

The Training of Artificial Neural Networks (ANN) involves a priori highly non-convex optimization, yet in practice one can obtain impressive convergence and generalization.

Question How to understand mathematically the dynamics of the training of an ANN as it becomes large?

Neural Tangent Kernel: Convergence and Generalization in Neural Networks

Arthur Jacot-Guillarmod* Franck Gabriel* Clément Hongler* *EPFL

Thanks to a new object the Neural Tangent Kernel, we can precisely describe the evolution, convergence, and generalization of ANNs of large width

Theorem In the infinite-width limit, the ANN function follows a kernel gradient descent with respect to the limiting NTK

1. ANN Training: Architecture and Optimization

- Fully connected neural net with $L+1$ layers of widths n_0, \dots, n_L and nonlinearity $\sigma: \mathbb{R} \rightarrow \mathbb{R}$

$$f_\theta: \mathbb{R}^{n_0} \rightarrow \mathbb{R}^{n_1} \rightarrow \mathbb{R}^{n_2} \rightarrow \dots \rightarrow \mathbb{R}^{n_L} \rightarrow \mathbb{R}^{n_L}$$

$$x \mapsto \sigma(W_{0,1}^{(1)} x) \mapsto \sigma(W_{1,2}^{(2)} \sigma(W_{0,1}^{(1)} x)) \mapsto \dots \mapsto \sigma(W_{L-1,L}^{(L)} \sigma(W_{L-2,L-1}^{(L-1)} \dots \sigma(W_{0,1}^{(1)} x) \dots)) = f_\theta(x)$$

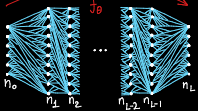
Network passes $(x, z) \in \mathbb{R}^{n_0} \times \mathbb{R}^{n_L}$ through the network to produce $f_\theta(x)$. The weights $W_{l,l+1}^{(l+1)}$ are initialized from Gaussian distributions.

- Gaussian Initialization: $\theta_p \sim \mathcal{N}(0, 1)$

Training Inputs $\leftarrow \{x_1, \dots, x_N\}$

- Training Set: $(x_i, z_i) \in \mathbb{R}^{n_0} \times \mathbb{R}^{n_L}$

Training Outputs $\leftarrow \{z_1, \dots, z_N\}$



- Training Loss: $C(\theta) = \frac{1}{N} \sum_{i=1}^N c(f_\theta(x_i), z_i)$ Goal: Minimize $C(\theta)$ with Gradient Descent

- Gradient Descent Step: $\theta \mapsto \theta - \eta \nabla C(\theta) \rightsquigarrow$ Flow: $\partial_t \theta(t) = -\nabla C(\theta)$

- Inference Flows: $f_\theta(x) \xrightarrow{t \rightarrow \infty} ?$ for $x \in \left\{ \begin{array}{l} \text{training set} \\ \text{test set} \end{array} \right.$

2. ANN Training: Functional Formulation

- Neural Realization Function: $F^{(L)}: \mathbb{R}^p \rightarrow \mathbb{R}^{n_L}$ $\theta \mapsto (f_\theta: \mathbb{R}^{n_0} \rightarrow \mathbb{R}^{n_L})$

- Neural Tangent Kernel: $\Theta_{B,i,j}(x,y) = \sum_{l=1}^L \partial_{\theta_{B,i}} f_{\theta_{B,i}}(x) \partial_{\theta_{B,j}} f_{\theta_{B,j}}(y)$

- Initialization: $F^{(L)}(\theta) \in \mathcal{F}$

- Non-Gaussian ∇C Level sets C level sets

- Functional Cost $\mathcal{C}: \mathcal{F} \rightarrow \mathbb{R}$ $f \mapsto \frac{1}{N} \sum_{i=1}^N c(f(x_i), z_i)$

- Training Loss: $C = \mathcal{C} \circ F^{(L)}$ Chain Rule

- Gradient Descent Step: $f_\theta \mapsto f_\theta - \eta \nabla C(\theta) \simeq f_\theta - \eta \sum_{l=1}^L \partial_{\theta_{B,i}} f_{\theta_{B,i}}(x) \partial_{\theta_{B,j}} f_{\theta_{B,j}}(y)$

- Flow: $\partial_t f_\theta(x) = -\nabla_{\theta(x)} \mathcal{C}|_{f_\theta}$

- Kernel Gradient: $\nabla_{\theta} \mathcal{C}|_{f_\theta} = \left(\frac{1}{N} \sum_{i=1}^N \sum_{j=1}^N c(x_i, z_i) \partial_{\theta_{B,i}} f_{\theta_{B,i}}(x) \partial_{\theta_{B,j}} f_{\theta_{B,j}}(y) \right)_{x,y \in \mathbb{R}^{n_0}}$

