

When do neural networks outperform kernel methods?

Song Mei

Stanford University

June 28, 2020

Joint work with Behrooz Ghorbani, Theodor Misiakiewicz, and Andrea Montanari

Neural tangent model

- ▶ Multi-layers NN: $f_N(\mathbf{x}; \boldsymbol{\theta})$, $\mathbf{x} \in \mathbb{R}^d$, $\boldsymbol{\theta} \in \mathbb{R}^N$

- ▶ Expanding around $\boldsymbol{\theta}_0$:

$$f_N(\mathbf{x}; \boldsymbol{\theta}) = f_N(\mathbf{x}; \boldsymbol{\theta}_0) + \langle \boldsymbol{\theta} - \boldsymbol{\theta}_0, \nabla_{\boldsymbol{\theta}} f_N(\mathbf{x}; \boldsymbol{\theta}_0) \rangle + o(\|\boldsymbol{\theta} - \boldsymbol{\theta}_0\|_2).$$

- ▶ Neural tangent model:

$$f_{\text{NT},N}(\mathbf{x}; \boldsymbol{\beta}, \boldsymbol{\theta}_0) = \langle \boldsymbol{\beta}, \nabla_{\boldsymbol{\theta}} f_N(\mathbf{x}; \boldsymbol{\theta}_0) \rangle.$$

- ▶ Coupled gradient flow:

$$\begin{aligned} \frac{d}{dt} \boldsymbol{\theta}^t &= -\nabla_{\boldsymbol{\theta}} \hat{\mathbb{E}}[(y - f_N(\mathbf{x}; \boldsymbol{\theta}^t))^2], & \boldsymbol{\theta}^0 &= \boldsymbol{\theta}_0, \\ \frac{d}{dt} \boldsymbol{\beta}^t &= -\nabla_{\boldsymbol{\beta}} \hat{\mathbb{E}}[(y - f_{\text{NT},N}(\mathbf{x}; \boldsymbol{\beta}^t, \boldsymbol{\theta}_0))^2], & \boldsymbol{\beta}^0 &= \mathbf{0}. \end{aligned}$$

- ▶ Under proper initialization and over-parameterization:

$$\lim_{N \rightarrow \infty} |f_N(\mathbf{x}; \boldsymbol{\theta}^t) - f_{\text{NT},N}(\mathbf{x}; \boldsymbol{\beta}^t)| = 0.$$

Neural tangent model

- ▶ Multi-layers NN: $f_N(\mathbf{x}; \boldsymbol{\theta})$, $\mathbf{x} \in \mathbb{R}^d$, $\boldsymbol{\theta} \in \mathbb{R}^N$

- ▶ Expanding around $\boldsymbol{\theta}_0$:

$$f_N(\mathbf{x}; \boldsymbol{\theta}) = f_N(\mathbf{x}; \boldsymbol{\theta}_0) + \langle \boldsymbol{\theta} - \boldsymbol{\theta}_0, \nabla_{\boldsymbol{\theta}} f_N(\mathbf{x}; \boldsymbol{\theta}_0) \rangle + o(\|\boldsymbol{\theta} - \boldsymbol{\theta}_0\|_2).$$

- ▶ Neural tangent model:

$$f_{\text{NT},N}(\mathbf{x}; \boldsymbol{\beta}, \boldsymbol{\theta}_0) = \langle \boldsymbol{\beta}, \nabla_{\boldsymbol{\theta}} f_N(\mathbf{x}; \boldsymbol{\theta}_0) \rangle.$$

- ▶ Coupled gradient flow:

$$\begin{aligned} \frac{d}{dt} \boldsymbol{\theta}^t &= -\nabla_{\boldsymbol{\theta}} \hat{\mathbb{E}}[(y - f_N(\mathbf{x}; \boldsymbol{\theta}^t))^2], & \boldsymbol{\theta}^0 &= \boldsymbol{\theta}_0, \\ \frac{d}{dt} \boldsymbol{\beta}^t &= -\nabla_{\boldsymbol{\beta}} \hat{\mathbb{E}}[(y - f_{\text{NT},N}(\mathbf{x}; \boldsymbol{\beta}^t, \boldsymbol{\theta}_0))^2], & \boldsymbol{\beta}^0 &= \mathbf{0}. \end{aligned}$$

