

# Neural Tangent Kernel

Setup:  $\mathcal{D} := \{x_i, y_i\}, i = 1, \dots, n$ ,  $y = \underbrace{f_\theta(x)}_{\text{MLP}} + \varepsilon$

Network outputs:  $f(x; \theta(t))$

Training via gradient go

$$\mathcal{L}(\theta) = \frac{1}{2} \sum_{i=1}^n [f(x_i; \theta(t)) - y_i]^2, \quad \frac{\partial \mathcal{L}}{\partial \theta} \rightarrow \underbrace{\theta^{n+1} = \theta^n - \eta \nabla_\theta \mathcal{L}(\theta^n)}$$

Gradient flow:  $\boxed{\frac{d\theta}{dt} = -\nabla_\theta \mathcal{L}(\theta)}$   $\xrightarrow{\text{forward Euler discretization}}$

• Derive the evolution of  $f(x; \theta(t))$ :

$$\frac{df(x; \theta(t))}{dt} = \frac{df(x; \theta(t))^T}{d\theta} \cdot \frac{d\theta}{dt}$$

$$= - \frac{df(x; \theta(t))^T}{d\theta} \underbrace{\nabla_\theta \mathcal{L}(\theta)}$$

$$= - \frac{df(x; \theta(t))^T}{d\theta} \sum_{i=1}^n [f(x_i; \theta(t)) - y_i] \frac{df(x_i; \theta(t))}{d\theta}$$

$$\Rightarrow \boxed{\frac{df(X; \theta(t))}{dt} = -\mathbb{K}_t(X, X) [f(X; \theta(t)) - y]}$$

where  $X := \begin{bmatrix} x_{11} & \dots & x_{1d} \\ \vdots & & \vdots \\ x_{n1} & \dots & x_{nd} \end{bmatrix}$   $y = \begin{bmatrix} y_1 \\ \vdots \\ y_n \end{bmatrix}$   $\nwarrow \textcircled{1}$

$$n \times d \quad \begin{bmatrix} x_{11} & \dots & x_{1d} \\ \vdots & & \vdots \\ x_{n1} & \dots & x_{nd} \end{bmatrix}, \quad \begin{bmatrix} y_1 \\ \vdots \\ y_n \end{bmatrix}$$

$$\underbrace{\mathbb{K}_t(x, x)_{ij} = \left\langle \frac{df(x_i; \theta(t))}{d\theta}, \frac{df(x_j; \theta(t))}{d\theta} \right\rangle}_{n \times n}$$

Neural Tangent Kernel

Remark #1: At the infinite width and infinitesimally small learning rate, the NTK  $\mathbb{K}_t$  converges to a deterministic kernel that remains constant during training.  $\mathbb{K}_t = \mathbb{K}(0) = \bar{\mathbb{K}}$

Remark #2:

$$\begin{aligned} f(x^*; \theta(t)) &\approx \mathbb{K}_t(x^*, x)^{-\frac{1}{2}} \mathbb{K}_t(x, x)^{-\frac{1}{2}} (\mathbb{I} - e^{-\frac{1}{2} \mathbb{K}_t t}) y \\ &\approx \underbrace{K^*(x^*, x) K^*(x, x)}_{\text{kernel regression}} y \end{aligned}$$

Remark #3:

$$\underbrace{\mathbb{K}^*}_{n \times n} = \underbrace{Q^T \Lambda Q}_{\text{SVD}} \quad \begin{cases} Q: \text{orthogonal} \\ \Lambda: \text{diagonal} \end{cases}$$

$$Q^T \underbrace{(f(x; \theta(t)) - y)}_{\text{training error}} = -e^{-\Lambda t} Q^T y$$

$$\Rightarrow f(x; \theta(t)) - y = \sum_{i=1}^n (e^{-\lambda_i t} q_i^T y) q_i \Rightarrow \text{Spectral Bias}$$

