
Scaling Neural Tangent Kernels via Sketching and Random Features

Insu Han

Yale University

Joint work with Amir Zandieh, Haim Avron,
Neta Shoham, Chaewon Kim, Jinwoo Shin

Outline

Introduction

- Kernel Regression
- Neural Tangent Kernel
- Existing Work

NTK Feature Map Construction

- Exact NTK Computation
- NTK Random Features
- NTK Sketch

Experiments

- Classification on MNIST datasets
- Classification on CIFAR-10 datasets

Conclusion

Outline

Introduction

- Kernel Regression
- Neural Tangent Kernel
- Existing Work

NTK Feature Map Construction

- Exact NTK Computation
- NTK Random Features
- NTK Sketch

Experiments

- Classification on MNIST datasets
- Classification on CIFAR-10 datasets

Conclusion

Kernel Regression

Kernel: a similarity function over pairs of data points in raw representation

- Mercer decomposition: for every kernel $K : \mathbb{R}^d \times \mathbb{R}^d \rightarrow \mathbb{R}$ and $\mathbf{x}, \mathbf{x}' \in \mathbb{R}^d$

$$K(\mathbf{x}, \mathbf{x}') = \langle \phi(\mathbf{x}), \phi(\mathbf{x}') \rangle$$

- ϕ is called a **feature map**

Kernel Regression

Kernel: a similarity function over pairs of data points in raw representation

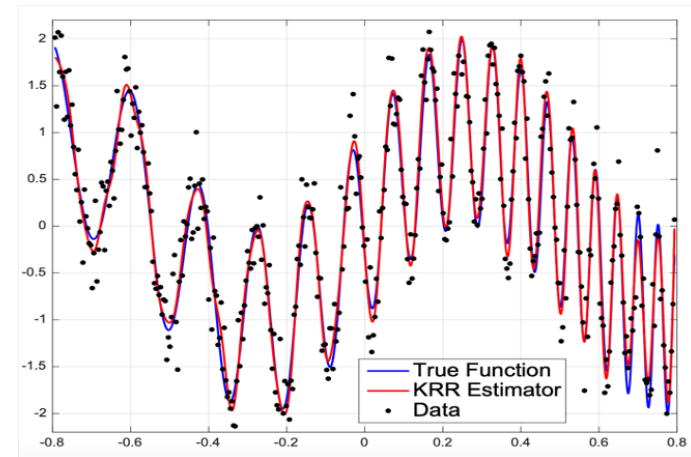
- Mercer decomposition: for every kernel $K : \mathbb{R}^d \times \mathbb{R}^d \rightarrow \mathbb{R}$ and $\mathbf{x}, \mathbf{x}' \in \mathbb{R}^d$

$$K(\mathbf{x}, \mathbf{x}') = \langle \phi(\mathbf{x}), \phi(\mathbf{x}') \rangle$$

- ϕ is called a **feature map**

Kernel ridge regression:

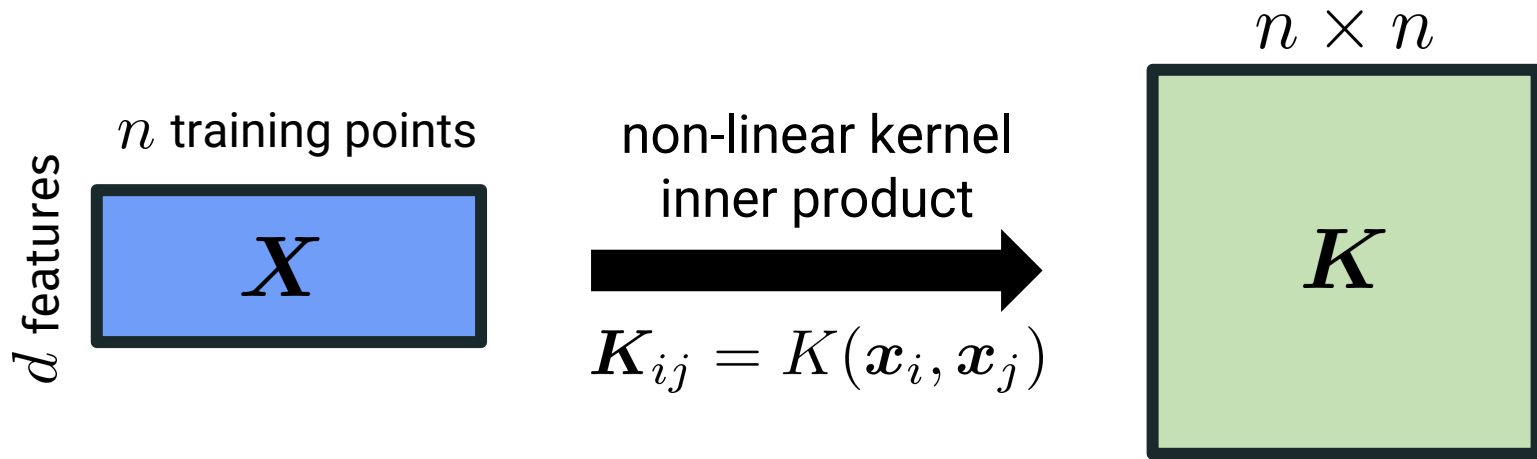
$$\mathbf{w}^* = \operatorname{argmin}_{\mathbf{w}} \frac{1}{n} \sum_{i=1}^n (y_i - \phi(\mathbf{x}_i)^\top \mathbf{w})^2 + \lambda \|\mathbf{w}\|_2^2$$



- Simple and yet powerful tool for learning non-linear relationships between data points

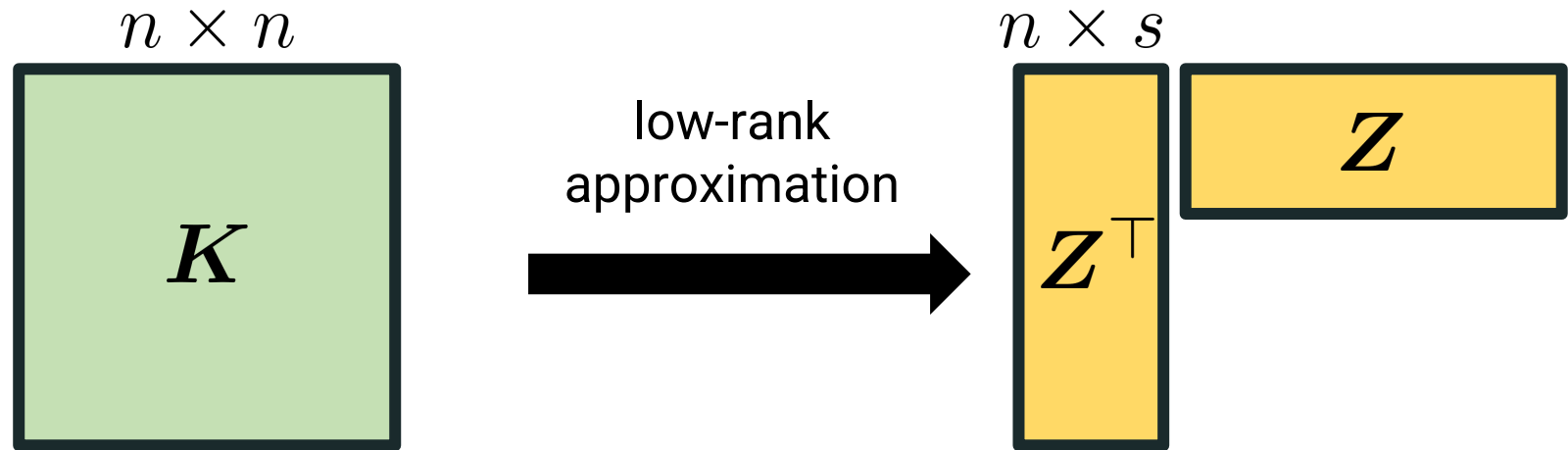
Scalability of Kernel Methods

Kernel methods are expensive



- Computing all kernel entries take $\Omega(n \cdot \text{nnz}(\mathbf{X}) + n^2)$ time
- Even writing it down takes $\Omega(n^2 d)$ time and $\Omega(n^2)$ memory
- A single iteration of a linear system solver takes $\Omega(n^2)$ time
- For $n = 100,000$, K has 10 billion entries. Take 80 GB of storage!

