
Learnability in the Context of Neural Tangent Kernels

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 Understanding the prioritization of certain samples over others during neural net-
2 work training is a fundamental challenge in deep learning. This prioritization is
3 intrinsically linked to the network’s inductive bias—the inherent assumptions that
4 enable generalization from training data to unseen data. In this study, we investigate
5 the role of the diagonal elements of the Neural Tangent Kernel (NTK), $k(x, x)$, in
6 determining sample learnability. Through theoretical analysis, we demonstrate that
7 higher values of $k(x, x)$ correlate with faster convergence rates of individual sam-
8 ple errors during training, indicating that such samples are learned more rapidly and
9 accurately. Conversely, lower $k(x, x)$ values are associated with slower learning
10 dynamics, classifying these samples as harder to learn. Empirical evaluations con-
11 ducted on standard datasets using convolutional neural networks (CNNs), validate
12 our theoretical predictions. We observe that samples with higher $k(x, x)$ values
13 consistently achieve higher accuracy in fewer training epochs compared to those
14 with lower values. Visual inspections further reveal that high- $k(x, x)$ samples are
15 typically clear and prototypical, whereas low- $k(x, x)$ samples often exhibit noise
16 or atypical characteristics.

17 1 Introduction

18 Understanding why neural networks prioritize learning certain samples over others is a fundamental
19 question in deep learning. This prioritization is closely tied to the network’s inductive bias—the
20 set of assumptions a model makes to generalize from training data to unseen data. The Neural
21 Tangent Kernel (NTK) has emerged as a powerful theoretical tool to analyze the training dynamics of
22 overparameterized neural networks [1, 2]. In this work, we focus on the diagonal elements of the
NTK, $k(x, x)$, and their relationship with sample learnability. Our empirical observations show that

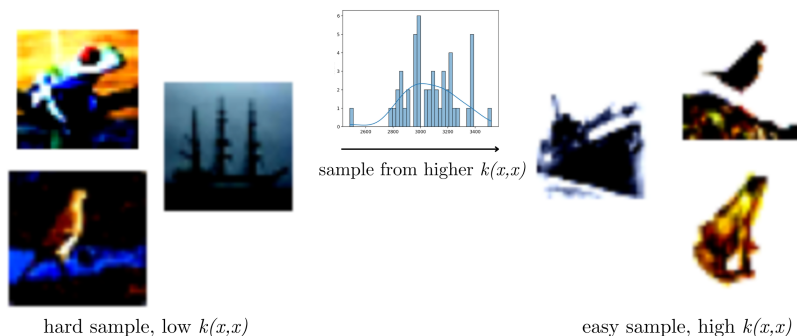


Figure 1: *Hard (left)* and *easy (right)* to learn CIFAR-10 samples for a CNN, picked by our method.

24 samples with higher $k(x, x)$ values tend to be learned faster and more accurately during training.
 25 These samples often correspond to "easy" examples with less complexity or ambiguity. Conversely,
 26 samples with lower $k(x, x)$ values are learned more slowly and are generally harder examples. We
 27 aim to theoretically justify this phenomenon by deriving convergence rates based on $k(x, x)$ and
 28 explaining how $k(x, x)$ affects the optimization dynamics of neural networks.

29 2 Theoretical Analysis

30 The NTK arises in the study of infinitely wide neural networks trained using gradient descent. For a
 31 neural network $f(\mathbf{x}; \boldsymbol{\theta})$ with parameters $\boldsymbol{\theta}$, the NTK is defined as [2][3]:

$$K(\mathbf{x}, \mathbf{x}') = \nabla_{\boldsymbol{\theta}} f(\mathbf{x}; \boldsymbol{\theta})^\top \nabla_{\boldsymbol{\theta}} f(\mathbf{x}'; \boldsymbol{\theta}). \quad (1)$$

32 The diagonal elements $k(x, x)$ represent the inner product of the gradient of the network's output
 33 with respect to its parameters at input \mathbf{x} :

$$k(x, x) = \|\nabla_{\boldsymbol{\theta}} f(\mathbf{x}; \boldsymbol{\theta})\|^2. \quad (2)$$

34 This quantity measures the sensitivity of the output at \mathbf{x} to changes in the parameters and thus can be
 35 interpreted as the influence of sample \mathbf{x} on the training dynamics. High $k(x, x)$ indicates that small
 36 parameter updates can significantly affect the output for \mathbf{x} , potentially leading to faster learning for
 37 that sample.

