# **A Neural Tangent Kernel Perspective of Infinite Tree Ensembles**

<sup>1</sup>National Institute of Informatics, <sup>2</sup>The Graduate University for Advanced Studies, SOKENDAI

The 10th International Conference on Learning Representations (ICLR 2022), April 25-29, 2022

### **Our Contribution**

- Theoretical support for empirical techniques

## Soft Tree Ensemble •



### Neural Tangent Kernel •

- Training behavior can be analyzed using kernel methods
- Many successes for infinitely wide neural networks

# Ryuichi Kanoh<sup>1,2</sup>, Mahito Sugiyama<sup>1,2</sup>

0.5

 $w_2$ 

 $(\boldsymbol{w}_3)$ 

#### • The TNTK induced by infinite trees converges to a deterministic kernel

 $\lim_{M \to \infty} \widehat{\Theta}_0^{(d)}(\boldsymbol{x}_i, \boldsymbol{x}_j) = \underbrace{2^d d \ \Sigma(\boldsymbol{x}_i, \boldsymbol{x}_j) (\mathcal{T}(\boldsymbol{x}_i, \boldsymbol{x}_j))^{d-1} \dot{\mathcal{T}}(\boldsymbol{x}_i, \boldsymbol{x}_j)}_{(d-1)} + \underbrace{2^d d \ \Sigma(\boldsymbol{x}_i, \boldsymbol{x}_j) (\mathcal{T}(\boldsymbol{x}_i, \boldsymbol{x}_j))^{d-1} \dot{\mathcal{T}}(\boldsymbol{x}_i, \boldsymbol{x}_j)}_{(d-1)} + \underbrace{2^d d \ \Sigma(\boldsymbol{x}_i, \boldsymbol{x}_j) (\mathcal{T}(\boldsymbol{x}_i, \boldsymbol{x}_j))^{d-1} \dot{\mathcal{T}}(\boldsymbol{x}_i, \boldsymbol{x}_j)}_{(d-1)} + \underbrace{2^d d \ \Sigma(\boldsymbol{x}_i, \boldsymbol{x}_j) (\mathcal{T}(\boldsymbol{x}_i, \boldsymbol{x}_j))^{d-1} \dot{\mathcal{T}}(\boldsymbol{x}_i, \boldsymbol{x}_j)}_{(d-1)} + \underbrace{2^d d \ \Sigma(\boldsymbol{x}_i, \boldsymbol{x}_j) (\mathcal{T}(\boldsymbol{x}_i, \boldsymbol{x}_j))^{d-1} \dot{\mathcal{T}}(\boldsymbol{x}_i, \boldsymbol{x}_j)}_{(d-1)} + \underbrace{2^d d \ \Sigma(\boldsymbol{x}_i, \boldsymbol{x}_j) (\mathcal{T}(\boldsymbol{x}_i, \boldsymbol{x}_j))^{d-1} \dot{\mathcal{T}}(\boldsymbol{x}_i, \boldsymbol{x}_j)}_{(d-1)} + \underbrace{2^d d \ \Sigma(\boldsymbol{x}_i, \boldsymbol{x}_j) (\mathcal{T}(\boldsymbol{x}_i, \boldsymbol{x}_j))^{d-1} \dot{\mathcal{T}}(\boldsymbol{x}_i, \boldsymbol{x}_j)}_{(d-1)} + \underbrace{2^d d \ \Sigma(\boldsymbol{x}_i, \boldsymbol{x}_j) (\mathcal{T}(\boldsymbol{x}_i, \boldsymbol{x}_j))^{d-1} \dot{\mathcal{T}}(\boldsymbol{x}_i, \boldsymbol{x}_j)}_{(d-1)} + \underbrace{2^d d \ \Sigma(\boldsymbol{x}_i, \boldsymbol{x}_j) (\mathcal{T}(\boldsymbol{x}_i, \boldsymbol{x}_j))^{d-1} \dot{\mathcal{T}}(\boldsymbol{x}_i, \boldsymbol{x}_j)}_{(d-1)} + \underbrace{2^d d \ \Sigma(\boldsymbol{x}_i, \boldsymbol{x}_j) (\mathcal{T}(\boldsymbol{x}_i, \boldsymbol{x}_j))^{d-1} \dot{\mathcal{T}}(\boldsymbol{x}_i, \boldsymbol{x}_j)}_{(d-1)} + \underbrace{2^d d \ \Sigma(\boldsymbol{x}_i, \boldsymbol{x}_j) (\mathcal{T}(\boldsymbol{x}_i, \boldsymbol{x}_j))^{d-1} \dot{\mathcal{T}}(\boldsymbol{x}_i, \boldsymbol{x}_j)}_{(d-1)} + \underbrace{2^d d \ \Sigma(\boldsymbol{x}_i, \boldsymbol{x}_j) (\mathcal{T}(\boldsymbol{x}_i, \boldsymbol{x}_j))^{d-1} \dot{\mathcal{T}}(\boldsymbol{x}_i, \boldsymbol{x}_j)}_{(d-1)} + \underbrace{2^d d \ \Sigma(\boldsymbol{x}_i, \boldsymbol{x}_j) (\mathcal{T}(\boldsymbol{x}_i, \boldsymbol{x}_j))^{d-1} \dot{\mathcal{T}}(\boldsymbol{x}_i, \boldsymbol{x}_j)}_{(d-1)} + \underbrace{2^d d \ \Sigma(\boldsymbol{x}_i, \boldsymbol{x}_j) (\mathcal{T}(\boldsymbol{x}_i, \boldsymbol{x}_j))^{d-1} \dot{\mathcal{T}}(\boldsymbol{x}_i, \boldsymbol{x}_j)}_{(d-1)} + \underbrace{2^d d \ \Sigma(\boldsymbol{x}_i, \boldsymbol{x}_j) (\mathcal{T}(\boldsymbol{x}_i, \boldsymbol{x}_j))^{d-1} \dot{\mathcal{T}}(\boldsymbol{x}_i, \boldsymbol{x}_j)}_{(d-1)} + \underbrace{2^d d \ \Sigma(\boldsymbol{x}_i, \boldsymbol{x}_j) (\mathcal{T}(\boldsymbol{x}_i, \boldsymbol{x}_j))^{d-1} \dot{\mathcal{T}}(\boldsymbol{x}_i, \boldsymbol{x}_j)}_{(d-1)} + \underbrace{2^d d \ \Sigma(\boldsymbol{x}_i, \boldsymbol{x}_j) (\mathcal{T}(\boldsymbol{x}_i, \boldsymbol{x}_j))^{d-1} \dot{\mathcal{T}}(\boldsymbol{x}_i, \boldsymbol{x}_j)}_{(d-1)} + \underbrace{2^d d \ \Sigma(\boldsymbol{x}_i, \boldsymbol{x}_j) (\mathcal{T}(\boldsymbol{x}_i, \boldsymbol{x}_j))^{d-1} \dot{\mathcal{T}}(\boldsymbol{x}_j, \boldsymbol{x}_j)}_{(d-1)} + \underbrace{2^d d \ \Sigma(\boldsymbol{x}_j, \boldsymbol{x}_j) (\mathcal{T}(\boldsymbol{x}_j, \boldsymbol{x}_j))^{d-1} \dot{\mathcal{T}}(\boldsymbol{x}_j, \boldsymbol{x}_j)}_{(d-1)} + \underbrace{2^d d \ \Sigma(\boldsymbol{x}_j, \boldsymbol{x}_j) (\mathcal{T}(\boldsymbol{x}_j, \boldsymbol{x}_j))^{d-1} \dot{\mathcal{T}}(\boldsymbol{x}_j, \boldsymbol{x}_j)}_{(d-1)} + \underbrace{2^d d \ \Sigma(\boldsymbol{x}_j, \boldsymbol{x}_j) (\mathcal{T}(\boldsymbol{x}_j, \boldsymbol{x}_j))^{d-1} \dot{\mathcal{T}}(\boldsymbol{x}_j, \boldsymbol{x}_j)}_{(d-1)} + \underbrace{2^d d \ \Sigma(\boldsymbol{x}_j, \boldsymbol{x}_j) (\mathcal{T}(\boldsymbol{x}_j, \boldsymbol{x}_j))^{d-1} \dot{\mathcal{T}}(\boldsymbol{x}_j$  $(2\mathcal{T}(oldsymbol{x}_i,oldsymbol{x}_j))^a$ contribution from inner nodes contribution from leave

#### As the #trees increases, the change from the initial value becomes smaller - Training behavior of infinite trees is analytically tractable





Inner product of the inputs