- ▶ Under proper initialization and over-parameterization:

$$\lim_{N \rightarrow \infty} |f_N(\mathbf{x}; \boldsymbol{\theta}^t) - f_{\text{NT},N}(\mathbf{x}; \boldsymbol{\beta}^t)| = 0.$$

Neural tangent model

- ▶ Multi-layers NN: $f_N(\mathbf{x}; \boldsymbol{\theta})$, $\mathbf{x} \in \mathbb{R}^d$, $\boldsymbol{\theta} \in \mathbb{R}^N$

- ▶ Expanding around $\boldsymbol{\theta}_0$:

$$f_N(\mathbf{x}; \boldsymbol{\theta}) = f_N(\mathbf{x}; \boldsymbol{\theta}_0) + \langle \boldsymbol{\theta} - \boldsymbol{\theta}_0, \nabla_{\boldsymbol{\theta}} f_N(\mathbf{x}; \boldsymbol{\theta}_0) \rangle + o(\|\boldsymbol{\theta} - \boldsymbol{\theta}_0\|_2).$$

- ▶ Neural tangent model:

$$f_{\text{NT},N}(\mathbf{x}; \boldsymbol{\beta}, \boldsymbol{\theta}_0) = \langle \boldsymbol{\beta}, \nabla_{\boldsymbol{\theta}} f_N(\mathbf{x}; \boldsymbol{\theta}_0) \rangle.$$

- ▶ Coupled gradient flow:

$$\begin{aligned} \frac{d}{dt} \boldsymbol{\theta}^t &= -\nabla_{\boldsymbol{\theta}} \hat{\mathbb{E}}[(y - f_N(\mathbf{x}; \boldsymbol{\theta}^t))^2], & \boldsymbol{\theta}^0 &= \boldsymbol{\theta}_0, \\ \frac{d}{dt} \boldsymbol{\beta}^t &= -\nabla_{\boldsymbol{\beta}} \hat{\mathbb{E}}[(y - f_{\text{NT},N}(\mathbf{x}; \boldsymbol{\beta}^t, \boldsymbol{\theta}_0))^2], & \boldsymbol{\beta}^0 &= \mathbf{0}. \end{aligned}$$

- ▶ Under proper initialization and over-parameterization:

$$\lim_{N \rightarrow \infty} |f_N(\mathbf{x}; \boldsymbol{\theta}^t) - f_{\text{NT},N}(\mathbf{x}; \boldsymbol{\beta}^t)| = 0.$$

Neural tangent model

- ▶ Multi-layers NN: $f_N(\mathbf{x}; \boldsymbol{\theta})$, $\mathbf{x} \in \mathbb{R}^d$, $\boldsymbol{\theta} \in \mathbb{R}^N$

- ▶ Expanding around $\boldsymbol{\theta}_0$:

$$f_N(\mathbf{x}; \boldsymbol{\theta}) = f_N(\mathbf{x}; \boldsymbol{\theta}_0) + \langle \boldsymbol{\theta} - \boldsymbol{\theta}_0, \nabla_{\boldsymbol{\theta}} f_N(\mathbf{x}; \boldsymbol{\theta}_0) \rangle + o(\|\boldsymbol{\theta} - \boldsymbol{\theta}_0\|_2).$$

- ▶ Neural tangent model:

$$f_{\text{NT},N}(\mathbf{x}; \boldsymbol{\beta}, \boldsymbol{\theta}_0) = \langle \boldsymbol{\beta}, \nabla_{\boldsymbol{\theta}} f_N(\mathbf{x}; \boldsymbol{\theta}_0) \rangle.$$

- ▶ Coupled gradient flow:

$$\begin{aligned} \frac{d}{dt} \boldsymbol{\theta}^t &= -\nabla_{\boldsymbol{\theta}} \hat{\mathbb{E}}[(y - f_N(\mathbf{x}; \boldsymbol{\theta}^t))^2], & \boldsymbol{\theta}^0 &= \boldsymbol{\theta}_0, \\ \frac{d}{dt} \boldsymbol{\beta}^t &= -\nabla_{\boldsymbol{\beta}} \hat{\mathbb{E}}[(y - f_{\text{NT},N}(\mathbf{x}; \boldsymbol{\beta}^t, \boldsymbol{\theta}_0))^2], & \boldsymbol{\beta}^0 &= \mathbf{0}. \end{aligned}$$