38 Consider training a neural network with mean squared error (MSE) loss on dataset $\{(\mathbf{x}_i, y_i)\}_{i=1}^N$. The
 39 network output at time t is $f_t(\mathbf{x})$. The dynamics of $f_t(\mathbf{x})$ during gradient descent are governed by:

$$\frac{df_t(\mathbf{x})}{dt} = -\eta \sum_{i=1}^N K(\mathbf{x}, \mathbf{x}_i) (f_t(\mathbf{x}_i) - y_i), \quad (3)$$

40 where η is the learning rate. This differential equation characterizes how the network's output evolves
 41 over time under gradient descent. In the infinite-width limit, $K(\mathbf{x}, \mathbf{x}')$ remains constant[2][1], and
 42 the solution can be expressed as:

$$f_t = f_0 - (I - e^{-\eta K t})(f_0 - y), \quad (4)$$

43 where f_0 is the initial network output, y is the vector of target values, and K is the NTK matrix.

44 2.1 Convergence Rate for Individual Samples

45 To understand how the error for each individual sample evolves during training, we derive the
 46 Ordinary Differential Equation (ODE) governing the error dynamics within the Neural Tangent
 47 Kernel (NTK) framework. This derivation is foundational for analyzing the convergence rates and
 48 establishing the relationship between the NTK's diagonal elements and sample learnability. Consider
 49 training a neural network $f(\mathbf{x}; \boldsymbol{\theta})$ with parameters $\boldsymbol{\theta}$ on a dataset $\{(\mathbf{x}_i, y_i)\}_{i=1}^N$ using gradient descent
 50 to minimize the mean squared error (MSE) loss:

$$\mathcal{L} = \frac{1}{2} \sum_{i=1}^N (f(\mathbf{x}_i; \boldsymbol{\theta}) - y_i)^2. \quad (5)$$

51 The continuous-time gradient descent (i.e., gradient flow) update rule is given by:

$$\frac{d\boldsymbol{\theta}(t)}{dt} = -\eta \nabla_{\boldsymbol{\theta}} \mathcal{L}, \quad (6)$$

52 where η is the learning rate. The time derivative of the network's output for sample \mathbf{x} is expressed as:

$$\frac{df(\mathbf{x}; \boldsymbol{\theta}(t))}{dt} = \nabla_{\boldsymbol{\theta}} f(\mathbf{x}; \boldsymbol{\theta}(t))^\top \frac{d\boldsymbol{\theta}(t)}{dt}. \quad (7)$$

53 Substituting the gradient descent update rule into this equation yields,

$$\frac{df(\mathbf{x}; \boldsymbol{\theta}(t))}{dt} = -\eta \sum_{i=1}^N (f(\mathbf{x}_i; \boldsymbol{\theta}(t)) - y_i) \nabla_{\boldsymbol{\theta}} f(\mathbf{x}; \boldsymbol{\theta}(t))^\top \nabla_{\boldsymbol{\theta}} f(\mathbf{x}_i; \boldsymbol{\theta}(t)) \quad (8)$$

$$= -\eta \sum_{i=1}^N K(\mathbf{x}, \mathbf{x}_i) (f(\mathbf{x}_i; \boldsymbol{\theta}(t)) - y_i), \quad (9)$$

