

Minimum complexity interpolation in random features models

Theodor Misiakiewicz

Stanford University

June 15th, 2021

Youth in High Dimensions 2021

Joint work with Michael Celentano and Andrea Montanari (Stanford)

RKHS methods (NNs in the 'lazy training' regime)

▶ Supervised learning setting: $\{(y_i, \mathbf{x}_i)\}_{i \in [n]}$, $\mathbf{x}_i \sim_{iid} (\mathcal{X}, \mathbb{P})$, $y_i = f_*(\mathbf{x}_i) + \varepsilon_i$.

▶ Kernel machines: starting from a weight space (Ω, μ)

▶ Featurization map: $\phi(\cdot; \mathbf{w}) : \mathcal{X} \rightarrow \mathbb{R}$, e.g., $\phi(\mathbf{x}; \mathbf{w}) = \sigma(\langle \mathbf{x}, \mathbf{w} \rangle)$,

$$\mathcal{F}_2 = \left\{ f(\mathbf{x}; a) = \int_{\Omega} a(\mathbf{w}) \phi(\mathbf{x}; \mathbf{w}) \mu(d\mathbf{w}) : \|a\|_{L^2}^2 = \int_{\Omega} |a(\mathbf{w})|^2 \mu(d\mathbf{w}) \right\}.$$

Associated kernel: $K(\mathbf{x}_1, \mathbf{x}_2) = \int_{\Omega} \phi(\mathbf{x}_1; \mathbf{w}) \phi(\mathbf{x}_2; \mathbf{w}) \mu(d\mathbf{w})$.

▶ Convex loss function: $\ell : \mathbb{R} \times \mathbb{R} \rightarrow \mathbb{R}$,

$$\hat{f}_n = \arg \min_{f \in \mathcal{F}_2} \left\{ \sum_{i=1}^n \ell(y_i, f(\mathbf{x}_i; a)) + \lambda \|a\|_{L^2}^2 \right\}.$$

▶ Can be solved efficiently despite \mathcal{F}_2 being infinite-dimensional.
The 'representer theorem':

$$\hat{a}(\mathbf{w}) = \sum_{i=1}^n c_i \phi(\mathbf{x}_i; \mathbf{w}).$$

RKHS methods (NNs in the 'lazy training' regime)

► **Supervised learning setting:** $\{(y_i, \mathbf{x}_i)\}_{i \in [n]}$, $\mathbf{x}_i \sim_{iid} (\mathcal{X}, \mathbb{P})$, $y_i = f_*(\mathbf{x}_i) + \varepsilon_i$.

► **Kernel machines:** starting from a weight space (Ω, μ)

► Featurization map: $\phi(\cdot; \mathbf{w}) : \mathcal{X} \rightarrow \mathbb{R}$, e.g., $\phi(\mathbf{x}; \mathbf{w}) = \sigma(\langle \mathbf{x}, \mathbf{w} \rangle)$,

$$\mathcal{F}_2 = \left\{ f(\mathbf{x}; a) = \int_{\Omega} a(\mathbf{w}) \phi(\mathbf{x}; \mathbf{w}) \mu(d\mathbf{w}) : \|a\|_{L^2}^2 = \int_{\Omega} |a(\mathbf{w})|^2 \mu(d\mathbf{w}) \right\}.$$

Associated kernel: $K(\mathbf{x}_1, \mathbf{x}_2) = \int_{\Omega} \phi(\mathbf{x}_1; \mathbf{w}) \phi(\mathbf{x}_2; \mathbf{w}) \mu(d\mathbf{w})$.

► Convex loss function: $\ell : \mathbb{R} \times \mathbb{R} \rightarrow \mathbb{R}$,

$$\hat{f}_n = \arg \min_{f \in \mathcal{F}_2} \left\{ \sum_{i=1}^n \ell(y_i, f(\mathbf{x}_i; a)) + \lambda \|a\|_{L^2}^2 \right\}.$$

► Can be solved efficiently despite \mathcal{F}_2 being infinite-dimensional.
The 'representer theorem':

$$\hat{a}(\mathbf{w}) = \sum_{i=1}^n c_i \phi(\mathbf{x}_i; \mathbf{w}).$$

RKHS methods (NNs in the 'lazy training' regime)

▶ **Supervised learning setting:** $\{(y_i, \mathbf{x}_i)\}_{i \in [n]}$, $\mathbf{x}_i \sim_{iid} (\mathcal{X}, \mathbb{P})$, $y_i = f_*(\mathbf{x}_i) + \varepsilon_i$.

▶ **Kernel machines:** starting from a weight space (Ω, μ)

▶ Featurization map: $\phi(\cdot; \mathbf{w}) : \mathcal{X} \rightarrow \mathbb{R}$, e.g., $\phi(\mathbf{x}; \mathbf{w}) = \sigma(\langle \mathbf{x}, \mathbf{w} \rangle)$,

$$\mathcal{F}_2 = \left\{ f(\mathbf{x}; a) = \int_{\Omega} a(\mathbf{w}) \phi(\mathbf{x}; \mathbf{w}) \mu(d\mathbf{w}) : \|a\|_{L^2}^2 = \int_{\Omega} |a(\mathbf{w})|^2 \mu(d\mathbf{w}) \right\}.$$

Associated kernel: $K(\mathbf{x}_1, \mathbf{x}_2) = \int_{\Omega} \phi(\mathbf{x}_1; \mathbf{w}) \phi(\mathbf{x}_2; \mathbf{w}) \mu(d\mathbf{w})$.

▶ Convex loss function: $\ell : \mathbb{R} \times \mathbb{R} \rightarrow \mathbb{R}$,

$$\hat{f}_n = \arg \min_{f \in \mathcal{F}_2} \left\{ \sum_{i=1}^n \ell(y_i, f(\mathbf{x}_i; a)) + \lambda \|a\|_{L^2}^2 \right\}.$$

▶ Can be solved efficiently despite \mathcal{F}_2 being infinite-dimensional.
The 'representer theorem':

$$\hat{a}(\mathbf{w}) = \sum_{i=1}^n c_i \phi(\mathbf{x}_i; \mathbf{w}).$$

RKHS methods (NNs in the 'lazy training' regime)

▶ **Supervised learning setting:** $\{(y_i, \mathbf{x}_i)\}_{i \in [n]}$, $\mathbf{x}_i \sim_{iid} (\mathcal{X}, \mathbb{P})$, $y_i = f_*(\mathbf{x}_i) + \varepsilon_i$.

▶ **Kernel machines:** starting from a weight space (Ω, μ)

▶ Featurization map: $\phi(\cdot; \mathbf{w}) : \mathcal{X} \rightarrow \mathbb{R}$, e.g., $\phi(\mathbf{x}; \mathbf{w}) = \sigma(\langle \mathbf{x}, \mathbf{w} \rangle)$,

$$\mathcal{F}_2 = \left\{ f(\mathbf{x}; a) = \int_{\Omega} a(\mathbf{w}) \phi(\mathbf{x}; \mathbf{w}) \mu(d\mathbf{w}) : \|a\|_{L^2}^2 = \int_{\Omega} |a(\mathbf{w})|^2 \mu(d\mathbf{w}) \right\}.$$

Associated kernel: $K(\mathbf{x}_1, \mathbf{x}_2) = \int_{\Omega} \phi(\mathbf{x}_1; \mathbf{w}) \phi(\mathbf{x}_2; \mathbf{w}) \mu(d\mathbf{w})$.

▶ Convex loss function: $\ell : \mathbb{R} \times \mathbb{R} \rightarrow \mathbb{R}$,

$$\hat{f}_n = \arg \min_{f \in \mathcal{F}_2} \left\{ \sum_{i=1}^n \ell(y_i, f(\mathbf{x}_i; a)) + \lambda \|a\|_{L^2}^2 \right\}.$$

▶ Can be solved efficiently despite \mathcal{F}_2 being infinite-dimensional.
The 'representer theorem':

$$\hat{a}(\mathbf{w}) = \sum_{i=1}^n c_i \phi(\mathbf{x}_i; \mathbf{w}).$$

RKHS methods (NNs in the 'lazy training' regime)

▶ **Supervised learning setting:** $\{(y_i, \mathbf{x}_i)\}_{i \in [n]}$, $\mathbf{x}_i \sim_{iid} (\mathcal{X}, \mathbb{P})$, $y_i = f_*(\mathbf{x}_i) + \varepsilon_i$.

▶ **Kernel machines:** starting from a weight space (Ω, μ)

▶ Featurization map: $\phi(\cdot; \mathbf{w}) : \mathcal{X} \rightarrow \mathbb{R}$, e.g., $\phi(\mathbf{x}; \mathbf{w}) = \sigma(\langle \mathbf{x}, \mathbf{w} \rangle)$,

$$\mathcal{F}_2 = \left\{ f(\mathbf{x}; a) = \int_{\Omega} a(\mathbf{w}) \phi(\mathbf{x}; \mathbf{w}) \mu(d\mathbf{w}) : \|a\|_{L^2}^2 = \int_{\Omega} |a(\mathbf{w})|^2 \mu(d\mathbf{w}) \right\}.$$

Associated kernel: $K(\mathbf{x}_1, \mathbf{x}_2) = \int_{\Omega} \phi(\mathbf{x}_1; \mathbf{w}) \phi(\mathbf{x}_2; \mathbf{w}) \mu(d\mathbf{w})$.