- ▶ Under proper initialization and over-parameterization:

$$\lim_{N \rightarrow \infty} |f_N(\mathbf{x}; \boldsymbol{\theta}^t) - f_{\text{NT},N}(\mathbf{x}; \boldsymbol{\beta}^t)| = 0.$$

Neural tangent model

- ▶ Multi-layers NN: $f_N(\mathbf{x}; \boldsymbol{\theta})$, $\mathbf{x} \in \mathbb{R}^d$, $\boldsymbol{\theta} \in \mathbb{R}^N$

- ▶ Expanding around $\boldsymbol{\theta}_0$:

$$f_N(\mathbf{x}; \boldsymbol{\theta}) = f_N(\mathbf{x}; \boldsymbol{\theta}_0) + \langle \boldsymbol{\theta} - \boldsymbol{\theta}_0, \nabla_{\boldsymbol{\theta}} f_N(\mathbf{x}; \boldsymbol{\theta}_0) \rangle + o(\|\boldsymbol{\theta} - \boldsymbol{\theta}_0\|_2).$$

- ▶ Neural tangent model:

$$f_{\text{NT},N}(\mathbf{x}; \boldsymbol{\beta}, \boldsymbol{\theta}_0) = \langle \boldsymbol{\beta}, \nabla_{\boldsymbol{\theta}} f_N(\mathbf{x}; \boldsymbol{\theta}_0) \rangle.$$

- ▶ Coupled gradient flow:

$$\begin{aligned} \frac{d}{dt} \boldsymbol{\theta}^t &= -\nabla_{\boldsymbol{\theta}} \hat{\mathbb{E}}[(y - f_N(\mathbf{x}; \boldsymbol{\theta}^t))^2], & \boldsymbol{\theta}^0 &= \boldsymbol{\theta}_0, \\ \frac{d}{dt} \boldsymbol{\beta}^t &= -\nabla_{\boldsymbol{\beta}} \hat{\mathbb{E}}[(y - f_{\text{NT},N}(\mathbf{x}; \boldsymbol{\beta}^t, \boldsymbol{\theta}_0))^2], & \boldsymbol{\beta}^0 &= \mathbf{0}. \end{aligned}$$

- ▶ Under proper initialization and over-parameterization:

$$\lim_{N \rightarrow \infty} |f_N(\mathbf{x}; \boldsymbol{\theta}^t) - f_{\text{NT},N}(\mathbf{x}; \boldsymbol{\beta}^t)| = 0.$$

How about generalization?

- ▶ [Arora, Du, Hu, Li, Salakhutdinov, Wang, 2019]:
Cifar10 experiments. NT: 23% test error. NN: less than 5% test error.
- ▶ [Arora, Du, Li, Salakhutdinov, Wang, Yu, 2019]:
Small dataset, NT sometimes generalize better than NN.
- ▶ [Shankar, Fang, Guo, Fridovich-Keil, Schmidt, Ragan-Kelley, Recht, 2020]
[Li, Wang, Yu, Du, Hu, Salakhutdinov, Arora, 2019]:
Smaller gap between NT and NN on Cifar10 (10% for NT).

Sometimes there is a large gap, while sometimes the gap is small.

How about generalization?

- ▶ [Arora, Du, Hu, Li, Salakhutdinov, Wang, 2019]:
Cifar10 experiments. NT: 23% test error. NN: less than 5% test error.
- ▶ [Arora, Du, Li, Salakhutdinov, Wang, Yu, 2019]:
Small dataset, NT sometimes generalize better than NN.
- ▶ [Shankar, Fang, Guo, Fridovich-Keil, Schmidt, Ragan-Kelley, Recht, 2020]
[Li, Wang, Yu, Du, Hu, Salakhutdinov, Arora, 2019]:
Smaller gap between NT and NN on Cifar10 (10% for NT).

Sometimes there is a large gap, while sometimes the gap is small.

Focus of this talk

When is there a large performance gap between NN and NT?