54 where the Neural Tangent Kernel (NTK) is defined as:

$$K(\mathbf{x}, \mathbf{x}') = \nabla_{\theta} f(\mathbf{x}; \theta)^{\top} \nabla_{\theta} f(\mathbf{x}'; \theta). \quad (10)$$

55 On defining the following vector quantities:

$$\mathbf{f}(t) = \begin{bmatrix} f(\mathbf{x}_1; \theta(t)) \\ \vdots \\ f(\mathbf{x}_N; \theta(t)) \end{bmatrix}, \quad \mathbf{y} = \begin{bmatrix} y_1 \\ \vdots \\ y_N \end{bmatrix}, \quad \mathbf{e}(t) = \mathbf{f}(t) - \mathbf{y}, \quad (11)$$

56 We may rewrite the evolution equation more compactly, as:

$$\frac{d\mathbf{f}(t)}{dt} = -\eta K \mathbf{e}(t), \quad (12)$$

57 where K is the NTK matrix with elements $K_{ij} = K(\mathbf{x}_i, \mathbf{x}_j)$. Substituting $\mathbf{e}(t)$ into the evolution
58 equation, we obtain a simple ordinary differential equation governing the error dynamics:

$$\frac{d\mathbf{e}(t)}{dt} = -\eta K \mathbf{e}(t). \quad (13)$$

59 Assuming that the NTK matrix K is *diagonal* or *approximately diagonal*—a reasonable assumption
60 in the infinite-width limit where off-diagonal elements become negligible [4, 2, 1]—the system
61 decouples into independent ODEs for each sample:

$$\frac{de_t(\mathbf{x}_i)}{dt} = -\eta k(x_i, x_i) e_t(\mathbf{x}_i), \implies e_t(\mathbf{x}_i) = e_0(\mathbf{x}_i) e^{-\eta k(x_i, x_i) t}, \quad (14)$$

62 where $e_0(\mathbf{x}_i) = e_{t=0}(\mathbf{x}_i)$ is the initial error for sample \mathbf{x}_i . The solution $e_t(\mathbf{x}_i) = e_0(\mathbf{x}_i) e^{-\eta k(x_i, x_i) t}$
63 indicates that the error for each sample decays *exponentially* over time. The *rate of convergence* is
64 governed by the product $\eta k(x_i, x_i)$.

65 2.2 Bounds on Convergence Rates

66 We can formalize the convergence rates by deriving bounds on the error $e_t(\mathbf{x}_i)$:

$$e^{-\eta k_{\max} t} e_0(\mathbf{x}_i) \leq e_t(\mathbf{x}_i) \leq e^{-\eta k_{\min} t} e_0(\mathbf{x}_i), \quad (15)$$

67 These bounds indicate that the error decay rate is bounded by the minimum and maximum diagonal
68 elements of the NTK. However, these bounds are loose as they do not capture the individual variability
69 of each $k(x_i, x_i)$. A more precise estimate considers each $k(x_i, x_i)$ individually:

$$e_t(\mathbf{x}_i) = e^{-\eta k(x_i, x_i) t} e_0(\mathbf{x}_i). \quad (16)$$

70 This relationship implies that for **higher** $k(x_i, x_i)$, **the error** $e_t(\mathbf{x}_i)$ **decreases faster**, leading to
71 quicker learning for sample \mathbf{x}_i . **Lower** $k(x_i, x_i)$ **leads to the error** $e_t(\mathbf{x}_i)$ **decreasing slower**,
72 indicating that sample \mathbf{x}_i is learned more gradually. Focusing on the error $e_t(\mathbf{x}) = f_t(\mathbf{x}) - y(\mathbf{x})$, we
73 analyze the convergence rate for each sample \mathbf{x} . Assuming K is positive definite, the error dynamics
74 for the i -th sample are:

$$e_t(\mathbf{x}_i) = e^{-\eta K t} e_0(\mathbf{x}_i). \quad (17)$$

75 If K is diagonal or approximately diagonal, the convergence rate for each sample simplifies to:

$$e_t(\mathbf{x}_i) \approx e^{-\eta k(x_i, x_i) t} e_0(\mathbf{x}_i). \quad (18)$$

76 This shows that the error for sample \mathbf{x}_i decreases exponentially with a rate proportional to its $k(x_i, x_i)$.
77 Thus, samples with higher $k(x_i, x_i)$ converge faster during training.

