

# An Introduction to the Neural Tangent Kernel (NTK)

---

Stephan Rabanser

[stephan@cs.toronto.edu](mailto:stephan@cs.toronto.edu)



University of Toronto  
Department of Computer Science



Vector Institute for  
Artificial Intelligence

February 2, 2023

## Recap: Linear Regression

- We are given a dataset  $\mathcal{D} = \{(\mathbf{x}_i, y_i)\}_{i=1}^n$  with  $\mathbf{x}_i \in \mathbb{R}^D$  and  $y_i \in \mathbb{R}$ .
- Goal: predict a response  $y$  from features  $\mathbf{x}$  using a linear function.
- Model the relationship using the predictive function:

$$f_{\mathbf{w}}(\mathbf{x}) = \langle \mathbf{w}, \mathbf{x} \rangle = \mathbf{w}^\top \mathbf{x}$$

- The quality of our fit is determined by a loss function:

$$\mathcal{L}_{\mathbf{w}} = \frac{1}{2} \sum_{i=1}^n (y_i - f_{\mathbf{w}}(\mathbf{x}_i))^2$$

- We want to minimize the loss:  $\min_{\mathbf{w}} \mathcal{L}_{\mathbf{w}}$ . We can do this using gradient descent:

$$\mathbf{w}_{t+1} = \mathbf{w}_t - \eta \nabla_{\mathbf{w}} \mathcal{L}_{\mathbf{w}}$$



## Recap: Kernel Methods

- Linear relations are restrictive, interesting datasets have non-linear interactions!
- Transform features into higher-dimensional space:

$$\mathbf{x} \in \mathbb{R}^D \rightarrow \phi(\mathbf{x}) \in \mathbb{R}^K, K \gg D$$

- Structure of prediction function stays the same:  $f_{\mathbf{w}}(\mathbf{x}) = \mathbf{w}^\top \phi(\mathbf{x})$ .
- The model is *linear* in  $\mathbf{w}$  but *non-linear* in  $\mathbf{x}$ .
- As a result, the objective  $\mathcal{L}_{\mathbf{w}}$  is still convex and solvable using gradient descent:

$$\mathcal{L}_{\mathbf{w}} = \frac{1}{2} \sum_{i=1}^n (y_i - f_{\mathbf{w}}(\mathbf{x}_i))^2 = \frac{1}{2} \sum_{i=1}^n (y_i - \mathbf{w}^\top \phi(\mathbf{x}_i))^2 \quad \min_{\mathbf{w}} \mathcal{L}_{\mathbf{w}}$$

- This is great! We have non-linear features but still a convex objective! However:
  - The transformation function  $\phi(\cdot)$  is fixed and needs to be tuned manually.
  - The computation of the feature map  $\phi(\mathbf{x}) \in \mathbb{R}^K$  with  $K \gg D$  can be expensive.

## Recap: Kernel Trick

- Many ML algorithms can be reformulated to feature space inner products:

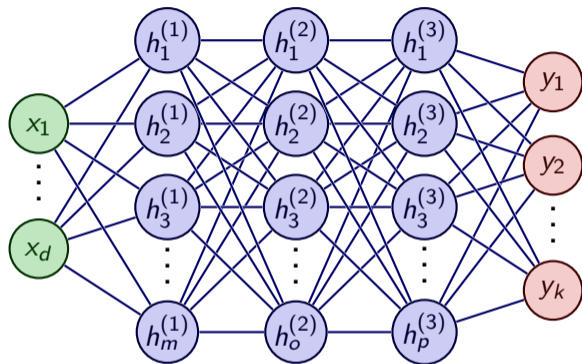
$$\kappa(\mathbf{x}, \mathbf{x}') = \langle \phi(\mathbf{x}), \phi(\mathbf{x}') \rangle$$

- This inner product is called a *kernel function*  $\kappa(\cdot, \cdot)$  and can be thought of as a similarity measure of two input vectors  $\mathbf{x}$  and  $\mathbf{x}'$ .
- Example kernels:

$$\kappa_{\text{Pol}}(\mathbf{x}, \mathbf{x}') = (\mathbf{x}^\top \mathbf{x}' + c)^d \qquad \kappa_{\text{SE}}(\mathbf{x}, \mathbf{x}') = \sigma^2 \exp\left(-\frac{(\mathbf{x} - \mathbf{x}')^2}{2\ell^2}\right)$$

- Collect all similarities in a positive semi-definite kernel matrix:  $\mathbf{K} \in \mathbb{R}^{n \times n} \succeq 0$ .
- Algorithms can be "kernelized": only rely on kernels instead of expl. feature maps.
- Are such kernels flexible enough? What about the computational concerns?

