

# 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 Tangent (NT) model

- Neural network:  $\text{NN}_p(x; \theta)$ ,  $x \in \mathbb{R}^d$ ,  $\theta \in \mathbb{R}^p$ .  
e.g., fully-connected neural network:

$$\text{NN}_p(x; \theta) = \text{h}(\theta; \text{vec}(W_2 (W_1 x)))$$

- Linearization around random initialization  $\theta_0$ :

$$\text{NN}_p(x; \theta) = \text{NN}_p(x; \theta_0) + \text{h}(\theta_0; \text{vec}(\text{NN}_p(x; \theta_0))) + o(\|\theta - \theta_0\|_2)$$

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

$$\text{NT}_p(x; y; \theta_0) = \text{h}(\theta_0; \text{vec}(\text{NN}_p(x; \theta_0)))$$

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

$$K_{\text{NT}}(x; y) = \mathbb{E}_{\theta_0}[\text{h}(\theta_0; \text{vec}(\text{NN}_p(x; \theta_0))) \text{h}(\theta_0; \text{vec}(\text{NN}_p(y; \theta_0)))]$$

## Neural Tangent (NT) model

- Neural network:  $\text{NN}_p(x; \theta)$ ,  $x \in \mathbb{R}^d$ ,  $\theta \in \mathbb{R}^p$ .  
e.g., fully-connected neural network:

$$\text{NN}_p(x; \theta) = \text{h} \text{a}; \left( \text{::: } W_2 \left( W_1 x \right) \right) i:$$

- Linearization around random initialization  $\theta_0$ :

$$\text{NN}_p(x; \theta) = \text{NN}_p(x; \theta_0) + \text{h} \text{r} \text{NN}_p(x; \theta_0) i + o(k \text{ } \theta_0 k_2):$$

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

$$\text{NT}_p(x; y; \theta_0) = \text{h} \text{r} \text{NN}_p(x; \theta_0) i:$$

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

$$K_{\text{NT}}(x; y) = \mathbb{E}_{\theta_0} [\text{h} \text{r} \text{NN}_p(x; \theta_0) \text{r} \text{NN}_p(y; \theta_0) i]:$$

## Neural Tangent (NT) model

- Neural network:  $\text{NN}_p(x; \theta)$ ,  $x \in \mathbb{R}^d$ ,  $\theta \in \mathbb{R}^p$ .  
e.g., fully-connected neural network:

$$\text{NN}_p(x; \theta) = h a; \left( \text{::: } W_2 (W_1 x) \right) i;$$

- Linearization around random initialization  $\theta_0$ :

$$\text{NN}_p(x; \theta) = \text{NN}_p(x; \theta_0) + \frac{h}{\sqrt{r}} \text{NN}_p(x; \theta_0) i + o(k \theta_0 k_2):$$

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

$$\text{NT}_p(x; y; \theta_0) = h ; r \text{NN}_p(x; \theta_0) i;$$

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

$$K_{\text{NT}}(x; y) = \mathbb{E}_{\theta_0} [h r \text{NN}_p(x; \theta_0) ; r \text{NN}_p(y; \theta_0) i];$$

## Neural Tangent (NT) model

- Neural network:  $\text{NN}_p(x; \theta)$ ,  $x \in \mathbb{R}^d$ ,  $\theta \in \mathbb{R}^p$ .  
e.g., fully-connected neural network:

$$\text{NN}_p(x; \theta) = \text{h}(\theta; \text{W}_2(\text{W}_1 x))$$

- Linearization around random initialization  $\theta_0$ :

$$\text{NN}_p(x; \theta) = \text{NN}_p(x; \theta_0) + \text{h}(\theta_0; r \text{NN}_p(x; \theta_0)) + o(k_0 k_2)$$

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

$$\text{NT}_p(x; y; \theta_0) = \text{h}(\theta_0; r \text{NN}_p(x; \theta_0))$$

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

$$K_{\text{NT}}(x; y) = \mathbb{E}_{\theta_0}[\text{h}(\theta_0; r \text{NN}_p(x; \theta_0)) \text{h}(\theta_0; r \text{NN}_p(y; \theta_0))]$$

Coupled gradient descent on empirical squared loss:

$$\begin{aligned} \frac{d}{dt} t &= -r \hat{E}[(y - \text{NN}_\rho(x; t))^2]; & t^0 &= t_0; \\ \frac{d}{dt} t &= -r \hat{E}[(y - \text{NT}_\rho(x; t; \mathbf{o}))^2]; & t^0 &= \mathbf{0}; \end{aligned}$$

### Theorem (informal)

For any  $\epsilon > 0$ , for large enough NNs and proper random initialization  $\mathbf{o}$  (Xavier initialization), we have with high probability

$$\sup_{t \geq 0} \sup_{\|x\|_2=1} |\text{NN}_\rho(x; t) - \text{NT}_\rho(x; t; \mathbf{o})| \leq \epsilon;$$

$\Rightarrow$  Intuition: *weights do not change much and  $\|t - \mathbf{o}\|_2$  remains small.*

[Jacot, Gabriel, Hongler, '18], [Du, Zhai, Póczos, 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], ...