78 3 Empirical Evaluation

79 We conducted experiments to validate the theoretical findings using standard deep learning datasets.
80 Specifically, we utilized a finite-width NTK [5][2][6] CNN trained on binary and multi-class classifi-
81 cation and probed learnability scores (i.e, the diagonal values on the NTK matrix; for the multi-class
82 case, we took the mean of the values across classes). Our findings on both MNIST and CIFAR-10
83 datasets were equivalent. We also looked at how samples with high and low learnability behave
84 during the training process, for which the finite-width approximation was necessary.

85 We divided the dataset into distinct groups based on their $k(x, x)$ values to analyze learning dynamics:

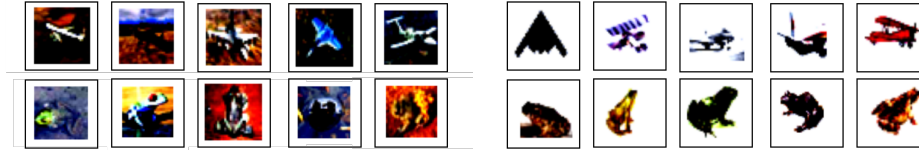


Figure 2: **Left:** Images sampled from CIFAR-10 with *low* learnability. **Right:** Images sampled from CIFAR-10 with *high* learnability. Sampling on the basis of a 3-layer CNN trained on binary classification to distinguish between frogs and planes.

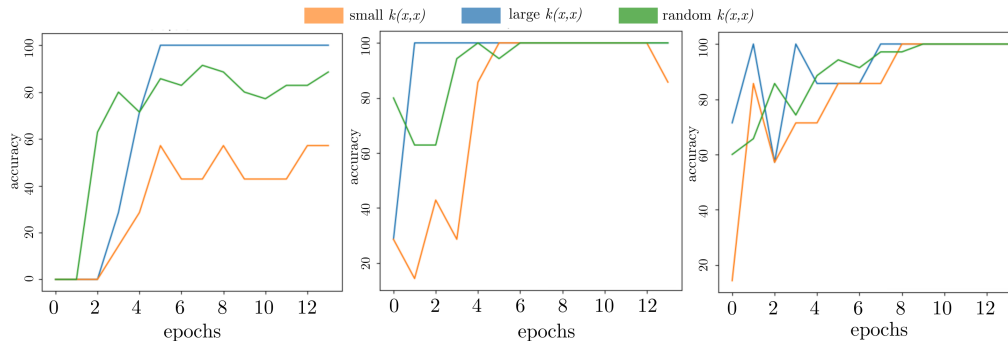


Figure 3: From *left to right*, accuracy vs. epochs for **ship | frog, horse | plane, bird | plane**.

- 86 • **High $k(x, x)$ Samples:** Top 10% of samples with the highest $k(x, x)$ values.
- 87 • **Low $k(x, x)$ Samples:** Bottom 10% of samples with the lowest $k(x, x)$ values.
- 88 • **Medium $k(x, x)$ Samples:** Remaining 80% of samples.

89 We tracked the training accuracy of each sample group over epochs. Figure 2 illustrates that the high
 90 $k(x, x)$ group achieves near-perfect accuracy within fewer epochs compared to the low $k(x, x)$ group.
 91 The medium group exhibits intermediate behavior, often *sheathed* by the easy and hard examples.

92 Visual inspection of samples from different groups reveals distinct characteristics. Figure 3 shows
 93 that high $k(x, x)$ samples are clear and prototypical, whereas low $k(x, x)$ samples often contain noise,
 94 distortions, or are atypical representations of their classes.

