# When do neural networks outperform kernel methods?

Song Mei

Stanford University

June 29, 2020

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

- ▶ Multi-layers NN:  $f_N(x; \theta), x \in \mathbb{R}^d, \theta \in \mathbb{R}^N$
- Expanding around  $\theta_0$ :

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

Neural tangent model:

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

Coupled gradient flow:

$$egin{aligned} &rac{\mathrm{d}}{\mathrm{d}t}m{ heta}^t = &-
abla_{m{ heta}}\hat{\mathbb{E}}[(y-f_N(m{x};m{ heta}^t))^2], &m{ heta}^0 = m{ heta}_0, \ &rac{\mathrm{d}}{\mathrm{d}t}m{ heta}^t = &-
abla_{m{ heta}}\hat{\mathbb{E}}[(y-f_{\mathsf{NT},N}(m{x};m{ heta}^t,m{ heta}_0))^2], &m{ heta}^0 = m{0}. \end{aligned}$$

Under proper initialization and over-parameterization:

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

- ▶ Multi-layers NN:  $f_N(x; \theta), x \in \mathbb{R}^d, \theta \in \mathbb{R}^N$
- Expanding around  $\theta_0$ :

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

Neural tangent model:

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

Coupled gradient flow:

$$egin{aligned} &rac{\mathrm{d}}{\mathrm{d}t}m{ heta}^t = & -
abla_{m{ heta}}\hat{\mathbb{E}}[(y-f_N(m{x};m{ heta}^t))^2], &m{ heta}^0 = m{ heta}_0, \ &rac{\mathrm{d}}{\mathrm{d}t}m{ heta}^t = & -
abla_{m{ heta}}\hat{\mathbb{E}}[(y-f_{\mathsf{NT},N}(m{x};m{ heta}^t,m{ heta}_0))^2], &m{ heta}^0 = m{0}. \end{aligned}$$

Under proper initialization and over-parameterization:

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

- ▶ Multi-layers NN:  $f_N(x; \theta), x \in \mathbb{R}^d, \theta \in \mathbb{R}^N$
- Expanding around  $\theta_0$ :

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

Neural tangent model:

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

Coupled gradient flow:

$$egin{aligned} &rac{\mathrm{d}}{\mathrm{d}t}m{ heta}^t = &-
abla_{m{ heta}}\hat{\mathbb{E}}[(y-f_N(m{x};m{ heta}^t))^2], &m{ heta}^0 = m{ heta}_0, \ &rac{\mathrm{d}}{\mathrm{d}t}m{ heta}^t = &-
abla_{m{ heta}}\hat{\mathbb{E}}[(y-f_{\mathsf{NT},N}(m{x};m{ heta}^t,m{ heta}_0))^2], &m{ heta}^0 = m{0}. \end{aligned}$$

Under proper initialization and over-parameterization:

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

- ▶ Multi-layers NN:  $f_N(x; \theta), x \in \mathbb{R}^d, \theta \in \mathbb{R}^N$
- Expanding around  $\theta_0$ :

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

Neural tangent model:

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

Coupled gradient flow:

$$\frac{\mathrm{d}}{\mathrm{d}t}\boldsymbol{\theta}^{t} = -\nabla_{\boldsymbol{\theta}}\hat{\mathbb{E}}[(\boldsymbol{y} - f_{N}(\boldsymbol{x};\boldsymbol{\theta}^{t}))^{2}], \qquad \boldsymbol{\theta}^{0} = \boldsymbol{\theta}_{0},$$
$$\frac{\mathrm{d}}{\mathrm{d}t}\boldsymbol{\beta}^{t} = -\nabla_{\boldsymbol{\beta}}\hat{\mathbb{E}}[(\boldsymbol{y} - f_{\mathsf{NT},N}(\boldsymbol{x};\boldsymbol{\beta}^{t},\boldsymbol{\theta}_{0}))^{2}], \qquad \boldsymbol{\beta}^{0} = \boldsymbol{0}.$$

Under proper initialization and over-parameterization:

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

- ▶ Multi-layers NN:  $f_N(x; \theta), x \in \mathbb{R}^d, \theta \in \mathbb{R}^N$
- Expanding around  $\theta_0$ :