Coupled gradient descent on empirical squared loss:

$$\begin{aligned} \frac{d}{dt} t &= -r \hat{E}[(y - \text{NN}_\rho(x; t))^2]; & t^0 &= t_0; \\ \frac{d}{dt} t &= -r \hat{E}[(y - \text{NT}_\rho(x; t; \theta))^2]; & t^0 &= 0; \end{aligned}$$

### Theorem (informal)

For any  $\epsilon > 0$ , for large enough NNs and **proper random initialization**  $\theta_0$  (Xavier initialization), we have with high probability

$$\sup_{t \geq 0} \sup_{\|x\|_2=1} |\text{NN}_\rho(x; t) - \text{NT}_\rho(x; t; \theta_0)| \leq \epsilon;$$

=> Intuition: *weights do not change much and  $\|t - \theta_0\|_2$  remains small.*

[Jacot, Gabriel, Hongler, '18], [Du, Zhai, Póczos, 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.

## 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.

## 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.

### 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.



### 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

## Setting

- Given  $n$  iid samples  $f(y_i; x_i)g_{i2[n]}$ :

$$y_i = f_*(x_i) + \epsilon_i; \quad x_i \text{ iid } \text{Unif}(S^{d-1}(\rho/d)); \quad \epsilon_i \text{ iid } N(0; \sigma^2):$$

- Two-layers neural networks:  $\text{NN}_N(x; \cdot) = \sum_{i=1}^N a_i (hw_i; x_i)$ .

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

$$\text{RF}_N(x; a; W) = \sum_{i=1}^N a_i (hw_i; x_i) = ha; \text{r}_a \text{NN}_p(x; \mathbf{0})i:$$

- Ridge regression with RF model (with fixed  $W = (w_i)_{i2[N]} \text{ iid } \text{Unif}(S^{d-1})$ )

$$\hat{a}(\cdot) = \arg \min_a \left\{ \frac{1}{n} \sum_{i=1}^n \left( y_i - \sum_{j=1}^N a_j (hw_j; x_i) \right)^2 + \lambda \|a\|_2^2 \right\}:$$

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

## Setting

- Given  $n$  iid samples  $f(y_i; x_i)_{i \in [n]}$ :

$$y_i = f_*(x_i) + \epsilon_i; \quad x_i \text{ iid } \text{Unif}(S^{d-1}(\rho^{-d})); \quad \epsilon_i \text{ iid } N(0; \sigma^2):$$

- Two-layers neural networks:  $\text{NN}_N(x; \cdot) = \sum_{i=1}^N a_i (h w_i; x_i)$ .

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

$$\text{RF}_N(x; a; W) = \sum_{i=1}^N a_i (h w_i; x_i) = h a; r_a \text{NN}_p(x; \mathbf{0})_i:$$

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

$$\hat{a}(\cdot) = \arg \min_a \left\{ \frac{1}{n} \sum_{i=1}^n \left( y_i - \sum_{j=1}^N a_j (h w_j; x_i) \right)^2 + k a k_2^2 \right\}:$$

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

## Setting

- Given  $n$  iid samples  $f(y_i; x_i)g_{i2[n]}$ :

$$y_i = f_*(x_i) + \epsilon_i; \quad x_i \text{ iid } \text{Unif}(S^{d-1}(\rho-d)); \quad \epsilon_i \text{ iid } N(0; \sigma^2):$$

- Two-layers neural networks:  $\text{NN}_N(x; \mathbf{a}) = \sum_{i=1}^N a_i (hW_i; x_i)$ .

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

$$\text{RF}_N(x; \mathbf{a}; W) = \sum_{i=1}^N a_i (hW_i; x_i) = \mathbf{h}\mathbf{a}; r_a \text{NN}_p(x; \mathbf{0})_i:$$

- Ridge regression with RF model (with fixed  $W = (w_i)_{i2[N]} \text{ iid } \text{Unif}(S^{d-1})$ )

$$\hat{\mathbf{a}}(\cdot) = \arg \min_{\mathbf{a}} \left\{ \frac{1}{n} \sum_{i=1}^n \left( y_i - \sum_{j=1}^N a_j (hW_j; x_{i,j}) \right)^2 + k \lambda k_2^2 \right\}:$$

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

## Setting

- Given  $n$  iid samples  $f(y_i; x_i)g_{i2[n]}$ :

$$y_i = f_*(x_i) + \epsilon_i; \quad x_i \text{ iid } \text{Unif}(S^{d-1}(\rho/d)); \quad \epsilon_i \text{ iid } N(0; \sigma^2):$$

- Two-layers neural networks:  $\text{NN}_N(x; \mathbf{a}) = \sum_{i=1}^N a_i (hW_i; x_i)$ .

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

$$\text{RF}_N(x; \mathbf{a}; W) = \sum_{i=1}^N a_i (hW_i; x_i) = \mathbf{h}\mathbf{a}; \mathbf{r}_a \text{NN}_p(x; \mathbf{0})_i:$$

- Ridge regression with RF model (with fixed  $W = (W_i)_{i2[N]} \text{ iid } \text{Unif}(S^{d-1})$ )

$$\hat{\mathbf{a}}(\cdot) = \arg \min_{\mathbf{a}} \left\{ \frac{1}{n} \sum_{i=1}^n \left( y_i - \sum_{j=1}^N a_j (hW_j; x_{ij}) \right)^2 + \lambda \|\mathbf{a}\|_2^2 \right\}:$$

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

## Test error of RF model

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

Assume  $d^{\ell+\delta} \min(n; N) \leq d^{\ell+1} \delta$  and  $j \log(n=N)j \leq \log d$ , satisfies some generic conditions and  $\delta$  is small enough. Then,

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

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

| For  $d^{\ell}$  parameters and samples, RF fits the best degree- $\ell$  polynomial approximation.

|  $n$  and  $N$  play a symmetric role:

$$\text{Test error}(n; N) = \max \left\{ \text{approx. error}(n = 1; N); \text{statistical error}(n; N = 1) \right\}:$$

| Statistical error  $(n; N = 1)$ : kernel method with an inner-product kernel  $K(x_1; x_2) = h(\langle x_1; x_2 \rangle / d)$  (e.g., NTK of fully-connected NNs of any depth).

