclic invariant f_* : given iid samples $\{(y_i, x_i)\}_{i \in [n]}$,

 $y_i = f_{\star}(\mathbf{x}_i) + \varepsilon_i, \qquad \mathbf{x}_i \sim_{iid} \text{Unif}(\mathbb{S}^{d-1}(\sqrt{d})), \quad \varepsilon_i \sim_{iid} N(0, \tau^2).$

Theorem (Mei, **Misiakiewicz**, Montanari, 2021)

Assume $d^{\ell-1+\delta} \leq \min(n, N) \leq d^{\ell-\delta}$ and $|\log(n/N)| \geq \delta \log d$, σ satisfies some generic conditions and λ is small enough. Then,

$$\|f_{\star} - \widehat{\mathsf{RF}}_{n,\boldsymbol{N},\lambda}^{\text{inv}}\|_{L^2}^2 = \|\mathsf{P}_{>\ell}f_{\star}\|_{L^2}^2 + o_{d,\mathbb{P}}(\cdot).$$

Standard RF model:

For $d^{\ell+\delta} \leq \min(n, \mathbb{N}) \leq d^{\ell+1-\delta}$, $\|f_{\star} - \widehat{\mathsf{RF}}_{n,\mathbb{N},\lambda}\|_{L^2}^2 = \|\mathsf{P}_{>\ell}f_{\star}\|_{L^2}^2 + o_d(\cdot)$.

Gain a factor *d* in features and sample complexity by using invariant RF.

• More generally, for \mathcal{G}_d with 'degeneracy α' , we gain a factor d^{α} .

Neural Networks and Kernel Methods

Test error of KRR with invariant kernel

• Cyclic invariant f_* : given iid samples $\{(y_i, x_i)\}_{i \in [n]}$,

 $y_i = f_*(\mathbf{x}_i) + \varepsilon_i, \qquad \mathbf{x}_i \sim_{iid} \text{Unif}(\mathbb{S}^{d-1}(\sqrt{d})), \quad \varepsilon_i \sim_{iid} N(0, \tau^2).$

Theorem (Mei, **Misiakiewicz**, Montanari, 2021)

Assume $d^{\ell-1+\delta} \leq \min(n, N) \leq d^{\ell-\delta}$ and $|\log(n/N)| \geq \delta \log d$, σ satisfies some generic conditions and λ is small enough. Then,

$$\|f_{\star} - \widehat{\mathsf{RF}}_{n,N,\lambda}^{\text{inv}}\|_{L^2}^2 = \|\mathsf{P}_{>\ell}f_{\star}\|_{L^2}^2 + o_{d,\mathbb{P}}(\cdot).$$

Standard RF model:

For $d^{\ell+\delta} \leq \min(n, N) \leq d^{\ell+1-\delta}$, $\|f_{\star} - \widehat{\mathsf{RF}}_{n,N,\lambda}\|_{L^2}^2 = \|\mathsf{P}_{>\ell}f_{\star}\|_{L^2}^2 + o_d(\cdot)$.

Gain a factor d in features and sample complexity by using invariant RF.

• More generally, for \mathcal{G}_d with 'degeneracy α ', we gain a factor d^{α} .

Neural Networks and Kernel Methods

Test error of KRR with invariant kernel

• Cyclic invariant f_* : given iid samples $\{(y_i, x_i)\}_{i \in [n]}$,

$$y_i = f_{\star}(\mathbf{x}_i) + \varepsilon_i, \qquad \mathbf{x}_i \sim_{iid} \text{Unif}(\mathbb{S}^{d-1}(\sqrt{d})), \quad \varepsilon_i \sim_{iid} N(0, \tau^2).$$

Theorem (Mei, **Misiakiewicz**, Montanari, 2021)

Assume $d^{\ell-1+\delta} \leq \min(n, N) \leq d^{\ell-\delta}$ and $|\log(n/N)| \geq \delta \log d$, σ satisfies some generic conditions and λ is small enough. Then,

$$\|f_{\star} - \widehat{\mathsf{RF}}_{n,N,\lambda}^{\text{inv}}\|_{L^2}^2 = \|\mathsf{P}_{>\ell}f_{\star}\|_{L^2}^2 + o_{d,\mathbb{P}}(\cdot).$$

Standard RF model:

For
$$d^{\ell+\delta} \leq \min(n, N) \leq d^{\ell+1-\delta}$$
, $\|f_\star - \widehat{\mathsf{RF}}_{n,N,\lambda}\|_{L^2}^2 = \|\mathsf{P}_{>\ell}f_\star\|_{L^2}^2 + o_d(\cdot)$.

Gain a factor d in features and sample complexity by using invariant RF.

• More generally, for \mathcal{G}_d with 'degeneracy α ', we gain a factor d^{α} .

Numerical simulations



Figure: Test error of KRR with invariant kernel and inner-product kernel.

Cyclic invariant MNIST



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

Theodor Misiakiewicz (Stanford)

Neural Networks and Kernel Methods

What I illustrated in this talk:

Isotropic covariates model: Kernel methods suffer from the *curse of dimensionality* and perform inferior to NNs due to lack of feature learning.

Spiked covariates model: For large covariate SNR, feature learning is unecessary and kernel methods can also perform well.

Invariant function estimation: Invariant kernels outperform standard kernels and we quantify the gain in statistical efficiency.

What I illustrated in this talk:

Isotropic covariates model:

Kernel methods suffer from the *curse of dimensionality* and perform inferior to NNs due to **lack of feature learning**.

Spiked covariates model: For large covariate SNR, feature learning is unecessary and kernel methods can also perform well.

Invariant function estimation: Invariant kernels outperform standard kernels and we quantify the gain in statistical efficiency.

What I illustrated in this talk:

Isotropic covariates model:

Kernel methods suffer from the *curse of dimensionality* and perform inferior to NNs due to **lack of feature learning**.

Spiked covariates model: For large covariate SNR, feature learning is unecessary and kernel methods can also perform well.

Invariant function estimation: Invariant kernels outperform standard kernels and we quantify the gain in statistical efficiency.

What I illustrated in this talk:

Isotropic covariates model:

Kernel methods suffer from the *curse of dimensionality* and perform inferior to NNs due to **lack of feature learning**.

Spiked covariates model: For large covariate SNR, feature learning is unecessary and kernel methods can also perform well.

Invariant function estimation: Invariant kernels outperform standard kernels and we quantify the gain in statistical efficiency.
Summary

What I illustrated in this talk:

Isotropic covariates model:

Kernel methods suffer from the *curse of dimensionality* and perform inferior to NNs due to **lack of feature learning**.

Spiked covariates model: For large covariate SNR, feature learning is unecessary and kernel methods can also perform well.

Invariant function estimation:

Invariant kernels outperform standard kernels and we quantify the gain in statistical efficiency.

All the results in this talk can be derived from a general framework for getting high-dimensional asymptotics of test error of random features and kernel methods. [Mei, Misiakiewicz, Montanari, '21a].

Open problems

- Show that NNs trained by gradient descent overcome the curse of dimensionality in the spiked covariates model (using mean-field dynamics, see [Chizat, Bach, '20]).
- What are other latent low-dimensional structure that we can study? (other than the spiked covariates model)
- Using our general framework to study more complicated kernels.
- Data dependent kernel methods?
- Practical advantages of using kernel methods vs NNs? Transfer learning, robustness...

Thank you!