Two-layers neural networks

Neural networks:

$$\mathcal{F}_{\text{NN},N} = \left\{ f_N(\mathbf{x}; \Theta) = \sum_{i=1}^N \mathbf{a}_i \sigma(\langle \mathbf{w}_i, \mathbf{x} \rangle) : \mathbf{a}_i \in \mathbb{R}, \mathbf{w}_i \in \mathbb{R}^d \right\}.$$

Linearization:

$$f_N(\mathbf{x}; \Theta) = f_N(\mathbf{x}; \Theta^0) + \underbrace{\sum_{i=1}^N \Delta \mathbf{a}_i \sigma(\langle \mathbf{w}_i^0, \mathbf{x} \rangle)}_{\text{Top layer linearization}} + \underbrace{\sum_{i=1}^N \mathbf{a}_i^0 \sigma'(\langle \mathbf{w}_i^0, \mathbf{x} \rangle) \langle \Delta \mathbf{w}_i, \mathbf{x} \rangle}_{\text{Bottom layer linearization}} + o(\cdot).$$

Linearized neural networks ($\mathbf{W} = (\mathbf{w}_i)_{i \in [N]} \sim_{iid} \text{Unif}(\mathbb{S}^{d-1})$):

$$\mathcal{F}_{\text{RF},N}(\mathbf{W}) = \left\{ f = \sum_{i=1}^N \mathbf{a}_i \sigma(\langle \mathbf{w}_i, \mathbf{x} \rangle) : \mathbf{a}_i \in \mathbb{R}, i \in [N] \right\},$$

$$\mathcal{F}_{\text{NT},N}(\mathbf{W}) = \left\{ f = \sum_{i=1}^N \sigma'(\langle \mathbf{w}_i, \mathbf{x} \rangle) \langle \mathbf{b}_i, \mathbf{x} \rangle : \mathbf{b}_i \in \mathbb{R}^d, i \in [N] \right\}.$$

Two-layers neural networks

Neural networks:

$$\mathcal{F}_{\text{NN},N} = \left\{ f_N(\mathbf{x}; \Theta) = \sum_{i=1}^N \mathbf{a}_i \sigma(\langle \mathbf{w}_i, \mathbf{x} \rangle) : \mathbf{a}_i \in \mathbb{R}, \mathbf{w}_i \in \mathbb{R}^d \right\}.$$

Linearization:

$$f_N(\mathbf{x}; \Theta) = f_N(\mathbf{x}; \Theta^0) + \underbrace{\sum_{i=1}^N \Delta \mathbf{a}_i \sigma(\langle \mathbf{w}_i^0, \mathbf{x} \rangle)}_{\text{Top layer linearization}} + \underbrace{\sum_{i=1}^N \mathbf{a}_i^0 \sigma'(\langle \mathbf{w}_i^0, \mathbf{x} \rangle) \langle \Delta \mathbf{w}_i, \mathbf{x} \rangle}_{\text{Bottom layer linearization}} + o(\cdot).$$

Linearized neural networks ($\mathbf{W} = (\mathbf{w}_i)_{i \in [N]} \sim_{iid} \text{Unif}(\mathbb{S}^{d-1})$):

$$\mathcal{F}_{\text{RF},N}(\mathbf{W}) = \left\{ f = \sum_{i=1}^N \mathbf{a}_i \sigma(\langle \mathbf{w}_i, \mathbf{x} \rangle) : \mathbf{a}_i \in \mathbb{R}, i \in [N] \right\},$$

$$\mathcal{F}_{\text{NT},N}(\mathbf{W}) = \left\{ f = \sum_{i=1}^N \sigma'(\langle \mathbf{w}_i, \mathbf{x} \rangle) \langle \mathbf{b}_i, \mathbf{x} \rangle : \mathbf{b}_i \in \mathbb{R}^d, i \in [N] \right\}.$$

Spiked features model

- ▶ Signal features and junk features

$$\mathbf{x} = (\mathbf{x}_1, \mathbf{x}_2) \in \mathbb{R}^d, \quad \mathbf{x}_1 \in \mathbb{R}^{d_s}, \quad \mathbf{x}_2 \in \mathbb{R}^{d-d_s},$$

$$d_s = d^\eta, \quad 0 \leq \eta \leq 1,$$

$$\text{Cov}(\mathbf{x}_1) = \text{snr}_f \cdot \mathbf{I}_{d_s}, \quad \text{Cov}(\mathbf{x}_2) = \mathbf{I}_{d-d_s},$$

$$\text{snr}_f = d^\kappa, \quad 0 \leq \kappa < \infty \quad (\text{feature SNR}).$$

- ▶ Response depend on signal features

$$y = f_\star(\mathbf{x}) + \varepsilon, \quad f_\star(\mathbf{x}) = \varphi(\mathbf{x}_1).$$

