# NEURAL TANGENT KERNEL (NTK) ANALYSIS
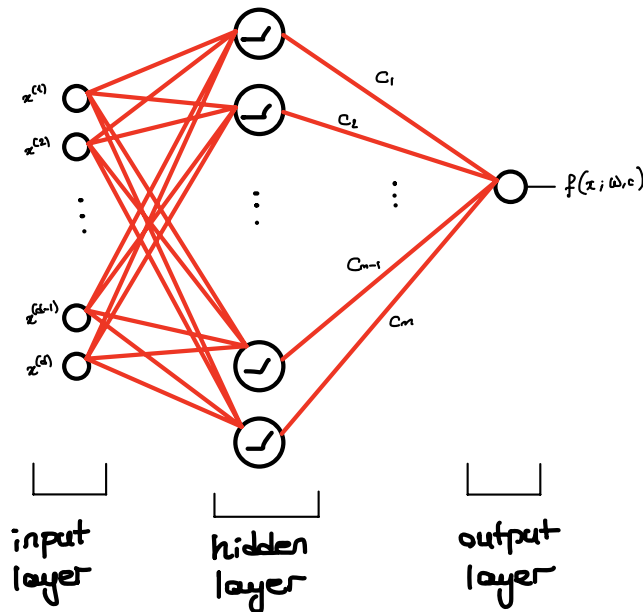
We will mainly consider feedforward ~~neural~~ neural neural networks.

$$f(x; W, c) = \sum_{i=1}^{m} c_i \sigma(w_i^T x) \quad \text{where} \quad w_i \in \mathbb{R}^d, \ c_i \in \mathbb{R} \quad \text{for all} \quad i \in \{1, 2, \ldots, m\}.$$

$m$ is the __width__ of the neural network.

neuron or hidden unit

$$W^T = \begin{bmatrix} w_1^T & \cdots & w_m^T \end{bmatrix} \in \mathbb{R}^{md}$$

$$c = \begin{bmatrix} c_1 & \cdots & c_m \end{bmatrix} \in \mathbb{R}^m$$



input layer    hidden layer    output layer

$\sigma : \mathbb{R} \to \mathbb{R}$ is the activation function.

Mainly used: $\sigma(z) = \max\{0, z\}$  (Rectified Linear Unit, ReLU)



$\sigma(z) = \max\{0, z\} = [z]_+$

Recall $\mathbb{1}\{z \geq 0\} := \begin{cases} 1, & \text{if } z \geq 0 \\ 0, & \text{otherwise.} \end{cases}$

Then, we also have $\sigma(z) = z \, \mathbb{1}\{z \geq 0\}$.

## Lemma 1 (Some elementary properties of ReLU networks)

① Given $x \in \mathbb{R}^d$, $x \, \mathbb{1}\{\theta^T x \geq 0\} \in \partial_\theta \sigma(\theta^T x)$ where $\sigma(z) = [z]_+$. ($\partial_\theta$ is subdifferential)

② Define $\nabla_{w_i} f(x; W, c) := c_i \, x \, \mathbb{1}\{w_i^T x \geq 0\}$ for $i = 1, 2, \ldots, m$,

$$\nabla_W^T f(x; W, c) := \begin{bmatrix} \nabla_{w_1}^T f(x; W, c) & \cdots & \nabla_{w_m}^T f(x; W, c) \end{bmatrix}.$$

Then,

$$\nabla_W^T f(x; W, c) \, W = f(x; W, c), \quad \forall x \in \mathbb{R}^d, \ (W, c) \in \mathbb{R}^{md} \times \mathbb{R}$$

**Pf:** ① 

$$\sigma(\theta^T x) + x^T \mathbb{1}\{\theta^T x \geq 0\}[\theta' - \theta]$$
$$= \sigma(\theta^T x) - \sigma(\theta^T x) + x^T \theta' \mathbb{1}\{\theta^T x \geq 0\}$$
$$= x^T \theta' \mathbb{1}\{\theta^T x \geq 0\}$$
$$\leq x^T \theta' \mathbb{1}\{x^T \theta' \geq 0\} = \sigma(x^T \theta').$$

② 

$$\nabla_\omega^T f(x; \omega, c)\omega = \sum_{i=1}^m \nabla_{\omega_i}^T f(x; \omega, c)\omega_i$$
$$= \sum_{i=1}^m c_i \omega_i^T x \,\mathbb{1}\{\omega_i^T x \geq 0\} = \sum_{i=1}^m c_i \sigma(\omega_i^T x).$$

**Note:** If $c_i \geq 0, \forall i$, then $f$ is convex and

$$\nabla_{\omega_i} f(x; \omega, c) \in \partial_{\omega_i} f(x; \omega, c).$$

However, $c$ can take on any value. Thus, we define

$\nabla_{\omega_i} f(x; \omega, c)$ as it may __not__ be a subgradient due to negative $c_i$.