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

Neural tangent model:

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

Coupled gradient flow:

$$\frac{\mathrm{d}}{\mathrm{d}t}\boldsymbol{\theta}^{t} = -\nabla_{\boldsymbol{\theta}}\hat{\mathbb{E}}[(\boldsymbol{y} - f_{N}(\boldsymbol{x};\boldsymbol{\theta}^{t}))^{2}], \qquad \boldsymbol{\theta}^{0} = \boldsymbol{\theta}_{0},$$
$$\frac{\mathrm{d}}{\mathrm{d}t}\boldsymbol{\beta}^{t} = -\nabla_{\boldsymbol{\beta}}\hat{\mathbb{E}}[(\boldsymbol{y} - f_{\mathsf{NT},N}(\boldsymbol{x};\boldsymbol{\beta}^{t},\boldsymbol{\theta}_{0}))^{2}], \qquad \boldsymbol{\beta}^{0} = \boldsymbol{0}.$$

Under proper initialization and over-parameterization:

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

# 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:

$$\mathcal{F}_{\mathsf{NN},N} = \left\{ f_N(oldsymbol{x};oldsymbol{\Theta}) = \sum_{i=1}^N oldsymbol{a}_i \sigma(\langleoldsymbol{w}_i,oldsymbol{x}
angle) : oldsymbol{a}_i \in \mathbb{R}, oldsymbol{w}_i \in \mathbb{R}^d 
ight\}.$$

Linearization:

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

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

$$egin{split} \mathcal{F}_{ ext{RF},N}(oldsymbol{W}) &= \Big\{f = \sum_{i=1}^N a_i \sigma(\langle oldsymbol{w}_i, oldsymbol{x} 
angle) : a_i \in \mathbb{R}, i \in [N] \Big\}, \ \mathcal{F}_{ ext{NT},N}(oldsymbol{W}) &= \Big\{f = \sum_{i=1}^N \sigma'(\langle oldsymbol{w}_i, oldsymbol{x} 
angle) : oldsymbol{b}_i \in \mathbb{R}^d, i \in [N] \Big\}. \end{split}$$

# Two-layers neural networks

$$\mathcal{F}_{\mathsf{NN},N} = \left\{ f_N(oldsymbol{x};oldsymbol{\Theta}) = \sum_{i=1}^N oldsymbol{a}_i \sigma(\langleoldsymbol{w}_i,oldsymbol{x}
angle) : oldsymbol{a}_i \in \mathbb{R}, oldsymbol{w}_i \in \mathbb{R}^d 
ight\}.$$

Linearization:

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

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

$$\mathcal{F}_{\mathsf{RF},N}(\boldsymbol{W}) = \Big\{ f = \sum_{i=1}^{N} \boldsymbol{a}_{i} \sigma(\langle \boldsymbol{w}_{i}, \boldsymbol{x} \rangle) : \boldsymbol{a}_{i} \in \mathbb{R}, i \in [N] \Big\},$$
$$\mathcal{F}_{\mathsf{NT},N}(\boldsymbol{W}) = \Big\{ f = \sum_{i=1}^{N} \sigma'(\langle \boldsymbol{w}_{i}, \boldsymbol{x} \rangle) \langle \boldsymbol{b}_{i}, \boldsymbol{x} \rangle : \boldsymbol{b}_{i} \in \mathbb{R}^{d}, i \in [N] \Big\}.$$

# Spiked features model

Signal features and junk features

$$egin{aligned} &oldsymbol{x} = (oldsymbol{x}_1, oldsymbol{x}_2) \in \mathbb{R}^d, \ oldsymbol{x}_1 \in \mathbb{R}^{d_{s}}, \ oldsymbol{x}_2 \in \mathbb{R}^{d-d_{s}}, \ d_{s} = d^{oldsymbol{\eta}}, \ 0 \leq oldsymbol{\eta} \leq 1, \ \mathrm{Cov}(oldsymbol{x}_1) = \mathrm{snr}_{\mathrm{f}} \cdot \mathbf{I}_{d_{s}}, \ \mathrm{Cov}(oldsymbol{x}_2) = \mathbf{I}_{d-d_{s}}, \ \mathrm{snr}_{\mathrm{f}} = d^{oldsymbol{\kappa}}, \ 0 \leq oldsymbol{\kappa} < \infty \ (\mathrm{feature\ SNR}). \end{aligned}$$