▶ Convex loss function: $\ell : \mathbb{R} \times \mathbb{R} \rightarrow \mathbb{R}$,

$$\hat{f}_n = \arg \min_{f \in \mathcal{F}_2} \left\{ \sum_{i=1}^n \ell(y_i, f(\mathbf{x}_i; a)) + \lambda \|a\|_{L^2}^2 \right\}.$$

▶ Can be solved efficiently despite \mathcal{F}_2 being infinite-dimensional.
The 'representer theorem':

$$\hat{a}(\mathbf{w}) = \sum_{i=1}^n c_i \phi(\mathbf{x}_i; \mathbf{w}).$$

RKHS methods (NNs in the 'lazy training' regime)

▶ **Supervised learning setting:** $\{(y_i, \mathbf{x}_i)\}_{i \in [n]}$, $\mathbf{x}_i \sim_{iid} (\mathcal{X}, \mathbb{P})$, $y_i = f_*(\mathbf{x}_i) + \varepsilon_i$.

▶ **Kernel machines:** starting from a weight space (Ω, μ)

▶ Featurization map: $\phi(\cdot; \mathbf{w}) : \mathcal{X} \rightarrow \mathbb{R}$, e.g., $\phi(\mathbf{x}; \mathbf{w}) = \sigma(\langle \mathbf{x}, \mathbf{w} \rangle)$,

$$\mathcal{F}_2 = \left\{ f(\mathbf{x}; a) = \int_{\Omega} a(\mathbf{w}) \phi(\mathbf{x}; \mathbf{w}) \mu(d\mathbf{w}) : \|a\|_{L^2}^2 = \int_{\Omega} |a(\mathbf{w})|^2 \mu(d\mathbf{w}) \right\}.$$

Associated kernel: $K(\mathbf{x}_1, \mathbf{x}_2) = \int_{\Omega} \phi(\mathbf{x}_1; \mathbf{w}) \phi(\mathbf{x}_2; \mathbf{w}) \mu(d\mathbf{w})$.

▶ Convex loss function: $\ell : \mathbb{R} \times \mathbb{R} \rightarrow \mathbb{R}$,

$$\hat{f}_n = \arg \min_{f \in \mathcal{F}_2} \left\{ \sum_{i=1}^n \ell(y_i, f(\mathbf{x}_i; a)) + \lambda \|a\|_{L^2}^2 \right\}.$$

▶ Can be solved efficiently despite \mathcal{F}_2 being infinite-dimensional.
The 'representer theorem':

$$\hat{a}(\mathbf{w}) = \sum_{i=1}^n c_i \phi(\mathbf{x}_i; \mathbf{w}).$$

Curse of dimensionality

▶ $\mathcal{F}_2(R) = \left\{ f(\mathbf{x}; \mathbf{a}) : \|\mathbf{a}\|_{L^2} \leq R \right\}, f_\star \in \mathcal{F}_2(R),$

Generalization error of learning f_\star from n samples $\leq C \frac{R}{\sqrt{n}}$.

▶ Kernel methods suffer from the **curse of dimensionality**.

▶ $\mathbf{x}_i \sim_{iid} \text{Unif}(\mathbb{S}^{d-1}(\sqrt{d})), \mathbf{w} \sim \text{Unif}(\mathbb{S}^{d-1}(1)), \phi(\mathbf{x}; \mathbf{w}) = \sigma(\langle \mathbf{x}, \mathbf{w} \rangle).$

Any target function $f_\star \in L^2(\mathcal{X})$, number of samples: $d^k \ll n \ll d^{k+1}$ [GMMM,'19]

Test error with squared loss: $\|f_\star - \hat{f}_n\|_{L^2}^2 \geq \|P_{>k} f_\star\|_{L^2}^2 + o_d(1).$

$\mathcal{F}_2(\sqrt{n}) \approx \{\text{degree-}k \text{ polynomials}\}.$

▶ $\mathcal{F}_2(R) = \left\{ f(\mathbf{x}; \mathbf{a}) : \|\mathbf{a}\|_{L^2} \leq R \right\}, f_* \in \mathcal{F}_2(R),$

Generalization error of learning f_* from n samples $\leq C \frac{R}{\sqrt{n}}$.

▶ Kernel methods suffer from the **curse of dimensionality**.

▶ $\mathbf{x}_i \sim_{iid} \text{Unif}(\mathbb{S}^{d-1}(\sqrt{d})), \mathbf{w} \sim \text{Unif}(\mathbb{S}^{d-1}(1)), \phi(\mathbf{x}; \mathbf{w}) = \sigma(\langle \mathbf{x}, \mathbf{w} \rangle).$

Any target function $f_* \in L^2(\mathcal{X})$, number of samples: $d^k \ll n \ll d^{k+1}$ [GMMM,'19]

Test error with squared loss: $\|f_* - \hat{f}_n\|_{L^2}^2 \geq \|P_{>k} f_*\|_{L^2}^2 + o_d(1).$

$\mathcal{F}_2(\sqrt{n}) \approx \{\text{degree-}k \text{ polynomials}\}.$

▶ $\mathcal{F}_2(R) = \left\{ f(\mathbf{x}; \mathbf{a}) : \|\mathbf{a}\|_{L^2} \leq R \right\}, f_* \in \mathcal{F}_2(R),$

$$\text{Generalization error of learning } f_* \text{ from } n \text{ samples} \leq C \frac{R}{\sqrt{n}}.$$

▶ Kernel methods suffer from the **curse of dimensionality**.

▶ $\mathbf{x}_i \sim_{iid} \text{Unif}(\mathbb{S}^{d-1}(\sqrt{d})), \mathbf{w} \sim \text{Unif}(\mathbb{S}^{d-1}(1)), \phi(\mathbf{x}; \mathbf{w}) = \sigma(\langle \mathbf{x}, \mathbf{w} \rangle).$

Any target function $f_* \in L^2(\mathcal{X})$, number of samples: $d^k \ll n \ll d^{k+1}$ [GMMM,'19]

$$\text{Test error with squared loss: } \|f_* - \hat{f}_n\|_{L^2}^2 \geq \|P_{>k} f_*\|_{L^2}^2 + o_d(1).$$

$$\mathcal{F}_2(\sqrt{n}) \approx \{\text{degree-}k \text{ polynomials}\}.$$

'Convex neural networks'

- ▶ $f_*(\mathbf{x}) = \phi(\mathbf{x}; \mathbf{w}_*)$. We have $a_*(\mathbf{w}) = \delta_{\mathbf{w}, \mathbf{w}_*}$ and $\|a_*\|_{L^2} = \infty$ (hence $f_* \notin \mathcal{F}_2$).

f_* is not learned efficiently by Kernel methods: f_* cannot be approximated by $f(\cdot; a)$ with $\|a\|_{L^2}$ bounded.

- ▶ 'Convex neural network' [Bengio et al., '06]: for $1 \leq p < 2$,

$$\mathcal{F}_p(R) = \left\{ f(\cdot; a) : \|a\|_{L^p} = \left(\int_{\Omega} |a(\mathbf{w})|^p \mu(d\mathbf{w}) \right)^{1/p} \leq R \right\}.$$

- ▶ By Jensen's inequality $\mathcal{F}_2(R) \subset \mathcal{F}_p(R) \subset \mathcal{F}_1(R)$.
- ▶ $a(\mathbf{w})$ may tend to singular distribution with $\|a\|_{L^1}$ bounded (not true for L^2 -norm).

[Bach, '17] \mathcal{F}_1 is adaptive and beat the curse of dimensionality for functions that depends on a low-dimensional projection of the covariates.

\mathcal{F}_2 is not adaptive.

'Convex neural networks'

- ▶ $f_*(\mathbf{x}) = \phi(\mathbf{x}; \mathbf{w}_*)$. We have $a_*(\mathbf{w}) = \delta_{\mathbf{w}, \mathbf{w}_*}$ and $\|a_*\|_{L^2} = \infty$ (hence $f_* \notin \mathcal{F}_2$).

f_* is not learned efficiently by Kernel methods: f_* cannot be approximated by $f(\cdot; a)$ with $\|a\|_{L^2}$ bounded.

- ▶ 'Convex neural network' [Bengio et al., '06]: for $1 \leq p < 2$,

$$\mathcal{F}_p(R) = \left\{ f(\cdot; a) : \|a\|_{L^p} = \left(\int_{\Omega} |a(\mathbf{w})|^p \mu(d\mathbf{w}) \right)^{1/p} \leq R \right\}.$$

- ▶ By Jensen's inequality $\mathcal{F}_2(R) \subset \mathcal{F}_p(R) \subset \mathcal{F}_1(R)$.
- ▶ $a(\mathbf{w})$ may tend to singular distribution with $\|a\|_{L^1}$ bounded (not true for L^2 -norm).

[Bach, '17] \mathcal{F}_1 is adaptive and beat the curse of dimensionality for functions that depends on a low-dimensional projection of the covariates.

\mathcal{F}_2 is not adaptive.

'Convex neural networks'

- ▶ $f_*(\mathbf{x}) = \phi(\mathbf{x}; \mathbf{w}_*)$. We have $a_*(\mathbf{w}) = \delta_{\mathbf{w}, \mathbf{w}_*}$ and $\|a_*\|_{L^2} = \infty$ (hence $f_* \notin \mathcal{F}_2$).

f_* is not learned efficiently by Kernel methods: f_* cannot be approximated by $f(\cdot; a)$ with $\|a\|_{L^2}$ bounded.

