

Neural Tangent Kernel

Comprehending Convergence & Generalization in Neural Networks¹

J. Setpal

March 6, 2025



**MACHINE LEARNING
@ PURDUE**

¹Jacot, Gabriel, Hongler. [NIPS 2018]

① Background & Intuition

② Neural Tangent Kernel

Outline

① Background & Intuition

② Neural Tangent Kernel

Problem Setting

Given network $f_\theta : \mathbb{R}^d \rightarrow \mathbb{R}$, we have two distinctive regimes:

Underparameterized Learning: $\mathbf{n} \geq \mathbf{p}$, $\theta \in \mathbb{R}^{\mathbf{p}}$, $\mathcal{D} = \{(x_i, y_i)\}_{i=1}^{\mathbf{n}}$

Overparameterized Learning: $\mathbf{p} \geq \mathbf{n}$, $\theta \in \mathbb{R}^{\mathbf{p}}$, $\mathcal{D} = \{(x_i, y_i)\}_{i=1}^{\mathbf{n}}$

Problem Setting

Given network $f_\theta : \mathbb{R}^d \rightarrow \mathbb{R}$, we have two distinctive regimes:

Underparameterized Learning: $\mathbf{n} \geq \mathbf{p}$, $\theta \in \mathbb{R}^{\mathbf{p}}$, $\mathcal{D} = \{(x_i, y_i)\}_{i=1}^n$

Overparameterized Learning: $\mathbf{p} \geq \mathbf{n}$, $\theta \in \mathbb{R}^{\mathbf{p}}$, $\mathcal{D} = \{(x_i, y_i)\}_{i=1}^n$

Either way, we optimize θ to minimize an empirical loss (here, quadratic):

$$\mathcal{L}(\theta, \mathcal{D}) = \frac{1}{n} \sum_{i=1}^n \ell(f_\theta(x_i), y_i) = \frac{1}{n} \|f_\theta(\mathbf{x}) - \mathbf{y}\|_2^2 \quad (1)$$

Problem Setting

Given network $f_\theta : \mathbb{R}^d \rightarrow \mathbb{R}$, we have two distinctive regimes:

Underparameterized Learning: $\mathbf{n} \geq \mathbf{p}$, $\theta \in \mathbb{R}^{\mathbf{p}}$, $\mathcal{D} = \{(x_i, y_i)\}_{i=1}^{\mathbf{n}}$

Overparameterized Learning: $\mathbf{p} \geq \mathbf{n}$, $\theta \in \mathbb{R}^{\mathbf{p}}$, $\mathcal{D} = \{(x_i, y_i)\}_{i=1}^{\mathbf{n}}$

Either way, we optimize θ to minimize an empirical loss (here, quadratic):

$$\mathcal{L}(\theta, \mathcal{D}) = \frac{1}{n} \sum_{i=1}^n \ell(f_\theta(x_i), y_i) = \frac{1}{n} \|f_\theta(\mathbf{x}) - \mathbf{y}\|_2^2 \quad (1)$$

Q: Why does overparameterized learning generalize?

Problem Setting

Given network $f_\theta : \mathbb{R}^d \rightarrow \mathbb{R}$, we have two distinctive regimes:

Underparameterized Learning: $\mathbf{n} \geq \mathbf{p}$, $\theta \in \mathbb{R}^{\mathbf{p}}$, $\mathcal{D} = \{(x_i, y_i)\}_{i=1}^{\mathbf{n}}$

Overparameterized Learning: $\mathbf{p} \geq \mathbf{n}$, $\theta \in \mathbb{R}^{\mathbf{p}}$, $\mathcal{D} = \{(x_i, y_i)\}_{i=1}^{\mathbf{n}}$

Either way, we optimize θ to minimize an empirical loss (here, quadratic):

$$\mathcal{L}(\theta, \mathcal{D}) = \frac{1}{n} \sum_{i=1}^n \ell(f_\theta(x_i), y_i) = \frac{1}{n} \|f_\theta(\mathbf{x}) - \mathbf{y}\|_2^2 \quad (1)$$

Q: Why does overparameterized learning generalize?

A: NTK's approach:

- a. Construct an analogy to a simpler paradigm (kernel methods).

