

## NTK Proof

$$f(x, \theta) = \sigma(wx + b) \quad w, b \in \theta$$

$$L(f, x, y, \theta) = (f(x, \theta) - y)^2$$

$$\theta_{t+1} = \theta_t + \eta \nabla_w L(f, x, y, \theta)$$

Denote  $\theta_t$  as dynamic  $\theta(t)$

$$\frac{d\theta(t)}{dt} = \nabla_{\theta} L(f, x, y, \theta)$$

$$= 2(f(x, \theta) - y) \nabla_{\theta} f(x, \theta)$$

$$= 2(f(x, \theta) - y) \nabla_{\theta} (f(x, \theta_0) + \nabla_{\theta} f(x, \theta_0)^{\top} (\theta - \theta_0))$$

$$= 2(f(x, \theta) - y) \nabla_{\theta} (f(x, \theta_0) + \nabla_{\theta} f(x, \theta_0)^{\top} \theta - \nabla_{\theta} f(x, \theta_0)^{\top} \theta_0)$$

$$= 2(f(x, \theta) - y) \nabla_{\theta} f(x, \theta_0)$$

---

## Abstract

To understand what is the neural tangent kernel(NTK), there are 3 important result that need to remember.

- When the width of the neural network goes to infinity, the network will be equivalent to a Gaussian process.
- When the width of the neural network goes to infinity, the weight of the network will remain almost unchanged during training. That is, the neural network will become a linear model.
- Since the network will become a linear model, the loss surface of the MSE of the infinite-width network will be a convex. As a result, we can optimize the infinite-width network just like optimizing a linear regression and solve it by ODE.

Combine these 3 points, NTK is a kernel that can kernelize the neural network architecture and it provides a closed-form solution of the kernel for anytime. Thus, We can compute the NTK at the end of training without actual training and compute the posterior and the prediction of the testing data with Bayesian inference.

## Infinite-Width Neural Network As Gaussian Process

### Define A Neural Network

First, given a training dataset  $\mathcal{D}$  including  $N = |\mathcal{D}|$  data points. We denote the data point as  $d_i = (x, y) \in \mathcal{D} \forall i$  where the feature vector  $x \in \mathbb{R}^{k \times 1}$  and the label  $y \in \mathbb{R}$ . The set of the feature vectors is  $\mathcal{X} = \{x : (x, y) \in \mathcal{D}\}$  and similarly, the set of the label  $\mathcal{Y} = \{y : (x, y) \in \mathcal{D}\}$ .

We represent the element of the neural network  $f(x, \theta)$ ,  $x \in \mathcal{X}$  as following

$$h^1 = W^1 x + b^1 h^{l+1} = W^{l+1} a^l + b^{l+1} = W^{l+1} \phi(h^l) + b^{l+1} \hat{y} = f(x, \theta) = h^L$$

where  $\phi$  is the activation function,  $h^{l+1} \in \mathbb{R}^{n_{l+1} \times 1}$  and  $a^{l+1} \in \mathbb{R}^{n_{l+1} \times 1}$  are the pre-activation and the post-activation respectively. The parameter  $\theta$  includes  $W^l, b^l \in \theta \forall l$ .  $W^{l+1} \in \mathbb{R}^{n_{l+1} \times n_l}$  and  $b^{l+1} \in \mathbb{R}^{n_{l+1} \times 1}$  are the weight and the bias respectively. They follow the LeCun and normal distribution respectively. Also, we denote the  $n_{l+1}$  as the network width of  $l+1$ -th layer and  $n_0 = k$  is the size of the feature vector. Finally, we denote  $\hat{y} \in \mathbb{R}$  is the prediction of the network based on input  $x$ . For convenience, we also denote  $\hat{\mathcal{Y}} = f(\mathcal{X}, \theta) \in \mathbb{R}^N$ .

$$W_{i,j}^{l+1} = \frac{\sigma_w}{\sqrt{n_{l+1}}} w_{i,j}^{l+1}, \quad b_i^{l+1} = \sigma_b \beta_i^{l+1} w_{i,j}^{l+1}, \beta_i^{l+1} \stackrel{i.i.d}{\sim} \mathcal{N}(0, 1) \forall l, i, j \quad 1 \leq l \leq L, \quad 1 \leq i \leq n_{l+1}, \quad 1 \leq j \leq n_l$$

where  $W_{i,j}^{l+1} \in \mathbb{R}$  is the entry of the matrix of  $l+1$ -th layer at the  $i$ -th row and  $j$ -th column.  $b_i^l \in \mathbb{R}$  is the  $i$ -th entry of the bias vector.

Combine them.

$$h_i^{l+1} = W_i^{l+1} a^l + b_i^{l+1} = \sum_{j=1}^{n_{l+1}} W_{i,j}^{l+1} a_j^l + b_i^{l+1} = \frac{\sigma_w}{\sqrt{n_{l+1}}} \sum_{j=1}^{n_{l+1}} w_{i,j}^{l+1} \phi(h_j^l) + \sigma_b \beta_i^{l+1}$$

where  $W_i^{l+1} \in \mathbb{R}^{1 \times n_l}$  is the  $i$ -th row of the matrix  $W^{l+1}$ . The post-activation vector is  $a^l \in \mathbb{R}^{n_l \times 1}$  and the  $j$ -th entry of the vector is  $a_j^l \in \mathbb{R}$ .

Let  $\theta^{l+1}$  denote as the all parameters of the layer  $l + 1$ .

$$\theta^{l+1} = \text{vec}(W^{l+1}, b^{l+1}) \in \mathbb{R}^{n_{l+1} \times (n_l + 1)} \theta = \text{vec}(\cup_{l=1}^L \theta^l) \in \mathbb{R}^{S \times 1}$$

where  $S = \sum_{l=1}^L n_l(n_{l-1} + 1)$  is the number of all parameters.

