

Optimization, Robustness and Attention in Deep Learning: Insights from Random and NTK Feature Models

Marco Mondelli

Institute of Science and Technology Austria (ISTA)

ISL Colloquium, April 11, 2024



Supervised learning

Data: $\{(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_n, y_n)\} \sim_{\text{i.i.d.}} \mathbb{P}(\mathbb{R}^d \times \mathbb{R})$

Goal: Given $(\mathbf{x}, y) \sim \mathbb{P}$, find $f : \mathbb{R}^d \rightarrow \mathbb{R}$ to predict y from \mathbf{x}

Supervised learning

Data: $\{(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_n, y_n)\} \sim_{\text{i.i.d.}} \mathbb{P}(\mathbb{R}^d \times \mathbb{R})$

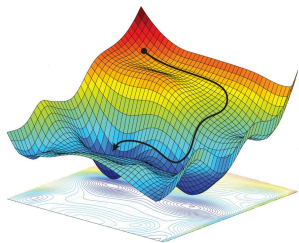
Goal: Minimize empirical risk $L_f(\boldsymbol{\theta}) = \frac{1}{n} \sum_{i=1}^n (y_i - f(\mathbf{x}_i; \boldsymbol{\theta}))^2$

Supervised learning

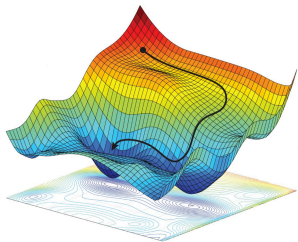
Data: $\{(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_n, y_n)\} \sim_{\text{i.i.d.}} \mathbb{P}(\mathbb{R}^d \times \mathbb{R})$

Goal: Minimize empirical risk $L_f(\boldsymbol{\theta}) = \frac{1}{n} \sum_{i=1}^n (y_i - f(\mathbf{x}_i; \boldsymbol{\theta}))^2$

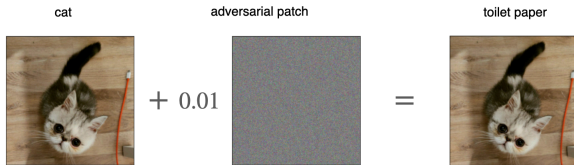
Gradient flow: $\dot{\boldsymbol{\theta}}(t) = -\nabla_{\boldsymbol{\theta}} L_f(\boldsymbol{\theta}(t))$



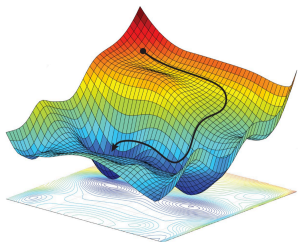
optimization



optimization



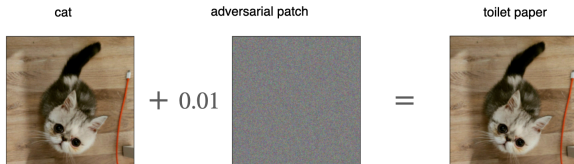
robustness



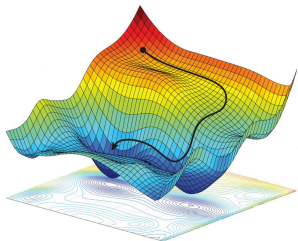
optimization



attention



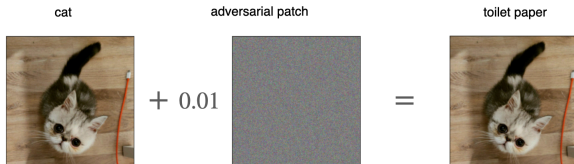
robustness



optimization



attention



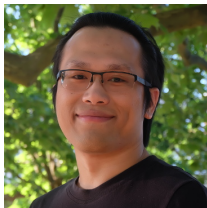
robustness



Simone Bombari (ISTA)



Mohammad Amani (ISTA → EPFL)



Quynh Nguyen (MPI)



Guido Montufar (MPI & UCLA)

Insights from the Neural Tangent Kernel (NTK)

$$\dot{\boldsymbol{\theta}}(t) = -\nabla_{\boldsymbol{\theta}} L_f(\boldsymbol{\theta}(t)), \quad \boldsymbol{\theta} \in \mathbb{R}^p, \quad p \gg n$$

Idea: $\boldsymbol{\theta}(0)$ not too far from an interpolator, so throughout the gradient flow trajectory we have

$$f(\mathbf{x}; \boldsymbol{\theta}(t)) \approx f(\mathbf{x}; \boldsymbol{\theta}(0)) + \langle \nabla_{\boldsymbol{\theta}} f(\mathbf{x}; \boldsymbol{\theta}(0)), \boldsymbol{\theta}(t) - \boldsymbol{\theta}(0) \rangle$$

Insights from the Neural Tangent Kernel (NTK)

$$\dot{\boldsymbol{\theta}}(t) = -\nabla_{\boldsymbol{\theta}} L_f(\boldsymbol{\theta}(t)), \quad \boldsymbol{\theta} \in \mathbb{R}^p, \quad p \gg n$$

Idea: $\boldsymbol{\theta}(0)$ not too far from an interpolator, so throughout the gradient flow trajectory we have

$$f(\mathbf{x}; \boldsymbol{\theta}(t)) \approx f(\mathbf{x}; \boldsymbol{\theta}(0)) + \langle \nabla_{\boldsymbol{\theta}} f(\mathbf{x}; \boldsymbol{\theta}(0)), \boldsymbol{\theta}(t) - \boldsymbol{\theta}(0) \rangle$$

$$L_f(\boldsymbol{\theta}(t)) \leq L_f(\boldsymbol{\theta}(0)) e^{-\lambda_{\min}(\mathbf{K})t/2}$$

- NTK: $\mathbf{K} = \mathbf{J}_f(\boldsymbol{\theta}_0) \mathbf{J}_f(\boldsymbol{\theta}_0)^\top \in \mathbb{R}^{n \times n}$
- Jacobian of f at initialization: $\mathbf{J}_f(\boldsymbol{\theta}_0) \in \mathbb{R}^{n \times p}$

[Jacot et al., 2018; Chizat et al., 2019; Du et al., 2019; Oymak et al., 2019; Bartlett et al., 2021; ...]

Convergence for (very) wide networks

	Deep?	Activation	Layer Width	# Wide Layers
[Oymak et al., '20]	No	Smooth	$\Omega(n^2 \lambda_0^{-2})$	×
[Montanari & Zhong, '22]	No	General	$\tilde{\Omega}(n/d)$	×
[Allen-Zhu et al., '19]	Yes	General	$\Omega(n^{24} L^{12} \phi^{-4})$	All
[Zou et al., '19]	Yes	ReLU	$\Omega(n^8 L^{12} \phi^{-4})$	All
[Du et al., '19]	Yes	Smooth	$\Omega\left(\frac{n^4 2^{\mathcal{O}(L)}}{\lambda_{\min}^4(\bar{\mathbf{K}}^{(L)})}\right)$	All

Convergence for (not so) wide networks

	Deep?	Activation	Layer Width	# Wide Layers
[Oymak et al., '20]	No	Smooth	$\Omega(n^2 \lambda_0^{-2})$	x
[Montanari & Zhong, '22]	No	General	$\tilde{\Omega}(n/d)$	x
[Allen-Zhu et al., '18]	Yes	General	$\Omega(n^{24} L^{12} \phi^{-4})$	All
[Zou et al., '19]	Yes	ReLU	$\Omega(n^8 L^{12} \phi^{-4})$	All
[Du et al., '19]	Yes	Smooth	$\Omega\left(\frac{n^4 2^{\mathcal{O}(L)}}{\lambda_{\min}^4(\mathbf{K}^{(L)})}\right)$	All
[Nguyen and M., '20]	Yes	Smooth	n	One

Need only **one wide layer + pyramidal topology**

Q. Nguyen and M. Mondelli, "Global Convergence of Deep Networks with One Wide Layer Followed by Pyramidal Topology", *NeurIPS*, 2020.