Classical Solution: Dimensionality Reduction



Storing Z uses $\mathcal{O}(ns)$ space and computing $Z^T Z x$ takes $\mathcal{O}(ns)$ time

Orthogonalization, eigen-decomposition and pseudo-inversion of $Z^T Z$ all take just $\mathcal{O}(ns^2)$ time

Neural Tangent Kernel (NTK)

Kernel: a similarity function over pairs of data points in raw representation

- Mercer kernel is a function $K : \mathbb{R}^d \times \mathbb{R}^d \rightarrow \mathbb{R}$ such that for $\mathbf{x}, \mathbf{x}' \in \mathbb{R}^d$

$$K(\mathbf{x}, \mathbf{x}') = \langle \phi(\mathbf{x}), \phi(\mathbf{x}') \rangle$$

- ϕ is called a **feature map** (can be a random function)

Tangent kernel:

- For a smooth function $f(\cdot, \theta)$, tangent kernel is an inner-product of gradient at a fixed θ_0 :

$$K_T(\mathbf{x}, \mathbf{x}'; \theta_0) = \left\langle \left. \frac{\partial f(\mathbf{x}, \theta)}{\partial \theta} \right|_{\theta=\theta_0}, \left. \frac{\partial f(\mathbf{x}', \theta)}{\partial \theta} \right|_{\theta=\theta_0} \right\rangle$$

Neural Tangent Kernel

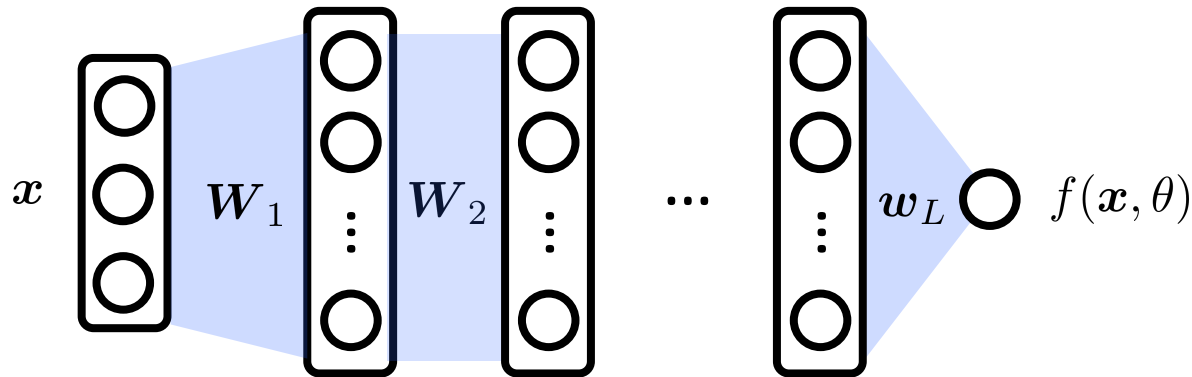
Fully-connected neural network:

$$f(\mathbf{x}, \theta) = \mathbf{h}_L^\top \mathbf{w}_L, \quad \mathbf{h}_\ell = \frac{c_\sigma}{\sqrt{d_\ell}} \sigma(\mathbf{h}_{\ell-1}^\top \mathbf{W}_\ell), \quad \mathbf{h}_0 = \mathbf{x}$$

- $\theta = (\mathbf{W}_1, \dots, \mathbf{w}_L)$: trainable parameters
- $c_\sigma = 1/\mathbb{E}_{z \sim \mathcal{N}}[\sigma(z)^2]$: normalization for activation

Neural Tangent Kernel [JGH18][LXS+19]:

$$K_{\text{NTK}}(\mathbf{x}, \mathbf{x}') = \mathbb{E}_{\theta \sim \mathcal{N}} \left\langle \frac{\partial f(\mathbf{x}, \theta)}{\partial \theta}, \frac{\partial f(\mathbf{x}', \theta)}{\partial \theta} \right\rangle$$



Neural Tangent Kernel

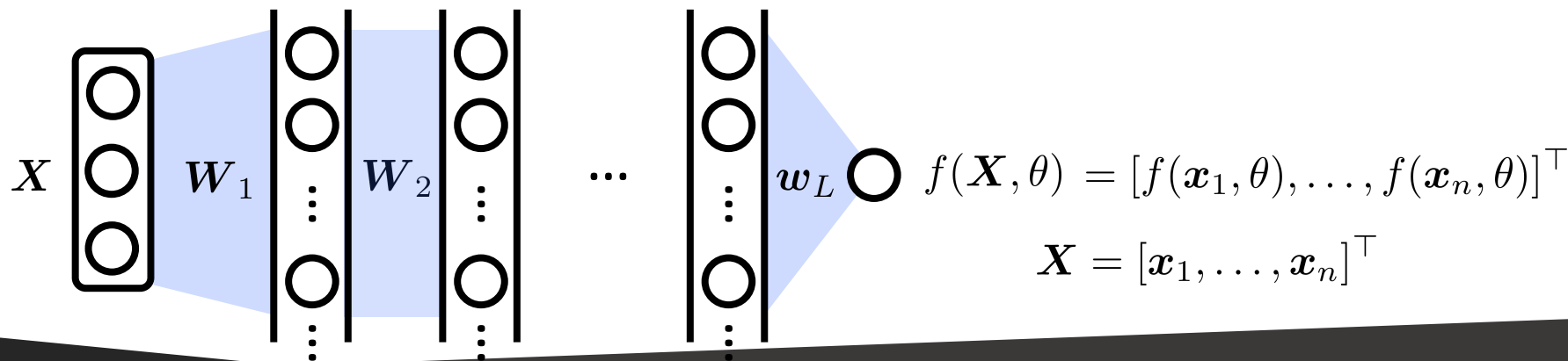
Under infinite width regime, neural network is **equivalent** to NTK

- Prediction of neural networks trained under ℓ_2 -loss by gradient flow
- NTK solved by kernel regression

$$y_{\text{test}} = K_{\text{NTK}}(\mathbf{x}_{\text{test}}, \mathbf{X}) K_{\text{NTK}}(\mathbf{X}, \mathbf{X})^{-1} \mathbf{y}$$

for training data $\mathbf{X} = [\mathbf{x}_1, \dots, \mathbf{x}_n]^\top$, $\mathbf{y} \in \mathbb{R}^n$

- No need to learn parameter θ if $K_{\text{NTK}}(\cdot)$ is given



Neural Tangent Kernel

NTK is useful both in theory and practice:

- Generalization [CG19], optimization [ALS19], regularization [HLY20]
- In practice, NTK can perform better regression results than training neural networks (NN) [ADL+19]