We also denote the empirical loss of all training dataset as  $\mathcal{L}(\mathcal{X}, \mathcal{Y})$ , and  $l(\hat{y}, y)$  as loss function.

$$l(\hat{y}, y) = l(f(x, \theta), y)$$

$$\mathcal{L}_{f(\cdot)}(\mathcal{X}, \mathcal{Y}) = \sum_{x \in \mathcal{X}, y \in \mathcal{Y}} l(f(x, \theta), y) \in \mathbb{R}$$

### The Mean Of The Layers

$$\begin{aligned} E[h_i^l(x)] &= E\left[\frac{\sigma_w}{\sqrt{n_1}} \sum_{j=1}^{n_1} w_{i,j}^1 a_j^{l-1} + \sigma_b \beta_i^1\right] \\ &= \frac{\sigma_w}{\sqrt{n_1}} \sum_{j=1}^{n_1} E[w_{i,j}^1] a_j^{l-1} + \sigma_b E[\beta_i^1] = 0 \end{aligned}$$

### The Covariance Of The First Layer

Consider the same entries.

$$\begin{aligned} k^1(x, x') &= \text{Cov}[h_i^1(x), h_i^1(x')] \\ &= E[(h_i^1(x) - E[h_i^1(x)])(h_i^1(x') - E[h_i^1(x')])] \\ &= E[h_i^1(x)h_i^1(x')] - E[h_i^1(x)]E[h_i^1(x')] \\ &= E\left[\left(\frac{\sigma_w}{\sqrt{n_1}} \sum_{j=1}^{n_1} w_{i,j}^1 x_j + \sigma_b \beta_i^1\right)\left(\frac{\sigma_w}{\sqrt{n_1}} \sum_{j=1}^{n_1} w_{i,j}^1 x'_j + \sigma_b \beta_i^1\right)\right] \\ &= E\left[\frac{\sigma_w^2}{n_1} \left(\sum_{j=1}^{n_1} w_{i,j}^1 x_j\right)\left(\sum_{k=1}^{n_1} w_{i,k}^1 x'_k\right) + (\sigma_b \beta_i^1)(\sigma_b \beta_i^1) + \frac{\sigma_w}{\sqrt{n_1}} (\sigma_b \beta_i^1 \sum_{j=1}^{n_1} w_{i,j}^1 x_j) + \frac{\sigma_w}{\sqrt{n_1}} (\sigma_b \beta_i^1 \sum_{j=1}^{n_1} w_{i,j}^1 x'_j)\right] \end{aligned}$$

$$= \frac{\sigma_w^2}{n_1} (E[\sum_{j=k}^{n_1} w_{i,j}^1 x_j w_{i,k}^1 x'_k] + E[\sum_{j \neq k}^{n_1} w_{i,j}^1 x_j w_{i,k}^1 x'_k]) + \sigma_b^2 E[(\beta_i^1)^2] + \frac{\sigma_w \sigma_b}{\sqrt{n_1}} (E[\beta_i^1] E[\sum_{j=1}^{n_1} w_{i,j}^1 x_j]) + \frac{\sigma_w \sigma_b}{\sqrt{n_1}} (E[\beta_i^1] E[\sum_{j=1}^{n_1} w_{i,j}^1 x'_j])$$

Since  $E[\beta_i^1] = E[w_{i,j}^1] = 0$ , thus we can derive the following formula.

$$= \frac{\sigma_w^2}{n_1} (\sum_{j=k}^{n_1} E[w_{i,j}^1 w_{i,k}^1] x_j x'_k + \sum_{j \neq k}^{n_1} E[w_{i,j}^1 w_{i,k}^1] x_j x'_k) + \sigma_b^2 (Var[\beta_i^1] + E[\beta_i^1]^2)$$

Since  $w_{i,j}^1, w_{i,k}^1 \stackrel{i.i.d}{\sim} \mathcal{N}(0, 1)$ , thus,  $E[w_{i,j}^1 w_{i,k}^1] = E[w_{i,j}^1] E[w_{i,k}^1] = 0$ . Also,  $Var[\beta_i^1] = 1$ .

$$\begin{aligned} &= \frac{\sigma_w^2}{n_1} (\sum_{j=1}^{n_1} (Var[w_{i,j}^1] + E[w_{i,j}^1]^2) x_j x'_j) + \sigma_b^2 \\ &= \frac{\sigma_w^2}{n_1} (\sum_{j=1}^{n_1} x_j x'_j) + \sigma_b^2 \\ &= \frac{\sigma_w^2}{n_1} x^\top x' + \sigma_b^2 \end{aligned}$$

**Consider the different entries.**

$$\begin{aligned} k_{i \neq j}^1(x, x') &= Cov[h_i^1(x), h_j^1(x')], \quad i \neq j \\ &= E[(\frac{\sigma_w}{\sqrt{n_1}} \sum_{k=1}^{n_1} w_{i,k}^1 x_k + \sigma_b \beta_i^1) (\frac{\sigma_w}{\sqrt{n_1}} \sum_{k'=1}^{n_1} w_{j,k'}^1 x'_{k'} + \sigma_b \beta_j^1)] \\ &= \frac{\sigma_w^2}{n_1} (E[\sum_{k=1}^{n_1} \sum_{k'=1}^{n_1} w_{i,k}^1 x_k w_{j,k'}^1 x'_{k'}]) + \sigma_b^2 E[\beta_i^1 \beta_j^1] + \frac{\sigma_w \sigma_b}{\sqrt{n_1}} (E[\beta_i^1] E[\sum_{k=1}^{n_1} w_{i,k}^1 x_k]) + \frac{\sigma_w \sigma_b}{\sqrt{n_1}} (E[\beta_j^1] E[\sum_{k'=1}^{n_1} w_{j,k'}^1 x'_{k'}]) \\ &= \frac{\sigma_w^2}{n_1} (\sum_{k=1}^{n_1} \sum_{k'=1}^{n_1} E[w_{i,k}^1] E[w_{j,k'}^1] x_k x'_{k'}) = 0 \end{aligned}$$