## Test error of RF model

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

Assume  $d^{\ell+\delta} \min(n; N) \leq d^{\ell+1} \delta$  and  $j \log(n=N)j \leq \log d$ , satisfies some generic conditions and  $\delta$  is small enough. Then,

$$\|kf_{\star} - \widehat{RF}_{n,N,\lambda}\|_{L^2}^2 = \|P_{>\ell} f_{\star}\|_{L^2}^2 + o_{d,P}(\delta)$$

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

| For  $d^{\ell}$  parameters and samples, RF fits the best degree- $\ell$  polynomial approximation.

|  $n$  and  $N$  play a symmetric role:

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

| Statistical error  $(n; N)$ : kernel method with an inner-product kernel  $K(x_1; x_2) = h(\langle x_1; x_2 \rangle; d)$  (e.g., NTK of fully-connected NNs of any depth).



## Test error of RF model

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

Assume  $d^{\ell+\delta} \min(n; N) \leq d^{\ell+1} \delta$  and  $j \log(n=N)j \leq \log d$ , satisfies some generic conditions and  $\delta$  is small enough. Then,

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

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

| For  $d^\ell$  parameters and samples, RF fits the best degree- $\ell$  polynomial approximation.

|  $n$  and  $N$  play a symmetric role:

$$\text{Test error}(n; N) = \max \left\{ \text{approx. error}(n=1; N); \text{statistical error}(n; N=1) \right\}$$

| Statistical error  $(n; N=1)$ : kernel method with an inner-product kernel  $K(x_1; x_2) = h(\langle x_1; x_2 \rangle / d)$  (e.g., NTK of fully-connected NNs of any depth).

## Test error of RF model

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

Assume  $d^{\ell+\delta} \min(n; N) \leq d^{\ell+1} \delta$  and  $j \log(n=N)j \leq \log d$ , satisfies some generic conditions and  $\delta$  is small enough. Then,

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

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

| For  $d^{\ell}$  parameters and samples, RF fits the best degree- $\ell$  polynomial approximation.

|  $n$  and  $N$  play a symmetric role:

$$\text{Test error}(n; N) = \max \left\{ \text{approx. error}(n=1; N); \text{statistical error}(n; N=1) \right\}:$$

| Statistical error  $(n; N=1)$ : kernel method with an inner-product kernel  $K(x_1; x_2) = h(hx_1; x_2)_{i=d}$  (e.g., NTK of fully-connected NNs of any depth).

## The truncated staircase decay

$$f = P_0 f + P_1 f + P_2 f + P_3 f + P_4 f :$$

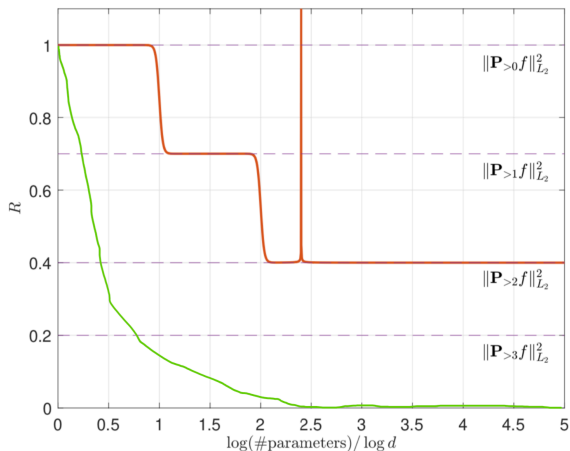


Figure: Test error (in red) v.s.  $\log(N)=\log(d)$  for  $n = d^{2.4}$  and  $\epsilon = 0^+$ . Train error in green.

## Performance gap between NN and kernel methods

- | We expect NNs to perform way better on 'low complexity' functions

E.g.,  $f_*(x) = (hw_*; xi)$ : if  $N \nearrow 1$  and  $n \nearrow d$ , then GD solution:

$$kf_* \quad \widehat{NN}_{N,n} k_{L^2}^2 = o_{d,P}() \quad [\text{Bai et al., '16}]$$

while for KRR:

$$kf_* \quad \widehat{KR}_{N,n} k_{L^2}^2 = kP_{>1} f_* k_{L^2}^2 + o_{d,P}():$$

- | 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  $f(x_1)$  using features  $(hw_i; xi)$

- | For RF,  $w_i$ 's are random,  $\text{corr}(x_1 e_1; hw_i; xi)^2 = hw_i; e_1 i^2 = kw_i k_2^2$  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.,  $f_*(x) = (hw_*; xi)$ : if  $N \nearrow 1$  and  $n \nearrow d$ , then GD solution:

$$kf_* \quad \widehat{NN}_{N,n} k_{L^2}^2 = o_{d,P}() \quad [\text{Bai et al., '16}]$$

while for KRR:

$$kf_* \quad \widehat{KR}_{N,n} k_{L^2}^2 = kP_{>1} f_* k_{L^2}^2 + o_{d,P}():$$

- | 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  $f_*(x_1)$  using features  $(hw_i; xi)$