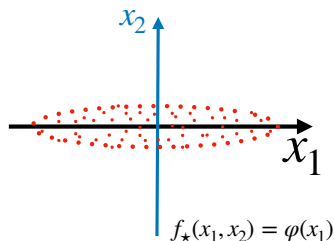


Figure: Anisotropic features:

$$\kappa > 0, \quad \text{snr}_f > 1$$

- ▶ Feature SNR: $\text{snr}_f = d^\kappa \geq 1$.
- ▶ Effective dimension: $d_{\text{eff}} = d_s \vee (d/\text{snr}_f)$. We have $d_s \leq d_{\text{eff}} \leq d$.
- ▶ Larger snr_f induces smaller d_{eff} .

More precisely: $\mathbf{x} \sim \text{Unif}(\mathbb{S}^{d_s}(r\sqrt{d_s})) \times \text{Unif}(\mathbb{S}^{d-d_s}(\sqrt{d}))$. Generalizable to multi-spheres.

Spiked features model

- ▶ Signal features and junk features

$$\mathbf{x} = (\mathbf{x}_1, \mathbf{x}_2) \in \mathbb{R}^d, \quad \mathbf{x}_1 \in \mathbb{R}^{d_s}, \quad \mathbf{x}_2 \in \mathbb{R}^{d-d_s},$$

$$d_s = d^\eta, \quad 0 \leq \eta \leq 1,$$

$$\text{Cov}(\mathbf{x}_1) = \text{snr}_f \cdot \mathbf{I}_{d_s}, \quad \text{Cov}(\mathbf{x}_2) = \mathbf{I}_{d-d_s},$$

$$\text{snr}_f = d^\kappa, \quad 0 \leq \kappa < \infty \quad (\text{feature SNR}).$$

- ▶ Response depend on signal features

$$y = f_\star(\mathbf{x}) + \varepsilon, \quad f_\star(\mathbf{x}) = \varphi(\mathbf{x}_1).$$

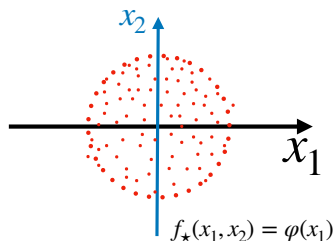


Figure: Isotropic features:

$$\kappa = 0, \quad \text{snr}_f = 1$$

- ▶ Feature SNR: $\text{snr}_f = d^\kappa \geq 1$.
- ▶ Effective dimension: $d_{\text{eff}} = d_s \vee (d/\text{snr}_f)$. We have $d_s \leq d_{\text{eff}} \leq d$.
- ▶ Larger snr_f induces smaller d_{eff} .

More precisely: $\mathbf{x} \sim \text{Unif}(\mathbb{S}^{d_s}(r\sqrt{d_s})) \times \text{Unif}(\mathbb{S}^{d-d_s}(\sqrt{d}))$. Generalizable to multi-spheres.

Spiked features model

- ▶ Signal features and junk features

$$\mathbf{x} = (\mathbf{x}_1, \mathbf{x}_2) \in \mathbb{R}^d, \quad \mathbf{x}_1 \in \mathbb{R}^{d_s}, \quad \mathbf{x}_2 \in \mathbb{R}^{d-d_s},$$

$$d_s = d^\eta, \quad 0 \leq \eta \leq 1,$$

$$\text{Cov}(\mathbf{x}_1) = \text{snr}_f \cdot \mathbf{I}_{d_s}, \quad \text{Cov}(\mathbf{x}_2) = \mathbf{I}_{d-d_s},$$

$$\text{snr}_f = d^\kappa, \quad 0 \leq \kappa < \infty \quad (\text{feature SNR}).$$

- ▶ Response depend on signal features

$$y = f_\star(\mathbf{x}) + \varepsilon, \quad f_\star(\mathbf{x}) = \varphi(\mathbf{x}_1).$$

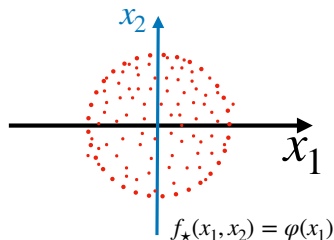


Figure: Isotropic features:

$$\kappa = 0, \quad \text{snr}_f = 1$$

- ▶ Feature SNR: $\text{snr}_f = d^\kappa \geq 1$.
- ▶ Effective dimension: $d_{\text{eff}} = d_s \vee (d/\text{snr}_f)$. We have $d_s \leq d_{\text{eff}} \leq d$.
- ▶ Larger snr_f induces smaller d_{eff} .