## The Covariance Of The Deeper Layers

Consider the different entries.

$$\begin{aligned}
k_{i \neq j}^l(x, x') &= Cov[h_i^l(x), h_j^l(x')], \quad i \neq j \\
&= E[(\frac{\sigma_w}{\sqrt{n_l}} \sum_{k=1}^{n_l} w_{i,k}^l \phi(h_k^l(x)) + \sigma_b \beta_i^l)(\frac{\sigma_w}{\sqrt{n_l}} \sum_{k'=1}^{n_l} w_{j,k'}^l \phi(h_{k'}^l(x')) + \sigma_b \beta_j^l)] \\
&= \frac{\sigma_w^2}{n_l} (E[\sum_{k=1}^{n_l} \sum_{k'=1}^{n_l} w_{i,k}^l \phi(h_k^l(x)) w_{j,k'}^l \phi(h_{k'}^l(x'))]) + \sigma_b^2 E[\beta_i^l \beta_j^l] + \frac{\sigma_w \sigma_b}{\sqrt{n_l}} (E[\beta_i^l] E[\sum_{k=1}^{n_l} w_{i,k}^l \phi(h_k^l(x))]) + \frac{\sigma_w \sigma_b}{\sqrt{n_l}} (E[\beta_j^l] E[\sum_{k'=1}^{n_l} w_{j,k'}^l \phi(h_{k'}^l(x'))]) \\
&= \frac{\sigma_w^2}{n_l} (\sum_{k=1}^{n_l} \sum_{k'=1}^{n_l} E[w_{i,k}^l] E[w_{j,k'}^l] \phi(h_k^l(x)) \phi(h_{k'}^l(x'))) = 0
\end{aligned}$$

We've got the covariance of different entries  $k_{i \neq j}^l(x, x')$  of sample  $x$  and  $x'$  is 0. Thus, we can conclude that different entries are independent.

Consider the same entries.

$$\begin{aligned}
k^l(x, x') &= Cov[h_i^l(x), h_i^l(x')] \\
&= E[(h_i^l(x) - E[h_i^l(x)])(h_i^l(x') - E[h_i^l(x')])] \\
&= E[h_i^l(x) h_i^l(x')] - E[h_i^l(x)] E[h_i^l(x')] \\
&= E[(\frac{\sigma_w}{\sqrt{n_1}} \sum_{j=1}^{n_1} w_{i,j}^1 \phi(h_j^{l-1}(x)) + \sigma_b \beta_i^1)(\frac{\sigma_w}{\sqrt{n_1}} \sum_{j=1}^{n_1} w_{i,j}^1 \phi(h_j^{l-1}(x')) + \sigma_b \beta_i^1)] \\
&= \frac{\sigma_w^2}{n_l} (E[\sum_{j=k}^{n_l} w_{i,j}^l \phi(h_j^{l-1}(x)) w_{i,k}^l \phi(h_j^{l-1}(x'))]) + E[\sum_{j \neq k}^{n_l} w_{i,j}^l \phi(h_j^{l-1}(x)) w_{i,k}^l \phi(h_j^{l-1}(x'))]) + \sigma_b^2 E[(\beta_i^l)^2] + \frac{\sigma_w \sigma_b}{\sqrt{n_l}} (E[\beta_i^l] E[\sum_{j=1}^{n_l} w_{i,j}^l \phi(h_j^{l-1}(x))]) \\
&= \frac{\sigma_w^2}{n_1} (\sum_{j=1}^{n_1} \phi(h_j^{l-1}(x)) \phi(h_j^{l-1}(x'))) + \sigma_b^2
\end{aligned}$$

Recall **central limit theorem (CLT)**, suppose we have random variables  $X_1, X_2, \dots, X_n$  are i.i.d, the expectation and the variance of the random variables are  $E[X_i] = \mu$  and  $Var[X_i] = \sigma^2$ . Denote the sum of the random variables is  $\bar{X} = \frac{1}{n} \sum_{i=1}^n X_i$ . The normalized  $\bar{X}$  is  $\zeta = \frac{\bar{X} - \mu}{\frac{\sigma}{\sqrt{n}}}$ .

$$\lim_{n \rightarrow \infty} P_\zeta(\zeta \leq z) = P_Z(Z \leq z)$$

where random variable  $Z \sim \mathcal{N}(0, 1)$  and  $z \in \mathbb{R}$  is a scalar. The functions  $P_\zeta$  and  $P_Z$  are C.D.F of the random variables  $\zeta$  and  $Z$  respectively.

**@ Need Proof: Connection between CLT and infinite width**

When  $n_l \rightarrow \infty$ , apply central limit theorem (CLT)

$$k^l(x, x') = \lim_{n_l \rightarrow \infty} \frac{\sigma_w^2}{n_1} \left( \sum_{j=1}^{n_1} \phi(h_j^{l-1}(x)) \phi(h_j^{l-1}(x')) \right) + \sigma_b^2 = \sigma_w^2 E_{h_j^{l-1}(x), h_j^{l-1}(x') \sim \mathcal{N}(0, K_{x, x'}^{l-1})} [\phi(h_j^{l-1}(x)) \phi(h_j^{l-1}(x'))] + \sigma_b^2$$

$$K_{xx'}^{l-1} = \begin{bmatrix} k^{l-1}(x, x) & k^{l-1}(x, x') \\ k^{l-1}(x', x) & k^{l-1}(x', x') \end{bmatrix}$$

$$k^l(x, x') = \frac{\sigma_w^2}{n_1} x^\top x' + \sigma_b^2$$

Define an operator  $\Gamma$ , for a positive semi-definite matrix  $\Sigma \in \mathbb{R}^{2 \times 2}$

$$\Gamma(\Sigma) = \sigma_w^2 E_{u, v \sim \mathcal{N}(0, \Sigma)} [\phi(u) \phi(v)] + \sigma_b^2$$

With  $\Gamma$ , the NNGP kernel  $\Sigma_{\mathcal{X}}^L$ .