3. Infinite-Width Limit: $n_1, \dots, n_{L-1} \rightarrow \infty$

- $(f_{\theta_{B,i}}(x))_{i=1, \dots, N} \xrightarrow[N \rightarrow \infty]{x \in \mathbb{R}^{n_0}} \mathbb{Z}$ centered Gaussian Process with $C_\theta(Z_i(x), Z_j(y)) = \sum_{l=1}^L \Theta_{B,i,j}^{(l)}(x,y)$ \leftarrow Neural, 1996 [2]

- $\Theta_{\theta_{B,i}} \xrightarrow[N \rightarrow \infty]{n_1, \dots, n_{L-1} \rightarrow \infty} \Theta \otimes \text{Id}_{n_L}$ \leftarrow Explicitly Computable

- $\Theta \otimes \text{Id}_{n_L} \xrightarrow[N \rightarrow \infty]{n_1, \dots, n_{L-1} \rightarrow \infty} \Theta \otimes \text{Id}_{n_L}$ \leftarrow At Initialization

- $\Theta \otimes \text{Id}_{n_L} \xrightarrow[N \rightarrow \infty]{n_1, \dots, n_{L-1} \rightarrow \infty} 0 \rightsquigarrow \partial_t f_\theta(x) = -\nabla_{\theta(x)} \mathcal{C}|_{f_\theta}$ \leftarrow During Training

4. Consequences as $n_1, \dots, n_{L-1} \rightarrow \infty$

Optimal Emergent Learning

- Individually, the weights evolve less and less: $|\theta_i(t) - \theta_i| \xrightarrow[t \rightarrow \infty]{} 0$

- ... yet learning occurs at each layer

- The influence of (x_i, z_i) on the prediction $f_\theta(x)$ for a test point x can be understood using the NTK

- Convergence to global min of \mathcal{C} is guaranteed if Θ_{∞} is positive-definite

- Turned for $n_1, \dots, n_{L-1} \rightarrow \infty$

Least-Squares Regression

- $f_\theta(t)$ is Gaussian for any $t \geq 0$

- $\mathbb{E}[f_\theta(t)] \xrightarrow[t \rightarrow \infty]{} \text{Bayesian Max A Posteriori}$

- w.r.t. $\mathcal{N}(0, \Theta \otimes \text{Id})$ prior

- = Kernel Ridge Regression

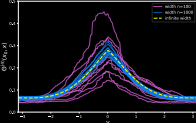
- with $\lambda \rightarrow 0$ regularization

- Learning governed by NTK PCA:

- conv. rate along $v_i \propto \lambda_i$

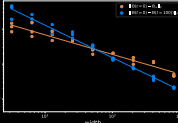
- PCA direction PCA eigenvalue

Neural Tangent Kernel at Initialization (L=4)



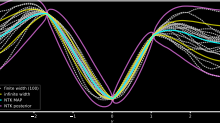
The NTK at initialization on the unit circle for 2 different widths (10 random initializations each), along with its infinite-width limit.

Convergence of the NTK in the Infinite-Width Limit



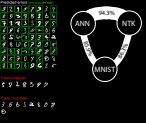
For a range of widths, the distances of the NTK at initialization to its infinite-width limit (orange) and to its limit at the end of training (blue). The distances are computed in Frobenius norm on a fixed batch.

Distribution of the Network Function at the End of Learning (L=3)



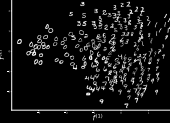
Distribution of f_θ at the end of training for an ANN of width 100, computed to the common mean (blue) and 80% confidence intervals of the infinite-width Gaussian distribution (green) and the posterior for NTK prior (pink).

Error Prediction on MNIST using the NTK



The NTK can be used to predict the misclassification of an ANN of width 500 with L=4 will make on MNIST. It yields a very good prediction on which weights will be chosen by the ANN.

Kernel PCA of MNIST with respect to the NTK



The 2nd and 3rd kernel principal components of MNIST with respect to the NTK (the first component is essentially constant with L=4). The three first eigenvalues are 41.03, 1.88, and 1.46.

References

- [1] A. Jacot, F. Gabriel, C. Hongler, *Neural Tangent Kernel: Convergence and Generalization in Neural Networks*, NeurIPS 2018, arXiv:1806.07572.
- [2] R. H. Neal, *Bayesian Learning for Neural Networks*, 1996.
- Acknowledgments: ERIC CG CRITICAL (Gabriel), ERIC SG CONSTAMIS, Blomqvist Family Foundation (Hongler), Labex Team Foundation.