More precisely: $\mathbf{x} \sim \text{Unif}(\mathbb{S}^{d_s}(r\sqrt{d_s})) \times \text{Unif}(\mathbb{S}^{d-d_s}(\sqrt{d}))$. Generalizable to multi-spheres.

Approximation error with N neurons

Approximation error: $R(f_\star, \mathcal{F}) = \inf_{f \in \mathcal{F}} \|f_\star - f\|_{L^2}^2$.

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

Assume $d_{\text{eff}}^{\ell+\delta} \leq N \leq d_{\text{eff}}^{\ell+1-\delta}$ and “generic condition” on σ , we have

$$R(f_\star, \mathcal{F}_{\text{RF}, N}(\mathbf{W})) = \|P_{>\ell} f_\star\|_{L^2}^2 + o_{d, \mathbb{P}}(\cdot),$$

$$R(f_\star, \mathcal{F}_{\text{NT}, N}(\mathbf{W})) = \|P_{>\ell+1} f_\star\|_{L^2}^2 + o_{d, \mathbb{P}}(\cdot).$$

On the contrary, assume $d_s^{\ell+\delta} \leq N \leq d_s^{\ell+1-\delta}$, we have

$$R(f_\star, \mathcal{F}_{\text{NN}, N}) \leq \|P_{>\ell+1} f_\star\|_{L^2}^2 + o_d(\cdot).$$

Moreover, $R(f_\star, \mathcal{F}_{\text{NN}, N})$ is independent of snr_f .

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

Approximation error with N neurons

Dim $d_{\text{eff}} \equiv d_s \vee (d/\text{snr}_f)$ and $d_s \leq d_{\text{eff}} \leq d$.

To approx. a degree- ℓ poly. in x_1 :

- ▶ NN need at most d_s^ℓ parameters*.
- ▶ RF need d_{eff}^ℓ parameters.
- ▶ NT need $d_{\text{eff}}^{\ell-1} \cdot d$ parameters.

Approximation power: $\text{NN} \geq \text{RF} \geq \text{NT}$.

* If we don't count parameters with value 0.

Extreme case: low feature SNR

Fix $0 < \eta < 1$, low snr_f : $\kappa = 0$.

To approx. a degree- ℓ poly. in x_1 :

- ▶ NN need at most $d^{\eta\ell}$ parameters*.
- ▶ RF need d^ℓ parameters.
- ▶ NT need d^ℓ parameters.

Approximation power: $\text{NN} > \text{RF} = \text{NT}$.

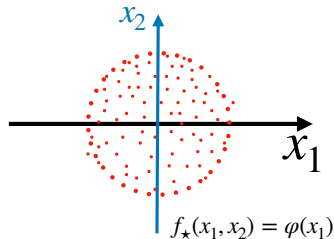


Figure: Isotropic features: $\kappa = 0$, $\text{snr}_f = 1$

* If we don't count parameters with value 0.

Extreme case: high feature SNR

Fix $0 < \eta < 1$, high snr_f : $\kappa \gg 1$.

To approx. a degree- ℓ poly. in x_1 :

- ▶ NN need at most $d^{\eta\ell}$ parameters*.
- ▶ RF need $d^{\eta\ell}$ parameters.
- ▶ NT need $d^{\eta(\ell-1)+1}$ parameters.

Approximation power: $\text{NN} \sim \text{RF} > \text{NT}$.

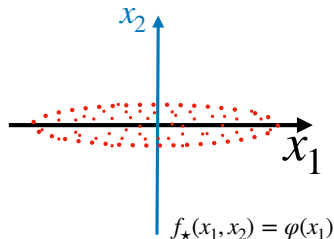
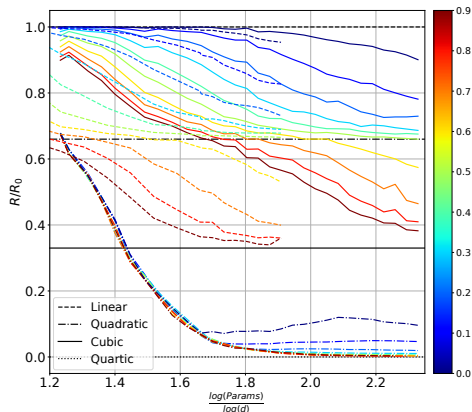


Figure: Anisotropic features:
 $\kappa > 0$, $\text{snr}_f > 1$

* If we don't count parameters with value 0.