- | For RF,  $w_i$ 's are random,  $\text{corr}(x_1 e_1; hw_i; xi)^2 = hw_i; e_1 i^2 = kw_i k_2^2$  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.,  $f_*(x) = (hw_*; xi)$ : if  $N \nearrow 1$  and  $n \nearrow d$ , then GD solution:

$$kf_* \quad \widehat{NN}_{N,n} k_{L^2}^2 = o_{d,P}() \quad [\text{Bai et al., '16}]$$

while for KRR:

$$kf_* \quad \widehat{KR}_{N,n} k_{L^2}^2 = kP_{>1} f_* k_{L^2}^2 + o_{d,P}():$$

- | 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  $f(x_1)$  using features  $(hw_i; xi)$

- | For RF,  $w_i$ 's are random,  $\text{corr}(x_1 e_1; hw_i; xi)^2 = hw_i; e_1 i^2 = kw_i k_2^2$  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.

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  $[U; U^\perp]$

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

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

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

$d_s$  = signal dimension.

$\text{snr}_c$  = covariate SNR.

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

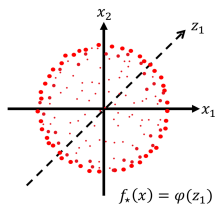


Figure: Isotropic covariates:  
 $\text{snr}_c = 1$ .

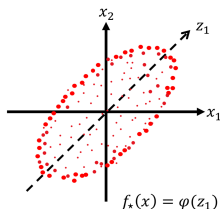


Figure: Anisotropic covariates:  
 $\text{snr}_c > 1$ .

## Test error of KRR in the spiked covariates model

- | Given iid samples  $(fX_i; y_i g)_{i \in [n]}$  from the spiked covariates model.
- | Effective dimension:  $d_{\text{eff}} = d_s \_ (d = \text{snr}_c)$ .

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

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

$$k f_{\star} \widehat{\text{RF}}_{n, N, \lambda} k_{L^2}^2 = k P_{> \ell} f_{\star} k_{L^2}^2 + o_{d, P}(\cdot)$$

- |  $d_{\text{eff}}$  capture the 'effective low-dimensionality' of the data.
- |  $d_{\text{eff}}$  decreases with  $\text{snr}_c$  (data more anisotropic) and kernel methods will perform better.
- | On the contrary, we expect NNs to learn features  $w_i$  aligned with  $z_1$  and to depend mildly on  $\text{snr}_c$ .

## Test error of KRR in the spiked covariates model

- | Given iid samples  $(fX_i; y_i g)_{i \in [n]}$  from the spiked covariates model.
- | Effective dimension:  $d_{\text{eff}} = d_s - (d = \text{snr}_c)$ .

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

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

$$k f_* \widehat{\text{RF}}_{n, N, \lambda} k_{L^2}^2 = k P_{> \ell} f_* k_{L^2}^2 + o_{d, P}(\cdot)$$

- |  $d_{\text{eff}}$  capture the 'effective low-dimensionality' of the data.
- |  $d_{\text{eff}}$  decreases with  $\text{snr}_c$  (data more anisotropic) and kernel methods will perform better.
- | On the contrary, we expect NNs to learn features  $w_i$  aligned with  $Z_1$  and to depend mildly on  $\text{snr}_c$ .

## Test error of KRR in the spiked covariates model

- | Given iid samples  $(fX_i; y_i g)_{i \in [n]}$  from the spiked covariates model.
- | Effective dimension:  $d_{\text{eff}} = d_s \_ (d = \text{snr}_c)$ .

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

Assume  $d_{\text{eff}}^{\ell + \delta} \min(n; N) d_{\text{eff}}^{\ell + 1} \delta$  and  $j \log(n = N) j \log d_{\text{eff}}$ , satisfies some generic conditions and is small enough. Then,

$$k f_{\star} \widehat{\text{RF}}_{n, N, \lambda} k_{L^2}^2 = k P_{> \ell} f_{\star} k_{L^2}^2 + o_{d, P}(\cdot):$$

- |  $d_{\text{eff}}$  capture the 'effective low-dimensionality' of the data.
- |  $d_{\text{eff}}$  decreases with  $\text{snr}_c$  (data more anisotropic) and kernel methods will perform better.
- | On the contrary, we expect NNs to learn features  $w_i$  aligned with  $Z_1$  and to depend mildly on  $\text{snr}_c$ .

## Test error of KRR in the spiked covariates model

- | Given iid samples  $(fX_i; y_i g)_{i \in [n]}$  from the spiked covariates model.
- | Effective dimension:  $d_{\text{eff}} = d_s - (d = \text{snr}_c)$ .