Problem Setting

Given network $f_\theta : \mathbb{R}^d \rightarrow \mathbb{R}$, we have two distinctive regimes:

Underparameterized Learning: $n \geq p$, $\theta \in \mathbb{R}^p$, $\mathcal{D} = \{(x_i, y_i)\}_{i=1}^n$

Overparameterized Learning: $p \geq n$, $\theta \in \mathbb{R}^p$, $\mathcal{D} = \{(x_i, y_i)\}_{i=1}^n$

Either way, we optimize θ to minimize an empirical loss (here, quadratic):

$$\mathcal{L}(\theta, \mathcal{D}) = \frac{1}{n} \sum_{i=1}^n \ell(f_\theta(x_i), y_i) = \frac{1}{n} \|f_\theta(\mathbf{x}) - \mathbf{y}\|_2^2 \quad (1)$$

Q: Why does overparameterized learning generalize?

A: NTK's approach:

- Construct an analogy to a simpler paradigm (kernel methods).
- Prove the analogy actually holds for sufficiently wide networks.

Problem Setting

Given network $f_\theta : \mathbb{R}^d \rightarrow \mathbb{R}$, we have two distinctive regimes:

Underparameterized Learning: $\mathbf{n} \geq \mathbf{p}$, $\theta \in \mathbb{R}^{\mathbf{p}}$, $\mathcal{D} = \{(x_i, y_i)\}_{i=1}^{\mathbf{n}}$

Overparameterized Learning: $\mathbf{p} \geq \mathbf{n}$, $\theta \in \mathbb{R}^{\mathbf{p}}$, $\mathcal{D} = \{(x_i, y_i)\}_{i=1}^{\mathbf{n}}$

Either way, we optimize θ to minimize an empirical loss (here, quadratic):

$$\mathcal{L}(\theta, \mathcal{D}) = \frac{1}{n} \sum_{i=1}^n \ell(f_\theta(x_i), y_i) = \frac{1}{n} \|f_\theta(\mathbf{x}) - \mathbf{y}\|_2^2 \quad (1)$$

Q: Why does overparameterized learning generalize?

A: NTK's approach:

- Construct an analogy to a simpler paradigm (kernel methods).
- Prove the analogy actually holds for sufficiently wide networks.
- Use the analysis from the kernel method to understand training dynamics of the neural network.

Pre-Requisite 1/3 – Taylor Series

We can approximate function g using a polynomial via the Taylor Series:

$$P_a(x) = \sum_{i=0}^{\infty} \frac{g^{(n)}(a)}{n!} (x - a)^n \quad (2)$$

Pre-Requisite 1/3 – Taylor Series

We can approximate function g using a polynomial via the Taylor Series:

$$P_a(x) = \sum_{i=0}^{\infty} \frac{g^{(i)}(a)}{i!} (x - a)^i \quad (2)$$

The first-order Taylor Expansion is a linear approximation of the function:

$$P_a(x) = g(a) + g'(a)(x - a) \quad (3)$$

Pre-Requisite 1/3 – Taylor Series

We can approximate function g using a polynomial via the Taylor Series:

$$P_a(x) = \sum_{i=0}^{\infty} \frac{g^{(n)}(a)}{n!} (x - a)^n \quad (2)$$

The first-order Taylor Expansion is a linear approximation of the function:

$$P_a(x) = g(a) + g'(a)(x - a) \quad (3)$$

We will use the first-order approximation to model **training dynamics in neural networks**.

Pre-Requisite 2/3 – Operator Norm

The operator norm of matrix A is the maximum amount of stretching that A can do to arbitrary vector x :

$$A \in \mathbb{R}^{m \times n}, \quad x \in \mathbb{R}^n \quad (4)$$

$$\|A\|_{\text{op}} = \max_{x \in \mathbb{R}^n} \frac{\|Ax\|_p}{\|x\|_q} \quad (5)$$

Usually $p = q = 2$.

Pre-Requisite 3/3 – Kernel Methods

$K : \mathbb{R}^m \times \mathbb{R}^m \rightarrow \mathbb{R}$ is a kernel if $\exists \phi : \mathcal{X} \rightarrow \mathcal{H}$ s.t:

$$K(x, x') = \phi(x)^T \phi(x') = \langle \phi(x), \phi(x') \rangle \quad (6)$$

Any PSD matrix defines a kernel.

$K : \mathbb{R}^m \times \mathbb{R}^m \rightarrow \mathbb{R}$ is a kernel if $\exists \phi : \mathcal{X} \rightarrow \mathcal{H}$ s.t:

$$K(x, x') = \phi(x)^T \phi(x') = \langle \phi(x), \phi(x') \rangle \quad (6)$$

Any PSD matrix defines a kernel.