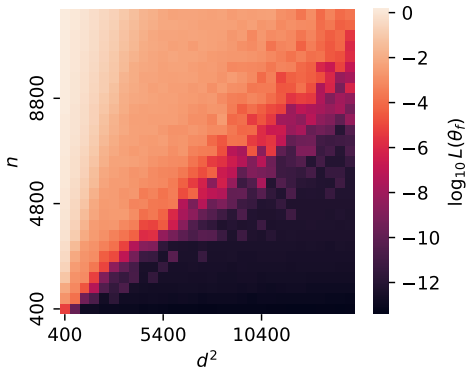
Convergence for (not so) wide networks

	Deep?	Activation	Layer Width	# Wide Layers
[Oymak et al., '20]	No	Smooth	$\Omega(n^2 \lambda_0^{-2})$	×
[Montanari & Zhong, '22]	No	General	$\tilde{\Omega}(n/d)$	×
[Allen-Zhu et al., '18]	Yes	General	$\Omega(n^{24} L^{12} \phi^{-4})$	All
[Zou et al., '19]	Yes	ReLU	$\Omega(n^8 L^{12} \phi^{-4})$	All
[Du et al., '19]	Yes	Smooth	$\Omega(\frac{n^4 2^{\mathcal{O}(L)}}{\lambda_{\min}^4(\bar{\mathbf{K}}(L))})$	All
[Nguyen and M., '20]	Yes	Smooth	n	One
[Bombari, Amani and M., '22]	Yes	Smooth	$\tilde{\Omega}(\sqrt{n})$	All

Need only **minimum over-parameterization**

S. Bombari, M. H. Amani, and M. Mondelli, "Memorization and Optimization in Deep Neural Networks with Minimum Over-parameterization", *NeurIPS*, 2022.

Optimization with minimum over-parameterization



- $\Omega(\sqrt{n})$ neurons **necessary** to interpolate (parameter counting or VC dimension bound [Bartlett et al., 2019])
- Scaling **close to practice** (back-of-the-envelope estimates on CIFAR-10, ImageNet)

Bounding $\lambda_{\min}(\mathbf{K})$

$$\mathbf{K} = \sum_{k=0}^{L-1} \mathbf{F}_k \mathbf{F}_k^\top \circ \mathbf{B}_{k+1} \mathbf{B}_{k+1}^\top$$

- $\mathbf{F}_k = [\mathbf{f}_k(\mathbf{x}_1), \dots, \mathbf{f}_k(\mathbf{x}_n)]^\top$, with $\mathbf{f}_k(\mathbf{x}_i)$ = feature vector at layer k with input \mathbf{x}_i
- $\mathbf{B}_{k+1} = [\mathbf{b}_{k+1}(\mathbf{x}_1), \dots, \mathbf{b}_{k+1}(\mathbf{x}_n)]^\top$, with $\mathbf{b}_{k+1}(\mathbf{x}_i)$ = back-propagation vector at layer $k + 1$ with input \mathbf{x}_i

Bounding $\lambda_{\min}(\mathbf{K})$

$$\mathbf{K} = \sum_{k=0}^{L-1} \mathbf{F}_k \mathbf{F}_k^{\top} \circ \mathbf{B}_{k+1} \mathbf{B}_{k+1}^{\top}$$

- $\mathbf{F}_k = [\mathbf{f}_k(\mathbf{x}_1), \dots, \mathbf{f}_k(\mathbf{x}_n)]^{\top}$, with $\mathbf{f}_k(\mathbf{x}_i)$ = feature vector at layer k with input \mathbf{x}_i
- $\mathbf{B}_{k+1} = [\mathbf{b}_{k+1}(\mathbf{x}_1), \dots, \mathbf{b}_{k+1}(\mathbf{x}_n)]^{\top}$, with $\mathbf{b}_{k+1}(\mathbf{x}_i)$ = back-propagation vector at layer $k + 1$ with input \mathbf{x}_i

First attempt: use matrix concentration on $\mathbf{F}_k \mathbf{F}_k^{\top}$

Q. Nguyen, M. Mondelli and G. Montufar, “Tight Bounds on the Smallest Eigenvalue of the Neural Tangent Kernel for Deep ReLU Networks”, *ICML*, 2021.

Bounding $\lambda_{\min}(\mathbf{K})$ with one wide layer

$$\mathbf{K} = \sum_{k=0}^{L-1} \mathbf{F}_k \mathbf{F}_k^\top \circ \mathbf{B}_{k+1} \mathbf{B}_{k+1}^\top$$

- $\mathbf{F}_k = [\mathbf{f}_k(\mathbf{x}_1), \dots, \mathbf{f}_k(\mathbf{x}_n)]^\top$, with $\mathbf{f}_k(\mathbf{x}_i)$ = feature vector at layer k with input \mathbf{x}_i
- $\mathbf{B}_{k+1} = [\mathbf{b}_{k+1}(\mathbf{x}_1), \dots, \mathbf{b}_{k+1}(\mathbf{x}_n)]^\top$, with $\mathbf{b}_{k+1}(\mathbf{x}_i)$ = back-propagation vector at layer $k+1$ with input \mathbf{x}_i

First attempt: use matrix concentration on $\mathbf{F}_k \mathbf{F}_k^\top$

Need **one wide layer** with $\Omega(n)$ neurons!

Q. Nguyen, M. Mondelli and G. Montufar, “Tight Bounds on the Smallest Eigenvalue of the Neural Tangent Kernel for Deep ReLU Networks”, *ICML*, 2021.

Bounding $\lambda_{\min}(\mathbf{K})$ with minimum over-parameterization

$$\mathbf{K} \succeq \mathbf{F}_{L-2}\mathbf{F}_{L-2}^{\top} \circ \mathbf{B}_{L-1}\mathbf{B}_{L-1}^{\top} := \mathbf{J}_{L-2}\mathbf{J}_{L-2}^{\top}$$

- $(\mathbf{J}_{L-2})_{i,:} = \mathbf{f}_{L-2}(\mathbf{x}_i) \otimes \mathbf{b}_{L-1}(\mathbf{x}_i)$
- $\mathbf{f}_{L-2}(\mathbf{x}_i)$ = feature vector at layer $L - 2$ with input \mathbf{x}_i
- $\mathbf{b}_{L-1}(\mathbf{x}_i)$ = back-propagation vector at layer $L - 1$ with input \mathbf{x}_i

Second attempt: directly center $\mathbf{J}_{L-2}\mathbf{J}_{L-2}^{\top}$

Bounding $\lambda_{\min}(\mathbf{K})$ with minimum over-parameterization

$$\mathbf{K} \succeq \mathbf{F}_{L-2}\mathbf{F}_{L-2}^T \circ \mathbf{B}_{L-1}\mathbf{B}_{L-1}^T := \mathbf{J}_{L-2}\mathbf{J}_{L-2}^T$$

- $(\mathbf{J}_{L-2})_{i,:} = \mathbf{f}_{L-2}(\mathbf{x}_i) \otimes \mathbf{b}_{L-1}(\mathbf{x}_i)$
- $\mathbf{f}_{L-2}(\mathbf{x}_i)$ = feature vector at layer $L - 2$ with input \mathbf{x}_i
- $\mathbf{b}_{L-1}(\mathbf{x}_i)$ = back-propagation vector at layer $L - 1$ with input \mathbf{x}_i

Second attempt: directly center $\mathbf{J}_{L-2}\mathbf{J}_{L-2}^T$

$$\mathbf{J}_{L-2}\mathbf{J}_{L-2}^T \succsim \mathbf{J}_{FB}\mathbf{J}_{FB}^T$$

- $(\mathbf{J}_{FB})_{i,:} = \tilde{\mathbf{f}}_{L-2}(\mathbf{x}_i) \otimes \tilde{\mathbf{b}}_{L-1}(\mathbf{x}_i)$
- $\tilde{\mathbf{f}}_{L-2}(\mathbf{x}_i) = \mathbf{f}_{L-2}(\mathbf{x}_i) - \mathbb{E}_{\mathbf{x}_i} \mathbf{f}_{L-2}(\mathbf{x}_i)$
- $\tilde{\mathbf{b}}_{L-1}(\mathbf{x}_i) = \mathbf{b}_{L-1}(\mathbf{x}_i) - \mathbb{E}_{\mathbf{x}_i} \mathbf{b}_{L-1}(\mathbf{x}_i)$

Features and back-propagations centered **together**



Bounding $\lambda_{\min}(\mathbf{K})$ with minimum over-parameterization

$$\mathbf{K} \succeq \mathbf{F}_{L-2}\mathbf{F}_{L-2}^T \circ \mathbf{B}_{L-1}\mathbf{B}_{L-1}^T := \mathbf{J}_{L-2}\mathbf{J}_{L-2}^T$$

- $(\mathbf{J}_{L-2})_{i,:} = \mathbf{f}_{L-2}(\mathbf{x}_i) \otimes \mathbf{b}_{L-1}(\mathbf{x}_i)$
- $\mathbf{f}_{L-2}(\mathbf{x}_i)$ = feature vector at layer $L - 2$ with input \mathbf{x}_i
- $\mathbf{b}_{L-1}(\mathbf{x}_i)$ = back-propagation vector at layer $L - 1$ with input \mathbf{x}_i

Second attempt: directly center $\mathbf{J}_{L-2}\mathbf{J}_{L-2}^T$

$$\mathbf{J}_{L-2}\mathbf{J}_{L-2}^T \succsim \mathbf{J}_{FB}\mathbf{J}_{FB}^T \succsim \tilde{\mathbf{J}}_{FB}\tilde{\mathbf{J}}_{FB}^T$$