Classifier	Friedman Rank	Average Accuracy	P90	P95	PMA
NTK	28.34	81.95%±14.10%	88.89%	72.22%	95.72% ±5.17%
NN (He init)	40.97	80.88%±14.96%	81.11%	65.56%	94.34% ±7.22%
NN (NTK init)	38.06	81.02%±14.47%	85.56%	60.00%	94.55% ±5.89%
RF	33.51	81.56% ±13.90%	85.56%	67.78%	95.25% ±5.30%
Gaussian Kernel	35.76	81.03% ± 15.09%	85.56%	72.22%	94.56% ±8.22%
Polynomial Kernel	38.44	78.21% ± 20.30%	80.00%	62.22%	91.29% ±18.05%

Comparisons of different classifiers on 90 UCI datasets [ADL+19]

Neural Tangent Kernel

NTK is useful both in theory and practice:

- Generalization [CG19], optimization [ALS19], regularization [HLY20]
- In practice, NTK can perform better regression results than training neural networks (NN) [ADL+19]

Computational issues:

- Even NTK can be computed exactly [ADH+19], solving kernel regression requires $\mathcal{O}(n^3)$ time \Rightarrow **infeasible** for large n
- Convolution NTK (CNTK) with $h \times w$ images takes $\Omega(n^2 h^2 w^2)$ time
E.g., Convolutional NTK with CIFAR-10 \geq **1000** GPU hours [SFG+20]

Neural Tangent Kernel

NTK is useful both in theory and practice:

- Generalization [CG19], optimization [ALS19], regularization [HLY20]
- In practice, NTK can perform better regression results than training neural networks (NN) [ADL+19]

Computational issues:

- Even NTK can be computed exactly [ADH+19], solving kernel regression requires $\mathcal{O}(n^3)$ time \Rightarrow **infeasible** for large n

Kernel approximation by feature map:

$$K_{\text{NTK}}(\mathbf{x}, \mathbf{x}') \approx \langle \phi(\mathbf{x}), \phi(\mathbf{x}') \rangle, \quad \phi : \mathbb{R}^d \rightarrow \mathbb{R}^m$$

- Kernel ridge regression can be solved in $\mathcal{O}(nm^2 + m^3)$ time
- Memory complexity can be $\mathcal{O}(nm + m^2)$
- Usually, $m \ll n$

Existing Work

Goal: feature map construction ϕ such that $K_{\text{NTK}}^{(L)}(\mathbf{x}, \mathbf{x}') \approx \langle \phi(\mathbf{x}), \phi(\mathbf{x}') \rangle$

Gradient features of finitely wide network:

$$K_{\text{NTK}}(\mathbf{x}, \mathbf{x}') := \mathbb{E}_{\theta' \sim \mathcal{N}} \left\langle \frac{\partial f(\mathbf{x}, \theta')}{\partial \theta}, \frac{\partial f(\mathbf{x}', \theta')}{\partial \theta} \right\rangle \quad (\text{Definition})$$
$$\approx \frac{1}{s} \sum_{i=1}^s \left\langle \frac{\partial f(\mathbf{x}, \theta_i)}{\partial \theta}, \frac{\partial f(\mathbf{x}', \theta_i)}{\partial \theta} \right\rangle$$

- $\theta_1, \dots, \theta_s$: random samples from Gaussian distribution
- Gradient features show **poor** performances in practice [ADH+19]

Contributions

Fast and accurate NTK approximation using Random Features/Sketching

	Feature dimension	Runtime
Gradient Features [ADH+19]	$\tilde{\Omega} \left(\frac{L^{13}}{\varepsilon^8} \right)$	$\tilde{\Omega} \left(\frac{L^{13}}{\varepsilon^8} \right)$
NTK Random Features (Our)	$\tilde{\mathcal{O}} \left(\frac{L^6}{\varepsilon^4} \right)$	$\tilde{\mathcal{O}} \left(\frac{L^{13}}{\varepsilon^8} \right)$
NTK Sketch (Our)	$\mathcal{O} \left(\frac{1}{\varepsilon^2} \log \frac{1}{\delta} \right)$	$\tilde{\mathcal{O}} \left(\frac{L^{11}}{\varepsilon^{6.7}} + \frac{L^3}{\varepsilon^2} d \right)$
CNTK Sketch (Our)	$\mathcal{O} \left(\frac{1}{\varepsilon^2} \log \frac{1}{\delta} \right)$	$\tilde{\mathcal{O}} \left(\frac{L^{11}}{\varepsilon^{6.7}} \cdot (d_1 d_2) \right)$

*CNTK with $d_1 \times d_2$ images

- For all methods, the runtime of kernel regression reduces from $\mathcal{O}(n^3 + n^2 d)$ to $\mathcal{O}(n(\text{dimension}^2 + \text{runtime}))$
- For n images, the exact CNTK can be computed in $\mathcal{O}(n^2 L d_1^2 d_2^2)$ time while CNTK Sketch runs in $\tilde{\mathcal{O}}(n L^{11} d_1 d_2 / \varepsilon^{6.7})$ time

Contributions

Fast and accurate NTK approximation using Random Features/Sketching

	Feature dimension	Runtime
Gradient Features [ADH+19]	$\tilde{\Omega} \left(\frac{L^{13}}{\varepsilon^8} \right)$	$\tilde{\Omega} \left(\frac{L^{13}}{\varepsilon^8} \right)$
NTK Random Features (Our)	$\tilde{\mathcal{O}} \left(\frac{L^6}{\varepsilon^4} \right)$	$\tilde{\mathcal{O}} \left(\frac{L^{13}}{\varepsilon^8} \right)$
NTK Sketch (Our)	$\mathcal{O} \left(\frac{1}{\varepsilon^2} \log \frac{1}{\delta} \right)$	$\tilde{\mathcal{O}} \left(\frac{L^{11}}{\varepsilon^{6.7}} + \frac{L^3}{\varepsilon^2} d \right)$
CNTK Sketch (Our)	$\mathcal{O} \left(\frac{1}{\varepsilon^2} \log \frac{1}{\delta} \right)$	$\tilde{\mathcal{O}} \left(\frac{L^{11}}{\varepsilon^{6.7}} \cdot (d_1 d_2) \right)$

*CNTK with $d_1 \times d_2$ images

Better performance than other heuristics under real-world datasets

- NTK Sketch with $L = 3$ achieves 72% of test accuracy of CIFAR-10 while the exact NTK, CNN achieve 70%, 64%, respectively
- CNTK Sketch is **190 times faster** than the exact CNTK

Outline

Introduction

- Kernel Regression
- Neural Tangent Kernel
- Existing Work

NTK Feature Map Construction

- Exact NTK Computation
- NTK Random Features
- NTK Sketch

Experiments

- Classification on MNIST datasets
- Classification on CIFAR-10 datasets

Conclusion

Exact NTK Computation

NTK with ReLU activation [ADH+19]: for layer $\ell = 1, \dots, L$ and $\mathbf{x} \in \mathbb{R}^d$

$$K_{\text{NTK}}^{(\ell)}(\mathbf{x}, \mathbf{x}') = K_{\text{NTK}}^{(\ell-1)}(\mathbf{x}, \mathbf{x}') \cdot \underbrace{\dot{\Sigma}^{(\ell)}(\mathbf{x}, \mathbf{x}') + \Sigma^{(\ell)}(\mathbf{x}, \mathbf{x}')}_{\text{function of } \Sigma^{(\ell-1)}}$$

where $K_{\text{NTK}}^{(0)}(\mathbf{x}, \mathbf{x}') = \langle \mathbf{x}, \mathbf{x}' \rangle$