**PART I: LINEARIZATION OF OVERPARAMETERIZED NETWORKS NEAR INITIALIZATION**

**DEF1** (Symmetric Xavier initialization)

Suppose that $m \in \mathbb{Z}_+$ is even. $\mathbb{R}^{m \times d} \times \mathbb{R}^m$-valued random variable $(\omega(0), c)$ is called a (symmetric) Xavier initialization if

$$\omega_i(0) = \omega_{i+\frac{m}{2}}(0) \overset{iid}{\sim} N(0, I_d)$$

$$c_i = -c_{i+\frac{m}{2}} \overset{iid}{\sim} \text{Rademacher},$$

for $i = 1, 2, \ldots, m.$

We consider $F(x; \omega, c) := \frac{1}{\sqrt{m}} \sum_{i=1}^m c_i \sigma(\omega_i^T x).$

Note that the scale $\frac{1}{\sqrt{m}}$ will be important.

**Lemma 2** Let $(\omega(0), c)$ be a symmetric Xavier initialization.

Then, show that $F(x; \omega(0), c) = 0$ w.p. 1 for any $x \in \mathbb{R}^d.$

Given $x \in \mathbb{R}^d$, let $(\omega(0), c)$ be a symmetric Xavier initialization and $\omega \in \mathbb{R}^{md}$ be an arbitrary weight vector. Then, w.p. 1,

$$\overbrace{F(x; \omega, c) = \nabla_\omega^T F(x; \omega(0), c)\left(\omega - \omega(0)\right)}^{\text{linear in } \omega} + \overbrace{\left(\nabla_\omega F(x; \omega, c) - \nabla_\omega F(x; \omega(0))\right)^T \omega}^{\text{nonlinearity}}$$

Pf: By Lemma 1. (2), $\nabla_\omega^T F(x; \omega(0), c)\, \omega(0) = F(x; \omega(0), c) = 0.$

Then, $\nabla_\omega^T F(x; \omega(0), c)\, \omega + \left(\nabla_\omega F(x; \omega, c) - \nabla_\omega F(x; \omega(0), c)\right)^T \omega$

$$= \nabla_\omega^T F(x; \omega, c)\, \omega = F(x; \omega, c). \quad \blacksquare$$

Remark: The idea will be to control the nonlinearity by large width $m$ (i.e., overparameterization) when $\|\omega - \omega(0)\|_2$ is small, i.e., lazy training or near-initialization or kernel regime.

We will see that $x \mapsto \nabla_\omega^T F(x; \omega(0), c)\, [\omega - \omega(0)]$ will be a powerful linear model (linear in $\omega$, nonlinear $x$).

The main result of this discussion: if $\|\omega_i - \omega_i(0)\|_2 \leq \frac{\alpha}{\sqrt{m}}$ for all $i$,

then $\left| F(x; \omega, c) - \nabla^T F(x; \omega(0), c)[\omega - \omega(0)] \right| = O\left(\frac{1}{\sqrt{m}}\right)$ with high prob.

Let $(\omega(0), c)$ be a symmetric random initialization, and $\omega \in \mathbb{R}^{md}$ be any vector s.t.

$$\max_{1 \leq i \leq m} \|\omega_i - \omega_i(0)\|_2 \leq \frac{\alpha}{\sqrt{m}},$$

for some $\alpha > 0$. Then, for any $x \in \mathbb{R}^d$, $\delta \in (0,1)$,

$$\left| F(x; \omega, c) - \nabla_\omega^T F(x; \omega(0), c)\left[\omega - \omega(0)\right]\right| \leq \frac{\alpha}{\sqrt{m}}\left(1 + \|x\|_2\right)\left[\alpha \|x\|_2 + \sqrt{\log(1/\delta)}\right],$$

with probability at least $1 - \delta$ over the random initialization.

**Proof**: By using Lemma 3, we can write:

$$F(x;\omega,c) = F(x;\omega(0),c) + \nabla_\omega^T F(x;\omega(0),c)[\omega - \omega(0)] + [\nabla F(x;\omega,c) - \nabla F(x;\omega(0),c)]^T \omega$$

By Lemma 2, we have $F(x;\omega(0),c) = 0$. Then, we can decompose $F(x;\omega,c)$ into linear and nonlinear parts (in $\omega$):

$$F(x;\omega,c) = \nabla_\omega^T F(x;\omega(0),c)[\omega - \omega(0)] + \Delta(\omega,m)$$

where

$$\Delta(\omega,c) := \left[\nabla_\omega F(x;\omega,c) - \nabla_\omega F(x;\omega(0),c)\right]^T \omega,$$

is the nonlinear term.