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

Assume  $d_{\text{eff}}^{\ell+\delta} \min(n; N) d_{\text{eff}}^{\ell+1} \delta$  and  $j \log(n=N) j \log d_{\text{eff}}$ , satisfies some generic conditions and is small enough. Then,

$$k f_{\star} \widehat{\text{RF}}_{n, N, \lambda} k_{L^2}^2 = k P_{> \ell} f_{\star} k_{L^2}^2 + o_{d, P}(\cdot):$$

- |  $d_{\text{eff}}$  capture the 'effective low-dimensionality' of the data.
- |  $d_{\text{eff}}$  decreases with  $\text{snr}_c$  (data more anisotropic) and kernel methods will perform better.
- | On the contrary, we expect NNs to learn features  $w_i$  aligned with  $z_1$  and to depend mildly on  $\text{snr}_c$ .

## Test error of KRR in the spiked covariates model

- | Given iid samples  $(fX_i; y_i g)_{i \in [n]}$  from the spiked covariates model.
- | Effective dimension:  $d_{\text{eff}} = d_s \_ (d = \text{snr}_c)$ .

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

Assume  $d_{\text{eff}}^{\ell + \delta} \min(n; N) d_{\text{eff}}^{\ell + 1} \delta$  and  $j \log(n = N) j \log d_{\text{eff}}$ , satisfies some generic conditions and is small enough. Then,

$$k f_{\star} \widehat{\text{RF}}_{n, N, \lambda} k_{L^2}^2 = k P_{> \ell} f_{\star} k_{L^2}^2 + o_{d, P}(\cdot):$$

- |  $d_{\text{eff}}$  capture the 'effective low-dimensionality' of the data.
- |  $d_{\text{eff}}$  decreases with  $\text{snr}_c$  (data more anisotropic) and kernel methods will perform better.
- | On the contrary, we expect NNs to learn features  $w_i$  aligned with  $Z_1$  and to depend mildly on  $\text{snr}_c$ .

## Test error of KRR in the spiked covariates model

- | Given iid samples  $(f_X^i; y_i g)_{i \in [n]}$  from the spiked covariates model.
- | Effective dimension:  $d_{\text{eff}} = d_s - (d = \text{snr}_c)$ .

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

Assume  $d_{\text{eff}}^{\ell+\delta} \min(n; N) d_{\text{eff}}^{\ell+1} \delta$  and  $j \log(n=N)^j \log d_{\text{eff}}$ , satisfies some generic conditions and is small enough. Then,

$$k f_{\star} \widehat{\text{RF}}_{n, N, \lambda} k_{L^2}^2 = k P_{> \ell} f_{\star} k_{L^2}^2 + o_{d, P}(\cdot):$$

- |  $d_{\text{eff}}$  capture the 'effective low-dimensionality' of the data.
- |  $d_{\text{eff}}$  decreases with  $\text{snr}_c$  (data more anisotropic) and kernel methods will perform better.
- | On the contrary, we expect NNs to learn features  $w_i$  aligned with  $z_1$  and to depend mildly on  $\text{snr}_c$ .

# NNs vs kernel methods in the spiked covariates model

Effective dimension:  $d_{\text{eff}} = d_s \_ (d = \text{snr}_c)$ .

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

- | Kernel methods to depend on  $d_{\text{eff}}$ ;
- | NNs to depend on  $d_s$ .

|  $d_s$   $d_{\text{eff}}$   $d$ :

Test error:      NN      Kernel methods

|  $d_{\text{eff}}$  decreases with  $\text{snr}_c$ :

- | Isotropic features:  $d_{\text{eff}} = d$  (low covariate SNR)

Test error:      NN      Kernel methods:

- | Highly anisotropic features:  $d_{\text{eff}} = d_s$  (high covariate SNR)

Test error:      NN      Kernel methods:



# NNs vs kernel methods in the spiked covariates model

Effective dimension:  $d_{\text{eff}} = d_s \_ (d = \text{snr}_c)$ .

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

- | Kernel methods to depend on  $d_{\text{eff}}$ ;
- | NNs to depend on  $d_s$ .

|  $d_s$   $d_{\text{eff}}$   $d$ :

Test error:      NN      Kernel methods

|  $d_{\text{eff}}$  decreases with  $\text{snr}_c$ :

- | Isotropic features:  $d_{\text{eff}} = d$  (low covariate SNR)

Test error:      NN      Kernel methods:

- | Highly anisotropic features:  $d_{\text{eff}} = d_s$  (high covariate SNR)

Test error:      NN      Kernel methods:

# NNs vs kernel methods in the spiked covariates model

Effective dimension:  $d_{\text{eff}} = d_s \_ (d = \text{snr}_c)$ .

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

- | Kernel methods to depend on  $d_{\text{eff}}$ ;
- | NNs to depend on  $d_s$ .

|  $d_s$   $d_{\text{eff}}$   $d$ :

Test error:      NN      Kernel methods

|  $d_{\text{eff}}$  decreases with  $\text{snr}_c$ :

- | Isotropic features:  $d_{\text{eff}} = d$  (low covariate SNR)

Test error:      NN      Kernel methods:

- | Highly anisotropic features:  $d_{\text{eff}} = d_s$  (high covariate SNR)

Test error:      NN      Kernel methods:

## Intuition for the gap: anisotropic case

Goal: fitting  $f(x_1)$  using features  $(hw_j; x_i)$ .

Consider  $E[xx^T] = \text{diag}(\text{snr}_c; \text{Id}_{d-1})$ .