Exact NTK Computation

NTK with ReLU activation [ADH+19]: for layer $\ell = 1, \dots, L$ and $\mathbf{x} \in \mathbb{R}^d$

$$K_{\text{NTK}}^{(\ell)}(\mathbf{x}, \mathbf{x}') = K_{\text{NTK}}^{(\ell-1)}(\mathbf{x}, \mathbf{x}') \cdot \underbrace{\dot{\Sigma}^{(\ell)}(\mathbf{x}, \mathbf{x}') + \Sigma^{(\ell)}(\mathbf{x}, \mathbf{x}')}_{\text{function of } \Sigma^{(\ell-1)}}$$

where $K_{\text{NTK}}^{(0)}(\mathbf{x}, \mathbf{x}') = \langle \mathbf{x}, \mathbf{x}' \rangle$

– $\dot{\Sigma}^{(\ell)}, \Sigma^{(\ell)}$ can be expressed as a **closed-form** formula using $\Sigma^{(\ell-1)}$

$$\dot{\Sigma}^{(\ell)}(\mathbf{x}, \mathbf{x}') = 1 - \frac{1}{\pi} \cos^{-1} \alpha^{(\ell)}$$

$$\Sigma^{(\ell)}(\mathbf{x}, \mathbf{x}') = \|\mathbf{x}\|_2 \|\mathbf{x}'\|_2 g(\alpha^{(\ell)})$$

$$\alpha^{(\ell)} = \frac{\Sigma^{(\ell-1)}(\mathbf{x}, \mathbf{x}')}{\sqrt{\Sigma^{(\ell-1)}(\mathbf{x}, \mathbf{x}) \Sigma^{(\ell-1)}(\mathbf{x}', \mathbf{x}')}}}$$

$$g(x) = \frac{\sqrt{1-x^2} + (\pi - \cos^{-1}(x))x}{\pi}$$

Efficient NTK Approximation

NTK with ReLU activation [ADH+19]: for layer $\ell = 1, \dots, L$ and $\mathbf{x} \in \mathbb{R}^d$

$$K_{\text{NTK}}^{(\ell)}(\mathbf{x}, \mathbf{x}') = K_{\text{NTK}}^{(\ell-1)}(\mathbf{x}, \mathbf{x}') \cdot \dot{\Sigma}^{(\ell)}(\mathbf{x}, \mathbf{x}') + \Sigma^{(\ell)}(\mathbf{x}, \mathbf{x}')$$

A key approach:

– If $\dot{\Sigma}^{(\ell)}(\mathbf{x}, \mathbf{x}') = \langle \mathbf{\Lambda}^{(\ell)}, \mathbf{\Lambda}^{(\ell)'} \rangle$ and $\Sigma^{(\ell)}(\mathbf{x}, \mathbf{x}') = \langle \mathbf{\Psi}^{(\ell)}, \mathbf{\Psi}^{(\ell)'} \rangle$

Efficient NTK Approximation

NTK with ReLU activation [ADH+19]: for layer $\ell = 1, \dots, L$ and $\mathbf{x} \in \mathbb{R}^d$

$$K_{\text{NTK}}^{(\ell)}(\mathbf{x}, \mathbf{x}') = K_{\text{NTK}}^{(\ell-1)}(\mathbf{x}, \mathbf{x}') \cdot \dot{\Sigma}^{(\ell)}(\mathbf{x}, \mathbf{x}') + \Sigma^{(\ell)}(\mathbf{x}, \mathbf{x}')$$

A key approach:

– If $\dot{\Sigma}^{(\ell)}(\mathbf{x}, \mathbf{x}') = \langle \mathbf{\Lambda}^{(\ell)}, \mathbf{\Lambda}^{(\ell)'} \rangle$ and $\Sigma^{(\ell)}(\mathbf{x}, \mathbf{x}') = \langle \mathbf{\Psi}^{(\ell)}, \mathbf{\Psi}^{(\ell)'} \rangle$

$$\begin{aligned} K_{\text{NTK}}^{(1)}(\mathbf{x}, \mathbf{x}') &= K_{\text{NTK}}^{(0)}(\mathbf{x}, \mathbf{x}') \cdot \langle \mathbf{\Lambda}^{(1)}, \mathbf{\Lambda}^{(1)'} \rangle + \langle \mathbf{\Psi}^{(1)}, \mathbf{\Psi}^{(1)'} \rangle \\ &= \langle \mathbf{x}, \mathbf{x}' \rangle \cdot \langle \mathbf{\Lambda}^{(1)}, \mathbf{\Lambda}^{(1)'} \rangle + \langle \mathbf{\Psi}^{(1)}, \mathbf{\Psi}^{(1)'} \rangle \\ &= \left\langle [\mathbf{x} \otimes \mathbf{\Lambda}^{(1)}, \mathbf{\Psi}^{(1)}], [\mathbf{x}' \otimes \mathbf{\Lambda}^{(1)'}, \mathbf{\Psi}^{(1)'}] \right\rangle \quad \otimes : \text{tensor product} \end{aligned}$$

Efficient NTK Approximation

NTK with ReLU activation [ADH+19]: for layer $\ell = 1, \dots, L$ and $\mathbf{x} \in \mathbb{R}^d$

$$K_{\text{NTK}}^{(\ell)}(\mathbf{x}, \mathbf{x}') = K_{\text{NTK}}^{(\ell-1)}(\mathbf{x}, \mathbf{x}') \cdot \dot{\Sigma}^{(\ell)}(\mathbf{x}, \mathbf{x}') + \Sigma^{(\ell)}(\mathbf{x}, \mathbf{x}')$$

A key approach:

– If $\dot{\Sigma}^{(\ell)}(\mathbf{x}, \mathbf{x}') = \langle \mathbf{\Lambda}^{(\ell)}, \mathbf{\Lambda}^{(\ell)'} \rangle$ and $\Sigma^{(\ell)}(\mathbf{x}, \mathbf{x}') = \langle \mathbf{\Psi}^{(\ell)}, \mathbf{\Psi}^{(\ell)'} \rangle$

$$\mathbf{\Phi}^{(\ell)} = [\mathbf{\Phi}^{(\ell-1)} \otimes \mathbf{\Lambda}^{(\ell)}, \mathbf{\Psi}^{(\ell)}], \quad \mathbf{\Phi}^{(0)} = \mathbf{x}$$

such that $K_{\text{NTK}}^{(\ell)}(\mathbf{x}, \mathbf{x}') = \langle \mathbf{\Phi}^{(\ell)}, \mathbf{\Phi}^{(\ell)'} \rangle$

\otimes : tensor product

Efficient NTK Approximation

NTK with ReLU activation [ADH+19]: for layer $\ell = 1, \dots, L$ and $\mathbf{x} \in \mathbb{R}^d$

$$K_{\text{NTK}}^{(\ell)}(\mathbf{x}, \mathbf{x}') = K_{\text{NTK}}^{(\ell-1)}(\mathbf{x}, \mathbf{x}') \cdot \dot{\Sigma}^{(\ell)}(\mathbf{x}, \mathbf{x}') + \Sigma^{(\ell)}(\mathbf{x}, \mathbf{x}')$$

A key approach:

– If $\dot{\Sigma}^{(\ell)}(\mathbf{x}, \mathbf{x}') = \langle \mathbf{\Lambda}^{(\ell)}, \mathbf{\Lambda}^{(\ell)'} \rangle$ and $\Sigma^{(\ell)}(\mathbf{x}, \mathbf{x}') = \langle \mathbf{\Psi}^{(\ell)}, \mathbf{\Psi}^{(\ell)'} \rangle$

$$\mathbf{\Phi}^{(\ell)} = [\mathbf{\Phi}^{(\ell-1)} \otimes \mathbf{\Lambda}^{(\ell)}, \mathbf{\Psi}^{(\ell)}], \quad \mathbf{\Phi}^{(0)} = \mathbf{x}$$

such that $K_{\text{NTK}}^{(\ell)}(\mathbf{x}, \mathbf{x}') = \langle \mathbf{\Phi}^{(\ell)}, \mathbf{\Phi}^{(\ell)'} \rangle$ \otimes : tensor product