A kernel methods are **instance-based** learners, instead of learning parameters over input features, we learn parameters over data-pairs, and interpolate using the kernel function to predict for the unseen sample.

Pre-Requisite 3/3 – Kernel Methods

$K : \mathbb{R}^m \times \mathbb{R}^m \rightarrow \mathbb{R}$ is a kernel if $\exists \phi : \mathcal{X} \rightarrow \mathcal{H}$ s.t:

$$K(x, x') = \phi(x)^T \phi(x') = \langle \phi(x), \phi(x') \rangle \quad (6)$$

Any PSD matrix defines a kernel.

A kernel methods are **instance-based** learners, instead of learning parameters over input features, we learn parameters over data-pairs, and interpolate using the kernel function to predict for the unseen sample.

As an example:

$$\hat{y} = \mathbf{sgn} \sum_{i=1}^n w_i y_i k(x, x') \quad (7)$$

Outline

① Background & Intuition

② Neural Tangent Kernel

Neural Network Setup

We have fully-connected network f_θ with n_0, \dots, n_L neurons in each layer, activation function $\sigma : \mathbb{R} \rightarrow \mathbb{R}$ with pre-activations $\tilde{\alpha}$ and activations α defined recursively as follows:

$$\alpha^{(0)}(x; \theta) = x \quad (8)$$

$$\tilde{\alpha}^{(\ell+1)}(x; \theta) = \frac{1}{\sqrt{n_\ell}} W^{(\ell)} \alpha^{(\ell)}(x; \theta) + \beta b^{(\ell)} \quad (9)$$

$$\alpha^{(\ell)}(x; \theta) = \sigma(\tilde{\alpha}^{(\ell)}(x; \theta)) \quad (10)$$

A **realization function** $F^{(L)} : \mathbb{R}^p \rightarrow \mathcal{F}$ maps parameters to a function.

Neural Network Setup

We have fully-connected network f_θ with n_0, \dots, n_L neurons in each layer, activation function $\sigma : \mathbb{R} \rightarrow \mathbb{R}$ with pre-activations $\tilde{\alpha}$ and activations α defined recursively as follows:

$$\alpha^{(0)}(x; \theta) = x \quad (8)$$

$$\tilde{\alpha}^{(\ell+1)}(x; \theta) = \frac{1}{\sqrt{n_\ell}} W^{(\ell)} \alpha^{(\ell)}(x; \theta) + \beta b^{(\ell)} \quad (9)$$

$$\alpha^{(\ell)}(x; \theta) = \sigma(\tilde{\alpha}^{(\ell)}(x; \theta)) \quad (10)$$

A **realization function** $F^{(L)} : \mathbb{R}^p \rightarrow \mathcal{F}$ maps parameters to a function.

At initialization, $\theta \sim \mathcal{N}(0, I_p)$.

Neural Network Setup

We have fully-connected network f_θ with n_0, \dots, n_L neurons in each layer, activation function $\sigma : \mathbb{R} \rightarrow \mathbb{R}$ with pre-activations $\tilde{\alpha}$ and activations α defined recursively as follows:

$$\alpha^{(0)}(x; \theta) = x \quad (8)$$

$$\tilde{\alpha}^{(\ell+1)}(x; \theta) = \frac{1}{\sqrt{n_\ell}} W^{(\ell)} \alpha^{(\ell)}(x; \theta) + \beta b^{(\ell)} \quad (9)$$

$$\alpha^{(\ell)}(x; \theta) = \sigma(\tilde{\alpha}^{(\ell)}(x; \theta)) \quad (10)$$

A **realization function** $F^{(L)} : \mathbb{R}^p \rightarrow \mathcal{F}$ maps parameters to a function.

At initialization, $\theta \sim \mathcal{N}(0, I_p)$. As a result at initialization, our network is equivalent to a gaussian process.

Neural Network Setup

We have fully-connected network f_θ with n_0, \dots, n_L neurons in each layer, activation function $\sigma : \mathbb{R} \rightarrow \mathbb{R}$ with pre-activations $\tilde{\alpha}$ and activations α defined recursively as follows:

$$\alpha^{(0)}(x; \theta) = x \quad (8)$$

$$\tilde{\alpha}^{(\ell+1)}(x; \theta) = \frac{1}{\sqrt{n_\ell}} W^{(\ell)} \alpha^{(\ell)}(x; \theta) + \beta b^{(\ell)} \quad (9)$$

$$\alpha^{(\ell)}(x; \theta) = \sigma(\tilde{\alpha}^{(\ell)}(x; \theta)) \quad (10)$$

A **realization function** $F^{(L)} : \mathbb{R}^p \rightarrow \mathcal{F}$ maps parameters to a function.