- ▶ 'Convex neural network' [Bengio et al., '06]: for $1 \leq p < 2$,

$$\mathcal{F}_p(R) = \left\{ f(\cdot; a) : \|a\|_{L^p} = \left(\int_{\Omega} |a(\mathbf{w})|^p \mu(d\mathbf{w}) \right)^{1/p} \leq R \right\}.$$

- ▶ By Jensen's inequality $\mathcal{F}_2(R) \subset \mathcal{F}_p(R) \subset \mathcal{F}_1(R)$.
- ▶ $a(\mathbf{w})$ may tend to singular distribution with $\|a\|_{L^1}$ bounded (not true for L^2 -norm).

[Bach, '17] \mathcal{F}_1 is adaptive and beat the curse of dimensionality for functions that depends on a low-dimensional projection of the covariates.

\mathcal{F}_2 is not adaptive.

'Convex neural networks'

- ▶ $f_*(\mathbf{x}) = \phi(\mathbf{x}; \mathbf{w}_*)$. We have $a_*(\mathbf{w}) = \delta_{\mathbf{w}, \mathbf{w}_*}$ and $\|a_*\|_{L^2} = \infty$ (hence $f_* \notin \mathcal{F}_2$).

f_* is not learned efficiently by Kernel methods: f_* cannot be approximated by $f(\cdot; a)$ with $\|a\|_{L^2}$ bounded.

- ▶ 'Convex neural network' [Bengio et al., '06]: for $1 \leq p < 2$,

$$\mathcal{F}_p(R) = \left\{ f(\cdot; a) : \|a\|_{L^p} = \left(\int_{\Omega} |a(\mathbf{w})|^p \mu(d\mathbf{w}) \right)^{1/p} \leq R \right\}.$$

- ▶ By Jensen's inequality $\mathcal{F}_2(R) \subset \mathcal{F}_p(R) \subset \mathcal{F}_1(R)$.
- ▶ $a(\mathbf{w})$ may tend to singular distribution with $\|a\|_{L^1}$ bounded (not true for L^2 -norm).

[Bach, '17] \mathcal{F}_1 is adaptive and beat the curse of dimensionality for functions that depends on a low-dimensional projection of the covariates.

\mathcal{F}_2 is not adaptive.

- ▶ **Minimum-norm interpolating solution:** ($\lambda \rightarrow 0^+$ in ERM)

$$\begin{aligned} & \text{minimize} && \int_{\Omega} |a(\mathbf{w})|^p \mu(d\mathbf{w}), \\ & \text{subj. to} && f(\mathbf{x}_i; a) = y_i, \quad \forall i \leq n. \end{aligned}$$

- ▶ Correspond to modern practice of training until interpolation.
- ▶ Infinite dimensional convex problem. Not clear if it is tractable for $p \neq 2$.
- ▶ For $p = 1$:
 - [Bengio et al., '06] incremental algorithm but no global optimality guarantees.
 - [Bach, '17] conditional gradient algorithm, but each step potentially hard.

- ▶ **Minimum-norm interpolating solution:** ($\lambda \rightarrow 0^+$ in ERM)

$$\begin{aligned} & \text{minimize} && \int_{\Omega} |a(\mathbf{w})|^p \mu(d\mathbf{w}), \\ & \text{subj. to} && f(\mathbf{x}_i; a) = y_i, \quad \forall i \leq n. \end{aligned}$$

- ▶ Correspond to modern practice of training until interpolation.
- ▶ Infinite dimensional convex problem. Not clear if it is tractable for $p \neq 2$.
- ▶ For $p = 1$:
 - [Bengio et al., '06] incremental algorithm but no global optimality guarantees.
 - [Bach, '17] conditional gradient algorithm, but each step potentially hard.

- ▶ **Minimum-norm interpolating solution:** ($\lambda \rightarrow 0^+$ in ERM)

$$\begin{aligned} & \text{minimize} && \int_{\Omega} |a(\mathbf{w})|^p \mu(d\mathbf{w}), \\ & \text{subj. to} && f(\mathbf{x}_i; a) = y_i, \quad \forall i \leq n. \end{aligned}$$

- ▶ Correspond to modern practice of training until interpolation.
- ▶ Infinite dimensional convex problem. Not clear if it is tractable for $p \neq 2$.
- ▶ For $p = 1$:
 - [Bengio et al., '06] incremental algorithm but no global optimality guarantees.
 - [Bach, '17] conditional gradient algorithm, but each step potentially hard.

- ▶ **Minimum-norm interpolating solution:** ($\lambda \rightarrow 0^+$ in ERM)

$$\begin{aligned} & \text{minimize} && \int_{\Omega} |a(\mathbf{w})|^p \mu(d\mathbf{w}), \\ & \text{subj. to} && f(\mathbf{x}_i; a) = y_i, \quad \forall i \leq n. \end{aligned}$$

- ▶ Correspond to modern practice of training until interpolation.
- ▶ Infinite dimensional convex problem. Not clear if it is tractable for $p \neq 2$.
- ▶ For $p = 1$:
 - [Bengio et al., '06] incremental algorithm but no global optimality guarantees.
 - [Bach, '17] conditional gradient algorithm, but each step potentially hard.

Random Features approximation

- ▶ Approximate μ by finitely supported $\hat{\mu}_M$: sample M weights $\mathbf{w}_j \sim_{iid} \mu$,

$$f(\mathbf{x}; \mathbf{a}) = \int_{\Omega} \mathbf{a}(\mathbf{w})\phi(\mathbf{x}; \mathbf{w})\mu(d\mathbf{w}) \longrightarrow f_M(\mathbf{x}; \mathbf{a}) = \frac{1}{M} \sum_{j=1}^M a_j\phi(\mathbf{x}; \mathbf{w}_j).$$

- ▶ 'Finite-width' problem is easy to solve: $\mathbf{a} \in \mathbb{R}^M$,

$$\text{minimize } \frac{1}{M} \sum_{j=1}^M |a_j|^p,$$

$$\text{subj. to } f_M(\mathbf{x}_i; \mathbf{a}) = y_i, \quad \forall i \leq n.$$

How large M needs to be for the finite-width solution $\hat{f}_{\text{RF},M,n}$ to approximate the infinite-width solution \hat{f}_n ?

- ▶ For $p = 2$, we have $\|\hat{f}_{\text{RF},M,n} - \hat{f}_n\|_{L^2} \approx 0$ when $M \geq n^{1+\delta}$ and $\|\hat{f}_{\text{RF},M,n} - \hat{f}_n\|_{L^2} > 0$ when $M \leq n^{1-\delta}$ [MMM,'21].

Random Features approximation

- ▶ Approximate μ by finitely supported $\hat{\mu}_M$: sample M weights $\mathbf{w}_j \sim_{iid} \mu$,

$$f(\mathbf{x}; \mathbf{a}) = \int_{\Omega} a(\mathbf{w})\phi(\mathbf{x}; \mathbf{w})\mu(d\mathbf{w}) \longrightarrow f_M(\mathbf{x}; \mathbf{a}) = \frac{1}{M} \sum_{j=1}^M a_j\phi(\mathbf{x}; \mathbf{w}_j).$$

- ▶ 'Finite-width' problem is easy to solve: $\mathbf{a} \in \mathbb{R}^M$,

$$\text{minimize } \frac{1}{M} \sum_{j=1}^M |a_j|^p,$$

$$\text{subj. to } f_M(\mathbf{x}_i; \mathbf{a}) = y_i, \quad \forall i \leq n.$$

How large M needs to be for the finite-width solution $\hat{f}_{\text{RF},M,n}$ to approximate the infinite-width solution \hat{f}_n ?

- ▶ For $p = 2$, we have $\|\hat{f}_{\text{RF},M,n} - \hat{f}_n\|_{L^2} \approx 0$ when $M \geq n^{1+\delta}$ and $\|\hat{f}_{\text{RF},M,n} - \hat{f}_n\|_{L^2} > 0$ when $M \leq n^{1-\delta}$ [MMM,'21].

Random Features approximation

- ▶ Approximate μ by finitely supported $\hat{\mu}_M$: sample M weights $\mathbf{w}_j \sim_{iid} \mu$,

$$f(\mathbf{x}; \mathbf{a}) = \int_{\Omega} a(\mathbf{w})\phi(\mathbf{x}; \mathbf{w})\mu(d\mathbf{w}) \longrightarrow f_M(\mathbf{x}; \mathbf{a}) = \frac{1}{M} \sum_{j=1}^M a_j\phi(\mathbf{x}; \mathbf{w}_j).$$

- ▶ 'Finite-width' problem is easy to solve: $\mathbf{a} \in \mathbb{R}^M$,

$$\text{minimize } \frac{1}{M} \sum_{j=1}^M |a_j|^p,$$

$$\text{subj. to } f_M(\mathbf{x}_i; \mathbf{a}) = y_i, \quad \forall i \leq n.$$

How large M needs to be for the finite-width solution $\hat{f}_{\text{RF},M,n}$ to approximate the infinite-width solution \hat{f}_n ?