95 4 Related Works and Discussion

96 The Neural Tangent Kernel (NTK), introduced by Jacot et al. (2018) [1], has become a fundamental
 97 tool for analyzing the training dynamics of overparameterized neural networks, offering insights into
 98 how models behave in the infinite-width limit. Early works focused on the connection between NTK
 99 and the generalization properties of deep networks, including Arora et al. (2019) [2], who demon-
 100 strated that the NTK matrix governs learning dynamics in function space, significantly influencing
 101 how networks fit data over time. Several studies, such as those by Novak et al. (2019) [5] and Yang
 102 et al. (2020) [7], have shown how NTK can be used to analyze the convergence behavior of neural
 103 networks across different architectures, linking specific NTK properties to model generalization and
 104 performance. More recent efforts by Du et al. (2018) [3] have investigated how NTK can predict
 105 training outcomes in different learning environments, especially for classification tasks, by leveraging
 106 its kernel structure to estimate the sample complexity. Robust learning techniques, such as Jacobian
 107 regularization explored by Hoffman et al. (2019) [8], have also utilized NTK concepts. Finally,
 108 Ilyas et al. (2019) [9] demonstrated how NTK theory can help understand adversarial examples,
 109 showing that adversarial attacks exploit features that NTK-based networks consider salient. Our work
 110 builds on these foundational studies, focusing on the empirical evaluation and theoretical analysis
 111 of individual samples and establishing a direct relationship between NTK diagonal values and the
 112 learnability of training data. This novel perspective provides new insights into optimizing training
 113 strategies by identifying "hard" and "easy" samples based on NTK properties, contributing to the
 114 broader understanding of neural network training dynamics.

115 **References**

- 116 [1] Arthur Jacot, Franck Gabriel, and Clément Hongler. Neural tangent kernel: Convergence and
117 generalization in neural networks. *Advances in neural information processing systems*, 31, 2018.
- 118 [2] Sanjeev Arora, Simon S Du, Wei Hu, Zhiyuan Li, Russ R Salakhutdinov, and Ruosong Wang. On
119 exact computation with an infinitely wide neural net. *Advances in neural information processing*
120 *systems*, 32, 2019.
- 121 [3] Simon S Du, Xiyu Zhai, Barnabas Poczos, and Aarti Singh. Gradient descent provably optimizes
122 over-parameterized neural networks. *arXiv preprint arXiv:1810.02054*, 2018.
- 123 [4] Boris Hanin and Mihai Nica. Finite depth and width corrections to the neural tangent kernel.
124 *arXiv preprint arXiv:1909.05989*, 2019.
- 125 [5] Roman Novak, Lechao Xiao, Jiri Hron, Jaehoon Lee, Alexander A Alemi, Jascha Sohl-Dickstein,
126 and Samuel S Schoenholz. Neural tangents: Fast and easy infinite neural networks in python.
127 *arXiv preprint arXiv:1912.02803*, 2019.
- 128 [6] Roman Novak, Jascha Sohl-Dickstein, and Samuel S Schoenholz. Fast finite width neural tangent
129 kernel. In *International Conference on Machine Learning*, pages 17018–17044. PMLR, 2022.
- 130 [7] Greg Yang and Edward J Hu. Feature learning in infinite-width neural networks. *arXiv preprint*
131 *arXiv:2011.14522*, 2020.
- 132 [8] Judy Hoffman, Daniel A Roberts, and Sho Yaida. Robust learning with jacobian regularization.
133 *arXiv preprint arXiv:1908.02729*, 5(6):7, 2019.
- 134 [9] Andrew Ilyas, Shibani Santurkar, Dimitris Tsipras, Logan Engstrom, Brandon Tran, and Alek-
135 sander Madry. Adversarial examples are not bugs, they are features. *Advances in neural*
136 *information processing systems*, 32, 2019.