At initialization, $\theta \sim \mathcal{N}(0, I_p)$. As a result at initialization, our network is equivalent to a gaussian process.

Aside

Factors $1/\sqrt{n_\ell}$, β scale gradients & enable consistent asymptotic behavior.

Studying Convergence in Function Space

Notice that quadratic loss² is **convex w.r.t. function space** but often **non-convex w.r.t. parameter space**.

$$C : \mathcal{F} \rightarrow \mathbb{R} \qquad \text{Convex :D} \qquad (11)$$

$$C \circ F^{(L)} : \mathbb{R}^p \rightarrow \mathbb{R} \qquad \text{Highly Non-Convex D:} \qquad (12)$$

²also most others

Studying Convergence in Function Space

Notice that quadratic loss² is **convex w.r.t. function space** but often **non-convex w.r.t. parameter space**.

$$C : \mathcal{F} \rightarrow \mathbb{R} \qquad \text{Convex :D} \qquad (11)$$

$$C \circ F^{(L)} : \mathbb{R}^p \rightarrow \mathbb{R} \qquad \text{Highly Non-Convex D:} \qquad (12)$$

Now, if we could analyze training in the convex setting, we can obtain meaningful insight about convergence of these networks.

²also most others

Studying Convergence in Function Space

Notice that quadratic loss² is **convex w.r.t. function space** but often **non-convex w.r.t. parameter space**.

$$C : \mathcal{F} \rightarrow \mathbb{R} \qquad \text{Convex :D} \qquad (11)$$

$$C \circ F^{(L)} : \mathbb{R}^p \rightarrow \mathbb{R} \qquad \text{Highly Non-Convex D:} \qquad (12)$$

Now, if we could analyze training in the convex setting, we can obtain meaningful insight about convergence of these networks.

Q₁: How can we do that?

A₁: By constructing a first-order approximation of training dynamics.

²also most others

Studying Convergence in Function Space

Notice that quadratic loss² is **convex w.r.t. function space** but often **non-convex w.r.t. parameter space**.

$$C : \mathcal{F} \rightarrow \mathbb{R} \qquad \text{Convex :D} \qquad (11)$$

$$C \circ F^{(L)} : \mathbb{R}^p \rightarrow \mathbb{R} \qquad \text{Highly Non-Convex D:} \qquad (12)$$

Now, if we could analyze training in the convex setting, we can obtain meaningful insight about convergence of these networks.

Q₁: How can we do that?

A₁: By constructing a first-order approximation of training dynamics.

Q₂: Why can we do that?

²also most others

Studying Convergence in Function Space

Notice that quadratic loss² is **convex w.r.t. function space** but often **non-convex w.r.t. parameter space**.

$$C : \mathcal{F} \rightarrow \mathbb{R} \quad \text{Convex :D} \quad (11)$$

$$C \circ F^{(L)} : \mathbb{R}^p \rightarrow \mathbb{R} \quad \text{Highly Non-Convex D:} \quad (12)$$

Now, if we could analyze training in the convex setting, we can obtain meaningful insight about convergence of these networks.

Q₁: How can we do that?

A₁: By constructing a first-order approximation of training dynamics.

Q₂: Why can we do that?

A₂: In the overparameterized regime, learning is *lazy*.

²also most others

Studying Convergence in Function Space

Notice that quadratic loss² is **convex w.r.t. function space** but often **non-convex w.r.t. parameter space**.

$$C : \mathcal{F} \rightarrow \mathbb{R} \quad \text{Convex :D} \quad (11)$$

$$C \circ F^{(L)} : \mathbb{R}^p \rightarrow \mathbb{R} \quad \text{Highly Non-Convex D:} \quad (12)$$

Now, if we could analyze training in the convex setting, we can obtain meaningful insight about convergence of these networks.

Q₁: How can we do that?

A₁: By constructing a first-order approximation of training dynamics.

Q₂: Why can we do that?

A₂: In the overparameterized regime, learning is *lazy*.