- ▶ For $p = 2$, we have $\|\hat{f}_{\text{RF},M,n} - \hat{f}_n\|_{L^2} \approx 0$ when $M \geq n^{1+\delta}$ and $\|\hat{f}_{\text{RF},M,n} - \hat{f}_n\|_{L^2} > 0$ when $M \leq n^{1-\delta}$ [MMM,'21].

Random Features approximation

- ▶ Approximate μ by finitely supported $\hat{\mu}_M$: sample M weights $\mathbf{w}_j \sim_{iid} \mu$,

$$f(\mathbf{x}; \mathbf{a}) = \int_{\Omega} a(\mathbf{w})\phi(\mathbf{x}; \mathbf{w})\mu(d\mathbf{w}) \longrightarrow f_M(\mathbf{x}; \mathbf{a}) = \frac{1}{M} \sum_{j=1}^M a_j\phi(\mathbf{x}; \mathbf{w}_j).$$

- ▶ 'Finite-width' problem is easy to solve: $\mathbf{a} \in \mathbb{R}^M$,

$$\text{minimize } \frac{1}{M} \sum_{j=1}^M |a_j|^p,$$

$$\text{subj. to } f_M(\mathbf{x}_i; \mathbf{a}) = y_i, \quad \forall i \leq n.$$

How large M needs to be for the finite-width solution $\hat{f}_{\text{RF},M,n}$ to approximate the infinite-width solution \hat{f}_n ?

- ▶ For $p = 2$, we have $\|\hat{f}_{\text{RF},M,n} - \hat{f}_n\|_{L^2} \approx 0$ when $M \geq n^{1+\delta}$ and $\|\hat{f}_{\text{RF},M,n} - \hat{f}_n\|_{L^2} > 0$ when $M \leq n^{1-\delta}$ [MMM,'21].

- ▶ **Data:** $\{y_i, \mathbf{x}_i\}_{i \leq n}$, bounded support: $\|\mathbf{x}\|_2 \leq C\sqrt{d}$.
- ▶ **Kernel matrix:** $K_n = (K(\mathbf{x}_i, \mathbf{x}_j))_{i,j \leq n}$, $K(\mathbf{x}_i, \mathbf{x}_j) = \int_{\Omega} \phi(\mathbf{x}_i; \mathbf{w})\phi(\mathbf{x}_j; \mathbf{w})\mu(d\mathbf{w})$.
- ▶ **Random Features vectors:** $\{\mathbf{w}_j\}_{j \leq M}$ fixed iid from μ ,
$$\phi_{n,j} = \phi_n(\mathbf{w}_j) := [\phi(\mathbf{x}_1; \mathbf{w}_j), \dots, \phi(\mathbf{x}_n; \mathbf{w}_j)] \in \mathbb{R}^n$$
.
- ▶ **Whitened features:** $\psi_n(\mathbf{w}) := K_n^{-1/2} \phi_n(\mathbf{w})$ (so that $\mathbb{E}_{\mathbf{w}}[\psi_n(\mathbf{w})\psi_n(\mathbf{w})^T] = I_n$).

Result will hold conditionally on realization of the data $\{y_i, \mathbf{x}_i\}$, and exploit randomness of weights $\{\mathbf{w}_j\}$ to show concentration.

E.g., $\phi_{n,j}$ conditional on \mathbf{X} : iid random vectors in \mathbb{R}^n .

- ▶ **Data:** $\{y_i, \mathbf{x}_i\}_{i \leq n}$, bounded support: $\|\mathbf{x}\|_2 \leq C\sqrt{d}$.
- ▶ **Kernel matrix:** $\mathbf{K}_n = (K(\mathbf{x}_i, \mathbf{x}_j))_{i,j \leq n}$, $K(\mathbf{x}_i, \mathbf{x}_j) = \int_{\Omega} \phi(\mathbf{x}_i; \mathbf{w})\phi(\mathbf{x}_j; \mathbf{w})\mu(d\mathbf{w})$.

- ▶ **Random Features vectors:** $\{\mathbf{w}_j\}_{j \leq M}$ fixed iid from μ ,

$$\phi_{n,j} = \phi_n(\mathbf{w}_j) := [\phi(\mathbf{x}_1; \mathbf{w}_j), \dots, \phi(\mathbf{x}_n; \mathbf{w}_j)] \in \mathbb{R}^n.$$

- ▶ **Whitened features:** $\psi_n(\mathbf{w}) := \mathbf{K}_n^{-1/2} \phi_n(\mathbf{w})$ (so that $\mathbb{E}_{\mathbf{w}}[\psi_n(\mathbf{w})\psi_n(\mathbf{w})^T] = \mathbf{I}_n$).

Result will hold conditionally on realization of the data $\{y_i, \mathbf{x}_i\}$, and exploit randomness of weights $\{\mathbf{w}_j\}$ to show concentration.

E.g., $\phi_{n,j}$ conditional on \mathbf{X} : iid random vectors in \mathbb{R}^n .

- ▶ **Data:** $\{y_i, \mathbf{x}_i\}_{i \leq n}$, bounded support: $\|\mathbf{x}\|_2 \leq C\sqrt{d}$.
- ▶ **Kernel matrix:** $\mathbf{K}_n = (K(\mathbf{x}_i, \mathbf{x}_j))_{i,j \leq n}$, $K(\mathbf{x}_i, \mathbf{x}_j) = \int_{\Omega} \phi(\mathbf{x}_i; \mathbf{w})\phi(\mathbf{x}_j; \mathbf{w})\mu(d\mathbf{w})$.

- ▶ **Random Features vectors:** $\{\mathbf{w}_j\}_{j \leq M}$ fixed iid from μ ,

$$\phi_{n,j} = \phi_n(\mathbf{w}_j) := [\phi(\mathbf{x}_1; \mathbf{w}_j), \dots, \phi(\mathbf{x}_n; \mathbf{w}_j)] \in \mathbb{R}^n.$$

- ▶ **Whitened features:** $\psi_n(\mathbf{w}) := \mathbf{K}_n^{-1/2} \phi_n(\mathbf{w})$ (so that $\mathbb{E}_{\mathbf{w}}[\psi_n(\mathbf{w})\psi_n(\mathbf{w})^T] = \mathbf{I}_n$).

Result will hold conditionally on realization of the data $\{y_i, \mathbf{x}_i\}$, and exploit randomness of weights $\{\mathbf{w}_j\}$ to show concentration.

E.g., $\phi_{n,j}$ conditional on \mathbf{X} : iid random vectors in \mathbb{R}^n .

- ▶ **Data:** $\{y_i, \mathbf{x}_i\}_{i \leq n}$, bounded support: $\|\mathbf{x}\|_2 \leq C\sqrt{d}$.
- ▶ **Kernel matrix:** $\mathbf{K}_n = (K(\mathbf{x}_i, \mathbf{x}_j))_{i,j \leq n}$, $K(\mathbf{x}_i, \mathbf{x}_j) = \int_{\Omega} \phi(\mathbf{x}_i; \mathbf{w})\phi(\mathbf{x}_j; \mathbf{w})\mu(d\mathbf{w})$.
- ▶ **Random Features vectors:** $\{\mathbf{w}_j\}_{j \leq M}$ fixed iid from μ ,
$$\phi_{n,j} = \phi_n(\mathbf{w}_j) := [\phi(\mathbf{x}_1; \mathbf{w}_j), \dots, \phi(\mathbf{x}_n; \mathbf{w}_j)] \in \mathbb{R}^n$$
.
- ▶ **Whitened features:** $\psi_n(\mathbf{w}) := \mathbf{K}_n^{-1/2} \phi_n(\mathbf{w})$ (so that $\mathbb{E}_{\mathbf{w}}[\psi_n(\mathbf{w})\psi_n(\mathbf{w})^T] = \mathbf{I}_n$).

Result will hold conditionally on realization of the data $\{y_i, \mathbf{x}_i\}$, and exploit randomness of weights $\{\mathbf{w}_j\}$ to show concentration.

E.g., $\phi_{n,j}$ conditional on \mathbf{X} : iid random vectors in \mathbb{R}^n .

- ▶ **Data:** $\{y_i, \mathbf{x}_i\}_{i \leq n}$, bounded support: $\|\mathbf{x}\|_2 \leq C\sqrt{d}$.
- ▶ **Kernel matrix:** $\mathbf{K}_n = (K(\mathbf{x}_i, \mathbf{x}_j))_{i,j \leq n}$, $K(\mathbf{x}_i, \mathbf{x}_j) = \int_{\Omega} \phi(\mathbf{x}_i; \mathbf{w})\phi(\mathbf{x}_j; \mathbf{w})\mu(d\mathbf{w})$.

- ▶ **Random Features vectors:** $\{\mathbf{w}_j\}_{j \leq M}$ fixed iid from μ ,

$$\phi_{n,j} = \phi_n(\mathbf{w}_j) := [\phi(\mathbf{x}_1; \mathbf{w}_j), \dots, \phi(\mathbf{x}_n; \mathbf{w}_j)] \in \mathbb{R}^n.$$

- ▶ **Whitened features:** $\psi_n(\mathbf{w}) := \mathbf{K}_n^{-1/2} \phi_n(\mathbf{w})$ (so that $\mathbb{E}_{\mathbf{w}}[\psi_n(\mathbf{w})\psi_n(\mathbf{w})^T] = \mathbf{I}_n$).