- To fit  $f(x_1)$ , we need to use features  $(hw_j; x_i)$  that are correlated to  $x_1$ , i.e.,

$$\text{corr}(x_1 e_1; hw_j; x_i) = \text{snr}_c \frac{hw_j; e_1 i^2}{k^{-1/2} w_j k_2^2} = \text{snr}_c O_{d,P}(d^{-1}) = O_{d,P}(d_{\text{eff}}^{-1}):$$

- For RF, when  $\text{snr}_c$  is large, all random  $w_j$ 's are good.
- For NN,  $w_j$ 's can be chosen to have a large correlation with  $e_1$ .

RF automatically have access to good  $w_j$ 's.

## Intuition for the gap: anisotropic case

Goal: fitting  $f(x_1)$  using features  $(hw_j; x_i)$ .

Consider  $E[xx^T] = \text{diag}(\text{snr}_c; \text{Id}_{d-1})$ .

- To fit  $f(x_1)$ , we need to use features  $(hw_j; x_i)$  that are correlated to  $x_1$ , i.e.,

$$\text{corr}(x_1 e_1; hw_j; x_i) = \text{snr}_c \frac{hw_j; e_1 j^2}{k^{-1/2} w_j k_2^2} = \text{snr}_c O_{d,P}(d^{-1}) = O_{d,P}(d_{\text{eff}}^{-1}):$$

- For RF, when  $\text{snr}_c$  is large, all random  $w_j$ 's are good.
- For NN,  $w_j$ 's can be chosen to have a large correlation with  $e_1$ .

RF automatically have access to good  $w_j$ 's.

## Intuition for the gap: anisotropic case

Goal: fitting  $y(x_1)$  using features  $(hw_j; x_i)$ .

Consider  $E[xx^T] = \text{diag}(\text{snr}_c; \text{Id}_{d-1})$ .

- To fit  $y(x_1)$ , we need to use features  $(hw_j; x_i)$  that are correlated to  $x_1$ , i.e.,

$$\text{corr}(x_1 e_1; hw_j; x_i) = \text{snr}_c \frac{hw_j; e_1 j^2}{k^{-1/2} w_j k_2^2} = \text{snr}_c O_{d,P}(d^{-1}) = O_{d,P}(d_{\text{eff}}^{-1}):$$

- For RF, when  $\text{snr}_c$  is large, all random  $w_j$ 's are good.
- For NN,  $w_j$ 's can be chosen to have a large correlation with  $e_1$ .

RF automatically have access to good  $w_j$ 's.

## Intuition for the gap: anisotropic case

Goal: fitting  $f(x_1)$  using features  $(hw_j; x_i)$ .

Consider  $E[xx^T] = \text{diag}(\text{snr}_c; \text{Id}_{d-1})$ .

- To fit  $f(x_1)$ , we need to use features  $(hw_j; x_i)$  that are correlated to  $x_1$ , i.e.,

$$\text{corr}(x_1 e_1; hw_j; x_i) = \text{snr}_c \frac{hw_j e_1^T j^2}{k^{-1/2} w_j k_2^2} = \text{snr}_c O_{d,P}(d^{-1}) = O_{d,P}(d_{\text{eff}}^{-1}):$$

- For RF, when  $\text{snr}_c$  is large, all random  $w_j$ 's are good.
- For NN,  $w_j$ 's can be chosen to have a large correlation with  $e_1$ .

RF automatically have access to good  $w_j$ 's.

## Intuition for the gap: anisotropic case

Goal: fitting  $y(x_1)$  using features  $(hw_i; x_i)$ .

Consider  $E[xx^T] = \text{diag}(\text{snr}_c; \text{Id}_{d-1})$ .

- To fit  $y(x_1)$ , we need to use features  $(hw_i; x_i)$  that are correlated to  $x_1$ , i.e.,

$$\text{corr}(x_1 e_1; hw_i; x_i) = \text{snr}_c \frac{hw_i; e_1 j^2}{k^{-1/2} w_i k_2^2} = \text{snr}_c O_{d,P}(d^{-1}) = O_{d,P}(d_{\text{eff}}^{-1}):$$

- For RF, when  $\text{snr}_c$  is large, all random  $w_i$ 's are good.
- For NN,  $w_i$ 's can be chosen to have a large correlation with  $e_1$ .

RF automatically have access to good  $w_i$ 's.

## 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 ( $Z_1$  part);
- | The labels to depend predominantly on the low-frequencies ( $y = f_*(x) = f(Z_1)$ ).

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

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 ( $Z_1$  part);
- | The labels to depend predominantly on the low-frequencies ( $y = f_*(x) = f(Z_1)$ ).

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

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 ( $z_1$  part);
- | The labels to depend predominantly on the low-frequencies ( $y = f_*(x) = f(z_1)$ ).

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

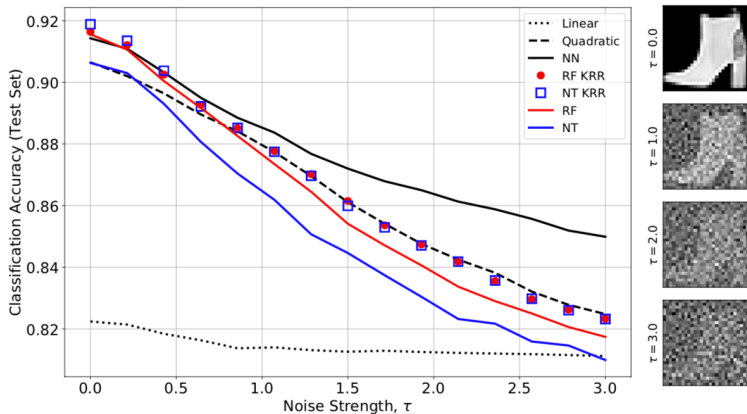


Figure: Test accuracy on FMNIST: adding noise to the high frequency components.