- $(\mathbf{J}_{FB})_{i,:} = \tilde{\mathbf{f}}_{L-2}(\mathbf{x}_i) \otimes \tilde{\mathbf{b}}_{L-1}(\mathbf{x}_i)$
- $\tilde{\mathbf{J}}_{FB} = \mathbf{J}_{FB} - \mathbb{E}_{\mathbf{X}}\mathbf{J}_{FB}$

Center **again** the whole matrix



Bounding $\lambda_{\min}(\mathbf{K})$ with minimum over-parameterization

$$\mathbf{K} \succeq \mathbf{F}_{L-2}\mathbf{F}_{L-2}^T \circ \mathbf{B}_{L-1}\mathbf{B}_{L-1}^T := \mathbf{J}_{L-2}\mathbf{J}_{L-2}^T$$

- $(\mathbf{J}_{L-2})_{i,:} = \mathbf{f}_{L-2}(\mathbf{x}_i) \otimes \mathbf{b}_{L-1}(\mathbf{x}_i)$
- $\mathbf{f}_{L-2}(\mathbf{x}_i)$ = feature vector at layer $L - 2$ with input \mathbf{x}_i
- $\mathbf{b}_{L-1}(\mathbf{x}_i)$ = back-propagation vector at layer $L - 1$ with input \mathbf{x}_i

Second attempt: directly center $\mathbf{J}_{L-2}\mathbf{J}_{L-2}^T$

$$\mathbf{J}_{L-2}\mathbf{J}_{L-2}^T \succsim \mathbf{J}_{FB}\mathbf{J}_{FB}^T \succsim \tilde{\mathbf{J}}_{FB}\tilde{\mathbf{J}}_{FB}^T \approx \mathbf{I}_n$$



- $(\mathbf{J}_{FB})_{i,:} = \tilde{\mathbf{f}}_{L-2}(\mathbf{x}_i) \otimes \tilde{\mathbf{b}}_{L-1}(\mathbf{x}_i)$
- $\tilde{\mathbf{J}}_{FB} = \mathbf{J}_{FB} - \mathbb{E}_{\mathbf{X}}\mathbf{J}_{FB}$

Concentration for i.i.d. rows with well-controlled $\|\cdot\|_{\psi_1}$
 [Adamczak et al., 2011]



Bounding $\lambda_{\min}(\mathbf{K})$ with minimum over-parameterization

$$\mathbf{K} \succeq \mathbf{F}_{L-2}\mathbf{F}_{L-2}^T \circ \mathbf{B}_{L-1}\mathbf{B}_{L-1}^T := \mathbf{J}_{L-2}\mathbf{J}_{L-2}^T$$

- $(\mathbf{J}_{L-2})_{i,:} = \mathbf{f}_{L-2}(\mathbf{x}_i) \otimes \mathbf{b}_{L-1}(\mathbf{x}_i)$
- $\mathbf{f}_{L-2}(\mathbf{x}_i)$ = feature vector at layer $L - 2$ with input \mathbf{x}_i
- $\mathbf{b}_{L-1}(\mathbf{x}_i)$ = back-propagation vector at layer $L - 1$ with input \mathbf{x}_i

Second attempt: directly center $\mathbf{J}_{L-2}\mathbf{J}_{L-2}^T$

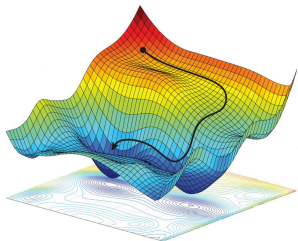
$$\mathbf{J}_{L-2}\mathbf{J}_{L-2}^T \succsim \mathbf{J}_{FB}\mathbf{J}_{FB}^T \succsim \tilde{\mathbf{J}}_{FB}\tilde{\mathbf{J}}_{FB}^T \approx \mathbf{I}_n$$



Concentration for i.i.d. rows with well-controlled $\|\cdot\|_{\psi_1}$
[Adamczak et al., 2011]



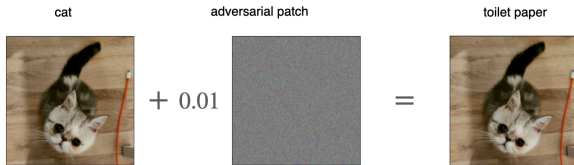
S. Bombari, M. H. Amani, and M. Mondelli, "Memorization and Optimization in Deep Neural Networks with Minimum Over-parameterization", *NeurIPS*, 2022.



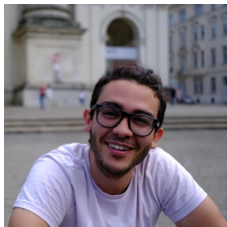
optimization



attention



robustness



Simone Bombari (ISTA)



Shayan Kiyani (ISTA \rightarrow UPenn)

S. Bombari, S. Kiyani, and M. Mondelli, "Beyond the Universal Law of Robustness: Sharper Laws for Random Features and Neural Tangent Kernels", *ICML*, 2023 (oral).

Robust interpolation needs more parameters

$p > n$ enough for interpolation...

Robust interpolation needs more parameters

$p > n$ enough for interpolation...
 but $p > nd$ **necessary** for **robust** interpolation

A Universal Law of Robustness via Isoperimetry

Sébastien Bubeck
 Microsoft Research
 sebubeck@microsoft.com

Mark Sellke
 Stanford University
 msellke@stanford.edu

[Bubeck &
 Sellke, 2021]

$$\frac{1}{n} \sum_{i=1}^n (f(x_i) - y_i)^2 \leq \sigma^2 - \epsilon \Rightarrow \text{Lip}(f) \geq \tilde{\Omega} \left(\epsilon \sqrt{\frac{nd}{p}} \right).$$

Robust interpolation needs more parameters

$p > n$ enough for interpolation...
 but $p > nd$ **necessary** for **robust** interpolation

A Universal Law of Robustness via Isoperimetry

Sébastien Bubeck
 Microsoft Research
 sebubeck@microsoft.com

Mark Sellke
 Stanford University
 msellke@stanford.edu

[Bubeck &
 Sellke, 2021]

$$\frac{1}{n} \sum_{i=1}^n (f(x_i) - y_i)^2 \leq \sigma^2 - \epsilon \Rightarrow \text{Lip}(f) \geq \tilde{\Omega} \left(\epsilon \sqrt{\frac{nd}{p}} \right).$$

[Bubeck et al., 2021] conjecture it is **sufficient** for two-layer networks

Supervised learning

Data: $\{(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_n, y_n)\} \sim_{\text{i.i.d.}} \mathbb{P}(\mathbb{R}^d \times \mathbb{R})$

Goal: Minimize empirical risk $L_f(\boldsymbol{\theta}) = \frac{1}{n} \sum_{i=1}^n (y_i - f(\mathbf{x}_i; \boldsymbol{\theta}))^2$

Gradient flow: $\dot{\boldsymbol{\theta}}(t) = -\nabla_{\boldsymbol{\theta}} L_f(\boldsymbol{\theta}(t))$

Generalized linear regression

Data: $\{(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_n, y_n)\} \sim_{\text{i.i.d.}} \mathbb{P}(\mathbb{R}^d \times \mathbb{R})$

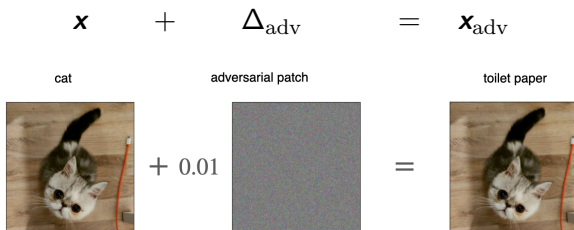
Goal: Minimize empirical risk $L_f(\boldsymbol{\theta}) = \frac{1}{n} \sum_{i=1}^n (y_i - \varphi(\mathbf{x}_i)^\top \boldsymbol{\theta})^2$

- $\varphi : \mathbb{R}^d \rightarrow \mathbb{R}^p$ feature map

Gradient flow solution: $\boldsymbol{\theta}^* = \Phi^\top (\Phi \Phi^\top)^{-1} \mathbf{y}$

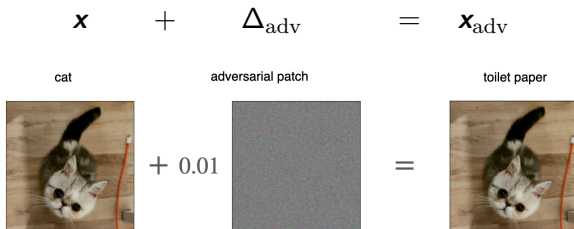
- $\Phi = [\varphi(\mathbf{x}_1), \dots, \varphi(\mathbf{x}_n)]^\top \in \mathbb{R}^{n \times p}$ feature matrix
- $\mathbf{y} = [y_1, \dots, y_n]^\top \in \mathbb{R}^n$ label vector