Result will hold conditionally on realization of the data $\{y_i, \mathbf{x}_i\}$, and exploit randomness of weights $\{\mathbf{w}_j\}$ to show concentration.

E.g., $\phi_{n,j}$ conditional on \mathbf{X} : iid random vectors in \mathbb{R}^n .

Assumptions

Conditionally on \mathbf{y}, \mathbf{X} :

A1 [Sub-gaussianity] For any $\|\mathbf{x}\|_2 \leq C\sqrt{d}$, $\phi(\mathbf{x}; \mathbf{w})$ is τ^2 -sub-Gaussian when $\mathbf{w} \sim \mu$.
Whitened feature vector $\psi_n(\mathbf{w})$ is τ^2 -sub-Gaussian when $\mathbf{w} \sim \mu$.

A2 [Lipschitz continuity] Feature $\phi(\mathbf{x}; \mathbf{w})$ is $L(\mathbf{w})$ -Lipschitz w.r.t \mathbf{x} and $L(\mathbf{w})$ is τ^2 -sub-Gaussian when $\mathbf{w} \sim \mu$.

A3 [Small ball property] There exist $\eta, c > 0$ such that

$$\inf_{\|\mathbf{u}\|_2 = \|\mathbf{v}\|_2 = 1} \mathbb{P}(|\langle \mathbf{u}, \psi_n(\mathbf{w}) \rangle| \geq \eta, |\langle \mathbf{v}, \psi_n(\mathbf{w}) \rangle| \geq \eta) \geq c.$$

Assumptions A1 and A3 can be hard to check in practice.

Assumptions

Conditionally on \mathbf{y}, \mathbf{X} :

A1 [Sub-gaussianity] For any $\|\mathbf{x}\|_2 \leq C\sqrt{d}$, $\phi(\mathbf{x}; \mathbf{w})$ is τ^2 -sub-Gaussian when $\mathbf{w} \sim \mu$.
Whitened feature vector $\psi_n(\mathbf{w})$ is τ^2 -sub-Gaussian when $\mathbf{w} \sim \mu$.

A2 [Lipschitz continuity] Feature $\phi(\mathbf{x}; \mathbf{w})$ is $L(\mathbf{w})$ -Lipschitz w.r.t \mathbf{x} and $L(\mathbf{w})$ is τ^2 -sub-Gaussian when $\mathbf{w} \sim \mu$.

A3 [Small ball property] There exist $\eta, c > 0$ such that

$$\inf_{\|\mathbf{u}\|_2 = \|\mathbf{v}\|_2 = 1} \mathbb{P}(|\langle \mathbf{u}, \psi_n(\mathbf{w}) \rangle| \geq \eta, |\langle \mathbf{v}, \psi_n(\mathbf{w}) \rangle| \geq \eta) \geq c.$$

Assumptions A1 and A3 can be hard to check in practice.

Assumptions

Conditionally on \mathbf{y}, \mathbf{X} :

A1 [Sub-gaussianity] For any $\|\mathbf{x}\|_2 \leq C\sqrt{d}$, $\phi(\mathbf{x}; \mathbf{w})$ is τ^2 -sub-Gaussian when $\mathbf{w} \sim \mu$.
Whitened feature vector $\psi_n(\mathbf{w})$ is τ^2 -sub-Gaussian when $\mathbf{w} \sim \mu$.

A2 [Lipschitz continuity] Feature $\phi(\mathbf{x}; \mathbf{w})$ is $L(\mathbf{w})$ -Lipschitz w.r.t \mathbf{x} and $L(\mathbf{w})$ is τ^2 -sub-Gaussian when $\mathbf{w} \sim \mu$.

A3 [Small ball property] There exist $\eta, c > 0$ such that

$$\inf_{\|\mathbf{u}\|_2 = \|\mathbf{v}\|_2 = 1} \mathbb{P}(|\langle \mathbf{u}, \psi_n(\mathbf{w}) \rangle| \geq \eta, |\langle \mathbf{v}, \psi_n(\mathbf{w}) \rangle| \geq \eta) \geq c.$$

Assumptions A1 and A3 can be hard to check in practice.

Assumptions

Conditionally on \mathbf{y}, \mathbf{X} :

A1 [Sub-gaussianity] For any $\|\mathbf{x}\|_2 \leq C\sqrt{d}$, $\phi(\mathbf{x}; \mathbf{w})$ is τ^2 -sub-Gaussian when $\mathbf{w} \sim \mu$.

Whitened feature vector $\psi_n(\mathbf{w})$ is τ^2 -sub-Gaussian when $\mathbf{w} \sim \mu$.

A2 [Lipschitz continuity] Feature $\phi(\mathbf{x}; \mathbf{w})$ is $L(\mathbf{w})$ -Lipschitz w.r.t \mathbf{x} and $L(\mathbf{w})$ is τ^2 -sub-Gaussian when $\mathbf{w} \sim \mu$.

A3 [Small ball property] There exist $\eta, c > 0$ such that

$$\inf_{\|\mathbf{u}\|_2 = \|\mathbf{v}\|_2 = 1} \mathbb{P}(|\langle \mathbf{u}, \psi_n(\mathbf{w}) \rangle| \geq \eta, |\langle \mathbf{v}, \psi_n(\mathbf{w}) \rangle| \geq \eta) \geq c.$$

Assumptions A1 and A3 can be hard to check in practice.

Theorem (Celentano, Misiakiewicz, Montanari, 2021)

For $p \in (1, 2]$. Assume A1, A2 and A3 hold conditionally on \mathbf{y}, \mathbf{X} . Then with probability at least $1 - CM^{-cn}$,

$$\|\hat{f}_{\text{RF}, M, n} - \hat{f}_n\|_{L^2}^2 \leq C \left(\frac{n \log M}{M} \vee \frac{(n \log M)^{p/(p-1)}}{M^2} \right) \|\mathbf{K}_n^{-1/2} \mathbf{y}\|_2^2.$$

- ▶ Typically, $\|\mathbf{K}_n^{-1/2} \mathbf{y}\|_2 \leq C\sqrt{n}$, hence we need $M \geq (n \log(n))^{2\sqrt{\frac{p-1}{p-1/2}}}$.
- ▶ Bound is not optimal: e.g., $p = 2$, we expect $M \geq n \log(n)$ to be sufficient.
- ▶ Bound diverges as $p \rightarrow 1$: we can't solve efficiently the infinite width problem \mathcal{F}_1 with the RF approach.

E.g., learning a single neuron:

[Bach, '17] can be learned with \mathcal{F}_1 with $n \leq d^{O(1)}$.

[GMMM, '19] need $M \geq e^{\Theta(d)}$ features to be approximated.

Theorem (Celentano, Misiakiewicz, Montanari, 2021)

For $p \in (1, 2]$. Assume A1, A2 and A3 hold conditionally on \mathbf{y}, \mathbf{X} . Then with probability at least $1 - CM^{-cn}$,

$$\|\hat{f}_{\text{RF}, M, n} - \hat{f}_n\|_{L^2}^2 \leq C \left(\frac{n \log M}{M} \vee \frac{(n \log M)^{p/(p-1)}}{M^2} \right) \|\mathbf{K}_n^{-1/2} \mathbf{y}\|_2^2.$$

- ▶ Typically, $\|\mathbf{K}_n^{-1/2} \mathbf{y}\|_2 \leq C\sqrt{n}$, hence we need $M \geq (n \log(n))^{2\vee\left(\frac{p-\frac{1}{2}}{p-1}\right)}$.
- ▶ Bound is not optimal: e.g., $p = 2$, we expect $M \geq n \log(n)$ to be sufficient.
- ▶ Bound diverges as $p \rightarrow 1$: we can't solve efficiently the infinite width problem \mathcal{F}_1 with the RF approach.

E.g., learning a single neuron:

[Bach, '17] can be learned with \mathcal{F}_1 with $n \leq d^{O(1)}$.

[GMMM, '19] need $M \geq e^{\Theta(d)}$ features to be approximated.

Theorem (Celentano, Misiakiewicz, Montanari, 2021)

For $p \in (1, 2]$. Assume A1, A2 and A3 hold conditionally on \mathbf{y}, \mathbf{X} . Then with probability at least $1 - CM^{-cn}$,

$$\|\hat{f}_{\text{RF}, M, n} - \hat{f}_n\|_{L^2}^2 \leq C \left(\frac{n \log M}{M} \vee \frac{(n \log M)^{p/(p-1)}}{M^2} \right) \|\mathbf{K}_n^{-1/2} \mathbf{y}\|_2^2.$$

- ▶ Typically, $\|\mathbf{K}_n^{-1/2} \mathbf{y}\|_2 \leq C\sqrt{n}$, hence we need $M \geq (n \log(n))^{2\sqrt{\frac{p-\frac{1}{2}}{p-1}}}$.
- ▶ Bound is not optimal: e.g., $p = 2$, we expect $M \geq n \log(n)$ to be sufficient.
- ▶ Bound diverges as $p \rightarrow 1$: we can't solve efficiently the infinite width problem \mathcal{F}_1 with the RF approach.

E.g., learning a single neuron:

[Bach, '17] can be learned with \mathcal{F}_1 with $n \leq d^{O(1)}$.

[GMMM, '19] need $M \geq e^{\Theta(d)}$ features to be approximated.