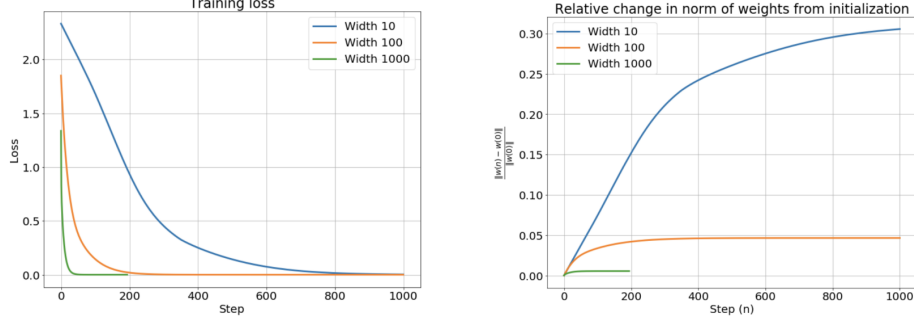
$$\Sigma_{\mathcal{X}}^L = \begin{bmatrix} k^L(x_1, x_1) & k^L(x_1, x_2) & \dots & k^L(x_1, x_N) \\ k^L(x_2, x_1) & k^L(x_2, x_2) & \dots & k^L(x_2, x_N) \\ \vdots & \vdots & \ddots & \vdots \\ k^L(x_N, x_1) & k^L(x_N, x_2) & \dots & k^L(x_N, x_N) \end{bmatrix}$$

$$k^L(x_1, x_2) = \Gamma(K_{x_1 x_2}^{L-1}) = \sigma_w^2 E_{u, v \sim \mathcal{N}(0, K^{L-1})} [\phi(u) \phi(v)] + \sigma_b^2$$

$$K_{x_1 x_2}^l = \begin{bmatrix} k^l(x_1, x_1) & k^l(x_1, x_2) \\ k^l(x_2, x_1) & k^l(x_2, x_2) \end{bmatrix}$$

## Infinite-Width Neural Network As A Linear Model

### Linear Model



As we've shown in the previous post, the parameters of the neural network change more slightly while the width of the network gets larger. In the other words, the neural network **remains almost unchanged during training**. As a result, the parameters  $\theta^{(T)} \in \mathbb{R}$  of neural network after training  $T$  steps will be very close to the initial parameters  $\theta^{(0)}$ .

$$\lim_{n_l \rightarrow \infty} f(x, \theta^{(t)}) = f(x, \theta^{(0)}), \quad 1 \leq l \leq L, \quad \forall t$$

where  $\bar{f}(x, \cdot)$  is the approximation function of neural network  $f(x, \cdot)$ .

### Taylor Expansion

Since the parameters of the infinite-width neural network only change slightly, thus, we can expand the neural network with Taylor expansion.

Taylor expansion

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

First-order Taylor expansion

$$g(x) \approx g(a) + \frac{dg(a)}{dx} (x - a)$$

We denote the parameters of the neural network at training step  $t$  as  $\theta^{(t)}$ . Then, expand the neural network  $f(x, \theta^{(t)})$  at training step  $t$  with data point  $x$ .

$$\hat{y}^{(t)} = f(x, \theta^{(t)}) \approx f_{lin}(x, \theta^{(t)}) = f(x, \theta^{(0)}) + \nabla_{\theta} f(x, \theta^{(0)}) (\theta^{(t)} - \theta^{(0)}), \quad \forall t$$

where  $\hat{y}^{(t)}$  is the prediction of the data point  $x$  from the network at training step  $t$ . The divergence  $\nabla_{\theta} f(x, \theta^{(0)}) \in \mathbb{R}^{N \times S}$  is a **Jacobian matrix**.

As for the whole dataset  $\mathcal{X}$ , the predictions  $\hat{\mathcal{Y}}^{(t)}$

$$\hat{\mathcal{Y}}^{(t)} = f(\mathcal{X}, \theta^{(t)}) \approx f_{lin}(x, \theta^{(t)}) = f(\mathcal{X}, \theta^{(0)}) + \nabla_{\theta} f(\mathcal{X}, \theta^{(0)})(\theta^{(t)} - \theta^{(0)}), \quad \forall t$$

$$f_{lin}(\mathcal{X}, \theta^{(t)}) - f(\mathcal{X}, \theta^{(0)}) = \nabla_{\theta} f(\mathcal{X}, \theta^{(0)})(\theta^{(t)} - \theta^{(0)})$$

Because  $f_{lin}(\mathcal{X}, \theta^{(0)}) = f(\mathcal{X}, \theta^{(0)}) + \nabla_{\theta} f(\mathcal{X}, \theta^{(0)})(\theta^{(0)} - \theta^{(0)}) = f(\mathcal{X}, \theta^{(0)})$ , replace  $f(\mathcal{X}, \theta^{(0)})$  with  $f_{lin}(\mathcal{X}, \theta^{(0)})$

$$f_{lin}(\mathcal{X}, \theta^{(t)}) - f_{lin}(\mathcal{X}, \theta^{(0)}) = \nabla_{\theta} f_{lin}(\mathcal{X}, \theta^{(0)})(\theta^{(t)} - \theta^{(0)})$$

where  $\nabla_{\theta} f_{lin}(x, \theta^{(0)}) \in \mathbb{R}^{N \times S}$  is also a **Jacobian matrix**.