**Lazy training or kernel regime:**

$$\omega_i = \omega_i(0) + u_i \quad , \quad \|u_i\|_2 \leq \frac{\alpha}{\sqrt{m}} \quad \text{for all} \quad i = 1,2,\ldots,m.$$

$$\Delta(\omega,m) \triangleq \left[\nabla_\omega F(x;\omega,c) - \nabla_\omega F(x;\omega(0),c)\right]^T \omega$$

$$= \Delta_1(\omega,m) + \Delta_2(\omega,m)$$

where

$$\Delta_1(\omega,m) := \left[\nabla_\omega F(x;\omega,c) - \nabla_\omega F(x;\omega(0),c)\right]^T \omega(0)$$

$$\Delta_2(\omega,m) := \left[\nabla_\omega F(x;\omega,c) - \nabla_\omega F(x;\omega(0),c)\right]^T u$$

Note that, by definition,

$$\left[\nabla_\omega F(x;\omega,c) - \nabla_\omega F(x;\omega(0),c)\right]^T v = \frac{1}{\sqrt{m}} \sum_i c_i \left[\mathbb{1}\{\omega_i^T x \geq 0\} - \mathbb{1}\{\omega_i^T(0)x \geq 0\}\right] v_i^T x$$

for any $v = \begin{bmatrix} v_1 \\ v_2 \\ \vdots \\ v_m \end{bmatrix} \in \mathbb{R}^{md}$.

We will need Lemma 4 to conclude the proof of Theorem 1.

**Lemma 4** | Let $\omega \in \mathbb{R}^{md}$ be s.t.

$$\max_{1 \leq i \leq m} \|\omega_i - \omega_i(0)\|_2 \leq \frac{\alpha}{\sqrt{m}} \quad \text{for some} \quad \alpha > 0. \text{ Then,}$$

for any $x \in \mathbb{R}^d$ and $\delta \in (0,1)$, $\exists A_\delta(x) \in \sigma(\omega(0), c)$ s.t.

$$\mathbb{P}(A_\delta(x)) \geq 1 - \delta, \text{ and}$$

$$|\Delta_1(\omega, m)| \leq \frac{\alpha}{\sqrt{m}} \cdot \left[ \alpha \cdot \|x\|_2 + \sqrt{\log(1/\delta)} \right] \text{ and}$$

$$|\Delta_2(\omega, m)| \leq \frac{\alpha \|x\|_2}{\sqrt{m}} \left[ \alpha \cdot \|x\|_2 + \sqrt{\log(1/\delta)} \right] \text{ in } A_\delta(x).$$

---

**Proof of Lemma 4**

Then,

$$|\Delta_1(\omega, m)| = \left| \frac{1}{m} \sum_{i=1}^{m} c_i \left[ \mathbb{1}\{\omega_i^T x \geq 0\} - \mathbb{1}\{\omega_i^T(0) x \geq 0\} \right] \omega_i^T(0) x \right|$$

$$\leq \frac{1}{\sqrt{m}} \sum_{i=1}^{m} \left| \mathbb{1}\{\omega_i^T x \geq 0\} - \mathbb{1}\{\omega_i^T(0) x \geq 0\} \right| \cdot |\omega_i^T(0) x|$$

Now, we start counting sign changes from $\omega_i^T(0) x$ to $\omega_i^T x$.

Let $$S(x) := \left\{ i \in [m] : \mathbb{1}\{\omega_i^T(0) x \geq 0\} \neq \mathbb{1}\{\omega_i^T x \geq 0\} \right\}.$$

$$\left| \mathbb{1}\{\omega_i^T x \geq 0\} - \mathbb{1}\{\omega_i^T(0) x \geq 0\} \right| = \mathbb{1}\{i \in S(x)\}.$$

Also,

$$i \in S(x) \implies |\omega_i^T(0) x| \leq |\omega_i^T(0) x - \omega_i^T x|$$

$$\leq \|\omega_i(0) - \omega_i\|_2 \cdot \|x\|_2$$

$$= \|u_i\|_2 \cdot \|x\|_2 \leq \frac{\alpha \|x\|_2}{\sqrt{m}}.$$

Then,

$$|\Delta_1(\omega, m)| \leq \frac{1}{\sqrt{m}} \sum_i \mathbb{1}\{i \in S(x)\} \cdot |\omega_i^T(0) x|$$

$$\leq \frac{\alpha}{m} \sum_{i=1}^{m} \mathbb{1}\{i \in S(x)\} \leq \alpha \cdot \frac{1}{m} \sum_{i=1}^{m} \mathbb{1}\{|\omega_i^T(0) x| \leq \frac{\alpha \|x\|_2}{\sqrt{m}}\}.$$

**Lemma 5** (Gaussian Anti-Concentration)

Let $u \sim N(0, I_d)$, $x \neq 0$. Then, for any $\varepsilon > 0$,

$$\mathbb{P}\left( |u^T x| \leq \|x\| \cdot \varepsilon \right) \leq \sqrt{\frac{2}{\pi}} \cdot \varepsilon.$$