Sensitivity to adversarial attacks



- $\|\Delta_{\text{adv}}\|_2 \leq \delta \|\mathbf{x}\|_2$ ($\delta = 0.01$)

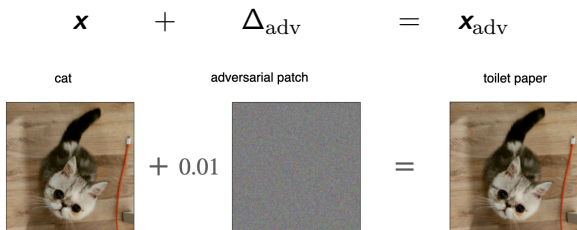
Sensitivity to adversarial attacks



- $\|\Delta_{\text{adv}}\|_2 \leq \delta \|\mathbf{x}\|_2$ ($\delta = 0.01$)

$$|f(\mathbf{x}_{\text{adv}}; \boldsymbol{\theta}) - f(\mathbf{x}; \boldsymbol{\theta})| \approx |\nabla_{\mathbf{x}} f(\mathbf{x}; \boldsymbol{\theta})^\top \Delta_{\text{adv}}|$$

Sensitivity to adversarial attacks



- $\|\Delta_{\text{adv}}\|_2 \leq \delta \|\mathbf{x}\|_2$ ($\delta = 0.01$)

$$|f(\mathbf{x}_{\text{adv}}; \boldsymbol{\theta}) - f(\mathbf{x}; \boldsymbol{\theta})| \approx |\nabla_{\mathbf{x}} f(\mathbf{x}; \boldsymbol{\theta})^\top \Delta_{\text{adv}}| \leq \delta \|\mathbf{x}\|_2 \|\nabla_{\mathbf{x}} f(\mathbf{x}; \boldsymbol{\theta})\|_2$$

- $f(\mathbf{x}; \boldsymbol{\theta}) = \varphi(\mathbf{x})^\top \boldsymbol{\theta}$

Sensitivity to adversarial attacks

$$|f(\mathbf{x}_{\text{adv}}; \boldsymbol{\theta}) - f(\mathbf{x}; \boldsymbol{\theta})| \approx |\nabla_{\mathbf{x}} f(\mathbf{x}; \boldsymbol{\theta})^{\top} \Delta_{\text{adv}}| \leq \delta \|\mathbf{x}\|_2 \|\nabla_{\mathbf{x}} f(\mathbf{x}; \boldsymbol{\theta})\|_2$$

- $f(\mathbf{x}; \boldsymbol{\theta}) = \varphi(\mathbf{x})^{\top} \boldsymbol{\theta}$

Sensitivity: $\mathcal{S}_{\varphi}(\mathbf{x}) = \|\mathbf{x}\|_2 \|\nabla_{\mathbf{x}} \varphi(\mathbf{x})^{\top} \boldsymbol{\theta}^*\|_2$

- $\mathcal{S}_{\varphi}(\mathbf{x}) = O(1) \implies$ model (at interpolation) is **robust**
- $\mathcal{S}_{\varphi}(\mathbf{x}) \gg 1 \implies$ model (at interpolation) is **not robust**

Sensitivity to adversarial attacks

$$|f(\mathbf{x}_{\text{adv}}; \boldsymbol{\theta}) - f(\mathbf{x}; \boldsymbol{\theta})| \approx |\nabla_{\mathbf{x}} f(\mathbf{x}; \boldsymbol{\theta})^{\top} \Delta_{\text{adv}}| \leq \delta \|\mathbf{x}\|_2 \|\nabla_{\mathbf{x}} f(\mathbf{x}; \boldsymbol{\theta})\|_2$$

- $f(\mathbf{x}; \boldsymbol{\theta}) = \varphi(\mathbf{x})^{\top} \boldsymbol{\theta}$

Sensitivity: $\mathcal{S}_{\varphi}(\mathbf{x}) = \|\mathbf{x}\|_2 \|\nabla_{\mathbf{x}} \varphi(\mathbf{x})^{\top} \boldsymbol{\theta}^*\|_2$

- $\mathcal{S}_{\varphi}(\mathbf{x}) = O(1) \implies$ model (at interpolation) is **robust**
- $\mathcal{S}_{\varphi}(\mathbf{x}) \gg 1 \implies$ model (at interpolation) is **not robust**

Having $\|\mathbf{x}\|_2$ on the RHS makes the sensitivity **scale-invariant**

Related work

- Adversarial training (instead of ERM) in linear models
[Donhauser et al., 2021; Javanmard et al., 2020, 2022; Taheri et al., 2020]
- [Bubeck & Sellke, 2021; Bubeck et al., 2021] consider Lipschitz constant
- [Dohmatob & Bietti, 2022; Dohmatob, 2022] consider $\mathbb{E}_{\mathbf{x}} \mathcal{S}_{\varphi}^2(\mathbf{x})$:
 - The former in the infinite-data regime ($n \rightarrow \infty$)
 - The latter in the infinite-width regime ($p \rightarrow \infty$) or proportional regime ($n = \Theta(p) = \Theta(d)$)
- [Zhu et al., 2022] consider average robustness
 $\mathbb{E}_{\mathbf{x}, \hat{\mathbf{x}}, \theta} \nabla_{\mathbf{x}} f(\mathbf{x}; \theta)^{\top} (\mathbf{x} - \hat{\mathbf{x}})$

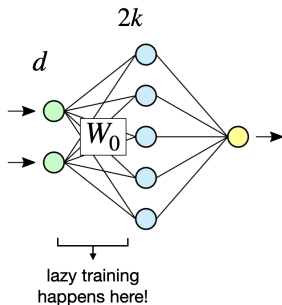
NTK features

$$f_{\text{NN}}(\mathbf{x}; \mathbf{W}) = \sum_{i=1}^k \phi(\mathbf{w}_i^T \mathbf{x}) - \sum_{i=k+1}^{2k} \phi(\mathbf{w}_i^T \mathbf{x})$$

NTK features

$$f_{\text{NN}}(\mathbf{x}; \mathbf{W}) = \sum_{i=1}^k \phi(\mathbf{w}_i^T \mathbf{x}) - \sum_{i=k+1}^{2k} \phi(\mathbf{w}_i^T \mathbf{x})$$

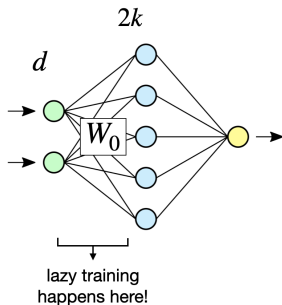
$$f_{\text{NTK}}(\mathbf{x}; \boldsymbol{\theta}) = \boldsymbol{\varphi}_{\text{NTK}}(\mathbf{x})^T \boldsymbol{\theta}, \quad \boldsymbol{\varphi}_{\text{NTK}}(\mathbf{x}) = \nabla_{\mathbf{W}} f_{\text{NN}}(\mathbf{x}; \mathbf{W}) \Big|_{\mathbf{W}=\mathbf{W}_0}$$



NTK features

$$f_{\text{NN}}(\mathbf{x}; \mathbf{W}) = \sum_{i=1}^k \phi(\mathbf{w}_i^T \mathbf{x}) - \sum_{i=k+1}^{2k} \phi(\mathbf{w}_i^T \mathbf{x})$$

$$f_{\text{NTK}}(\mathbf{x}; \boldsymbol{\theta}) = \varphi_{\text{NTK}}(\mathbf{x})^T \boldsymbol{\theta}, \quad \varphi_{\text{NTK}}(\mathbf{x}) = \mathbf{x} \otimes \phi'(\mathbf{W}_0 \mathbf{x})$$



NTK features are robust

Theorem [Bombari, Kiyani, and M., 2023]

Let $\varphi_{\text{NTK}}(\mathbf{x}) = \mathbf{x} \otimes \phi'(\mathbf{W}_0 \mathbf{x})$. Assume $p \gg n$, $k = O(d)$, $n = O(k)$, ϕ even and smooth. Then, with high probability,

$$\mathcal{S}_{\text{NTK}}(\mathbf{x}) = \tilde{O} \left(\sqrt{\frac{nd}{p}} \right).$$

NTK features are robust

Theorem [Bombari, Kiyani, and M., 2023]

Let $\varphi_{\text{NTK}}(\mathbf{x}) = \mathbf{x} \otimes \phi'(\mathbf{W}_0 \mathbf{x})$. Assume $p \gg n$, $k = O(d)$, $n = O(k)$, ϕ even and smooth. Then, with high probability,

$$\mathcal{S}_{\text{NTK}}(\mathbf{x}) = \tilde{O} \left(\sqrt{\frac{nd}{p}} \right).$$

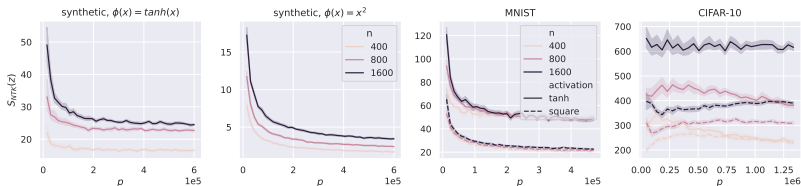
Saturates lower bound!