Q: How to find $\mathbf{\Lambda}^{(\ell)}, \mathbf{\Psi}^{(\ell)}$ that approximates $\dot{\Sigma}^{(\ell)}, \Sigma^{(\ell)}$

\Rightarrow Random Features [CS09], Polynomial Sketch [AKK+21]

NTK Random Features

Goal: $\dot{\Sigma}^{(\ell)}(\mathbf{x}, \mathbf{x}') \approx \langle \mathbf{\Lambda}^{(\ell)}, \mathbf{\Lambda}^{(\ell)'} \rangle$ and $\Sigma^{(\ell)}(\mathbf{x}, \mathbf{x}') \approx \langle \mathbf{\Psi}^{(\ell)}, \mathbf{\Psi}^{(\ell)'} \rangle$ $\ell = 1, \dots, L$

$$K_{\text{NTK}}^{(\ell)}(\mathbf{x}, \mathbf{x}') = K_{\text{NTK}}^{(\ell-1)}(\mathbf{x}, \mathbf{x}') \cdot \dot{\Sigma}^{(\ell)}(\mathbf{x}, \mathbf{x}') + \Sigma^{(\ell)}(\mathbf{x}, \mathbf{x}')$$

What is Random Features?

– A feature map $\phi : \mathbb{R}^d \times \mathbb{R}^d \rightarrow \mathbb{R}$ such that

$$f(\mathbf{x}, \mathbf{x}') = \mathbb{E}_{\mathbf{w} \sim \mathcal{P}} [\phi(\mathbf{x}, \mathbf{w}) \cdot \phi(\mathbf{x}', \mathbf{w})]$$

– Sample m independent vectors $\mathbf{w}_1, \dots, \mathbf{w}_m$ and construct

$$\mathbf{\Phi} = \frac{1}{\sqrt{m}} [\phi(\mathbf{x}, \mathbf{w}_1), \dots, \phi(\mathbf{x}, \mathbf{w}_m)]^\top \in \mathbb{R}^m \quad (\text{Random Features})$$

such that $f(\mathbf{x}, \mathbf{x}') = \mathbb{E}[\langle \mathbf{\Phi}, \mathbf{\Phi}' \rangle]$

NTK Random Features

Goal: $\dot{\Sigma}^{(\ell)}(\mathbf{x}, \mathbf{x}') \approx \langle \mathbf{\Lambda}^{(\ell)}, \mathbf{\Lambda}^{(\ell)'} \rangle$ and $\Sigma^{(\ell)}(\mathbf{x}, \mathbf{x}') \approx \langle \mathbf{\Psi}^{(\ell)}, \mathbf{\Psi}^{(\ell)'} \rangle$ $\ell = 1, \dots, L$

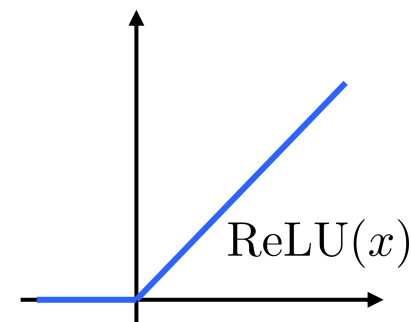
$$K_{\text{NTK}}^{(\ell)}(\mathbf{x}, \mathbf{x}') = K_{\text{NTK}}^{(\ell-1)}(\mathbf{x}, \mathbf{x}') \cdot \dot{\Sigma}^{(\ell)}(\mathbf{x}, \mathbf{x}') + \Sigma^{(\ell)}(\mathbf{x}, \mathbf{x}')$$

Random Features of Arc-cosine Kernel [CS09]:

$$\Sigma^{(1)}(\mathbf{x}, \mathbf{x}') = 2 \mathbb{E}_{\mathbf{w} \sim \mathcal{N}} [\text{ReLU}(\langle \mathbf{w}, \mathbf{x} \rangle) \cdot \text{ReLU}(\langle \mathbf{w}, \mathbf{x}' \rangle)] = \mathbb{E} [\langle \mathbf{\Psi}, \mathbf{\Psi}' \rangle]$$

– m_1 independent samples:

$$\mathbf{\Psi} = \sqrt{\frac{2}{m_1}} \text{ReLU}(\langle \mathbf{w}_1, \mathbf{x} \rangle, \dots, \langle \mathbf{w}_{m_1}, \mathbf{x} \rangle)^\top \in \mathbb{R}^{m_1}$$



NTK Random Features

Goal: $\dot{\Sigma}^{(\ell)}(\mathbf{x}, \mathbf{x}') \approx \langle \mathbf{\Lambda}^{(\ell)}, \mathbf{\Lambda}^{(\ell)'} \rangle$ and $\Sigma^{(\ell)}(\mathbf{x}, \mathbf{x}') \approx \langle \mathbf{\Psi}^{(\ell)}, \mathbf{\Psi}^{(\ell)'} \rangle$ $\ell = 1, \dots, L$

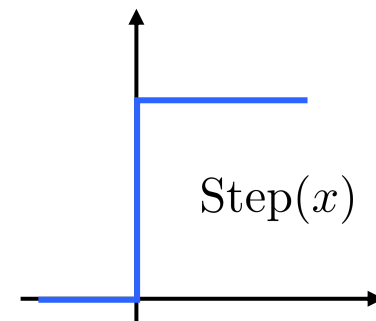
$$K_{\text{NTK}}^{(\ell)}(\mathbf{x}, \mathbf{x}') = K_{\text{NTK}}^{(\ell-1)}(\mathbf{x}, \mathbf{x}') \cdot \dot{\Sigma}^{(\ell)}(\mathbf{x}, \mathbf{x}') + \Sigma^{(\ell)}(\mathbf{x}, \mathbf{x}')$$

Random Features of Arc-cosine Kernel [CS09]:

$$\dot{\Sigma}^{(1)}(\mathbf{x}, \mathbf{x}') = 2 \mathbb{E}_{\mathbf{w} \sim \mathcal{N}} [\text{Step}(\langle \mathbf{w}, \mathbf{x} \rangle) \cdot \text{Step}(\langle \mathbf{w}, \mathbf{x}' \rangle)] = \mathbb{E} [\langle \mathbf{\Lambda}, \mathbf{\Lambda}' \rangle]$$

– m_0 independent samples:

$$\mathbf{\Lambda} = \sqrt{\frac{2}{m_0}} \text{Step}(\langle \mathbf{w}_1, \mathbf{x} \rangle, \dots, \langle \mathbf{w}_{m_1}, \mathbf{x} \rangle)^\top \in \mathbb{R}^{m_0}$$



NTK Random Features

NTK Random Features.

Initialize $\Phi^{(0)} = \Psi^{(0)} = \mathbf{x} \in \mathbb{R}^d$

For $\ell = 1, \dots, L$

$$\Lambda^{(\ell)} = \sqrt{\frac{2}{m_0}} \text{Step}(\mathbf{W}' \Psi^{(\ell-1)}) \in \mathbb{R}^{m_0}$$

$$\Psi^{(\ell)} = \sqrt{\frac{2}{m_1}} \text{ReLU}(\mathbf{W} \Psi^{(\ell-1)}) \in \mathbb{R}^{m_1}$$

$$\Phi^{(\ell)} = [\Phi^{(\ell-1)} \otimes \Lambda^{(\ell)}, \Psi^{(\ell)}]$$

Return $\Phi^{(L)}$ such that $K_{\text{NTK}}^{(L)}(\mathbf{x}, \mathbf{x}') \approx \langle \Phi^{(L)}, \Phi^{(L)'} \rangle$

Random Features

$$\mathbf{W}'_{ij} \sim \mathcal{N}(0, 1)$$

$$\mathbf{W}_{ij} \sim \mathcal{N}(0, 1)$$

NTK Random Features

NTK Random Features.