Theorem (Celentano, Misiakiewicz, Montanari, 2021)

For $p \in (1, 2]$. Assume A1, A2 and A3 hold conditionally on \mathbf{y}, \mathbf{X} . Then with probability at least $1 - CM^{-cn}$,

$$\|\hat{f}_{\text{RF}, M, n} - \hat{f}_n\|_{L^2}^2 \leq C \left(\frac{n \log M}{M} \vee \frac{(n \log M)^{p/(p-1)}}{M^2} \right) \|\mathbf{K}_n^{-1/2} \mathbf{y}\|_2^2.$$

- ▶ Typically, $\|\mathbf{K}_n^{-1/2} \mathbf{y}\|_2 \leq C\sqrt{n}$, hence we need $M \geq (n \log(n))^{2\sqrt{\frac{p-\frac{1}{2}}{p-1}}}$.
- ▶ Bound is not optimal: e.g., $p = 2$, we expect $M \geq n \log(n)$ to be sufficient.
- ▶ Bound diverges as $p \rightarrow 1$: we can't solve efficiently the infinite width problem \mathcal{F}_1 with the RF approach.

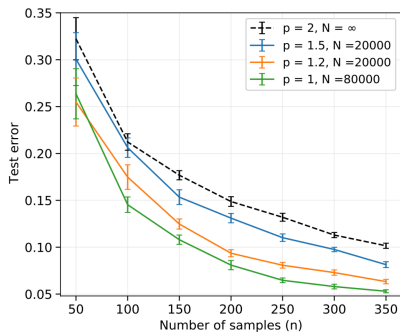
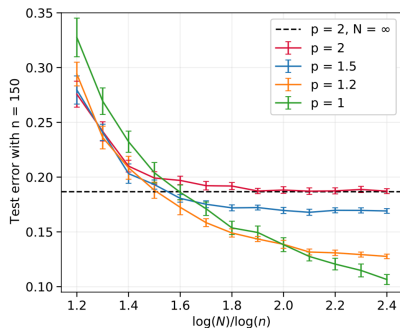
E.g., learning a single neuron:

[Bach, '17] can be learned with \mathcal{F}_1 with $n \leq d^{O(1)}$.

[GMMM, '19] need $M \geq e^{\Theta(d)}$ features to be approximated.

Illustration

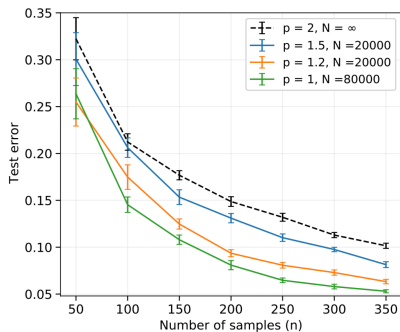
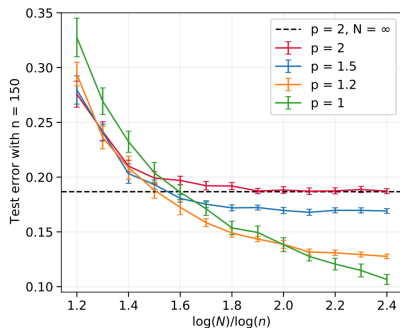
Learning single neuron: $\sigma(\langle \mathbf{x}, \mathbf{w}_* \rangle)$, $\mathbf{x} \sim \text{Unif}(\mathbb{S}^{d-1}(\sqrt{d}))$, $\mathbf{w} \sim N(0, I_d/d)$, $d = 30$.



- ▶ Left: test error settles on 'infinite-width' solution error when $M \geq M_*(n, p)$.
- ▶ Left: $M_*(n, p)$ increases as p decreases ($p = 1$, unable to reach $M = \infty$ error).
- ▶ Right: test error decreases when p decreases. $\mathcal{F}_2(R) \subset \mathcal{F}_p(R) \subset \mathcal{F}_1(R)$ capture better and better functions highly dependent on low-dimensional projection of \mathbf{x} .

Illustration

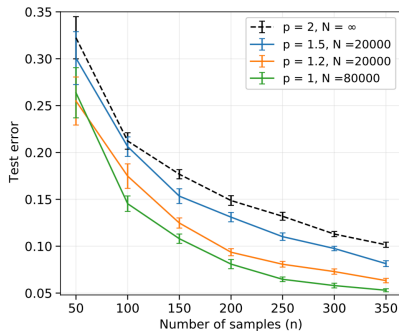
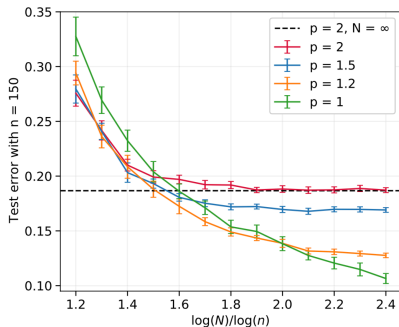
Learning single neuron: $\sigma(\langle \mathbf{x}, \mathbf{w}_* \rangle)$, $\mathbf{x} \sim \text{Unif}(\mathbb{S}^{d-1}(\sqrt{d}))$, $\mathbf{w} \sim N(0, I_d/d)$, $d = 30$.



- ▶ **Left:** test error settles on 'infinite-width' solution error when $M \geq M_*(n, p)$.
- ▶ **Left:** $M_*(n, p)$ increases as p decreases ($p = 1$, unable to reach $M = \infty$ error).
- ▶ **Right:** test error decreases when p decreases. $\mathcal{F}_2(R) \subset \mathcal{F}_p(R) \subset \mathcal{F}_1(R)$ capture better and better functions highly dependent on low-dimensional projection of \mathbf{x} .

Illustration

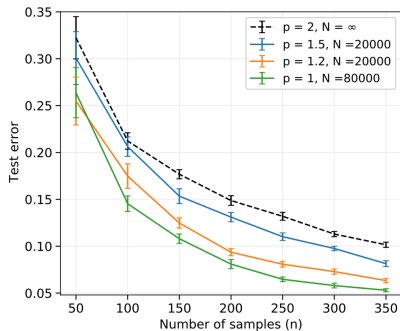
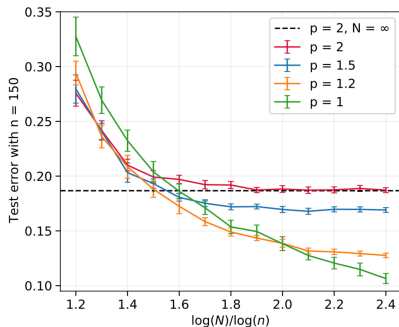
Learning single neuron: $\sigma(\langle \mathbf{x}, \mathbf{w}_* \rangle)$, $\mathbf{x} \sim \text{Unif}(\mathbb{S}^{d-1}(\sqrt{d}))$, $\mathbf{w} \sim N(0, I_d/d)$, $d = 30$.



- ▶ **Left:** test error settles on 'infinite-width' solution error when $M \geq M_*(n, p)$.
- ▶ **Left:** $M_*(n, p)$ increases as p decreases ($p = 1$, unable to reach $M = \infty$ error).
- ▶ **Right:** test error decreases when p decreases. $\mathcal{F}_2(R) \subset \mathcal{F}_p(R) \subset \mathcal{F}_1(R)$ capture better and better functions highly dependent on low-dimensional projection of \mathbf{x} .

Illustration

Learning single neuron: $\sigma(\langle \mathbf{x}, \mathbf{w}_* \rangle)$, $\mathbf{x} \sim \text{Unif}(\mathbb{S}^{d-1}(\sqrt{d}))$, $\mathbf{w} \sim N(0, I_d/d)$, $d = 30$.



- ▶ **Left:** test error settles on 'infinite-width' solution error when $M \geq M_*(n, p)$.
- ▶ **Left:** $M_*(n, p)$ increases as p decreases ($p = 1$, unable to reach $M = \infty$ error).
- ▶ **Right:** test error decreases when p decreases. $\mathcal{F}_2(R) \subset \mathcal{F}_p(R) \subset \mathcal{F}_1(R)$ capture better and better functions highly dependent on low-dimensional projection of \mathbf{x} .

Idea of the proof (I): dual problem

$$\text{Infinite-width: } \min_{a: \Omega \rightarrow \mathbb{R}} \left\{ \int_{\Omega} \rho(a(\mathbf{w})) \mu(d\mathbf{w}) : \forall i \leq n, \int_{\Omega} a(\mathbf{w}) \phi(\mathbf{x}_i; \mathbf{w}) \mu(d\mathbf{w}) = y_i \right\},$$

$$\text{Finite-width: } \min_{\mathbf{a} \in \mathbb{R}^M} \left\{ \frac{1}{M} \sum_{j \leq M} \rho(a_j) : \forall i \leq n, \frac{1}{M} \sum_{j \leq M} a_j \phi(\mathbf{x}_i; \mathbf{w}_j) = y_i \right\}.$$

Dual problems: $\rho(x) = \frac{1}{p}|x|^p$, $\rho^* = \frac{1}{q}|x|^q$, $q = p/(p-1)$.

$$\text{Infinite-width: } \max_{\lambda \in \mathbb{R}^n} F(\lambda) := \langle \mathbf{y}, \lambda \rangle - \int_{\Omega} \rho^*(\langle \phi_n(\mathbf{w}), \lambda \rangle) \mu(d\mathbf{w}),$$

$$\text{Finite-width: } \max_{\lambda \in \mathbb{R}^n} F_M(\lambda) := \langle \mathbf{y}, \lambda \rangle - \frac{1}{M} \sum_{j \leq M} \rho^*(\langle \phi_{n,j}, \lambda \rangle).$$

► **Primal:** overparametrized, difficult to control.