We'll answer the questions in detail, but in reverse order.

²also most others

Lazy Learning

For L -smooth f , gradient descent using small η guarantees convergent loss.

$$\mathcal{L}(\theta_{t+1}, \mathcal{D}) \leq \mathcal{L}(\theta_t, \mathcal{D}), \quad \eta \leq 1/L \quad (13)$$

$$\|\theta - \theta_0\| \approx \frac{\|f_{\theta_0}(x) - y\|}{\|\nabla_{\theta} f_{\theta_0}(x)\|_{\text{op}}} \quad (14)$$

Lazy Learning

For L -smooth f , gradient descent using small η guarantees convergent loss.

$$\mathcal{L}(\theta_{t+1}, \mathcal{D}) \leq \mathcal{L}(\theta_t, \mathcal{D}), \quad \eta \leq 1/L \quad (13)$$

$$\|\theta - \theta_0\| \approx \frac{\|f_{\theta_0}(x) - y\|}{\|\nabla_{\theta} f_{\theta_0}(x)\|_{\text{op}}} \quad (14)$$

Takeaway: With gaussian initialization, as $\theta \in \mathbb{R}^p$ s.t. $p \rightarrow \infty$

$$\exists \theta^* \quad \text{s.t.} \quad \mathcal{L}(f_{\theta^*}, \mathcal{D}) = \min_{\theta} \mathcal{L}(f_{\theta}, \mathcal{D}) \quad \text{and} \quad \|\theta^* - \theta_0\| < B \ll \infty \quad (15)$$

Lazy Learning

For L -smooth f , gradient descent using small η guarantees convergent loss.

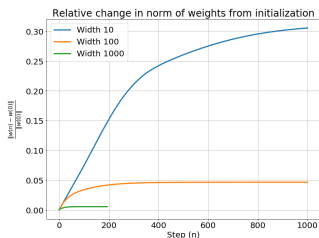
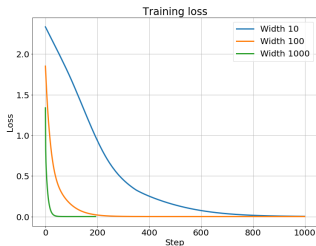
$$\mathcal{L}(\theta_{t+1}, \mathcal{D}) \leq \mathcal{L}(\theta_t, \mathcal{D}), \quad \eta \leq 1/L \quad (13)$$

$$\|\theta - \theta_0\| \approx \frac{\|f_{\theta_0}(x) - y\|}{\|\nabla_{\theta} f_{\theta_0}(x)\|_{\text{op}}} \quad (14)$$

Takeaway: With gaussian initialization, as $\theta \in \mathbb{R}^p$ s.t. $p \rightarrow \infty$

$$\exists \theta^* \text{ s.t. } \mathcal{L}(f_{\theta^*}, \mathcal{D}) = \min_{\theta} \mathcal{L}(f_{\theta}, \mathcal{D}) \quad \text{and} \quad \|\theta^* - \theta_0\| < B \ll \infty \quad (15)$$

We can see this by the change of $\|\theta\|$ as $\theta \in \mathbb{R}^{10}, \mathbb{R}^{100}, \mathbb{R}^{1000}$:



Linear Approximation of f_θ

Applying first-order Taylor expansion of f with respect to the evolution of parameters θ , we have:

$$g_\theta(x) := f_{\theta_0}(x) + \langle \nabla_\theta f_{\theta_0}(x), \underbrace{\theta - \theta_0}_{\Delta\theta} \rangle \quad (16)$$

Linear Approximation of f_θ

Applying first-order Taylor expansion of f with respect to the evolution of parameters θ , we have:

$$g_\theta(x) := f_{\theta_0}(x) + \langle \nabla_\theta f_{\theta_0}(x), \underbrace{\theta - \theta_0}_{\Delta\theta} \rangle \quad (16)$$

For convenience of analysis, we can set θ_0 s.t. $f_{\theta_0}(x) = 0 \ \forall x$.