# Recap: Neural Networks



- We don't want to define non-linear transformations ourselves: we want to *learn* them!
- Discover multiple (hierarchical) feature spaces:

$$\mathbb{R}^d \rightarrow \mathbb{R}^m \rightarrow \mathbb{R}^o \rightarrow \mathbb{R}^p \rightarrow \dots \rightarrow \mathbb{R}^k$$

- The objective  $\mathcal{L}_{\mathbf{w}}$  is not convex, yet we still use gradient descent:

$$\mathcal{L}_{\mathbf{w}} = \frac{1}{2} \sum_{i=1}^n (y_i - f_{\mathbf{w}}(\mathbf{x}_i))^2 = \frac{1}{2} \sum_{i=1}^n \left( y_i - \sigma(\mathbf{W}_l^\top \dots \sigma(\mathbf{W}_2^\top \sigma(\mathbf{W}_1^\top \mathbf{x}_i))) \right)^2 \quad \mathbf{w} = (\mathbf{W}_1, \dots, \mathbf{W}_l)$$

# Understanding the Performance of Neural Nets: Central Questions

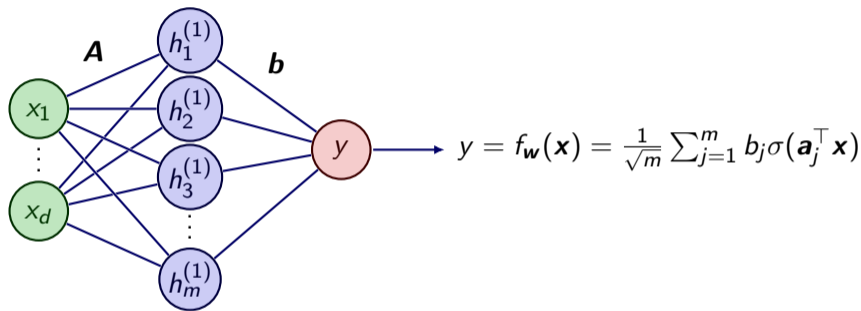
**Neural networks empirically perform well, but their convergence and generalization properties are hard to analyze.**

Deep learning is poorly understood. Kernel methods on the other hand are based on solid mathematical theory. *Can we distill a neural network into a kernel?*

It is a known result that at initialization time, a wide neural network is a Gaussian process. *Can we also describe the training process using a kernel?*

# Deriving a Kernel from a Neural Network

- Assume a simple neural network with a single hidden layer and params  $\mathbf{w} = (\mathbf{A}, \mathbf{b})$ :



- As discussed, the objective  $\mathcal{L}_{\mathbf{w}}$  is no longer convex:

$$\mathcal{L}_{\mathbf{w}} = \frac{1}{2} \sum_{i=1}^n (y_i - f_{\mathbf{w}}(\mathbf{x}_i))^2 = \frac{1}{2} \sum_{i=1}^n \left( y_i - \frac{1}{\sqrt{m}} \sum_{j=1}^m b_j \sigma(\mathbf{a}_j^{\top} \mathbf{x}_i) \right)^2$$

## Deriving a Kernel from a Neural Network (cont'd)

$$\mathcal{L}_{\mathbf{w}} = \frac{1}{2} \sum_{i=1}^n (y_i - f_{\mathbf{w}}(\mathbf{x}_i))^2 = \frac{1}{2} \sum_{i=1}^n \left( y_i - \frac{1}{\sqrt{m}} \sum_{j=1}^m b_j \sigma(\mathbf{a}_j^\top \mathbf{x}_i) \right)^2$$

- We can still minimize the objective  $\mathcal{L}_{\mathbf{w}}$  using *full-batch* gradient descent:

$$\begin{aligned} \mathbf{w}_{t+1} &= \mathbf{w}_t - \eta \nabla_{\mathbf{w}} \mathcal{L}_{\mathbf{w}_t} \\ &= \mathbf{w}_t - \eta \sum_{i=1}^n (y_i - f_{\mathbf{w}_t}(\mathbf{x}_i)) \underbrace{\nabla_{\mathbf{w}} f_{\mathbf{w}_t}(\mathbf{x}_i)}_{\substack{\text{variable during training} \\ \text{changes depending on } \mathbf{w}_t}} \end{aligned}$$

Under what circumstances does  $\nabla_{\mathbf{w}} f_{\mathbf{w}_t}(\mathbf{x}_i)$  not change much during training?