Response depend on signal features

$$y=f_\star(x)+arepsilon, \ \ f_\star(x)=arphi(x_1).$$



Figure: Anisotropic features:  $\kappa > 0, \operatorname{snr}_{f} > 1$ 

Feature SNR:  $\operatorname{snr}_{\mathbf{f}} = d^{\kappa} \geq 1$ .

▶ Effective dimension:  $d_{\text{eff}} = d_{\text{s}} \lor (d/\text{snr}_{\text{f}})$ . We have  $d_{\text{s}} \le d_{\text{eff}} \le d$ .

• Larger  $\operatorname{snr}_{\mathbf{f}}$  induces smaller  $d_{\operatorname{eff}}$ .

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

,

# Spiked features model

Signal features and junk features

$$egin{aligned} &oldsymbol{x} = (oldsymbol{x}_1, oldsymbol{x}_2) \in \mathbb{R}^d, \ oldsymbol{x}_1 \in \mathbb{R}^{d_{\mathrm{s}}}, \ oldsymbol{x}_2 \in \mathbb{R}^{d-d_{\mathrm{s}}}, \ d_{\mathrm{s}} = d^{\eta}, \ 0 \leq \eta \leq 1, \ \mathrm{Cov}(oldsymbol{x}_1) = \mathrm{snr}_{\mathrm{f}} \cdot \mathbf{I}_{d_{\mathrm{s}}}, \ \mathrm{Cov}(oldsymbol{x}_2) = \mathbf{I}_{d-d_{\mathrm{s}}}, \ \mathrm{snr}_{\mathrm{f}} = d^{\kappa}, \ 0 \leq \kappa < \infty \ (\mathrm{feature\ SNR}). \end{aligned}$$

Response depend on signal features

$$y=f_\star(x)+arepsilon, \ \ f_\star(x)=arphi(x_1).$$



Figure: Isotropic features:  $\kappa = 0, \ \operatorname{snr}_{f} = 1$ 

Feature SNR:  $\operatorname{snr}_{\mathbf{f}} = d^{\kappa} \geq 1$ .

▶ Effective dimension:  $d_{\text{eff}} = d_{\text{s}} \lor (d/\text{snr}_{\text{f}})$ . We have  $d_{\text{s}} \le d_{\text{eff}} \le d$ .

Larger  $\operatorname{snr}_{\mathbf{f}}$  induces smaller  $d_{\operatorname{eff}}$ .

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

,

# Spiked features model

► Signal features and junk features

$$egin{aligned} &oldsymbol{x} = (oldsymbol{x}_1, oldsymbol{x}_2) \in \mathbb{R}^d, \ oldsymbol{x}_1 \in \mathbb{R}^{d_{\mathrm{s}}}, \ oldsymbol{x}_2 \in \mathbb{R}^{d-d_{\mathrm{s}}}, \ d_{\mathrm{s}} = d^{\eta}, \ 0 \leq \eta \leq 1, \ \mathrm{Cov}(oldsymbol{x}_1) = \mathrm{snr}_{\mathrm{f}} \cdot \mathbf{I}_{d_{\mathrm{s}}}, \ \mathrm{Cov}(oldsymbol{x}_2) = \mathbf{I}_{d-d_{\mathrm{s}}}, \ \mathrm{snr}_{\mathrm{f}} = d^{\kappa}, \ 0 \leq \kappa < \infty \ (\mathrm{feature\ SNR}). \end{aligned}$$

Response depend on signal features

$$y=f_\star(x)+arepsilon, \ \ f_\star(x)=arphi(x_1).$$



Figure: Isotropic features:  $\kappa = 0, \ \operatorname{snr}_{f} = 1$ 

Feature SNR:  $\operatorname{snr}_{\mathbf{f}} = d^{\kappa} \geq 1$ .