Linear Approximation of f_θ

Applying first-order Taylor expansion of f with respect to the evolution of parameters θ , we have:

$$g_\theta(x) := f_{\theta_0}(x) + \underbrace{\langle \nabla_\theta f_{\theta_0}(x), \theta - \theta_0 \rangle}_{\Delta\theta} \quad (16)$$

For convenience of analysis, we can set θ_0 s.t. $f_{\theta_0}(x) = 0 \forall x$. Now,

$$g_\theta(x) := \underbrace{\langle \nabla_\theta f_{\theta_0}(x), \Delta\theta \rangle}_{\phi(x)} \quad (17)$$

Linear Approximation of f_θ

Applying first-order Taylor expansion of f with respect to the evolution of parameters θ , we have:

$$g_\theta(x) := f_{\theta_0}(x) + \underbrace{\langle \nabla_\theta f_{\theta_0}(x), \theta - \theta_0 \rangle}_{\Delta\theta} \quad (16)$$

For convenience of analysis, we can set θ_0 s.t. $f_{\theta_0}(x) = 0 \forall x$. Now,

$$g_\theta(x) := \underbrace{\langle \nabla_\theta f_{\theta_0}(x), \Delta\theta \rangle}_{\phi(x)} \quad (17)$$

Since we are in the overparameterized lazy learning regime, we have $\theta \approx \theta_0$, and $g_\theta \approx f_\theta$, which means:

$$\mathcal{L}(f_\theta, \mathcal{D}) \approx \mathcal{L}(g_\theta, \mathcal{D}) \quad (18)$$

Linear Approximation of f_θ

Applying first-order Taylor expansion of f with respect to the evolution of parameters θ , we have:

$$g_\theta(x) := f_{\theta_0}(x) + \underbrace{\langle \nabla_\theta f_{\theta_0}(x), \theta - \theta_0 \rangle}_{\Delta\theta} \quad (16)$$

For convenience of analysis, we can set θ_0 s.t. $f_{\theta_0}(x) = 0 \forall x$. Now,

$$g_\theta(x) := \underbrace{\langle \nabla_\theta f_{\theta_0}(x), \Delta\theta \rangle}_{\phi(x)} \quad (17)$$

Since we are in the overparameterized lazy learning regime, we have $\theta \approx \theta_0$, and $g_\theta \approx f_\theta$, which means:

$$\mathcal{L}(f_\theta, \mathcal{D}) \approx \mathcal{L}(g_\theta, \mathcal{D}) \quad (18)$$

Now for $G^{(L)} : \mathbb{R}^p \rightarrow \mathcal{G}$ cost-composition $C \circ G^L : \mathbb{R}^p \rightarrow \mathbb{R}$ is convex!!

Connection to Convergence

Next, notice $\phi(x)$ **does not** depend on post training parameters.

Connection to Convergence

Next, notice $\phi(x)$ **does not** depend on post training parameters.

We can use this to define feature map Φ over the entire dataset:

$$\Phi = [\phi(x_i)^T]_n = [\nabla f_{\theta_0}(x_i)^T]_n \in \mathbb{R}^{n \times p} \quad (19)$$

Connection to Convergence

Next, notice $\phi(x)$ **does not** depend on post training parameters.

We can use this to define feature map Φ over the entire dataset:

$$\Phi = [\phi(x_i)^T]_n = [\nabla f_{\theta_0}(x_i)^T]_n \in \mathbb{R}^{n \times p} \quad (19)$$

Optimizing $\mathcal{L}(f_\theta) \approx$ optimizing $\mathcal{L}(g_\theta)$ which is convex in parameter space.
So this is basically just linear regression:

$$\min_{g_\theta} \|y - g_\theta(x)\|_2^2 = \min_{\Delta\theta} \|y - \Phi \cdot \Delta\theta\|_2^2 \quad (20)$$

Comparing in output space, we can show that decay of error is exponential.

Connection to Convergence

Next, notice $\phi(x)$ **does not** depend on post training parameters.

We can use this to define feature map Φ over the entire dataset:

$$\Phi = [\phi(x_i)^T]_n = [\nabla f_{\theta_0}(x_i)^T]_n \in \mathbb{R}^{n \times p} \quad (19)$$

Optimizing $\mathcal{L}(f_\theta) \approx$ optimizing $\mathcal{L}(g_\theta)$ which is convex in parameter space. So this is basically just linear regression:

$$\min_{g_\theta} \|y - g_\theta(x)\|_2^2 = \min_{\Delta\theta} \|y - \Phi \cdot \Delta\theta\|_2^2 \quad (20)$$

Comparing in output space, we can show that decay of error is exponential.

Caveat: We're discussing gradient descent, not *stochastic* gradient descent. In practice, we see performance dropoffs when optimizing over g_θ .