#### @ Need Hessian Proof

With Lipschitz continuity, we can guarantee that  $\theta^{(t+\Delta t)} - \theta^{(0)} \geq \theta^{(t+\Delta t)} - \theta^{(t)}$  and  $f(\mathcal{X}, \theta^{(t+\Delta t)}) - f(\mathcal{X}, \theta^{(0)}) \geq f(\mathcal{X}, \theta^{(t+\Delta t)}) - f(\mathcal{X}, \theta^{(t)})$ , and we can also write down the formula.

$$\lim_{\Delta t \rightarrow 0} \frac{f_{lin}(\mathcal{X}, \theta^{(t+\Delta t)}) - f_{lin}(\mathcal{X}, \theta^{(t)})}{\Delta t} = \nabla_{\theta} f_{lin}(\mathcal{X}, \theta^{(0)}) \lim_{\Delta t \rightarrow 0} \frac{(\theta^{(t+\Delta t)} - \theta^{(t)})}{\Delta t}$$

$$\frac{\partial f_{lin}(\mathcal{X}, \theta^{(t)})}{\partial t} = \nabla_{\theta} f_{lin}(\mathcal{X}, \theta^{(0)}) \frac{\partial \theta^{(t)}}{\partial t}$$

In the form of dynamic

$$\dot{f}_{lin}(\mathcal{X}, \theta^{(t)}) = \nabla_{\theta} f_{lin}(\mathcal{X}, \theta^{(0)}) \dot{\theta}(t)$$

#### Combine With Gradient Descent

Denote the loss of linearized neural network  $f_{lin}(\cdot, \cdot)$  on the training dataset  $(\mathcal{X}, \mathcal{Y})$  at the training step  $t$  as  $\mathcal{L}_{lin}^{(t)}(\mathcal{X}, \mathcal{Y}) = \sum_{x \in \mathcal{X}, y \in \mathcal{Y}} l(f_{lin}(x, \theta^{(t)}), y)$ . The gradient descent can be represented as the following formula.

$$\theta^{(t+1)} = \theta^{(t)} + \eta \nabla_{\theta} \mathcal{L}_{lin}^{(t)}(\mathcal{X}, \mathcal{Y}) = \theta^{(t)} + \eta \nabla_{\theta} f_{lin}(\mathcal{X}, \theta^{(t)})^{\top} \nabla_{f_{lin}(\mathcal{X}, \theta^{(t)})} \mathcal{L}_{lin}^{(t)}(\mathcal{X}, \mathcal{Y})$$



$$\theta^{(t+1)} - \theta^{(t)} = \eta \nabla_{\theta} f_{lin}(\mathcal{X}, \theta^{(t)})^{\top} \nabla_{f_{lin}(\mathcal{X}, \theta^{(t)})} \mathcal{L}_{lin}^{(t)}(\mathcal{X}, \mathcal{Y})$$

where  $\nabla_{f_{lin}(\mathcal{X}, \theta^{(t)})} \mathcal{L}_{lin}^{(t)}(\mathcal{X}, \mathcal{Y}) \in \mathbb{R}^{N \times 1}$

In the form of gradient flow dynamic

$$\dot{\theta}(t) = \eta \nabla_{\theta} f_{lin}(\mathcal{X}, \theta^{(t)})^{\top} \nabla_{f_{lin}(\mathcal{X}, \theta^{(t)})} \mathcal{L}_{lin}^{(t)}(\mathcal{X}, \mathcal{Y})$$

Replace  $\dot{\theta}(t)$  with  $\eta \nabla_{\theta} f_{lin}(\mathcal{X}, \theta^{(t)})^{\top} \nabla_{f_{lin}(\mathcal{X}, \theta^{(t)})} \mathcal{L}_{lin}^{(t)}(\mathcal{X}, \mathcal{Y})$

$$\dot{f}_{lin}(\mathcal{X}, \theta^{(t)}) = \eta \nabla_{\theta} f_{lin}(\mathcal{X}, \theta^{(t)}) \nabla_{\theta} f_{lin}(\mathcal{X}, \theta^{(t)})^{\top} \nabla_{f_{lin}(\mathcal{X}, \theta^{(t)})} \mathcal{L}_{lin}^{(t)}(\mathcal{X}, \mathcal{Y})$$

Then, we can expand  $\nabla_{\theta} f_{lin}(\mathcal{X}, \theta^{(t)})$  with Taylor expansion.

$$\nabla_{\theta} f_{lin}(\mathcal{X}, \theta^{(t)}) = \nabla_{\theta} (f(\mathcal{X}, \theta^{(0)}) + \nabla_{\theta} f(\mathcal{X}, \theta^{(0)}) (\theta^{(t)} - \theta^{(0)})) = \nabla_{\theta} (f(\mathcal{X}, \theta^{(0)}) + \nabla_{\theta} f(\mathcal{X}, \theta^{(0)}) \theta^{(t)} - \nabla_{\theta} f(\mathcal{X}, \theta^{(0)}) \theta^{(0)})$$

Finally, we can replace  $\nabla_{\theta} f_{lin}(\mathcal{X}, \theta^{(t)})$  with  $\nabla_{\theta} f(\mathcal{X}, \theta^{(0)})$ , since they are identical due to the linearity.

$$\dot{f}_{lin}(\mathcal{X}, \theta^{(t)}) = \eta \nabla_{\theta} f(\mathcal{X}, \theta^{(0)}) \nabla_{\theta} f(\mathcal{X}, \theta^{(0)})^{\top} \nabla_{f_{lin}(\mathcal{X}, \theta^{(t)})} \mathcal{L}_{lin}^{(t)}(\mathcal{X}, \mathcal{Y})$$

Let  $T_{\mathcal{X}\mathcal{X}}^{(0)} = \nabla_{\theta} f(\mathcal{X}, \theta^{(0)}) \nabla_{\theta} f(\mathcal{X}, \theta^{(0)})^{\top}$

$$\dot{f}_{lin}(\mathcal{X}, \theta^{(t)}) = \eta T_{\mathcal{X}\mathcal{X}}^{(0)} \nabla_{f_{lin}(\mathcal{X}, \theta^{(t)})} \mathcal{L}_{lin}^{(t)}(\mathcal{X}, \mathcal{Y})$$

where  $T_{\mathcal{X}\mathcal{X}}^{(0)} \in \mathbb{R}^{|\mathcal{D}| \times |\mathcal{D}|}$  is the **Neural Tangent Kernel (NTK)**