# Lazy Training

- Assume we initialize the weights randomly using a standard Gaussian:  $\mathcal{N}(0, 1)$ .
- We can observe the trajectory of the weights during training:

$$\mathbf{w}_0 \longrightarrow \mathbf{w}_1 \longrightarrow \mathbf{w}_2 \longrightarrow \dots \longrightarrow \mathbf{w}_T$$

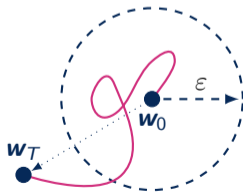
## Empirical Observation

When  $m$  is large ( $m \rightarrow \infty$ ), parameters show stable evolution patterns ( $\varepsilon \rightarrow 0$ ).

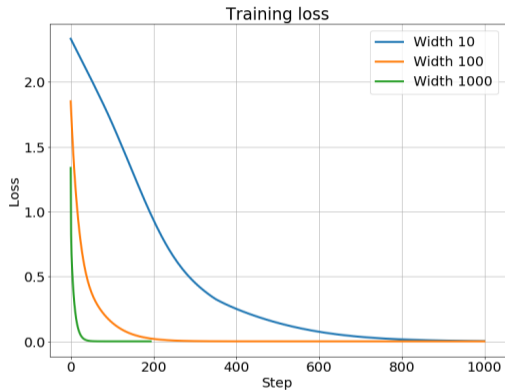
Lazy training



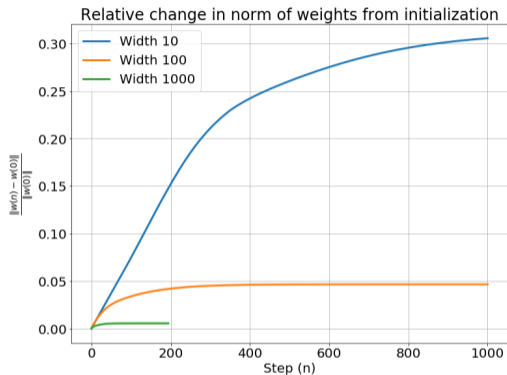
Non-Lazy training



# Depiction of Lazy Training for Wide Neural Nets



(a) Loss evolution



(b) Weight distance from origin

<https://rajatvd.github.io/NTK/>



# Neural Tangent Kernel: Approximating a Lazy Trajectory

- If lazy training holds, then a first-order Taylor approximation of the function around its initialization  $\mathbf{w}_0$  might be helpful:

$$\begin{aligned} f_{\mathbf{w}}(\mathbf{x}) &\approx f_{\mathbf{w}_0}(\mathbf{x}) + \underbrace{\nabla_{\mathbf{w}} f_{\mathbf{w}_0}(\mathbf{x})^\top}_{\phi(\mathbf{x})^\top} (\mathbf{w} - \mathbf{w}_0) + \underbrace{\dots}_{\text{higher order Taylor terms}} \\ &\approx \underbrace{c + \phi(\mathbf{x})^\top (\mathbf{w} - \mathbf{w}_0)}_{\text{model is an affine func in } \mathbf{w}} \end{aligned}$$

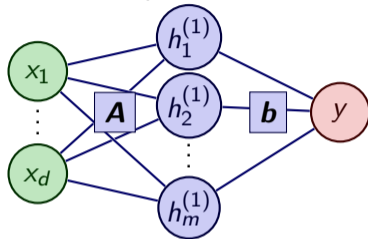
## Neural Tangent Kernel (NTK)

For the standard kernel definition  $\kappa(\mathbf{x}, \mathbf{x}') = \langle \phi(\mathbf{x}), \phi(\mathbf{x}') \rangle$ , use the gradient of the neural network's output evaluated at  $\mathbf{w}_0$  as the kernel function:

$$\phi(\mathbf{x}) = \nabla_{\mathbf{w}} f_{\mathbf{w}_0}(\mathbf{x})$$

# Computing the NTK for Our Toy Example

- Recall our simple neural network from before:



$$y = f_{\mathbf{w}}(\mathbf{x}) = \frac{1}{\sqrt{m}} \sum_{j=1}^m b_j \sigma(\mathbf{a}_j^\top \mathbf{x})$$

$$\begin{cases} \nabla_{\mathbf{a}_j} f_{\mathbf{w}}(\mathbf{x}) = \frac{1}{\sqrt{m}} b_j \sigma'(\mathbf{a}_j^\top \mathbf{x}) \mathbf{x} \\ \nabla_{b_j} f_{\mathbf{w}}(\mathbf{x}) = \frac{1}{\sqrt{m}} \sigma(\mathbf{a}_j^\top \mathbf{x}) \end{cases}$$