**Proof:** For any $x \neq 0$, $\dfrac{u^T x}{\|x\|_2} \sim N(0,1)$ since

$$\mathbb{E}\left[ \frac{u^T x}{\|x\|_2} \right] = \left( \mathbb{E}[u] \right)^T \frac{x}{\|x\|_2} = 0, \text{ and}$$

$$\mathrm{Var}\left( u^T \frac{x}{\|x\|_2} \right) = \mathbb{E}\left[ \frac{x^T u}{\|x\|_2} \cdot \frac{u^T x}{\|x\|_2} \right] = \frac{x^T \mathbb{E}[u\, u^T] x}{\|x\|_2^2} = 1.$$

Then,

$$\mathbb{P}\left( |u^T x| < \|x\| \cdot \varepsilon \right) = \mathbb{P}\left( -\varepsilon \leq \frac{u^T z}{\|x\|_2} \leq \varepsilon \right)$$

$$= \int_{-\varepsilon}^{\varepsilon} \frac{1}{\sqrt{2\pi}} \underbrace{e^{-\frac{z^2}{2}}}_{\leq 1} dy \leq \int_{-\varepsilon}^{\varepsilon} \frac{1}{\sqrt{2\pi}} dy$$

$$\leq \frac{\sqrt{2}\varepsilon}{\sqrt{\pi}}. \quad \blacksquare$$

Also,

$$\frac{1}{m} \sum_{i=1}^{m} \mathbb{1}\left\{ |w_i^T(0) x| \leq \frac{\alpha}{\sqrt{m}} \cdot \|x\|_2 \right\} - \mathbb{P}\left( |u^T x| \leq \frac{\alpha}{\sqrt{m}} \|x\|_2 \right) \leq \sqrt{\frac{\log(1/\delta)}{m}}$$

with probability at least $1 - \delta$.

Thus, let

$$A_1(x) := \left\{ \frac{1}{m} \sum_{i=1}^{m} \mathbb{1}\left\{ |w_i^T(0) x| \leq \frac{\alpha \|x\|_2}{\sqrt{m}} \right\} \leq \frac{1}{\sqrt{m}}\left[ \alpha \|x\|_2 + \sqrt{\log(1/\delta)} \right] \right\}.$$

Then, $\mathbb{P}(A_1(x)) \geq 1 - \delta$, and

$$|\Delta_1(w, x)| \leq \frac{\alpha}{\sqrt{m}}\left[ \alpha \cdot \|x\|_2 + \sqrt{\log(1/\delta)} \right] \text{ in } A_1(x).$$

$$\Delta_2(\omega, z) = \left[ \nabla F(z; \omega, c) - \nabla F(z, \omega(0), c) \right]^T u$$

$$= \frac{1}{\sqrt{m}} \sum_{i=1}^{n} c_i \left[ \mathbf{1}\{\omega_i^T z \geq 0\} - \mathbf{1}\{\omega_i^T(0) z \geq 0\} \right] u_i^T x$$

Recall that $\max_i \|u_i\|_2 \leq \frac{\alpha}{\sqrt{m}}$. Then,

$$|\Delta_2(\omega, z)| \leq \frac{1}{\sqrt{m}} \sum_{i=1}^{m} \left| \mathbf{1}\{\omega_i^T z \geq 0\} - \mathbf{1}\{\omega_i^T(0) z \geq 0\} \right| \cdot \underbrace{|u_i^T x|}_{\text{Cauchy-Schwarz}}$$

$$\color{blue}{\text{Cauchy-Schwarz}}$$
$$\color{blue}{\leq \|u_i\|_2 \cdot \|x\|_2}$$
$$\color{blue}{\leq \frac{\alpha}{\sqrt{m}} \cdot \|x\|_2}$$

$$\leq \frac{\alpha \|x\|_2}{m} \sum_i \mathbf{1}\{ i \in S(z)\}$$

$$\leq \frac{\beta}{\sqrt{m}} \left[ \alpha \cdot \|x\|_2 + \sqrt{\log(1/\delta)} \right] \quad \text{in} \quad A_1(x).$$

Thus, we conclude the proof of Lemma 4. ∎

Coming back to the proof of Theorem 1,

$$\left| F(x; \omega, c) - \nabla^T F(z; \omega(0), c) \left[ \omega - \omega(0) \right] \right| = |\Delta(\omega, m)|$$

$$= |\Delta_1(\omega, m) + \Delta_2(\omega, m)|$$

$$\leq |\Delta_1(\omega, m)| + |\Delta_2(\omega, m)|.$$

Substituting the bounds on $|\Delta_1|$ and $|\Delta_2|$ that we found in Lemma 4, we conclude the proof of Theorem 1. ∎