## Gradient Flow Of MSE As A Linear Regression

### Mean Square Error(MSE)

In the previous section, we've derive the relation between the prediction of the neural network and the gradient descent at time step  $t$ . In this section, we'll dive into the point of view of gradient flow. Now, we've known that the linear approximation of the neural network regression.

$$\dot{f}_{lin}(\mathcal{X}, \theta^{(t)}) = \eta T_{\mathcal{X}\mathcal{X}}^{(0)} \nabla_{f_{lin}(\mathcal{X}, \theta^{(t)})} \mathcal{L}_{lin}^{(t)}(\mathcal{X}, \mathcal{Y})$$

Usually, in regression tasks, we use **Mean Square Error(MSE)** as error function.

$$l(f_{lin}(x, \theta), y) = \frac{1}{2} \|f_{lin}(x, \theta) - y\|_2^2$$

Plug in the MSE into the neural network regression

$$\begin{aligned} \dot{f}_{lin}(\mathcal{X}, \theta^{(t)}) &= \eta T_{\mathcal{X}\mathcal{X}}^{(0)} \nabla_{f_{lin}(\mathcal{X}, \theta^{(t)})} \mathcal{L}_{lin}^{(t)}(\mathcal{X}, \mathcal{Y}) \\ &= \eta T_{\mathcal{X}\mathcal{X}}^{(0)} \nabla_{f_{lin}(\mathcal{X}, \theta^{(t)})} \sum_{x \in \mathcal{X}, y \in \mathcal{Y}} \frac{1}{2} \|f_{lin}(x, \theta^{(t)}) - y\|_2^2 \\ &= \eta T_{\mathcal{X}\mathcal{X}}^{(0)} \sum_{x \in \mathcal{X}, y \in \mathcal{Y}} \|f_{lin}(x, \theta^{(t)}) - y\|_2 \\ &= \eta T_{\mathcal{X}\mathcal{X}}^{(0)} (f_{lin}(\mathcal{X}, \theta^{(t)}) - \mathcal{Y}) \end{aligned}$$

### The Dynamic

In the previous post, we've shown that we can continue the discrete step and take the discrete trajectory as a continuous trajectory as if the step is small enough. For instance, Malthus' population growth model says that "the population growth rate is proportional to the total population".

$$\begin{cases} \frac{dP(t)}{dt} = \lambda P(t), \lambda > 0 \\ P(t_0) = P_0 \end{cases}$$

where the dynamic  $P(t)$  is the population at time  $t$  and the  $P_0$  is the initial population. Usually, we also denote  $\frac{dP(t)}{dt} = \dot{P}(t)$ .

Now, we can solve the dynamic  $P(t)$ .

$$\begin{aligned} \int_{t_0}^{t_1} \frac{dP(t)}{dt} \frac{1}{P(t)} dt &= \int_{t_0}^{t_1} \lambda dt \\ \int_{t_0}^{t_1} \frac{d \ln(P(t))}{dt} dt &= \lambda(t_1 - t_0) \\ \ln(P(t_1)) - \ln(P(t_0)) &= \lambda(t_1 - t_0) \end{aligned}$$

$$\ln\left(\frac{P(t_1)}{P(t_0)}\right) = \lambda(t_1 - t_0)$$

$$P(t_1) = P(t_0)e^{\lambda(t_1-t_0)}$$

The solution is

$$P(t) = P_0 e^{\lambda(t-t_0)}$$

### Prediction Dynamic

$$\dot{f}_{lin}(\mathcal{X}, \theta^{(t)}) = \eta T_{\mathcal{X}\mathcal{X}}^{(0)}(f_{lin}(\mathcal{X}, \theta^{(t)}) - \mathcal{Y})$$

### @ Need Proof

Thus, we can write down the ODE.

$$\dot{f}_{lin}(\mathcal{X}, \theta^{(t)}) = \eta T_{\mathcal{X}\mathcal{X}}^{(0)} f_{lin}(\mathcal{X}, \theta^{(t)})$$

It's trivial that since the term's derivation is itself, thus the solution of the term must contain natural exponential  $e$ . Usually, we guess the solution in the form of  $Ae^{\lambda t}$

$$f_{lin}(\mathcal{X}, \theta^{(t)}) = Ae^{\eta T_{\mathcal{X}\mathcal{X}}^{(0)} t}$$

Consider the training time step  $t = 0$ .

$$f_{lin}(\mathcal{X}, \theta^{(0)}) = Ae^{\eta T_{\mathcal{X}\mathcal{X}}^{(0)} 0} = A$$

Plug into the  $Ae^{\eta T_{\mathcal{X}\mathcal{X}}^{(0)} t}$

$$f_{lin}(\mathcal{X}, \theta^{(t)}) = f_{lin}(\mathcal{X}, \theta^{(0)})e^{\eta T_{\mathcal{X}\mathcal{X}}^{(0)} t}$$

### Weight Dynamic

$$\theta^{(t+1)} - \theta^{(t)} = \eta \nabla_{\theta} f(\mathcal{X}, \theta^{(t)}) \nabla_{f(\mathcal{X}, \theta^{(t)})} \mathcal{L}^{(t)}(\mathcal{X}, \mathcal{Y})$$

We can write down the ODE of the dynamic  $\dot{\theta}(t)$

$$\lim_{\Delta t \rightarrow 0} \frac{\theta^{(t+\Delta t)} - \theta^{(t)}}{\Delta t} = \frac{\partial \theta^{(t)}}{\partial t} = \dot{\theta}(t) = \eta \nabla_{\theta} f(\mathcal{X}, \theta^{(t)}) \nabla_{f(\mathcal{X}, \theta^{(t)})} \mathcal{L}^{(t)}(\mathcal{X}, \mathcal{Y})$$