- Computing the Neural Tangent Kernel for this network:

$$\begin{aligned} \kappa_{\text{NTK}}(\mathbf{x}, \mathbf{x}') &= \kappa_{\text{NTK}}^{(\mathbf{A})}(\mathbf{x}, \mathbf{x}') + \kappa_{\text{NTK}}^{(\mathbf{b})}(\mathbf{x}, \mathbf{x}') \\ &= \frac{1}{m} \sum_{i=1}^m b_i^2 \sigma'(\mathbf{a}_i^\top \mathbf{x}) \sigma'(\mathbf{a}_i^\top \mathbf{x}') \mathbf{x} \mathbf{x}' + \frac{1}{m} \sum_{j=1}^m \sigma(\mathbf{a}_j^\top \mathbf{x}) \sigma(\mathbf{a}_j^\top \mathbf{x}') \\ &\stackrel{m \rightarrow \infty}{\equiv} \mathbb{E}[\mathbf{b}^2 \sigma'(\mathbf{A}^\top \mathbf{x}) \sigma'(\mathbf{A}^\top \mathbf{x}') \mathbf{x} \mathbf{x}'] + \mathbb{E}[\sigma(\mathbf{A}^\top \mathbf{x}) \sigma(\mathbf{A}^\top \mathbf{x}')] \end{aligned}$$

# Analyzing Gradient Descent

- This kernel can directly be used in a kernel machine!
- We can analyze further properties of neural network training using flow dynamics.
- Example: parameter gradient flow dynamics ( $\eta \rightarrow 0$ ):

$$\begin{aligned}\mathbf{w}_{t+1} &= \mathbf{w}_t - \eta \nabla_{\mathbf{w}} \mathcal{L}_{\mathbf{w}_t} \\ \frac{\mathbf{w}_{t+1} - \mathbf{w}_t}{\eta} &= -\nabla_{\mathbf{w}} \mathcal{L}_{\mathbf{w}_t} \\ \frac{d\mathbf{w}_t}{dt} &\stackrel{\eta \rightarrow 0}{=} -\nabla_{\mathbf{w}} \mathcal{L}_{\mathbf{w}_t}\end{aligned}$$

Can we use differential equations to analyze the evolution of various properties of neural networks (weights, predictions, losses, etc. over time)?

## Parameter Dynamics

How do the parameters change over the course of training?

Assume  $\mathcal{L}_{\mathbf{w}_t} = \frac{1}{2} \|\mathbf{f}_{\mathbf{w}_t} - \mathbf{y}\|_2^2$ . Then:

$$\nabla_{\mathbf{w}} \mathcal{L}_{\mathbf{w}_t} = \nabla_{\mathbf{w}} \mathbf{f}_{\mathbf{w}_t} (\mathbf{f}_{\mathbf{w}_t} - \mathbf{y})$$

$$\begin{aligned} \frac{d\mathbf{w}_t}{dt} &= -\nabla_{\mathbf{w}} \mathcal{L}_{\mathbf{w}_t} \\ &= -\nabla_{\mathbf{w}} \mathbf{f}_{\mathbf{w}_t} (\mathbf{f}_{\mathbf{w}_t} - \mathbf{y}) \end{aligned}$$

## Predictions Dynamics

How do the predictions change over the course of training?

Approximate w/ NTK matrix  $\mathbf{K}_{\text{NTK}}$ :

$$\begin{aligned} \frac{d\mathbf{f}_{\mathbf{w}_t}}{dt} &= -\nabla_{\mathbf{w}} \mathbf{f}_{\mathbf{w}_t}^\top \frac{d\mathbf{w}_t}{dt} \\ &= -\underbrace{\nabla_{\mathbf{w}} \mathbf{f}_{\mathbf{w}_t}^\top \nabla_{\mathbf{w}} \mathbf{f}_{\mathbf{w}_t}}_{\text{NTK evaluated at } \mathbf{w}_t} (\mathbf{f}_{\mathbf{w}_t} - \mathbf{y}) \\ &\approx -\mathbf{K}_{\text{NTK}} (\mathbf{f}_{\mathbf{w}_t} - \mathbf{y}) \end{aligned}$$