Initialize $\Phi^{(0)} = \Psi^{(0)} = \mathbf{x} \in \mathbb{R}^d$

For $\ell = 1, \dots, L$

Random Features

$$\Lambda^{(\ell)} = \sqrt{\frac{2}{m_0}} \text{Step}(\mathbf{W}' \Psi^{(\ell-1)}) \in \mathbb{R}^{m_0}$$

$$\mathbf{W}'_{ij} \sim \mathcal{N}(0, 1)$$

$$\Psi^{(\ell)} = \sqrt{\frac{2}{m_1}} \text{ReLU}(\mathbf{W} \Psi^{(\ell-1)}) \in \mathbb{R}^{m_1}$$

$$\mathbf{W}_{ij} \sim \mathcal{N}(0, 1)$$

$$\Phi^{(\ell)} = [\Phi^{(\ell-1)} \otimes \Lambda^{(\ell)}, \Psi^{(\ell)}]$$

Return $\Phi^{(L)}$ such that $K_{\text{NTK}}^{(L)}(\mathbf{x}, \mathbf{x}') \approx \langle \Phi^{(L)}, \Phi^{(L)'} \rangle$

Dimension of tensor product: $\dim(\mathbf{x} \otimes \mathbf{y}) = \dim(\mathbf{x}) \cdot \dim(\mathbf{y})$

Dimension of $\Phi^{(L)}$ is $\mathcal{O}(m_0^L (m_1 + d))$ 🙄

When L is large, the output feature can be very large \Rightarrow no efficiency

NTK Random Features

Computational bottleneck:

$$\Lambda^{(\ell)} = \sqrt{\frac{2}{m_0}} \text{Step}(\mathbf{W}' \Psi^{(\ell-1)}) \in \mathbb{R}^{m_0} \quad \mathbf{W}'_{ij} \sim \mathcal{N}(0, 1)$$

$$\Psi^{(\ell)} = \sqrt{\frac{2}{m_1}} \text{ReLU}(\mathbf{W} \Psi^{(\ell-1)}) \in \mathbb{R}^{m_1} \quad \mathbf{W}_{ij} \sim \mathcal{N}(0, 1)$$

$$\Phi^{(\ell)} = [\Phi^{(\ell-1)} \otimes \Lambda^{(\ell)}, \Psi^{(\ell)}]$$

TensorSketch [PP13; AKK+20]:

$$\langle \mathbf{x} \otimes \mathbf{y}, \mathbf{z} \otimes \mathbf{w} \rangle = \mathbb{E} \langle \mathcal{T}(\mathbf{x}, \mathbf{y}), \mathcal{T}(\mathbf{z}, \mathbf{w}) \rangle, \quad \mathbf{x}, \mathbf{y}, \mathbf{z}, \mathbf{w} \in \mathbb{R}^d$$

- Inner-product preserving dimensionality reduction
- $\mathcal{T} : \mathbb{R}^{d^2} \rightarrow \mathbb{R}^{m_s}$ can be computed in time $\mathcal{O}(d + m_s \log m_s)$ using FFT
- m_s is a tradeoff between runtime and error bound
- $m_s = \mathcal{O}\left(\frac{1}{\varepsilon^2} \log^3 \frac{1}{\varepsilon \delta}\right)$ for given accuracy ε and failure probability δ

NTK Random Features

NTK Random Features.

Initialize $\Phi^{(0)} = \Psi^{(0)} = \mathbf{x} \in \mathbb{R}^d$

For $\ell = 1, \dots, L$

Random Features

$$\Lambda^{(\ell)} = \sqrt{\frac{2}{m_0}} \text{Step}(\mathbf{W}' \Psi^{(\ell-1)}) \in \mathbb{R}^{m_0}$$

$$\mathbf{W}'_{ij} \sim \mathcal{N}(0, 1)$$

$$\Psi^{(\ell)} = \sqrt{\frac{2}{m_1}} \text{ReLU}(\mathbf{W} \Psi^{(\ell-1)}) \in \mathbb{R}^{m_1}$$

$$\mathbf{W}_{ij} \sim \mathcal{N}(0, 1)$$

$$\Phi^{(\ell)} = [\mathcal{T}(\Phi^{(\ell-1)}, \Lambda^{(\ell)}), \Psi^{(\ell)}] \quad \mathcal{T} : \text{TensorSketch to } \mathbb{R}^{m_s}$$

Return $\Phi^{(L)}$ such that $K_{\text{NTK}}^{(L)}(\mathbf{x}, \mathbf{x}') \approx \langle \Phi^{(L)}, \Phi^{(L)'} \rangle$

Dimension of $\Phi^{(L)}$ is $m_1 + m_s$ due to **Tensor Sketch** 😊

Larger m_1, m_0, m_s guarantee more accurate NTK approximation

Question: what is the error bound in terms of m_1, m_0, m_s ?

NTK Random Features

Theorem 1. Given $\mathbf{x}, \mathbf{x}' \in \mathbb{R}^d$ and $\delta \in (0, 1), \varepsilon \in (0, 1/L)$, if

$$m_0 = \mathcal{O}\left(\frac{L^2}{\varepsilon^2} \log \frac{L}{\delta}\right), \quad m_1 = \mathcal{O}\left(\frac{L^6}{\varepsilon^4} \log \frac{L}{\delta}\right), \quad m_s = \mathcal{O}\left(\frac{L^2}{\varepsilon^2} \log^3 \frac{L}{\varepsilon\delta}\right),$$

then

$$\Pr \left[\left| \left\langle \Phi^{(L)}, \Phi^{(L)'} \right\rangle - K_{\text{NTK}}^{(L)}(\mathbf{x}, \mathbf{x}') \right| \leq (L+1)\varepsilon \right] \geq 1 - \delta$$

NTK Random Features

Theorem 1. Given $\mathbf{x}, \mathbf{x}' \in \mathbb{R}^d$ and $\delta \in (0, 1), \varepsilon \in (0, 1/L)$, if

$$m_0 = \mathcal{O}\left(\frac{L^2}{\varepsilon^2} \log \frac{L}{\delta}\right), \quad m_1 = \mathcal{O}\left(\frac{L^6}{\varepsilon^4} \log \frac{L}{\delta}\right), \quad m_s = \mathcal{O}\left(\frac{L^2}{\varepsilon^2} \log^3 \frac{L}{\varepsilon\delta}\right),$$

then

$$\Pr \left[\left| \left\langle \Phi^{(L)}, \Phi^{(L)'} \right\rangle - K_{\text{NTK}}^{(L)}(\mathbf{x}, \mathbf{x}') \right| \leq (L+1)\varepsilon \right] \geq 1 - \delta$$

[ADH+19] showed that the gradient random features can guarantee

$$\Pr \left[\left| \left\langle \frac{\partial f(\mathbf{x}, \theta)}{\partial \theta}, \frac{\partial f(\mathbf{x}', \theta)}{\partial \theta} \right\rangle - K_{\text{NTK}}^{(L)}(\mathbf{x}, \mathbf{x}') \right| \leq (L+1)\varepsilon \right] \geq 1 - \delta$$

- Their feature dimension can be $\|\theta\|_0 = \tilde{\Omega}\left(\frac{L^{13}}{\varepsilon^8}\right)$
- Our feature dimension is $m_1 + m_s = \tilde{\mathcal{O}}\left(\frac{L^6}{\varepsilon^4}\right)$ 😄