NTK features are robust

Theorem [Bombari, Kiyani, and M., 2023]

Let $\varphi_{\text{NTK}}(\mathbf{x}) = \mathbf{x} \otimes \phi'(\mathbf{W}_0 \mathbf{x})$. Assume $p \gg n$, $k = O(d)$, $n = O(k)$, ϕ even and smooth. Then, with high probability,

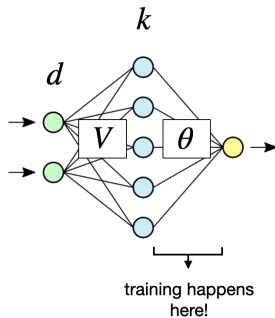
$$\mathcal{S}_{\text{NTK}}(\mathbf{x}) = \tilde{O}\left(\sqrt{\frac{nd}{p}}\right).$$



- Even activations more robust than odd ones.

Random features

$$f_{\text{RF}}(\mathbf{x}; \boldsymbol{\theta}) = \varphi_{\text{RF}}(\mathbf{x})^\top \boldsymbol{\theta}, \quad \varphi_{\text{RF}}(\mathbf{x}) = \phi(\mathbf{V}\mathbf{x})$$



Random features are not robust

Theorem [Bombari, Kiyani, and M., 2023]

Let $\varphi_{\text{RF}}(\mathbf{x}) = \phi(\mathbf{V}\mathbf{x})$. Assume $p \gg n$, $p \gg d$, $d \gg n^{2/3}$, ϕ smooth and $\mathbb{E}_{\rho \sim \mathcal{N}(0,1)} \phi'(\rho) \neq 0$. Then, with high probability,

$$\mathcal{S}_{\text{RF}}(\mathbf{x}) = \Omega\left(n^{1/6}\right) \gg 1.$$

Random features are not robust

Theorem [Bombari, Kiyani, and M., 2023]

Let $\varphi_{\text{RF}}(\mathbf{x}) = \phi(\mathbf{V}\mathbf{x})$. Assume $p \gg n$, $p \gg d$, $d \gg n^{2/3}$, ϕ smooth and $\mathbb{E}_{\rho \sim \mathcal{N}(0,1)} \phi'(\rho) \neq 0$. Then, with high probability,

$$\mathcal{S}_{\text{RF}}(\mathbf{x}) = \Omega\left(n^{1/6}\right) \gg 1.$$

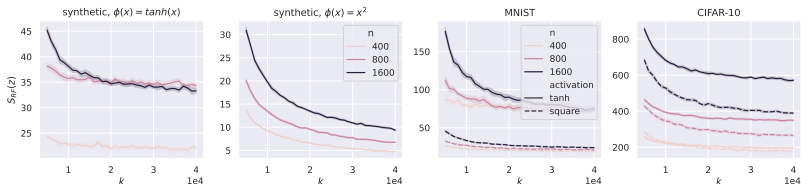
Never robust, regardless of over-parameterization!

Random features are not robust

Theorem [Bombari, Kiyani, and M., 2023]

Let $\varphi_{\text{RF}}(\mathbf{x}) = \phi(\mathbf{V}\mathbf{x})$. Assume $p \gg n$, $p \gg d$, $d \gg n^{2/3}$, ϕ smooth and $\mathbb{E}_{\rho \sim \mathcal{N}(0,1)} \phi'(\rho) \neq 0$. Then, with high probability,

$$\mathcal{S}_{\text{RF}}(\mathbf{x}) = \Omega\left(n^{1/6}\right) \gg 1.$$



- Having $\mathbb{E}_{\rho \sim \mathcal{N}(0,1)} \phi'(\rho) = 0$ improves robustness.

Proof ideas for NTK

$$\mathcal{S}_{\text{NTK}}(\mathbf{x}) = \|\mathbf{x}\|_2 \|\nabla_{\mathbf{x}} \varphi_{\text{NTK}}(\mathbf{x})^\top \boldsymbol{\theta}_{\text{NTK}}^*\|_2$$

Proof ideas for NTK

$$\begin{aligned}\mathcal{S}_{\text{NTK}}(\mathbf{x}) &= \|\mathbf{x}\|_2 \|\nabla_{\mathbf{x}} \varphi_{\text{NTK}}(\mathbf{x})^\top \boldsymbol{\theta}_{\text{NTK}}^*\|_2 \\ &= \|\mathbf{x}\|_2 \|\nabla_{\mathbf{x}} \varphi_{\text{NTK}}(\mathbf{x})^\top \Phi_{\text{NTK}}^\top (\Phi_{\text{NTK}} \Phi_{\text{NTK}}^\top)^{-1} \mathbf{y}\|_2\end{aligned}$$

- $\Phi_{\text{NTK}} = [\varphi_{\text{NTK}}(\mathbf{x}_1), \dots, \varphi_{\text{NTK}}(\mathbf{x}_n)]^\top$

Proof ideas for NTK

$$\begin{aligned}
 \mathcal{S}_{\text{NTK}}(\mathbf{x}) &= \|\mathbf{x}\|_2 \|\nabla_{\mathbf{x}} \varphi_{\text{NTK}}(\mathbf{x})^\top \boldsymbol{\theta}_{\text{NTK}}^*\|_2 \\
 &= \|\mathbf{x}\|_2 \|\nabla_{\mathbf{x}} \varphi_{\text{NTK}}(\mathbf{x})^\top \Phi_{\text{NTK}}^\top (\Phi_{\text{NTK}} \Phi_{\text{NTK}}^\top)^{-1} \mathbf{y}\|_2 \\
 &\leq \|\mathbf{x}\|_2 \|\mathcal{I}_{\text{NTK}}\|_{\text{op}} \lambda_{\min}^{-1} \left(\Phi_{\text{NTK}} \Phi_{\text{NTK}}^\top \right) \|\mathbf{y}\|_2
 \end{aligned}$$

- $\Phi_{\text{NTK}} = [\varphi_{\text{NTK}}(\mathbf{x}_1), \dots, \varphi_{\text{NTK}}(\mathbf{x}_n)]^\top$
- $\mathcal{I}_{\text{NTK}} = \nabla_{\mathbf{x}} \varphi_{\text{NTK}}(\mathbf{x})^\top \Phi_{\text{NTK}}^\top$ **interaction matrix**

Proof ideas for NTK

$$\begin{aligned}
 \mathcal{S}_{\text{NTK}}(\mathbf{x}) &= \|\mathbf{x}\|_2 \|\nabla_{\mathbf{x}} \varphi_{\text{NTK}}(\mathbf{x})^\top \boldsymbol{\theta}_{\text{NTK}}^*\|_2 \\
 &= \|\mathbf{x}\|_2 \|\nabla_{\mathbf{x}} \varphi_{\text{NTK}}(\mathbf{x})^\top \Phi_{\text{NTK}}^\top (\Phi_{\text{NTK}} \Phi_{\text{NTK}}^\top)^{-1} \mathbf{y}\|_2 \\
 &\leq \|\mathbf{x}\|_2 \|\mathcal{I}_{\text{NTK}}\|_{\text{op}} \lambda_{\min}^{-1} \left(\Phi_{\text{NTK}} \Phi_{\text{NTK}}^\top \right) \|\mathbf{y}\|_2
 \end{aligned}$$

- $\Phi_{\text{NTK}} = [\varphi_{\text{NTK}}(\mathbf{x}_1), \dots, \varphi_{\text{NTK}}(\mathbf{x}_n)]^\top$
- $\mathcal{I}_{\text{NTK}} = \nabla_{\mathbf{x}} \varphi_{\text{NTK}}(\mathbf{x})^\top \Phi_{\text{NTK}}^\top$ **interaction matrix**

$\|\mathcal{I}_{\text{NTK}}\|_{\text{op}}$ computed explicitly for even ϕ

$\lambda_{\min} (\Phi_{\text{NTK}} \Phi_{\text{NTK}}^\top)$ is the smallest NTK eigenvalue