Now, since the ODE  $\dot{\theta}(t)$  is linear, it has closed form solution

$$\begin{aligned}
\dot{\theta}(t) &= \eta \nabla_{\theta} f(\mathcal{X}, \theta^{(t)}) \nabla_{f(\mathcal{X}, \theta^{(t)})} \sum_{x \in \mathcal{X}, y \in \mathcal{Y}} \frac{1}{2} \|f(x, \theta^{(t)}) - y\|_2^2 \\
&= \eta \nabla_{\theta} f(\mathcal{X}, \theta^{(t)}) \sum_{x \in \mathcal{X}, y \in \mathcal{Y}} \|f(x, \theta^{(t)}) - y\|_2 \\
&= \eta \nabla_{\theta} f(\mathcal{X}, \theta^{(t)}) (\hat{\mathcal{Y}}^{(t)} - \mathcal{Y})
\end{aligned}$$

## Prediction of Trained NN

### Gradient Kernel

Recall the NTK  $T_{\mathcal{X}\mathcal{X}}^{(0)} = \nabla_{\theta} f(\mathcal{X}, \theta^{(0)}) \nabla_{\theta} f(\mathcal{X}, \theta^{(0)})^{\top} \in \mathbb{R}^{N \times N}$

$$T_{\mathcal{X}\mathcal{X}}^{(0)} = \nabla_{\theta} f(\mathcal{X}, \theta^{(0)}) \nabla_{\theta} f(\mathcal{X}, \theta^{(0)})^{\top}$$

Since  $\nabla_{\theta} f(\mathcal{X}, \theta^{(0)}) \in \mathbb{R}^{N \times S}$  is a Jacobian matrix, the matrix multiplication can be written as the following matrix.

$$= \begin{bmatrix} \dot{k}^L(x_1, x_1) & \dot{k}^L(x_1, x_2) & \dots & \dot{k}^L(x_1, x_N) \\ \dot{k}^L(x_2, x_1) & \dot{k}^L(x_2, x_2) & \dots & \dot{k}^L(x_2, x_N) \\ \vdots & \vdots & & \vdots \\ \dot{k}^L(x_N, x_1) & \dot{k}^L(x_N, x_2) & \dots & \dot{k}^L(x_N, x_N) \end{bmatrix}$$

where the function  $\dot{k}^L(x, x') = \nabla_{\theta} f(x, \theta^{(0)})^{\top} \nabla_{\theta} f(x', \theta^{(0)})$  is the kernel function of NTK.

$$\dot{k}^L(x, x') = \nabla_{\theta} f(x, \theta^{(0)})^{\top} \nabla_{\theta} f(x', \theta^{(0)})$$

$$= \nabla_{\theta \leq L} h^L(x)^{\top} \nabla_{\theta \leq L} h^L(x')$$

Since  $\nabla_{\theta^L} h^L(x)^{\top} = [\nabla_{\theta_1^L} h_1^L(x), \nabla_{\theta_{\leq L-1}^L} h^L(x)]$ ,

$$= \nabla_{\theta^L} h^L(x)^{\top} \nabla_{\theta^L} h^L(x') + \nabla_{\theta \leq L-1} h^L(x)^{\top} \nabla_{\theta \leq L-1} h^L(x')$$

Since  $\nabla_{\theta^L} h^L(x)^{\top} = [\nabla_{\theta_1^L} h_1^L(x), \nabla_{\theta_2^L} h_2^L(x), \dots, \nabla_{\theta_{n_L}^L} h_{n_L}^L(x)] \in \mathbb{R}^{1 \times n_L}$ , the inner product of the vector can be written as a form of summation. Note that  $\theta_i^{\leq l}$  represents the  $i$ -th entry of the  $l$ -th layer.

$$\begin{aligned}
&= \sum_{i=1}^{n_L} \nabla_{\theta_i^L} h_i^L(x)^\top \nabla_{\theta_i^L} h_i^L(x') + (\nabla_{\theta^{\leq L-1}} h^{L-1}(x) \nabla_{h^{L-1}(x)} h^L(x))^\top (\nabla_{\theta^{\leq L-1}} h^{L-1}(x') \nabla_{h^{L-1}(x')} h^L(x')) \\
&= \sum_{i=1}^{n_L} \nabla_{\theta_i^L} h_i^L(x)^\top \nabla_{\theta_i^L} h_i^L(x') + \nabla_{h^{L-1}(x)} h^L(x)^\top \nabla_{\theta^{\leq L-1}} h^{L-1}(x)^\top \nabla_{\theta^{\leq L-1}} h^{L-1}(x') \nabla_{h^{L-1}(x')} h^L(x') \\
&= \sum_{i=1}^{n_L} \nabla_{\theta_i^L} h_i^L(x)^\top \nabla_{\theta_i^L} h_i^L(x') + \nabla_{h^{L-1}(x)} h^L(x)^\top \dot{k}^{L-1}(x, x') \nabla_{h^{L-1}(x')} h^L(x') \\
&= n_L \Gamma(K_{xx'}^{L-1}) + \dot{k}^{L-1}(x, x') \dot{\Gamma}(K_{xx'}^{L-1})
\end{aligned}$$


---

### The First Term

Derive the inner product of gradient at  $i$ -th entry.

$$\begin{aligned}
&\nabla_{\theta_i^L} h_i^L(x)^\top \nabla_{\theta_i^L} h_i^L(x') \\
&= \nabla_{\theta_i^L} \left( \frac{\sigma_w}{\sqrt{n_l}} \sum_{j=1}^{n_l} w_{i,j}^L \phi(h_j^{L-1}(x)) + \sigma_b \beta_i^L \right)^\top \nabla_{\theta_i^L} \left( \frac{\sigma_w}{\sqrt{n_l}} \sum_{j=1}^{n_l} w_{i,j}^L \phi(h_j^{L-1}(x')) + \sigma_b \beta_i^L \right) \\
&= \left[ \frac{\sigma_w}{\sqrt{n_l}} \sum_{j=1}^{n_l} \phi(h_j^{L-1}(x)), \sigma_b \right] \left[ \frac{\sigma_w}{\sqrt{n_l}} \sum_{j=1}^{n_l} \phi(h_j^{L-1}(x')), \sigma_b \right]^\top
\end{aligned}$$