NTK Random Features

Theorem 1. Given $\mathbf{x}, \mathbf{x}' \in \mathbb{R}^d$ and $\delta \in (0, 1), \varepsilon \in (0, 1/L)$, if

$$m_0 = \mathcal{O}\left(\frac{L^2}{\varepsilon^2} \log \frac{L}{\delta}\right), \quad m_1 = \mathcal{O}\left(\frac{L^6}{\varepsilon^4} \log \frac{L}{\delta}\right), \quad m_s = \mathcal{O}\left(\frac{L^2}{\varepsilon^2} \log^3 \frac{L}{\varepsilon\delta}\right),$$

then

$$\Pr \left[\left| \left\langle \Phi^{(L)}, \Phi^{(L)'} \right\rangle - K_{\text{NTK}}^{(L)}(\mathbf{x}, \mathbf{x}') \right| \leq (L+1)\varepsilon \right] \geq 1 - \delta$$

[ADH+19] showed that the gradient random features can guarantee

$$\Pr \left[\left| \left\langle \frac{\partial f(\mathbf{x}, \theta)}{\partial \theta}, \frac{\partial f(\mathbf{x}', \theta)}{\partial \theta} \right\rangle - K_{\text{NTK}}^{(L)}(\mathbf{x}, \mathbf{x}') \right| \leq (L+1)\varepsilon \right] \geq 1 - \delta$$

- Their feature dimension can be $\|\theta\|_0 = \tilde{\Omega}\left(\frac{L^{13}}{\varepsilon^8}\right)$
- Our feature dimension is $m_1 + m_s = \tilde{\mathcal{O}}\left(\frac{L^6}{\varepsilon^4}\right)$ 😊

We can improve the error bound using [importance sampling](#)

NTK Sketch

Goal: $\dot{\Sigma}^{(\ell)}(\mathbf{x}, \mathbf{x}') \approx \langle \mathbf{\Lambda}^{(\ell)}, \mathbf{\Lambda}^{(\ell)'} \rangle$ and $\Sigma^{(\ell)}(\mathbf{x}, \mathbf{x}') \approx \langle \mathbf{\Psi}^{(\ell)}, \mathbf{\Psi}^{(\ell)'} \rangle \quad \ell = 1, \dots, L$

$$K_{\text{NTK}}^{(\ell)}(\mathbf{x}, \mathbf{x}') = K_{\text{NTK}}^{(\ell-1)}(\mathbf{x}, \mathbf{x}') \cdot \dot{\Sigma}^{(\ell)}(\mathbf{x}, \mathbf{x}') + \Sigma^{(\ell)}(\mathbf{x}, \mathbf{x}')$$

PolySketch [AKK+20]:

– $\dot{\Sigma}^{(\ell)}, \Sigma^{(\ell)}$ are dot-product kernels & approximated by a polynomial

$$\Sigma^{(1)}(\mathbf{x}, \mathbf{x}') = 1 - \frac{1}{\pi} \cos^{-1} \alpha \approx \sum_{j=0}^p c_j \alpha^j \quad \alpha := \frac{\langle \mathbf{x}, \mathbf{x}' \rangle}{\|\mathbf{x}\|_2 \|\mathbf{x}'\|_2}$$

– Each monomial term can be approximated by TensorSketch [PP13]

$$\alpha^j \approx \langle \mathcal{T}(\mathbf{x}, j), \mathcal{T}(\mathbf{x}', j) \rangle \quad \mathcal{T}(\cdot, j) : \text{TensorSketch of degree } j$$

$$\mathbf{\Lambda}^{(1)} \approx [\sqrt{c_0}, \sqrt{c_1} \mathcal{T}(\mathbf{x}, 1), \dots, \sqrt{c_p} \mathcal{T}(\mathbf{x}, p)]$$

$$\Sigma^{(1)}(\mathbf{x}, \mathbf{x}') \approx \langle \mathbf{\Lambda}^{(1)}, \mathbf{\Lambda}^{(1)'} \rangle$$

NTK Sketch

NTK Sketch.

Initialize $\Phi^{(0)} = \Lambda^{(0)} = \mathbf{x} \in \mathbb{R}^d$

For $\ell = 1, \dots, L$

Sketching method

$$\Lambda^{(\ell)} = \text{PolySketch}(\Psi^{(\ell-1)}, p) \quad \text{for } 1 - \frac{1}{\pi} \cos^{-1}(x)$$

$$\Psi^{(\ell)} = \text{PolySketch}(\Psi^{(\ell-1)}, p') \quad \text{for } \frac{1}{\pi} (\sqrt{1-x^2} + (\pi - \cos^{-1}(x))x)$$

$$\Phi^{(\ell)} = [\mathcal{T}(\Phi^{(\ell-1)}, \Lambda^{(\ell)}), \Psi^{(\ell)}] \quad \mathcal{T} : \text{TensorSketch to } \mathbb{R}^{m_s}$$

Return $\Phi^{(L)}$ such that $K_{\text{NTK}}^{(L)}(\mathbf{x}, \mathbf{x}') \approx \langle \Phi^{(L)}, \Phi^{(L)'} \rangle$

$\dot{\Sigma}^{(\ell)}, \Sigma^{(\ell)}$ can be approximated by PolySketch [AKK+20]

PolySketch is faster than Random Features (i.e., matrix-matrix products)

With degree p and sketch dimension m , it can run in $\mathcal{O}(p(d + m \log m))$

NTK Sketch

Theorem 2. Given $\mathbf{x}, \mathbf{x}' \in \mathbb{R}^d$ and $\delta \in (0, 1), \varepsilon \in (0, 1/L)$, the NTK Sketch computes $\Phi^{(L)} \in \mathbb{R}^m, m = \mathcal{O}(\frac{1}{\varepsilon^2} \log \frac{1}{\delta})$ in time

$$\tilde{\mathcal{O}} \left(\frac{L^{11}}{\varepsilon^{6.7}} + \frac{L^3}{\varepsilon^2} d \right)$$

such that

$$\Pr \left[\left| \langle \Phi^{(L)}, \Phi^{(L)'} \rangle - K_{\text{NTK}}^{(L)}(\mathbf{x}, \mathbf{x}') \right| \leq \varepsilon \cdot K_{\text{NTK}}^{(L)}(\mathbf{x}, \mathbf{x}') \right] \geq 1 - \delta$$

Feature dimension m does not depend on L

Compared to **NTK Random Features**, NTK Sketch is fast and efficient

	Feature dimension	Running time
NTK Random Features	$\tilde{\mathcal{O}} \left(\frac{L^6}{\varepsilon^4} \right)$	$\tilde{\mathcal{O}} \left(\frac{L^{13}}{\varepsilon^8} \right)$

NTK RF/Sketch methods can be extended to Convolutional Neural Networks (CNNs)