Proof ideas for NTK

$$\begin{aligned}
 \mathcal{S}_{\text{NTK}}(\mathbf{x}) &= \|\mathbf{x}\|_2 \|\nabla_{\mathbf{x}} \varphi_{\text{NTK}}(\mathbf{x})^\top \boldsymbol{\theta}_{\text{NTK}}^*\|_2 \\
 &= \|\mathbf{x}\|_2 \|\nabla_{\mathbf{x}} \varphi_{\text{NTK}}(\mathbf{x})^\top \Phi_{\text{NTK}}^\top (\Phi_{\text{NTK}} \Phi_{\text{NTK}}^\top)^{-1} \mathbf{y}\|_2 \\
 &\leq \|\mathbf{x}\|_2 \|\mathcal{I}_{\text{NTK}}\|_{\text{op}} \lambda_{\min}^{-1} \left(\Phi_{\text{NTK}} \Phi_{\text{NTK}}^\top \right) \|\mathbf{y}\|_2 \\
 &= \tilde{O} \left(\sqrt{\frac{nd}{p}} \right)
 \end{aligned}$$

- $\Phi_{\text{NTK}} = [\varphi_{\text{NTK}}(\mathbf{x}_1), \dots, \varphi_{\text{NTK}}(\mathbf{x}_n)]^\top$
- $\mathcal{I}_{\text{NTK}} = \nabla_{\mathbf{x}} \varphi_{\text{NTK}}(\mathbf{x})^\top \Phi_{\text{NTK}}^\top$ **interaction matrix**

$\|\mathcal{I}_{\text{NTK}}\|_{\text{op}}$ computed explicitly for even ϕ

$\lambda_{\min} (\Phi_{\text{NTK}} \Phi_{\text{NTK}}^\top)$ is the smallest NTK eigenvalue



Proof ideas (RF)

A bit more involved...

$$\mathcal{S}_{\text{RF}}(\mathbf{x}) = \Omega \left(\|\mathbf{x}\|_2 \|\tilde{\mathcal{L}}_{\text{RF}}\|_F \lambda_{\max}^{-1} \left(\tilde{\Phi}_{\text{RF}} \tilde{\Phi}_{\text{RF}}^T \right) \right)$$

- Remove low-rank components by centering:
 $\tilde{\Phi}_{\text{RF}} = \Phi_{\text{RF}} - \mathbb{E}_{\mathbf{X}} [\Phi_{\text{RF}}]$
- Interaction matrix captures the effect of the activation

$$\|\tilde{\mathcal{L}}_{\text{RF}}\|_F = \frac{k\sqrt{n}}{\sqrt{d}} \left(\mathbb{E}_{\rho \sim \mathcal{N}(0,1)}^2 \phi'(\rho) + o(1) \right)$$



Proof ideas (RF)

A bit more involved...

$$\mathcal{S}_{\text{RF}}(\mathbf{x}) = \Omega \left(\|\mathbf{x}\|_2 \|\tilde{\mathcal{L}}_{\text{RF}}\|_F \lambda_{\max}^{-1} \left(\tilde{\Phi}_{\text{RF}} \tilde{\Phi}_{\text{RF}}^T \right) \right) = \Omega(n^{1/6})$$

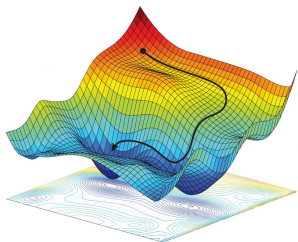
- Remove low-rank components by centering:

$$\tilde{\Phi}_{\text{RF}} = \Phi_{\text{RF}} - \mathbb{E}_{\mathbf{x}} [\Phi_{\text{RF}}]$$

- Interaction matrix captures the effect of the activation

$$\|\tilde{\mathcal{L}}_{\text{RF}}\|_F = \frac{k\sqrt{n}}{\sqrt{d}} \left(\mathbb{E}_{\rho \sim \mathcal{N}(0,1)} \phi'(\rho)^2 + o(1) \right)$$

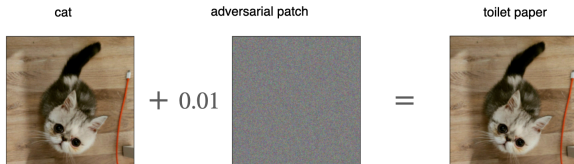




optimization



attention



robustness



Simone Bombari (ISTA)

S. Bombari and M. Mondelli, "Towards Understanding the Word Sensitivity of Attention Layers: A Study via Random Features", *arXiv preprint*, 2024.

Changing a word changes the meaning

Prompt	Output
Reply with "Yes" if the review I will provide you is positive , and "No" otherwise. Review: Sorry, gave it a 1, which is the rating I give to movies on which I walk out or fall asleep.	No
Reply with "Yes" if the review I will provide you is negative , and "No" otherwise. Review: Sorry, gave it a 1, which is the rating I give to movies on which I walk out or fall asleep.	Yes

- Different output of the Llama2-7b-chat model

Changing a word changes the scores



- Different attention score pattern of BERT-Base model

Changing a word changes the scores



- Different attention score pattern of BERT-Base model

Transformers capture the effect of changing a single word in a sentence.

Insights from random features

Data: $\mathbf{X} = [\mathbf{x}_1, \dots, \mathbf{x}_N]^T \in \mathbb{R}^{N \times d}$

- Tokens $\{\mathbf{x}_i\}_{i=1}^N$
- $N =$ context length, $d =$ embedding dimension

Insights from random features

Data: $\mathbf{X} = [\mathbf{x}_1, \dots, \mathbf{x}_N]^T \in \mathbb{R}^{N \times d}$

- Tokens $\{\mathbf{x}_i\}_{i=1}^N$
- $N =$ context length, $d =$ embedding dimension

Random features: $\varphi_{\text{RF}} : \mathbb{R}^{N \times d} \rightarrow \mathbb{R}^k$

$$\varphi_{\text{RF}}(\mathbf{X}) = \phi(\mathbf{V} \text{flat}(\mathbf{X}))$$

Insights from random features

Data: $\mathbf{X} = [\mathbf{x}_1, \dots, \mathbf{x}_N]^T \in \mathbb{R}^{N \times d}$

- Tokens $\{\mathbf{x}_i\}_{i=1}^N$
- $N =$ context length, $d =$ embedding dimension

Random features: $\varphi_{\text{RF}} : \mathbb{R}^{N \times d} \rightarrow \mathbb{R}^k$

$$\varphi_{\text{RF}}(\mathbf{X}) = \phi(\mathbf{V} \text{flat}(\mathbf{X}))$$

Random attention features: $\varphi_{\text{QKV}} : \mathbb{R}^{N \times d} \rightarrow \mathbb{R}^{N \times d'}$

$$\varphi_{\text{QKV}}(\mathbf{X}) = \text{softmax} \left(\frac{\mathbf{X} \mathbf{W}_Q^T \mathbf{W}_K \mathbf{X}^T}{\sqrt{d'}} \right) \mathbf{X} \mathbf{W}_V^T$$

- $\text{softmax}(\mathbf{s})_i = e^{s_i} / \sum_j e^{s_j}$
- $\mathbf{W}_Q, \mathbf{W}_K, \mathbf{W}_V \in \mathbb{R}^{d' \times d}$ queries, keys and values matrices

Insights from random features

Data: $\mathbf{X} = [\mathbf{x}_1, \dots, \mathbf{x}_N]^T \in \mathbb{R}^{N \times d}$

- Tokens $\{\mathbf{x}_i\}_{i=1}^N$
- $N =$ context length, $d =$ embedding dimension

Random features: $\varphi_{\text{RF}} : \mathbb{R}^{N \times d} \rightarrow \mathbb{R}^k$

$$\varphi_{\text{RF}}(\mathbf{X}) = \phi(\mathbf{V} \text{flat}(\mathbf{X}))$$

Random attention features: $\varphi_{\text{RAF}} : \mathbb{R}^{N \times d} \rightarrow \mathbb{R}^{N \times d}$

$$\varphi_{\text{RAF}}(\mathbf{X}) = \text{softmax} \left(\frac{\mathbf{X} \mathbf{W} \mathbf{X}^T}{\sqrt{d}} \right) \mathbf{X}$$

- $\text{softmax}(\mathbf{s})_i = e^{s_i} / \sum_j e^{s_j}$

Sample complexity comparison between RF and RAF in (Fu et al., 2023)

Sensitivity to changing a word

Word sensitivity: $\mathcal{WS}_\varphi(\mathbf{X}) = \sup_{j \in [N], \|\Delta\|_2 \leq \sqrt{d}} \frac{\|\varphi(\mathbf{X}^j(\Delta)) - \varphi(\mathbf{X})\|_2}{\|\varphi(\mathbf{X})\|_2}$

- $\varphi(\mathbf{X}^j(\Delta)) = \mathbf{X} + \mathbf{e}_j \Delta^\top$ (only j -th token changed)
- $\|\Delta\|_2 \leq \sqrt{d}$ (perturbation size bounded by token size)

Low WS of random features

Theorem [Bombari and M., 2024]