# ODE Solution for Residual Loss

- Under a simple residual loss, we can model the dynamics w/ linear ODE:

$$\mathbf{u} = \mathbf{f}_{\mathbf{w}_t} - \mathbf{y} \quad \implies \quad \frac{d\mathbf{u}}{dt} \approx -\mathbf{K}_{\text{NTK}}\mathbf{u} \quad \implies \quad \mathbf{u}_t = \mathbf{u}_0 \exp(-\mathbf{K}_{\text{NTK}}t)$$

- As  $t \rightarrow \infty$ ,  $\mathbf{u} \rightarrow 0$  and  $\mathbf{f}_{\mathbf{w}_t} \rightarrow \mathbf{y}$ .
- In over-parameterized networks:  $\mathbf{K}_{\text{NTK}} \succ 0$ , i.e. smallest eigenvalue larger than 0!
- We can factorize the kernel matrix  $\mathbf{K}_{\text{NTK}} = \sum_{i=1}^k \lambda_i \mathbf{v}_i \mathbf{v}_i^\top$  with  $0 < \lambda_1 < \dots < \lambda_k$ .
- Substituting this factorization back into the ODE:

$$\mathbf{u}_t = \mathbf{u}_0 \prod_{i=1}^k \exp(-\lambda_i \mathbf{v}_i \mathbf{v}_i^\top t)$$

- $\lambda_1$  governs the rate of convergence.

## Summary of Neural Tangent Kernel

- Parameters hardly move from their initialization for  $m \rightarrow \infty$ .
- NTK is defined as the gradient of the NN evaluated at init:  $\nabla_{\mathbf{w}} \mathbf{f}_{\mathbf{w}_0}^\top \nabla_{\mathbf{w}} \mathbf{f}_{\mathbf{w}_0}$ .
- NTK is deterministic at initialization and constant during training.
- Can be used in ODEs to study evolution of various quantities.





## Open Questions

- Results from the NTK limit are not SOTA. Where does the gap come from?
- Results only hold for full-batch GD. What is the role of SGD?
- NTK results talk about convergence? What about performance?
- How can we use these training dynamics insights for trust?

Thanks! :)

# References

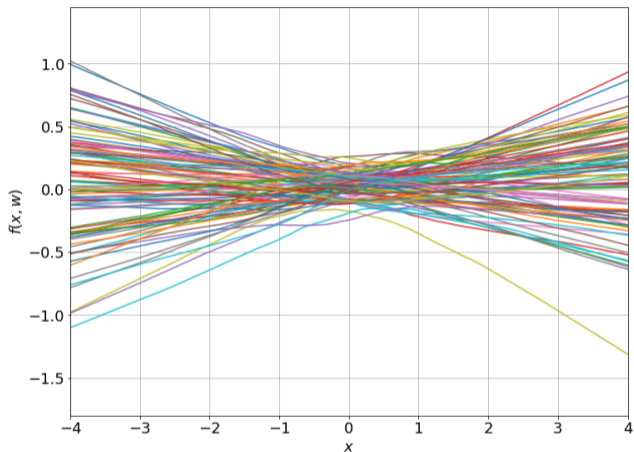
- <https://rajatvd.github.io/NTK/>
- <https://lilianweng.github.io/posts/2022-09-08-ntk/>
- <https://www.youtube.com/watch?v=D0bobAnELkU>

-  Sanjeev Arora, Simon S Du, Wei Hu, Zhiyuan Li, Russ R Salakhutdinov, and Ruosong Wang, *On exact computation with an infinitely wide neural net*, NeurIPS **32** (2019).
-  Lenaic Chizat and Francis Bach, *On the global convergence of gradient descent for over-parameterized models using optimal transport*, NeurIPS **31** (2018).
-  Lenaic Chizat, Edouard Oyallon, and Francis Bach, *On lazy training in differentiable programming*, NeurIPS **32** (2019).
-  Arthur Jacot, Franck Gabriel, and Clément Hongler, *Neural tangent kernel: Convergence and generalization in neural networks*, NeurIPS **31** (2018).



Backup

# Over-Parameterization At Initialization Time $\rightarrow$ GP



<https://rajatvd.github.io/NTK/>

# Depiction of Lazy Training for Wide Neural Nets

---

(a) 10 layer

(b) 100 layer

(c) 1000 layer

<https://rajatvd.github.io/NTK/>