Dual: underparametrized, can use uniform convergence argument.

► ‘Representer’ theorem for $p \in (1, 2]$: with $s(x) = (\rho^*)'(x) = \text{sign}(x)|x|^{q-1}$, $\lambda \in \mathbb{R}^n$

$$\hat{f}_n(x; \lambda) = \int_{\Omega} s(\langle \phi_n(\mathbf{w}), \lambda \rangle) \phi(x; \mathbf{w}) \mu(d\mathbf{w}),$$

$$\hat{f}_{\text{RF}, M, n}(x; \lambda) = \frac{1}{M} \sum_{j \leq M} s(\langle \phi_n(\mathbf{w}_j), \lambda \rangle) \phi(x; \mathbf{w}_j).$$

Idea of the proof (I): dual problem

$$\text{Infinite-width: } \min_{a: \Omega \rightarrow \mathbb{R}} \left\{ \int_{\Omega} \rho(a(\mathbf{w})) \mu(d\mathbf{w}) : \forall i \leq n, \int_{\Omega} a(\mathbf{w}) \phi(\mathbf{x}_i; \mathbf{w}) \mu(d\mathbf{w}) = y_i \right\},$$

$$\text{Finite-width: } \min_{\mathbf{a} \in \mathbb{R}^M} \left\{ \frac{1}{M} \sum_{j \leq M} \rho(a_j) : \forall i \leq n, \frac{1}{M} \sum_{j \leq M} a_j \phi(\mathbf{x}_i; \mathbf{w}_j) = y_i \right\}.$$

Dual problems: $\rho(x) = \frac{1}{p}|x|^p$, $\rho^* = \frac{1}{q}|x|^q$, $q = p/(p-1)$.

$$\text{Infinite-width: } \max_{\boldsymbol{\lambda} \in \mathbb{R}^n} F(\boldsymbol{\lambda}) := \langle \mathbf{y}, \boldsymbol{\lambda} \rangle - \int_{\Omega} \rho^*(\langle \phi_n(\mathbf{w}), \boldsymbol{\lambda} \rangle) \mu(d\mathbf{w}),$$

$$\text{Finite-width: } \max_{\boldsymbol{\lambda} \in \mathbb{R}^n} F_M(\boldsymbol{\lambda}) := \langle \mathbf{y}, \boldsymbol{\lambda} \rangle - \frac{1}{M} \sum_{j \leq M} \rho^*(\langle \phi_{n,j}, \boldsymbol{\lambda} \rangle).$$

► **Primal:** overparametrized, difficult to control.

Dual: underparametrized, can use uniform convergence argument.

► ‘Representer’ theorem for $p \in (1, 2]$: with $s(x) = (\rho^*)'(x) = \text{sign}(x)|x|^{q-1}$, $\boldsymbol{\lambda} \in \mathbb{R}^n$

$$\hat{f}_n(x; \boldsymbol{\lambda}) = \int_{\Omega} s(\langle \phi_n(\mathbf{w}), \boldsymbol{\lambda} \rangle) \phi(x; \mathbf{w}) \mu(d\mathbf{w}),$$

$$\hat{f}_{\text{RF}, M, n}(x; \boldsymbol{\lambda}) = \frac{1}{M} \sum_{j \leq M} s(\langle \phi_n(\mathbf{w}_j), \boldsymbol{\lambda} \rangle) \phi(x; \mathbf{w}_j).$$

Idea of the proof (I): dual problem

$$\text{Infinite-width: } \min_{a: \Omega \rightarrow \mathbb{R}} \left\{ \int_{\Omega} \rho(a(\mathbf{w})) \mu(d\mathbf{w}) : \forall i \leq n, \int_{\Omega} a(\mathbf{w}) \phi(\mathbf{x}_i; \mathbf{w}) \mu(d\mathbf{w}) = y_i \right\},$$

$$\text{Finite-width: } \min_{\mathbf{a} \in \mathbb{R}^M} \left\{ \frac{1}{M} \sum_{j \leq M} \rho(a_j) : \forall i \leq n, \frac{1}{M} \sum_{j \leq M} a_j \phi(\mathbf{x}_i; \mathbf{w}_j) = y_i \right\}.$$

Dual problems: $\rho(x) = \frac{1}{p}|x|^p$, $\rho^* = \frac{1}{q}|x|^q$, $q = p/(p-1)$.

$$\text{Infinite-width: } \max_{\boldsymbol{\lambda} \in \mathbb{R}^n} F(\boldsymbol{\lambda}) := \langle \mathbf{y}, \boldsymbol{\lambda} \rangle - \int_{\Omega} \rho^*(\langle \phi_n(\mathbf{w}), \boldsymbol{\lambda} \rangle) \mu(d\mathbf{w}),$$

$$\text{Finite-width: } \max_{\boldsymbol{\lambda} \in \mathbb{R}^n} F_M(\boldsymbol{\lambda}) := \langle \mathbf{y}, \boldsymbol{\lambda} \rangle - \frac{1}{M} \sum_{j \leq M} \rho^*(\langle \phi_{n,j}, \boldsymbol{\lambda} \rangle).$$

► **Primal:** overparametrized, difficult to control.

Dual: underparametrized, can use uniform convergence argument.

► ‘Representer’ theorem for $p \in (1, 2]$: with $s(x) = (\rho^*)'(x) = \text{sign}(x)|x|^{q-1}$, $\boldsymbol{\lambda} \in \mathbb{R}^n$

$$\hat{f}_n(x; \boldsymbol{\lambda}) = \int_{\Omega} s(\langle \phi_n(\mathbf{w}), \boldsymbol{\lambda} \rangle) \phi(x; \mathbf{w}) \mu(d\mathbf{w}),$$

$$\hat{f}_{\text{RF}, M, n}(x; \boldsymbol{\lambda}) = \frac{1}{M} \sum_{j \leq M} s(\langle \phi_n(\mathbf{w}_j), \boldsymbol{\lambda} \rangle) \phi(x; \mathbf{w}_j).$$

Idea of the proof (I): dual problem

$$\text{Infinite-width: } \min_{a: \Omega \rightarrow \mathbb{R}} \left\{ \int_{\Omega} \rho(a(\mathbf{w})) \mu(d\mathbf{w}) : \forall i \leq n, \int_{\Omega} a(\mathbf{w}) \phi(\mathbf{x}_i; \mathbf{w}) \mu(d\mathbf{w}) = y_i \right\},$$

$$\text{Finite-width: } \min_{\mathbf{a} \in \mathbb{R}^M} \left\{ \frac{1}{M} \sum_{j \leq M} \rho(a_j) : \forall i \leq n, \frac{1}{M} \sum_{j \leq M} a_j \phi(\mathbf{x}_i; \mathbf{w}_j) = y_i \right\}.$$

Dual problems: $\rho(x) = \frac{1}{p}|x|^p$, $\rho^* = \frac{1}{q}|x|^q$, $q = p/(p-1)$.

$$\text{Infinite-width: } \max_{\boldsymbol{\lambda} \in \mathbb{R}^n} F(\boldsymbol{\lambda}) := \langle \mathbf{y}, \boldsymbol{\lambda} \rangle - \int_{\Omega} \rho^*(\langle \phi_n(\mathbf{w}), \boldsymbol{\lambda} \rangle) \mu(d\mathbf{w}),$$

$$\text{Finite-width: } \max_{\boldsymbol{\lambda} \in \mathbb{R}^n} F_M(\boldsymbol{\lambda}) := \langle \mathbf{y}, \boldsymbol{\lambda} \rangle - \frac{1}{M} \sum_{j \leq M} \rho^*(\langle \phi_{n,j}, \boldsymbol{\lambda} \rangle).$$

► **Primal:** overparametrized, difficult to control.

Dual: underparametrized, can use uniform convergence argument.

► ‘Representer’ theorem for $p \in (1, 2]$: with $s(x) = (\rho^*)'(x) = \text{sign}(x)|x|^{q-1}$, $\boldsymbol{\lambda} \in \mathbb{R}^n$

$$\hat{f}_n(\mathbf{x}; \boldsymbol{\lambda}) = \int_{\Omega} s(\langle \phi_n(\mathbf{w}), \boldsymbol{\lambda} \rangle) \phi(\mathbf{x}; \mathbf{w}) \mu(d\mathbf{w}),$$

$$\hat{f}_{\text{RF}, M, n}(\mathbf{x}; \boldsymbol{\lambda}) = \frac{1}{M} \sum_{i < M} s(\langle \phi_n(\mathbf{w}_i), \boldsymbol{\lambda} \rangle) \phi(\mathbf{x}; \mathbf{w}_i).$$

Idea of the proof (II): uniform convergence

$$\hat{\lambda} = \arg \max_{\lambda \in \mathbb{R}^n} F(\lambda) := \langle \mathbf{y}, \lambda \rangle - \int_{\Omega} \rho^*(\langle \phi_n(\mathbf{w}), \lambda \rangle) \mu(d\mathbf{w}),$$

$$\hat{\lambda}_M = \arg \max_{\lambda \in \mathbb{R}^n} F_M(\lambda) := \langle \mathbf{y}, \lambda \rangle - \frac{1}{M} \sum_{j \leq M} \rho^*(\langle \phi_{n,j}, \lambda \rangle).$$

- ▶ Concentration of the landscape uniformly over balls in λ :

$$\nabla F_M(\lambda) \rightarrow \nabla F(\lambda), \quad \lambda_{\max}(\nabla^2 F_M(\lambda)) \rightarrow \lambda_{\max}(\nabla^2 F(\lambda)),$$

which shows $\|\hat{\lambda}_M - \hat{\lambda}\|_2 \leq \varepsilon_1(n, M, \rho)$.