Since they are identical entries( $i$ -th), we can combine them into single term.

$$= \frac{\sigma_w^2}{n_l} \left( \sum_{j=1}^{n_l} \phi(h_j^{L-1}(x)) \phi(h_j^{L-1}(x')) + \sum_{j \neq j'}^{n_l} \phi(h_j^{L-1}(x)) \phi(h_{j'}^{L-1}(x')) \right) + \sigma_b^2$$

By central limit theorem(CLT), if  $n_l \rightarrow \infty$ , it will converge.

$$= \sigma_w^2 E_{h_j^{L-1}(x), h_j^{L-1}(x') \sim \mathcal{N}(0, K_{xx'}^{L-1})} [\phi(h_j^{L-1}(x)) \phi(h_j^{L-1}(x'))] + \sigma_b^2$$

$$\begin{aligned}
&= \sigma_w^2 E_{u,v \sim \mathcal{N}(0, K_{xx'}^{l-1})} [\phi(u)\phi(v)] + \sigma_b^2 \\
&= \Gamma(K_{xx'}^{l-1})
\end{aligned}$$

$$K_{xx'}^{l-1} = \begin{bmatrix} k^{l-1}(x, x) & k^{l-1}(x, x') \\ k^{l-1}(x', x) & k^{l-1}(x', x') \end{bmatrix}$$

### The Second Term

$$\begin{aligned}
&\nabla_{h^{l-1}(x)} h^l(x)^\top \dot{k}^{l-1}(x, x') \nabla_{h^{l-1}(x')} h^l(x') \\
&= \dot{k}^{l-1}(x, x') \nabla_{h^{l-1}(x)} h^l(x)^\top \nabla_{h^{l-1}(x')} h^l(x') \\
&= \dot{k}^{l-1}(x, x') \nabla_{h^{l-1}(x)} \left( \frac{\sigma_w}{\sqrt{n_l}} \sum_{j=1}^{n_l} w_{i,j}^l \phi(h_j^{l-1}(x)) + \sigma_b \beta_i^l \right)^\top \nabla_{h^{l-1}(x')} \left( \frac{\sigma_w}{\sqrt{n_l}} \sum_{j=1}^{n_l} w_{i,j}^l \phi(h_j^{l-1}(x')) + \sigma_b \beta_i^l \right)
\end{aligned}$$

Due to  $\nabla_{h^{l-1}(x)} h^l(x)^\top = [\nabla_{h_1^{l-1}(x)} h^l(x), \nabla_{h_2^{l-1}(x)} h^l(x), \dots, \nabla_{h_{n_l}^{l-1}(x)} h^l(x)] \in \mathbb{R}^{1 \times n_l}$ ,  $\nabla_{h^{l-1}(x)} h^l(x)^\top \nabla_{h^{l-1}(x')} h^l(x')$  is an inner product.

$$\begin{aligned}
&= \dot{k}^{l-1}(x, x') \frac{\sigma_w^2}{n_l} \left( \sum_{j=1}^{n_l} w_{i,j}^l \nabla_{h^{l-1}(x)} \phi(h_j^{l-1}(x)) w_{i,j}^l \nabla_{h^{l-1}(x')} \phi(h_j^{l-1}(x')) \right) \\
&= \dot{k}^{l-1}(x, x') \frac{\sigma_w^2}{n_l} \left( \sum_{j=1}^{n_l} w_{i,j}^l \nabla_{h^{l-1}(x)} \phi(h_j^{l-1}(x)) w_{i,j}^l \nabla_{h^{l-1}(x')} \phi(h_j^{l-1}(x')) \right)
\end{aligned}$$

With central limit theorem (CLT), we can take an expectation over the equation. Since  $w_{i,j}^l$  is independent to  $h_j^{l-1}(x)$  and  $h_j^{l-1}(x')$ , we can separate the expectation  $E[w_{i,j}^l w_{i,j}^l]$ .

$$\begin{aligned}
&= \dot{k}^{l-1}(x, x') \sigma_w^2 E_{w_{i,j}^l \sim \mathcal{N}(0,1)} [w_{i,j}^l w_{i,j}^l] E_{h_j^{l-1}(x), h_j^{l-1}(x') \sim \mathcal{N}(0, K_{xx'}^{l-1})} [\nabla_{h^{l-1}(x)} \phi(h_j^{l-1}(x)) \nabla_{h^{l-1}(x')} \phi(h_j^{l-1}(x')))] \\
&= \dot{k}^{l-1}(x, x') \sigma_w^2 E_{h_j^{l-1}(x), h_j^{l-1}(x') \sim \mathcal{N}(0, K_{xx'}^{l-1})} [\nabla_{h^{l-1}(x)} \phi(h_j^{l-1}(x)) \nabla_{h^{l-1}(x')} \phi(h_j^{l-1}(x')))]
\end{aligned}$$

$$= \dot{k}^{l-1}(x, x') \sigma_w^2 E_{u, v \sim \mathcal{N}(0, K_{xx'}^{l-1})} [\nabla_u \phi(u) \nabla_v \phi(v)]$$

$$= \dot{k}^{l-1}(x, x') \dot{\Gamma}(K_{xx'}^{l-1})$$

$$K_{xx'}^{l-1} = \begin{bmatrix} k^{l-1}(x, x) & k^{l-1}(x, x') \\ k^{l-1}(x', x) & k^{l-1}(x', x') \end{bmatrix}$$