▶ Effective dimension:  $d_{\text{eff}} = d_{\text{s}} \lor (d/\operatorname{snr}_{\text{f}})$ . We have  $d_{\text{s}} \le d_{\text{eff}} \le d$ .

• Larger  $\operatorname{snr}_{f}$  induces smaller  $d_{eff}$ .

More precisely:  $\boldsymbol{x} \sim \mathrm{Unif}(\mathbb{S}^{d_{\mathbb{S}}}(r\sqrt{d_{\mathbb{S}}})) \times \mathrm{Unif}(\mathbb{S}^{d-d_{\mathbb{S}}}(\sqrt{d}))$ . 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  $\sigma$ , we have

$$egin{aligned} &R(f_\star,\mathcal{F}_{\mathsf{RF},oldsymbol{N}}(oldsymbol{W})) = \|\mathsf{P}_{>\ell}f_\star\|_{L^2}^2 + o_{d,\mathbb{P}}(\cdot), \ &R(f_\star,\mathcal{F}_{\mathsf{NT},oldsymbol{N}}(oldsymbol{W})) = \|\mathsf{P}_{>\ell+1}f_\star\|_{L^2}^2 + o_{d,\mathbb{P}}(\cdot) \end{aligned}$$

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

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

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

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

# Approximation error with N neurons

Dim  $d_{\text{eff}} \equiv d_{\text{s}} \lor (d/\text{snr}_{\text{f}})$  and  $d_{\text{s}} \leq d_{\text{eff}} \leq d$ .

To approx. a degree- $\ell$  poly. in  $x_1$ :

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

Approximation power:  $NN \ge RF \ge NT$ .

<sup>\*</sup> If we don't count parameters with value 0.

# Extreme case: low feature SNR

Fix  $0 < \eta < 1$ , low  $\operatorname{snr}_{f}$ :  $\kappa = 0$ .

To approx. a degree- $\ell$  poly. in  $x_1$ :

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

Approximation power: NN > RF = NT.



Figure: Isotropic features:  $\kappa = 0$ ,  $\operatorname{snr}_{f} = 1$ 

<sup>\*</sup> If we don't count parameters with value 0.

# Extreme case: high feature SNR

Fix  $0 < \eta < 1$ , high  $\operatorname{snr}_{f}: \kappa \gg 1$ .

To approx. a degree- $\ell$  poly. in  $x_1$ :

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

Approximation power:  $NN \sim RF > NT$ .



Figure: Anisotropic features:  $\kappa > 0, \ \operatorname{snr}_{f} > 1$ 

<sup>\*</sup> 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:  $NN \ge RF \ge NT$ . (b) Risk of NN independent of  $snr_f$ . (c) Larger  $snr_f$  induces larger power of {RF, NT}.

Similar results for generalization error with finite samples n

# Extreme case: low feature SNR

Fix  $0 < \eta < 1$ , low  $\operatorname{snr}_{f}$ :  $\kappa = 0$ .

To fit a degree- $\ell$  poly. in  $x_1$ :

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

 $\begin{array}{l} \mbox{Potential generalization power:} \\ \mbox{NN} > \mbox{Kernel methods.} \end{array}$ 



Figure: Isotropic features:  $\kappa = 0$ ,  $\operatorname{snr}_{f} = 1$ 

# Extreme case: high feature SNR

Fix  $0 < \eta < 1$ , high  $\operatorname{snr}_{f}: \kappa \gg 1$ .

To fit a degree- $\ell$  poly. in  $x_1$ :

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

 $\begin{array}{l} \mbox{Potential generalization power:} \\ \mbox{NN} \sim \mbox{Kernel methods.} \end{array}$ 



Figure: Anisotropic features:  $\kappa > 0, \operatorname{snr}_{f} > 1$ 

# Implications

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

 $snr_{f} = Feature SNR = \frac{Signal \ features \ variance}{Junk \ features \ variance}$ 

Small snr<sub>f</sub>, there is a large separation.

▶ Large snr<sub>f</sub>, {RF, NT} performs closer to NN.

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

$$ext{snr_f} 
eq ext{SNR} = \|f_\star\|_{L^2}^2 / \mathbb{E}[arepsilon^2]$$

Song Mei (Stanford University)

Thank you!