- ▶ Uniform concentration of the predictor over balls in λ :

$$\max_{\lambda \in B} \|\hat{f}_{\text{RF},M,n}(\cdot; \lambda) - \hat{f}_n(\cdot; \lambda)\|_{L^2} \leq \varepsilon_2(n, M, \rho).$$

- ▶ Combining the two + Lipschitzness of \hat{f}_n w.r.t λ :

$$\begin{aligned} \|\hat{f}_{\text{RF},M,n}(\cdot; \hat{\lambda}_M) - \hat{f}_n(\cdot; \hat{\lambda})\|_{L^2} &\leq \|\hat{f}_{\text{RF},M,n}(\cdot; \hat{\lambda}_M) - \hat{f}_n(\cdot; \hat{\lambda}_M)\|_{L^2} + \|\hat{f}_n(\cdot; \hat{\lambda}_M) - \hat{f}_n(\cdot; \hat{\lambda})\|_{L^2} \\ &\leq \varepsilon_2(n, M, \rho) + \varepsilon_1(n, M, \rho). \end{aligned}$$

Idea of the proof (II): uniform convergence

$$\hat{\lambda} = \arg \max_{\lambda \in \mathbb{R}^n} F(\lambda) := \langle \mathbf{y}, \lambda \rangle - \int_{\Omega} \rho^*(\langle \phi_n(\mathbf{w}), \lambda \rangle) \mu(d\mathbf{w}),$$

$$\hat{\lambda}_M = \arg \max_{\lambda \in \mathbb{R}^n} F_M(\lambda) := \langle \mathbf{y}, \lambda \rangle - \frac{1}{M} \sum_{j \leq M} \rho^*(\langle \phi_{n,j}, \lambda \rangle).$$

- ▶ Concentration of the landscape uniformly over balls in λ :

$$\nabla F_M(\lambda) \rightarrow \nabla F(\lambda), \quad \lambda_{\max}(\nabla^2 F_M(\lambda)) \rightarrow \lambda_{\max}(\nabla^2 F(\lambda)),$$

which shows $\|\hat{\lambda}_M - \hat{\lambda}\|_2 \leq \varepsilon_1(n, M, p)$.

- ▶ Uniform concentration of the predictor over balls in λ :

$$\max_{\lambda \in \mathcal{B}} \|\hat{f}_{\text{RF},M,n}(\cdot; \lambda) - \hat{f}_n(\cdot; \lambda)\|_{L^2} \leq \varepsilon_2(n, M, p).$$

- ▶ Combining the two + Lipschitzness of \hat{f}_n w.r.t λ :

$$\begin{aligned} \|\hat{f}_{\text{RF},M,n}(\cdot; \hat{\lambda}_M) - \hat{f}_n(\cdot; \hat{\lambda})\|_{L^2} &\leq \|\hat{f}_{\text{RF},M,n}(\cdot; \hat{\lambda}_M) - \hat{f}_n(\cdot; \hat{\lambda}_M)\|_{L^2} + \|\hat{f}_n(\cdot; \hat{\lambda}_M) - \hat{f}_n(\cdot; \hat{\lambda})\|_{L^2} \\ &\leq \varepsilon_2(n, M, p) + \varepsilon_1(n, M, p). \end{aligned}$$

Idea of the proof (II): uniform convergence

$$\hat{\lambda} = \arg \max_{\lambda \in \mathbb{R}^n} F(\lambda) := \langle \mathbf{y}, \lambda \rangle - \int_{\Omega} \rho^*(\langle \phi_n(\mathbf{w}), \lambda \rangle) \mu(d\mathbf{w}),$$

$$\hat{\lambda}_M = \arg \max_{\lambda \in \mathbb{R}^n} F_M(\lambda) := \langle \mathbf{y}, \lambda \rangle - \frac{1}{M} \sum_{j \leq M} \rho^*(\langle \phi_{n,j}, \lambda \rangle).$$

- ▶ Concentration of the landscape uniformly over balls in λ :

$$\nabla F_M(\lambda) \rightarrow \nabla F(\lambda), \quad \lambda_{\max}(\nabla^2 F_M(\lambda)) \rightarrow \lambda_{\max}(\nabla^2 F(\lambda)),$$

which shows $\|\hat{\lambda}_M - \hat{\lambda}\|_2 \leq \varepsilon_1(n, M, \rho)$.

- ▶ Uniform concentration of the predictor over balls in λ :

$$\max_{\lambda \in B} \|\hat{f}_{\text{RF},M,n}(\cdot; \lambda) - \hat{f}_n(\cdot; \lambda)\|_{L^2} \leq \varepsilon_2(n, M, \rho).$$

- ▶ Combining the two + Lipschitzness of \hat{f}_n w.r.t λ :

$$\begin{aligned} \|\hat{f}_{\text{RF},M,n}(\cdot; \hat{\lambda}_M) - \hat{f}_n(\cdot; \hat{\lambda})\|_{L^2} &\leq \|\hat{f}_{\text{RF},M,n}(\cdot; \hat{\lambda}_M) - \hat{f}_n(\cdot; \hat{\lambda}_M)\|_{L^2} + \|\hat{f}_n(\cdot; \hat{\lambda}_M) - \hat{f}_n(\cdot; \hat{\lambda})\|_{L^2} \\ &\leq \varepsilon_2(n, M, \rho) + \varepsilon_1(n, M, \rho). \end{aligned}$$

Conclusion

- ▶ \mathcal{F}_p ($1 < p < 2$): function spaces of 2-layers NNs that are tractable and not RKHS.
- ▶ Studying effect of regularization in RF models: reduces to studying the corresponding infinite width model when sufficiently overparametrized.
- ▶ New proof of double descent phenomenon, that does not rely on strong assumptions. Simple mechanism: uniform concentration of the dual to an infinite-width problem.

Future directions:

- ▶ Generalization properties of \mathcal{F}_p .
- ▶ Non-zero regularization (non-interpolating solution).
- ▶ The most interesting case $p = 1$. Is there an efficient algorithm?

Thank you!

Conclusion

- ▶ \mathcal{F}_p ($1 < p < 2$): function spaces of 2-layers NNs that are tractable and not RKHS.
- ▶ Studying effect of regularization in RF models: reduces to studying the corresponding infinite width model when sufficiently overparametrized.
- ▶ New proof of double descent phenomenon, that does not rely on strong assumptions. Simple mechanism: uniform concentration of the dual to an infinite-width problem.

Future directions:

- ▶ Generalization properties of \mathcal{F}_p .
- ▶ Non-zero regularization (non-interpolating solution).
- ▶ The most interesting case $p = 1$. Is there an efficient algorithm?

Thank you!

Conclusion

- ▶ \mathcal{F}_p ($1 < p < 2$): function spaces of 2-layers NNs that are tractable and not RKHS.
- ▶ Studying effect of regularization in RF models: reduces to studying the corresponding infinite width model when sufficiently overparametrized.
- ▶ New proof of double descent phenomenon, that does not rely on strong assumptions. Simple mechanism: uniform concentration of the dual to an infinite-width problem.

Future directions:

- ▶ Generalization properties of \mathcal{F}_p .
- ▶ Non-zero regularization (non-interpolating solution).
- ▶ The most interesting case $p = 1$. Is there an efficient algorithm?

Thank you!

Conclusion

- ▶ \mathcal{F}_p ($1 < p < 2$): function spaces of 2-layers NNs that are tractable and not RKHS.
- ▶ Studying effect of regularization in RF models: reduces to studying the corresponding infinite width model when sufficiently overparametrized.
- ▶ New proof of double descent phenomenon, that does not rely on strong assumptions. Simple mechanism: uniform concentration of the dual to an infinite-width problem.

Future directions:

- ▶ Generalization properties of \mathcal{F}_p .
- ▶ Non-zero regularization (non-interpolating solution).
- ▶ The most interesting case $p = 1$. Is there an efficient algorithm?

Thank you!

Conclusion

- ▶ \mathcal{F}_p ($1 < p < 2$): function spaces of 2-layers NNs that are tractable and not RKHS.
- ▶ Studying effect of regularization in RF models: reduces to studying the corresponding infinite width model when sufficiently overparametrized.
- ▶ New proof of double descent phenomenon, that does not rely on strong assumptions. Simple mechanism: uniform concentration of the dual to an infinite-width problem.

Future directions:

- ▶ Generalization properties of \mathcal{F}_p .
- ▶ Non-zero regularization (non-interpolating solution).
- ▶ The most interesting case $p = 1$. Is there an efficient algorithm?

Thank you!