Spiked covariates model:

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

$$\text{snr}_c = \text{Covariate SNR} = \frac{\text{Signal covariates variance}}{\text{Noise covariates variance}};$$

- | Small  $\text{snr}_c$ : large separation.
- | Large  $\text{snr}_c$ : kernel methods perform closer to NN.

Intuition:

- | When  $\text{snr}_c$  is small, RF fail to find the signal covariates.
- | When  $\text{snr}_c$  is big, RF automatically find the signal covariates.
- | NNs always look for the signal covariates.

## C. Invariant function estimation on the sphere

## Symmetries in the data

- | 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?

## Symmetries in the data

- | 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?

## Symmetries in the data

- | 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?



## Symmetries in the data

- | 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?

## Cyclic function estimation on the sphere

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

| Consider  $G_d = \langle g_0; g_1; \dots; g_{d-1} \rangle$  the cyclic group:

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

| Goal: learn a cyclic invariant function  $f_*$ :

$$\text{i.e., } f_*(g(x)) = f_*(x) \text{ for all } g \in G_d.$$

$$\text{E.g., } f(x) = \sum_{i=1}^d x_i x_{i+1}.$$

Stylized model for an image label  $y = f_*(x)$  invariant by translation of image  $x$ .

## Cyclic function estimation on the sphere

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

| Consider  $G_d = \langle g_0; g_1; \dots; g_{d-1} \rangle$  the cyclic group:

$$g_i x = (x_{d-i+1}; x_{d-i+2}; \dots; x_d; x_1; x_2; \dots; x_{d-i}):$$

| Goal: learn a cyclic invariant function  $f_*$ :

$$\text{i.e., } f_*(g x) = f_*(x) \text{ for all } g \in G_d.$$

$$\text{E.g., } f(x) = \sum_{i=1}^d x_i x_{i+1}.$$

Stylized model for an image label  $y = f_*(x)$  invariant by translation of image  $x$ .

## Cyclic function estimation on the sphere

| Data  $x \sim \text{Unif}(S^{d-1}(\rho^{-d}))$ .

| Consider  $G_d = \langle g_0; g_1; \dots; g_{d-1} \rangle$  the cyclic group:

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

| Goal: learn a cyclic invariant function  $f_\star$ :

$$\text{i.e., } f_\star(g x) = f_\star(x) \text{ for all } g \in G_d.$$

$$\text{E.g., } f(x) = \sum_{i=1}^d x_i x_{i+1}.$$

Stylized model for an image label  $y = f_\star(x)$  invariant by translation of image  $x$ .

## Cyclic function estimation on the sphere

| Data  $x \sim \text{Unif}(S^{d-1}(\rho^{-d}))$ .

| Consider  $G_d = \langle g_0; g_1; \dots; g_{d-1} \rangle$  the cyclic group:

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

| Goal: learn a cyclic invariant function  $f_\star$ :

$$\text{i.e., } f_\star(g x) = f_\star(x) \text{ for all } g \in G_d.$$

$$\text{E.g., } f(x) = \sum_{i=1}^d x_i x_{i+1}.$$

Stylized model for an image label  $y = f_\star(x)$  invariant by translation of image  $x$ .

## Invariant random features and kernels

| Invariant RF model:

$$\text{RF}_N^{\text{inv}}(x; a; W) = \sum_{i=1}^N a_i \int_{g \in G_d} (h(w_i; g \cdot x)) \, (dg) = \frac{1}{d} \sum_{i=1}^N a_i \sum_{k=1}^d (h(w_i; g_k \cdot x)):$$

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

Associated invariant kernel:

$$K_{\text{inv}}(x; y) = \int_{g \in G_d} h(g \cdot x; g \cdot y) \, (dg):$$

| Compared to standard RF model:

$$\text{RF}_N(x; a; W) = \sum_{i=1}^N a_i (h(x; w_i));$$

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

$$K(x; y) = h(x; y):$$

## Invariant random features and kernels

| Invariant RF model:

$$\text{RF}_N^{\text{inv}}(x; a; W) = \sum_{i=1}^N a_i \int_{g \in G_d} (h(w_i; g \cdot x)) \, (dg) = \frac{1}{d} \sum_{i=1}^N a_i \sum_{k=1}^d (h(w_i; g_k \cdot x)):$$

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

Associated invariant kernel:

$$K_{\text{inv}}(x; y) = \int_{g \in G_d} h(hg \cdot x; y) \, (dg):$$

| Compared to standard RF model:

$$\text{RF}_N(x; a; W) = \sum_{i=1}^N a_i (hx; w_i);$$

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

$$K(x; y) = h(hx; y):$$

## Test error of KRR with invariant kernel

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

$$y_i = f_*(x_i) + \epsilon_i; \quad x_i \stackrel{iid}{\sim} \text{Unif}(S^{d-1}(\sqrt{d})); \quad \epsilon_i \stackrel{iid}{\sim} N(0, \sigma^2):$$

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

Assume  $d^{\ell+1+\delta} \ll \min(n; N)$ ,  $d^{\ell+\delta} \ll j \log(n=N)$ ,  $\log d$ , satisfies some generic conditions and  $\sigma$  is small enough. Then,

$$k_{f_*} \widehat{\text{RF}}_{n, N, \lambda}^{\text{inv}} k_{L^2}^2 = k_{P_{>\ell} f_*} k_{L^2}^2 + o_{d, P}():$$

- Standard RF model:

$$\text{For } d^{\ell+\delta} \ll \min(n; N), \quad d^{\ell+1+\delta} \ll j \log(n=N), \quad k_{f_*} \widehat{\text{RF}}_{n, N, \lambda} k_{L^2}^2 = k_{P_{>\ell} f_*} k_{L^2}^2 + o_d():$$

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

- More generally, for  $G_d$  with 'degeneracy', we gain a factor  $d^\alpha$ .