Outline

Introduction

- Kernel Regression
- Neural Tangent Kernel
- Existing Work

NTK Feature Map Construction

- Exact NTK Computation
- NTK Random Features
- NTK Sketch

Experiments

- Classification on MNIST datasets
- Classification on CIFAR-10 datasets

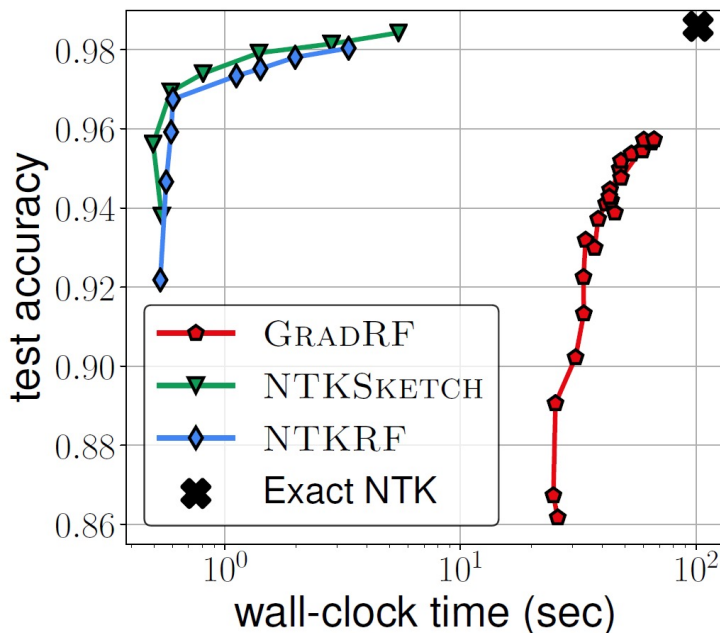
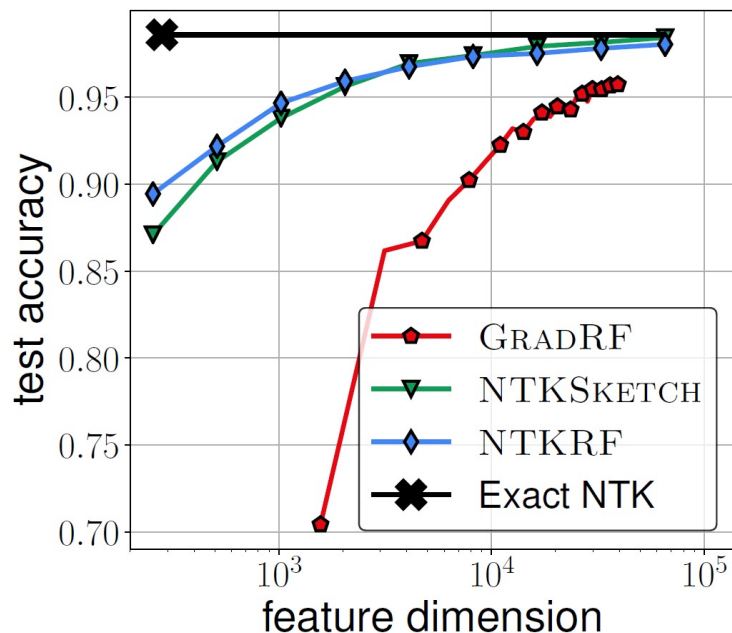
Conclusion

Experiments: MNIST Classification

Comparison to other NTK approximation methods:

- Kernel regression by the exact NTK (Exact NTK)
- Gradient features by Monte-Carlo sampling (GradRF)
- NTK Random Features (NTKRF) / Sketch (NTKSketch)

Test accuracy versus feature dimension / computing time ($L = 1$)

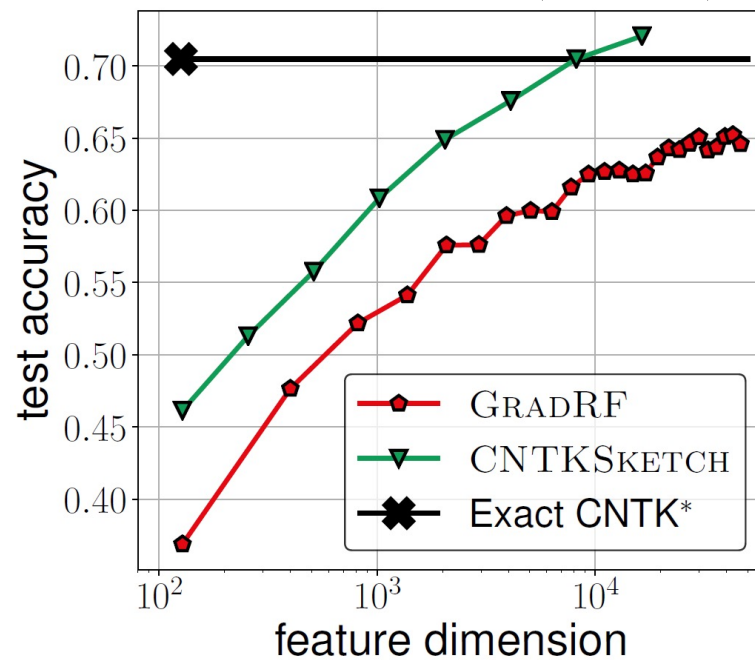
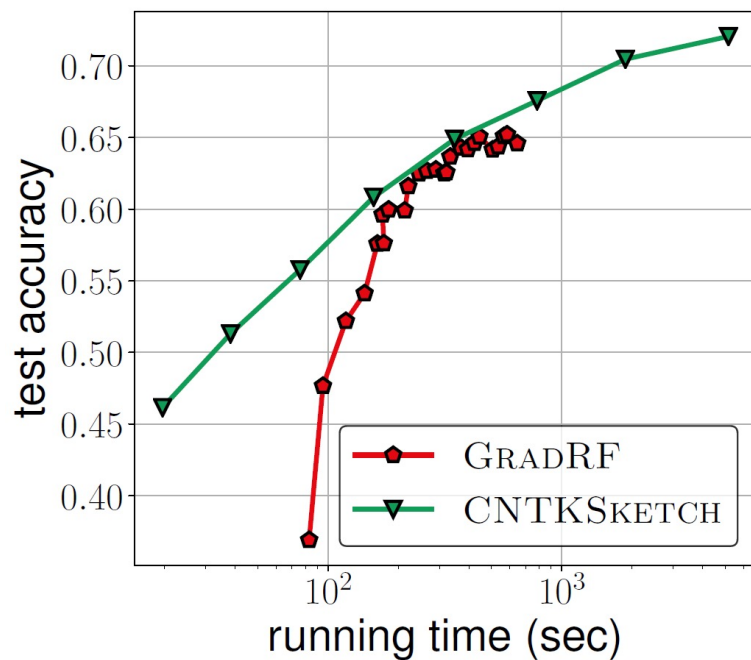


Experiments: CIFAR-10 Classification

Comparison to other NTK approximation methods:

- Kernel regression by the exact CNTK (Exact CNTK)
- Gradient features by Monte-Carlo sampling (GradRF)
- CNTK Sketch (CNTKSketch)

Test accuracy versus feature dimension / computing time ($L = 3$)



Experiments: CIFAR-10 Classification

Comparison to other NTK approximation methods:

- Kernel regression by the exact CNTK (Exact CNTK)
- Gradient features by Monte-Carlo sampling (GradRF)
- CNTK Sketch (CNTKSketch)

Test accuracy versus feature dimension / computing time ($L = 3$)

	CNTKSKETCH (ours)			GRADRF			Exact CNTK	CNN
Feature dimension	4,096	8,192	16,384	9,328	17,040	42,816		
Test accuracy (%)	67.58	70.46	72.06	62.49	62.57	65.21	70.47*	63.81*
Time (s)	780	1,870	5,160	300	360	580	> 1,000,000	

* means that the result is from [ADH+19]

Conclusion

Summary:

- We propose efficient feature maps of NTK / CNTK
- We design two approaches utilize sketching algorithm and arc-cosine random features, respectively
- We provide an entry-wise error bound for both algorithms
- We additionally provide a spectral approximation bound of random features approach with leverage score sampling
- The proposed methods outperform other baselines

Future work:

- Spectral approximation guarantee for deeper layers
- Applying our method to convolutional NTK, attention network, etc