Numerical simulations



Colorbar: $\kappa \in [0, 1]$.

Dot-dashed: NN.

Dashed lines: RF;

Continuous lines: NT;

Dimension: $d = 1024$.

Eff. dim: $d_s = 16$.

Conclusion

- (a) Power: $\text{NN} \geq \text{RF} \geq \text{NT}$.
- (b) Risk of NN independent of snr_f .
- (c) Larger snr_f induces larger power of $\{\text{RF}, \text{NT}\}$.

Similar results for generalization error with finite samples n

Extreme case: low feature SNR

Fix $0 < \eta < 1$, low snr_f : $\kappa = 0$.

To fit a degree- ℓ poly. in x_1 :

- ▶ $\exists \nu$, NN need at most $d^{\eta\ell}$ samples.
- ▶ {RF, NT} kernel need d^ℓ samples.

Potential generalization power:

NN > Kernel methods.

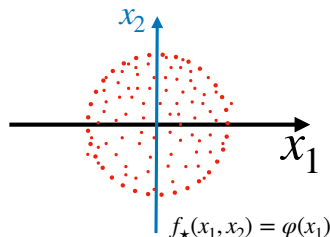


Figure: Isotropic features: $\kappa = 0$,
 $\text{snr}_f = 1$

Extreme case: high feature SNR

Fix $0 < \eta < 1$, high snr_f : $\kappa \gg 1$.

To fit a degree- ℓ poly. in x_1 :

- ▶ $\exists \nu$, NN need at most $d^{\eta \ell}$ samples.
- ▶ {RF, NT} kernel need $d^{\eta \ell}$ samples.

Potential generalization power:

NN \sim Kernel methods.

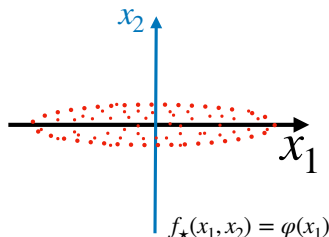


Figure: Anisotropic features:

$$\kappa > 0, \text{snr}_f > 1$$

Implications

Adding isotropic noise in features (i.e., decreasing snr_f), performance gap between NN and $\{\text{RF}, \text{NT}\}$ becomes larger.

Numerical simulations

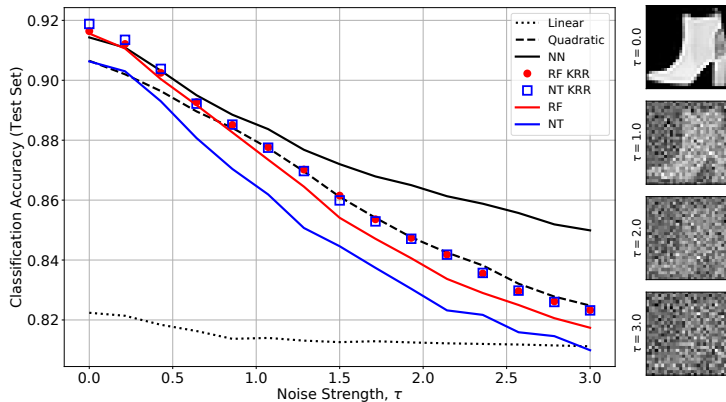


Figure: Underlying assumption: labels depend on low frequency components of images.

Message

In spiked features model, a controlling parameter of the performance gap between NN and $\{\text{RF}, \text{NT}\}$ is

$$\text{snr}_f = \text{Feature SNR} = \frac{\text{Signal features variance}}{\text{Junk features variance}}.$$

- ▶ Small snr_f , there is a large separation.
- ▶ Large snr_f , $\{\text{RF}, \text{NT}\}$ performs closer to NN.

Somewhat implicitly, NN first finds the signal features (PCA), and then perform kernel methods on these features.

$$\text{snr}_f \neq \text{SNR} = \|f_\star\|_{L^2}^2 / \mathbb{E}[\epsilon^2]$$

Thank you!