Assume ϕ Lipschitz and $k = \Omega(Nd)$. Then, with high probability over \mathbf{V} ,

$$\mathcal{WS}_{\text{RF}}(\mathbf{X}) = O\left(\frac{1}{\sqrt{N}}\right)$$

Word sensitivity vanishes as context length N grows

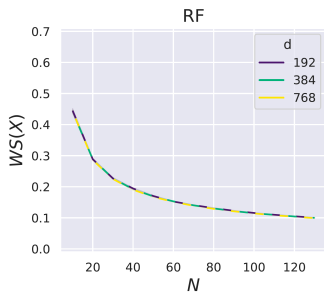
Low WS of random features

Theorem [Bombari and M., 2024]

Assume ϕ Lipschitz and $k = \Omega(Nd)$. Then, with high probability over \mathbf{V} ,

$$\mathcal{WS}_{\text{RF}}(\mathbf{X}) = O\left(\frac{1}{\sqrt{N}}\right)$$

Word sensitivity vanishes as context length N grows



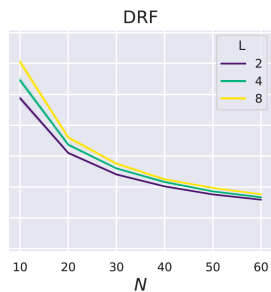
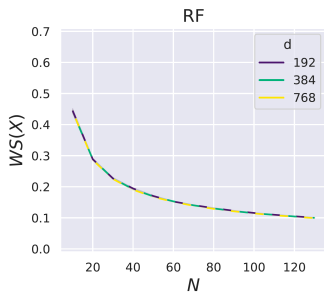
Low WS of random features

Theorem [Bombari and M., 2024]

Assume ϕ Lipschitz and $k = \Theta(Nd)$. Then, with high probability over $\mathbf{V}_1, \dots, \mathbf{V}_L$,

$$WS_{\text{DRF}}(\mathbf{X}) = O\left(\frac{e^{CL}}{\sqrt{N}}\right)$$

Word sensitivity vanishes as context length N grows



High WS of random attention features

Theorem [Bombari and M., 2024]

Assume $d = \tilde{\Omega}(N)$. Then, with high probability over \mathbf{W} ,

$$\mathcal{WS}_{\text{RAF}}(\mathbf{X}) = \Omega(1)$$

High word sensitivity regardless of the context length N

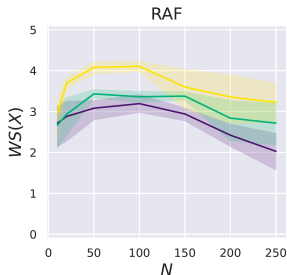
High WS of random attention features

Theorem [Bombari and M., 2024]

Assume $d = \tilde{\Omega}(N)$. Then, with high probability over \mathbf{W} ,

$$\mathcal{WS}_{\text{RAF}}(\mathbf{X}) = \Omega(1)$$

High word sensitivity regardless of the context length N



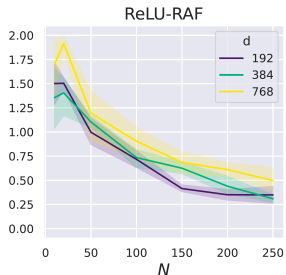
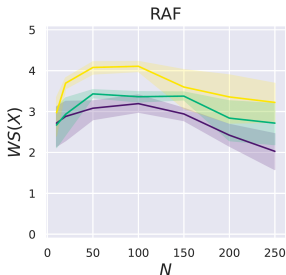
High WS of random attention features

Theorem [Bombari and M., 2024]

Assume $d = \tilde{\Omega}(N)$. Then, with high probability over \mathbf{W} ,

$$\mathcal{WS}_{\text{RAF}}(\mathbf{X}) = \Omega(1)$$

High word sensitivity regardless of the context length N



Word sensitivity decreases when replacing softmax with ReLU

Proof ideas for RAF

1. Find a direction δ^* aligned with many words \mathbf{x}_i 's.
- Constant fraction of the entries of $\mathbf{X}\delta^*$ is $\Omega(d/\sqrt{N})$ in modulus

Proof ideas for RAF

- Find a direction δ^* aligned with many words \mathbf{x}_i 's.
 - Constant fraction of the entries of $\mathbf{X}\delta^*$ is $\Omega(d/\sqrt{N})$ in modulus
- Exhibit two different directions Δ_1^* and Δ_2^* both aligned with many words in the feature space $\{\mathbf{W}^\top \mathbf{x}_i\}_{i=1}^N$.
 - $\|\Delta_1^* - \Delta_2^*\|_2 = \Omega(\sqrt{d})$
 - Constant fraction of the entries of $\mathbf{XW}\Delta_k^*/\sqrt{d}$ is $\Omega(\log^2 d)$

Proof ideas for RAF

- Find a direction δ^* aligned with many words \mathbf{x}_i 's.
 - Constant fraction of the entries of $\mathbf{X}\delta^*$ is $\Omega(d/\sqrt{N})$ in modulus
- Exhibit two different directions Δ_1^* and Δ_2^* both aligned with many words in the feature space $\{\mathbf{W}^\top \mathbf{x}_i\}_{i=1}^N$.
 - $\|\Delta_1^* - \Delta_2^*\|_2 = \Omega(\sqrt{d})$
 - Constant fraction of the entries of $\mathbf{XW}\Delta_k^*/\sqrt{d}$ is $\Omega(\log^2 d)$
- Attention concentrates towards the perturbed word.
 - Constant fraction of rows of $\text{softmax}(\mathbf{X}^j(\Delta_k^*)\mathbf{W}(\mathbf{X}^j(\Delta_k^*))^\top/\sqrt{d}) \approx \mathbf{e}_j$

Key role of softmax

Proof ideas for RAF

- Find a direction δ^* aligned with many words \mathbf{x}_i 's.
 - Constant fraction of the entries of $\mathbf{X}\delta^*$ is $\Omega(d/\sqrt{N})$ in modulus
- Exhibit two different directions Δ_1^* and Δ_2^* both aligned with many words in the feature space $\{\mathbf{W}^\top \mathbf{x}_i\}_{i=1}^N$.
 - $\|\Delta_1^* - \Delta_2^*\|_2 = \Omega(\sqrt{d})$
 - Constant fraction of the entries of $\mathbf{XW}\Delta_k^*/\sqrt{d}$ is $\Omega(\log^2 d)$
- Attention concentrates towards the perturbed word.
 - Constant fraction of rows of $\text{softmax}(\mathbf{X}^j(\Delta_k^*)\mathbf{W}(\mathbf{X}^j(\Delta_k^*))^\top/\sqrt{d}) \approx \mathbf{e}_j$

Key role of softmax

- Conclude with at least one perturbation between Δ_1^* and Δ_2^* .
 - $\|\varphi_{\text{RAF}}(\mathbf{X}) - \varphi_{\text{RAF}}(\mathbf{X}^j(\Delta_k^*))\|_F = o(\sqrt{dN})$ for $k = 1, 2 \Rightarrow \|\Delta_1^* - \Delta_2^*\|_2 = o(\sqrt{d})$, which is a contradiction

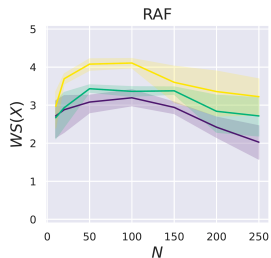
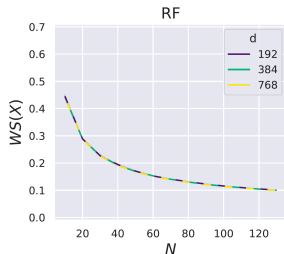
Generalization on context modification

1. Random features have low word sensitivity:

$$\mathcal{WS}_{\text{RF}}(\mathbf{X}) = O(1/\sqrt{N}).$$

2. Random attention features have high word sensitivity:

$$\mathcal{WS}_{\text{RAF}}(\mathbf{X}) = \Omega(1).$$



Idea: random features cannot learn to distinguish \mathbf{X} and $\mathbf{X}^j(\Delta)$, while random attention features can!

Generalized linear regression

Data: $(\mathbf{X}_1, y_1), \dots, (\mathbf{X}_n, y_n) \in \mathbb{R}^{N \times d} \times \{-1, 1\}$

Goal: Minimize empirical risk $L(\boldsymbol{\theta}) = \frac{1}{n} \sum_{i=1}^n \left(y_i - \varphi(\mathbf{X}_i)^\top \boldsymbol{\theta} \right)^2$

- $\varphi : \mathbb{R}^{N \times d} \rightarrow \mathbb{R}^p$ feature map

Gradient descent solution: $\boldsymbol{\theta}^* = \boldsymbol{\theta}_0 + \Phi^\top (\Phi \Phi^\top)^{-1} (\mathbf{y} - \Phi \boldsymbol{\theta}_0)$

- $\boldsymbol{\theta}_0$ initialization
- $\Phi = [\varphi(\mathbf{X}_1), \dots, \varphi(\mathbf{X}_n)]^\top$ feature matrix
- $\mathbf{y} = [y_1, \dots, y_n]^\top$ label vector

More training and generalization

Does further training on (\mathbf{X}, \tilde{y}) allow to generalize on $(\mathbf{X}^j(\Delta), -\tilde{y})$?