## Test error of KRR with invariant kernel

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

$$y_i = f_*(x_i) + \epsilon_i; \quad x_i \stackrel{iid}{\sim} \text{Unif}(S^{d-1}(\sqrt{d})); \quad \epsilon_i \stackrel{iid}{\sim} N(0, \sigma^2):$$

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

Assume  $d^{\ell+1+\delta} \min(n; N) \gg d^{\ell+\delta}$  and  $j \log(n=N) \ll \log d$ , satisfies some generic conditions and  $\sigma$  is small enough. Then,

$$k_{f_*} \widehat{\text{RF}}_{n, N, \lambda}^{\text{inv}} k_{L^2}^2 = k_{P_{>\ell} f_*} k_{L^2}^2 + o_{d, P}():$$

- Standard RF model:

$$\text{For } d^{\ell+\delta} \min(n; N) \gg d^{\ell+1+\delta}; \quad k_{f_*} \widehat{\text{RF}}_{n, N, \lambda} k_{L^2}^2 = k_{P_{>\ell} f_*} k_{L^2}^2 + o_d():$$

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

- More generally, for  $G_d$  with 'degeneracy', we gain a factor  $d^\alpha$ .

## Test error of KRR with invariant kernel

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

$$y_i = f_*(x_i) + \epsilon_i; \quad x_i \stackrel{iid}{\sim} \text{Unif}(S^{d-1}(\sqrt{d})); \quad \epsilon_i \stackrel{iid}{\sim} N(0; \sigma^2):$$

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

Assume  $d^{\ell+1+\delta} \min(n; N) \gg d^{\ell+\delta}$  and  $j \log(n=N) \ll \log d$ , satisfies some generic conditions and  $\sigma$  is small enough. Then,

$$k_{f_*} \widehat{\text{RF}}_{n, N, \lambda}^{\text{inv}} k_{L^2}^2 = k_{P_{>\ell} f_*} k_{L^2}^2 + o_{d, P}():$$

- Standard RF model:

$$\text{For } d^{\ell+\delta} \min(n; N) \gg d^{\ell+1+\delta}; \quad k_{f_*} \widehat{\text{RF}}_{n, N, \lambda} k_{L^2}^2 = k_{P_{>\ell} f_*} k_{L^2}^2 + o_d():$$

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

- More generally, for  $G_d$  with 'degeneracy  $\alpha$ ', we gain a factor  $d^\alpha$ .

# Numerical simulations

$$f_{\text{lin}} = \frac{1}{d} \sum_{i=1}^d X_i; \quad f_{\text{quad}} = \frac{1}{d} \sum_{i=1}^d X_i X_{i+1}; \quad f_{\text{cube}} = \frac{1}{d} \sum_{i=1}^d X_i X_{i+1} X_{i+2}.$$

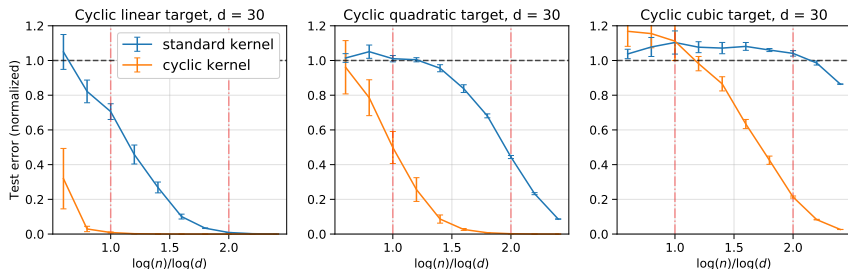


Figure: Test error of KRR with invariant kernel and inner-product kernel.

# Cyclic invariant MNIST

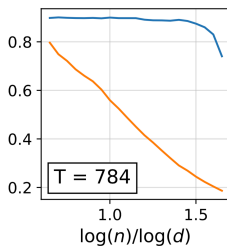
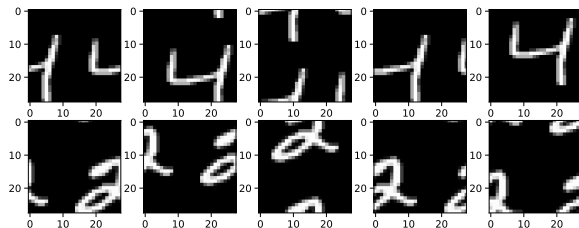


Figure: Test accuracy against number of samples (orange: cyclic kernel, blue: standard kernel).

## 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 unnecessary 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].

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 unnecessary 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].

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 unnecessary 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].

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 unnecessary 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].



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 unnecessary 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!