Defining NTK

We can define a kernel using our feature map ϕ :

$$\kappa(x, x') = \phi(x)^T \phi(x') = \langle \nabla_{\theta} f_{\theta_0}(x), \nabla_{\theta} f_{\theta_0}(x') \rangle \quad (21)$$

This kernel is the **neural tangent kernel**.

Defining NTK

We can define a kernel using our feature map ϕ :

$$\kappa(x, x') = \phi(x)^T \phi(x') = \langle \nabla_{\theta} f_{\theta_0}(x), \nabla_{\theta} f_{\theta_0}(x') \rangle \quad (21)$$

This kernel is the **neural tangent kernel**.

The NTK effectively exists as part of gradient descent:

$$\theta_{t+1} = \theta_t - \nabla_{\theta} \mathcal{L}(f_{\theta_t}, \mathcal{D}) \quad (22)$$

$$\frac{\partial \theta}{\partial t} = -\eta \nabla_{\theta} \mathcal{L}(\theta, \mathcal{D}) = -\frac{\eta}{n} \sum_{i=1}^n \nabla_{\theta} f_{\theta}(x_i) \nabla_{f_{\theta}} \ell(f_{\theta}, y_i) \quad (23)$$

$$\frac{\partial f_{\theta}(x)}{\partial t} = \frac{\partial f_{\theta}(x)}{\partial \theta} \frac{\partial \theta}{\partial t} = -\frac{\eta}{n} \sum_{i=1}^n \underbrace{\nabla_{\theta} f_{\theta}(x) \nabla_{\theta} f_{\theta}(x_i)}_{\text{NTK}} \nabla_{f_{\theta}} \ell(f_{\theta}, y_i) \quad (24)$$

Defining NTK

We can define a kernel using our feature map ϕ :

$$\kappa(x, x') = \phi(x)^T \phi(x') = \langle \nabla_{\theta} f_{\theta_0}(x), \nabla_{\theta} f_{\theta_0}(x') \rangle \quad (21)$$

This kernel is the **neural tangent kernel**.

The NTK effectively exists as part of gradient descent:

$$\theta_{t+1} = \theta_t - \nabla_{\theta} \mathcal{L}(f_{\theta_t}, \mathcal{D}) \quad (22)$$

$$\frac{\partial \theta}{\partial t} = -\eta \nabla_{\theta} \mathcal{L}(\theta, \mathcal{D}) = -\frac{\eta}{n} \sum_{i=1}^n \nabla_{\theta} f_{\theta}(x_i) \nabla_{f_{\theta}} \ell(f_{\theta}, y_i) \quad (23)$$

$$\frac{\partial f_{\theta}(x)}{\partial t} = \frac{\partial f_{\theta}(x)}{\partial \theta} \frac{\partial \theta}{\partial t} = -\frac{\eta}{n} \sum_{i=1}^n \underbrace{\nabla_{\theta} f_{\theta}(x) \nabla_{\theta} f_{\theta}(x_i)}_{\text{NTK}} \nabla_{f_{\theta}} \ell(f_{\theta}, y_i) \quad (24)$$

The NTK depends on θ_0 , is random at initialization and varies during training, but in the limit, this changes.

Limiting Behavior of NTK

When the width tends to infinity, the NTK is deterministic at initialization and doesn't change through training:

$$k \rightarrow \infty \implies f_{\theta,k} \rightarrow \mathcal{N}(0, \Sigma_k) \quad \forall k \in \{1, \dots, n_L\} \quad (25)$$

Limiting Behavior of NTK

When the width tends to infinity, the NTK is deterministic at initialization and doesn't change through training:

$$k \rightarrow \infty \implies f_{\theta,k} \rightarrow \mathcal{N}(0, \Sigma_k) \quad \forall k \in \{1, \dots, n_L\} \quad (25)$$

Where covariance matrices Σ_k are defined recursively as follows:

$$\Sigma_1(x, x') = \frac{1}{n_0} x^T x' + \beta^2 \quad (26)$$

$$\Sigma_{k+1}(x, x') = \mathbb{E}_{f \sim \mathcal{N}(0, \Sigma_k)} [\sigma(f(x))\sigma(f(x'))] + \beta^2 \quad (27)$$

Limiting Behavior of NTK

When the width tends to infinity, the NTK is deterministic at initialization and doesn't change through training:

$$k \rightarrow \infty \implies f_{\theta,k} \rightarrow \mathcal{N}(0, \Sigma_k) \quad \forall k \in \{1, \dots, n_L\} \quad (25)$$