	Prompt		Output
\mathbf{X}	<hr/> Reply with "Yes" if the review I will provide you is positive , and "No" otherwise. Review: Sorry, gave it a 1, which is the rating I give to movies on which I walk out or fall asleep. <hr/>	\tilde{y}	<hr/> No <hr/>
$\mathbf{X}^j(\Delta)$	<hr/> Reply with "Yes" if the review I will provide you is negative , and "No" otherwise. Review: Sorry, gave it a 1, which is the rating I give to movies on which I walk out or fall asleep. <hr/>	$-\tilde{y}$	<hr/> Yes <hr/>

More training and generalization

Does further training on (\mathbf{X}, \tilde{y}) allow to generalize on $(\mathbf{X}^j(\Delta), -\tilde{y})$?

Fine-tuning. Initialize with θ^* and train only the extra sample:

$$\theta_f^* = \theta^* + \frac{\varphi(\mathbf{X})}{\|\varphi(\mathbf{X})\|_2^2} \left(\tilde{y} - \varphi(\mathbf{X})^\top \theta^* \right).$$

More training and generalization

Does further training on (\mathbf{X}, \tilde{y}) allow to generalize on $(\mathbf{X}^j(\Delta), -\tilde{y})$?

Fine-tuning. Initialize with θ^* and train only the extra sample:

$$\theta_f^* = \theta^* + \frac{\varphi(\mathbf{X})}{\|\varphi(\mathbf{X})\|_2^2} \left(\tilde{y} - \varphi(\mathbf{X})^\top \theta^* \right).$$

Re-training. Add (\mathbf{X}, \tilde{y}) to training set and train from scratch:

$$\theta_r^* = \theta_0 + \Phi_r^\top (\Phi_r \Phi_r^\top)^{-1} (\mathbf{y}_r - \Phi_r \theta_0).$$

- θ_0 initialization
- $\Phi_r = [\varphi(\mathbf{X}_1), \dots, \varphi(\mathbf{X}_n), \varphi(\mathbf{X})]^\top$ feature matrix
- $\mathbf{y}_r = [y_1, \dots, y_n, \tilde{y}]^\top$ label vector

Random features do not generalize...

Theorem [Bombari and M., 2024]

Let $|\varphi(\mathbf{X}^j(\Delta))^T \boldsymbol{\theta}^* - \varphi(\mathbf{X})^T \boldsymbol{\theta}^*| \leq \gamma$ for $\gamma \in [0, 2)$.

Random features do not generalize...

Theorem [Bombari and M., 2024]

Let $|\varphi(\mathbf{X}^j(\Delta))^T \boldsymbol{\theta}^* - \varphi(\mathbf{X})^T \boldsymbol{\theta}^*| \leq \gamma$ for $\gamma \in [0, 2)$. Then, under some technical assumptions, with high probability,

$$\text{Err}_{\text{RF}}(\mathbf{X}^j(\Delta), \boldsymbol{\theta}_{f/r}^*) := \left(\varphi(\mathbf{X}^j(\Delta))^T \boldsymbol{\theta}_{f/r}^* + \tilde{y} \right)^2 > (2-\gamma)^2 - \mathcal{O}\left(\frac{1}{\sqrt{N}}\right)$$

Random features do not generalize...

Theorem [Bombari and M., 2024]

Let $|\varphi(\mathbf{X}^j(\Delta))^T \boldsymbol{\theta}^* - \varphi(\mathbf{X})^T \boldsymbol{\theta}^*| \leq \gamma$ for $\gamma \in [0, 2)$. Then, under some technical assumptions, with high probability,

$$\text{Err}_{\text{RF}}(\mathbf{X}^j(\Delta), \boldsymbol{\theta}_{f/r}^*) := \left(\varphi(\mathbf{X}^j(\Delta))^T \boldsymbol{\theta}_{f/r}^* + \tilde{y} \right)^2 > (2-\gamma)^2 - O\left(\frac{1}{\sqrt{N}}\right)$$

Unless the correct label is already known ($\gamma = 2$),
fine-tuning/re-training does not help much.

Idea: after perturbing the j -th token, the model cannot move more than its WS, which is $O(1/\sqrt{N})$.

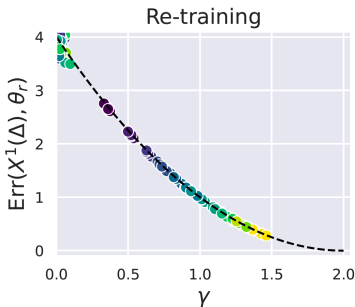
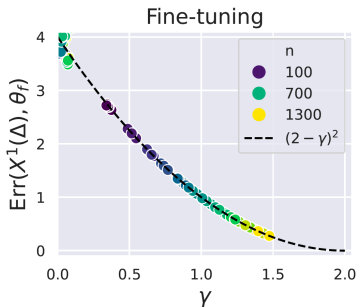


Random features do not generalize...

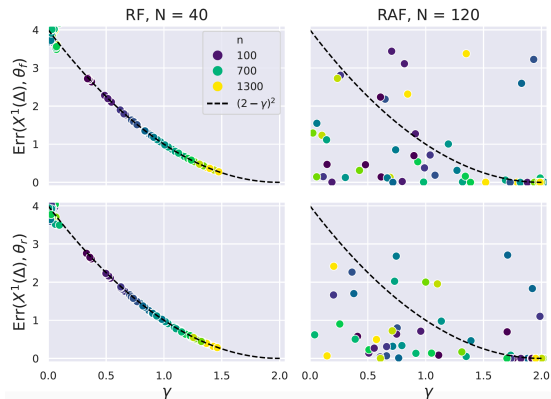
Theorem [Bombari and M., 2024]

Let $|\varphi(\mathbf{X}^j(\Delta))^T \boldsymbol{\theta}^* - \varphi(\mathbf{X})^T \boldsymbol{\theta}^*| \leq \gamma$ for $\gamma \in [0, 2)$. Then, under some technical assumptions, with high probability,

$$\text{Err}_{\text{RF}}(\mathbf{X}^j(\Delta), \boldsymbol{\theta}_{f/r}^*) := \left(\varphi(\mathbf{X}^j(\Delta))^T \boldsymbol{\theta}_{f/r}^* + \tilde{y} \right)^2 > (2-\gamma)^2 - O\left(\frac{1}{\sqrt{N}}\right)$$

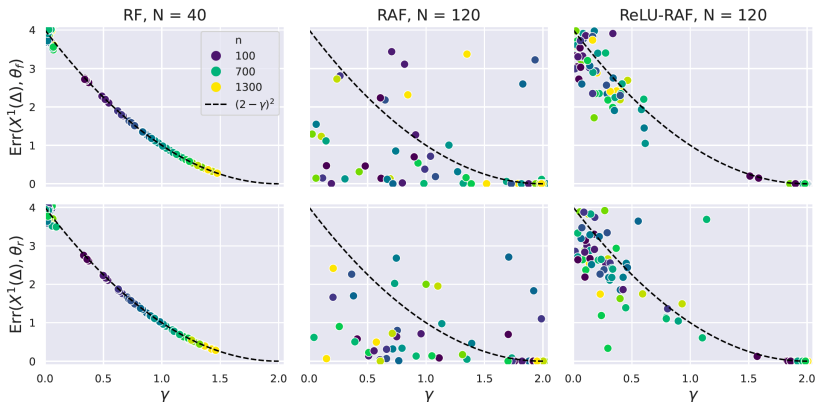


... but random attention features do generalize



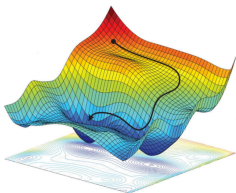
- Random attention features generalize even when \mathbf{X} and $\mathbf{X}^j(\Delta)$ were indistinguishable before the extra training ($\gamma \approx 0$).

... but random attention features do generalize



- Random attention features generalize even when \mathbf{X} and $\mathbf{X}^j(\Delta)$ were indistinguishable before the extra training ($\gamma \approx 0$).
- Replacing softmax with ReLU increases error.

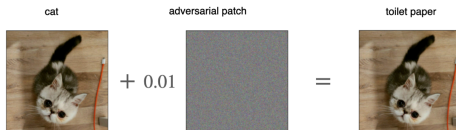
Take home



optimization



attention



robustness

Quantitative understanding via random and NTK features