Where covariance matrices Σ_k are defined recursively as follows:

$$\Sigma_1(x, x') = \frac{1}{n_0} x^T x' + \beta^2 \quad (26)$$

$$\Sigma_{k+1}(x, x') = \mathbb{E}_{f \sim \mathcal{N}(0, \Sigma_k)} [\sigma(f(x))\sigma(f(x'))] + \beta^2 \quad (27)$$

Which we further defines the limiting NTK:

$$\kappa_1(x, x') = \Sigma_1(x, x') \quad (28)$$

$$\kappa_{k+1}(x, x') = \kappa_k(x, x') \Sigma'_{k+1}(x, x') + \Sigma_k(x, x') \quad (29)$$

$$\Sigma'_{k+1}(x, x') = \mathbb{E}_{f \sim \mathcal{N}(0, \Sigma_k)} [\sigma'(f(x))\sigma'(f(x'))] + \beta^2 \quad (30)$$

Which is independent of initialization.

Training within the NTK Regime

During training f_θ follows a descent along the negative kernel gradient:

$$\partial_t f_{\theta_t} = -\nabla_{\Phi_{(L)}} C|_{f_{\theta_t}} \quad (31)$$

The Limiting Kernel is always *positive definite*. By performing an eigendecomposition, we can decouple the gradient flow into eigenvectors, that decay at the rate of their eigenvalues.

Training within the NTK Regime

During training f_θ follows a descent along the negative kernel gradient:

$$\partial_t f_{\theta_t} = -\nabla_{\Phi_{(L)}} C|_{f_{\theta_t}} \quad (31)$$

The Limiting Kernel is always *positive definite*. By performing an eigendecomposition, we can decouple the gradient flow into eigenvectors, that decay at the rate of their eigenvalues.

Since they all decay, NTK guarantees that infinite width neural networks converge to a global minimum when trained to minimize empirical loss.

Training within the NTK Regime

During training f_θ follows a descent along the negative kernel gradient:

$$\partial_t f_{\theta_t} = -\nabla_{\Phi_{(L)}} C|_{f_{\theta_t}} \quad (31)$$

The Limiting Kernel is always *positive definite*. By performing an eigendecomposition, we can decouple the gradient flow into eigenvectors, that decay at the rate of their eigenvalues.

Since they all decay, NTK guarantees that infinite width neural networks converge to a global minimum when trained to minimize empirical loss.

Bonus: We have theoretical motivations for early-stopping; once the k largest eigenvectors decay beyond a set threshold, stop training.

Training within the NTK Regime

During training f_θ follows a descent along the negative kernel gradient:

$$\partial_t f_{\theta_t} = -\nabla_{\Phi_{(L)}} C|_{f_{\theta_t}} \quad (31)$$

The Limiting Kernel is always *positive definite*. By performing an eigendecomposition, we can decouple the gradient flow into eigenvectors, that decay at the rate of their eigenvalues.

Since they all decay, NTK guarantees that infinite width neural networks converge to a global minimum when trained to minimize empirical loss.

Bonus: We have theoretical motivations for early-stopping; once the k largest eigenvectors decay beyond a set threshold, stop training.

Caveat: Computing NTK is expensive $\Omega(n^2)$ while regular GD uses $\mathcal{O}(d)$ samples.

Thank you!

Have an awesome rest of your day!

Slides: <https://jinen.setpal.net/slides/ntk.pdf>

References:

1. Dwaraknath, Rajat. (Nov 2019). Understanding the Neural Tangent Kernel. <https://www.eigntales.com/NTK/>
2. Ma, Tengyu. (Nov 2022). Stanford CS229M - Lecture 13: Neural Tangent Kernel. <https://youtu.be/btphvvnad0A>
3. Ma, Tengyu. (Nov 2022). Stanford CS229M - Lecture 14: Neural Tangent Kernel, Implicit regularization of gradient descent. <https://youtu.be/xpT1ymwCk9w>
4. Weng, Lilian. (Sep 2022). Some math behind neural tangent kernel. <https://lilianweng.github.io/posts/2022-09-08-ntk/>
5. Pietraho, Thomas. Math 2805: Mathematical principles of machine learning. https://web.bowdoin.edu/~tpietrah/hugo/math2805/docs/unsupervised_learning/dim_reduction/hw_operator_norms/