

Limitations of Lazy Training of Two-layers Neural Networks

Behrooz Ghorbani ^{1,*} Song Mei ^{2,*} Theodor Misiakiewicz ^{3,*} Andrea Montanari ^{1,3}

¹Department of Electrical Engineering, Stanford University

²ICME, Stanford University

³Department of Statistics, Stanford University

*Equal contributions

Introduction

Consider the function class of two-layers neural networks

$$F_{NN,N} = \{f(\mathbf{x}) = \sum_{i=1}^N a_i \phi(\mathbf{w}_i, \mathbf{x}) : a_i \in \mathbb{R}, \mathbf{w}_i \in \mathbb{R}^d\}$$

- Linearization around (random) parameter $\mathbf{w}_i^0 = (a_i^0, \mathbf{w}_i^0)$

$$\hat{f}_{NN}(\mathbf{x}; \mathbf{w}^0) = \hat{f}_{NN}(\mathbf{x}; \mathbf{w}^0) + \sum_{i=1}^N a_i^0 \phi(\mathbf{w}_i^0, \mathbf{x})$$

- Lazy training [1]: under certain initialization and for a large number of parameters N , the parameters learned by SGD stay close to the initialization \mathbf{w}^0 and the above approximation is accurate [2].

- In this regime, learning the neural network is essentially the same as learning the linearized part:

$$\hat{f}_{NN}(\mathbf{x}; \mathbf{w}^0) = \sum_{i=1}^N a_i^0 \phi(\mathbf{w}_i^0, \mathbf{x}) + \sum_{i=1}^N a_i^0 \phi(\mathbf{w}_i^0, \mathbf{x})$$

Second layer linearization
First layer linearization

We consider the following two function classes which we will refer to as the random feature model (RF) [6], and the neural tangent model (NT) [4]: for $\mathbf{w}_i \stackrel{i.i.d.}{\sim} \mathcal{N}(0, I_d)$,

$$F_{RF,N}(\mathcal{W}) = \{f_N(\mathbf{x}) = \sum_{i=1}^N a_i \phi(\mathbf{w}_i, \mathbf{x}) : a_i \in \mathbb{R}\}$$

$$F_{NT,N}(\mathcal{W}) = \{f_N(\mathbf{x}) = \sum_{i=1}^N a_i \phi(\mathbf{w}_i, \mathbf{x}) : a_i \in \mathbb{R}^d\}$$

Blue: random and fixed. Red: parameters to be optimized.

Questions

- Do RF/NT models provide a good approximation to effectively trained NN (e.g. by SGD)?
- Do RF/NT learn good representations of the data?

We provide two simple settings where we can fully characterize the behavior of RF, NT and SGD-trained NN. In these settings, these two questions admit negative answers.

The prediction risk achieved within any of the regimes $\mathcal{M} \in \{\text{RF}, \text{NT}, \text{NN}\}$ is defined by

$$R_{M,N}(f) = \min_{\hat{f} \in F_{M,N}(\mathcal{W})} \mathbb{E} \int (f(\mathbf{x}) - \hat{f}(\mathbf{x}))^2$$

$$R_{NN,N}(f; \mathbf{w}^0) = \mathbb{E} \int (f(\mathbf{x}) - \hat{f}(\mathbf{x}; \mathbf{w}^0))^2$$

where $\hat{f}(\cdot; \mathbf{w}^0)$ is the neural network produced by steps of stochastic gradient descent (SGD) where each sample is used once, and the stepsize is set to

Quadratic Functions (QF)

Setting: $\mathbf{x}_i \sim \mathcal{N}(0, I_d)$ and responses

$$y_i = f(\mathbf{x}_i) = b_0 + \mathbf{x}_i^T \mathbf{B} \mathbf{x}_i, \text{ with } \mathbf{B} \succeq 0.$$

We take a quadratic activation $\phi(u) = u^2 + c_0$ and consider the high-dimensional regime: $N, d \rightarrow \infty, N/d \rightarrow \gamma \in (0, \infty)$.

Results [5]:

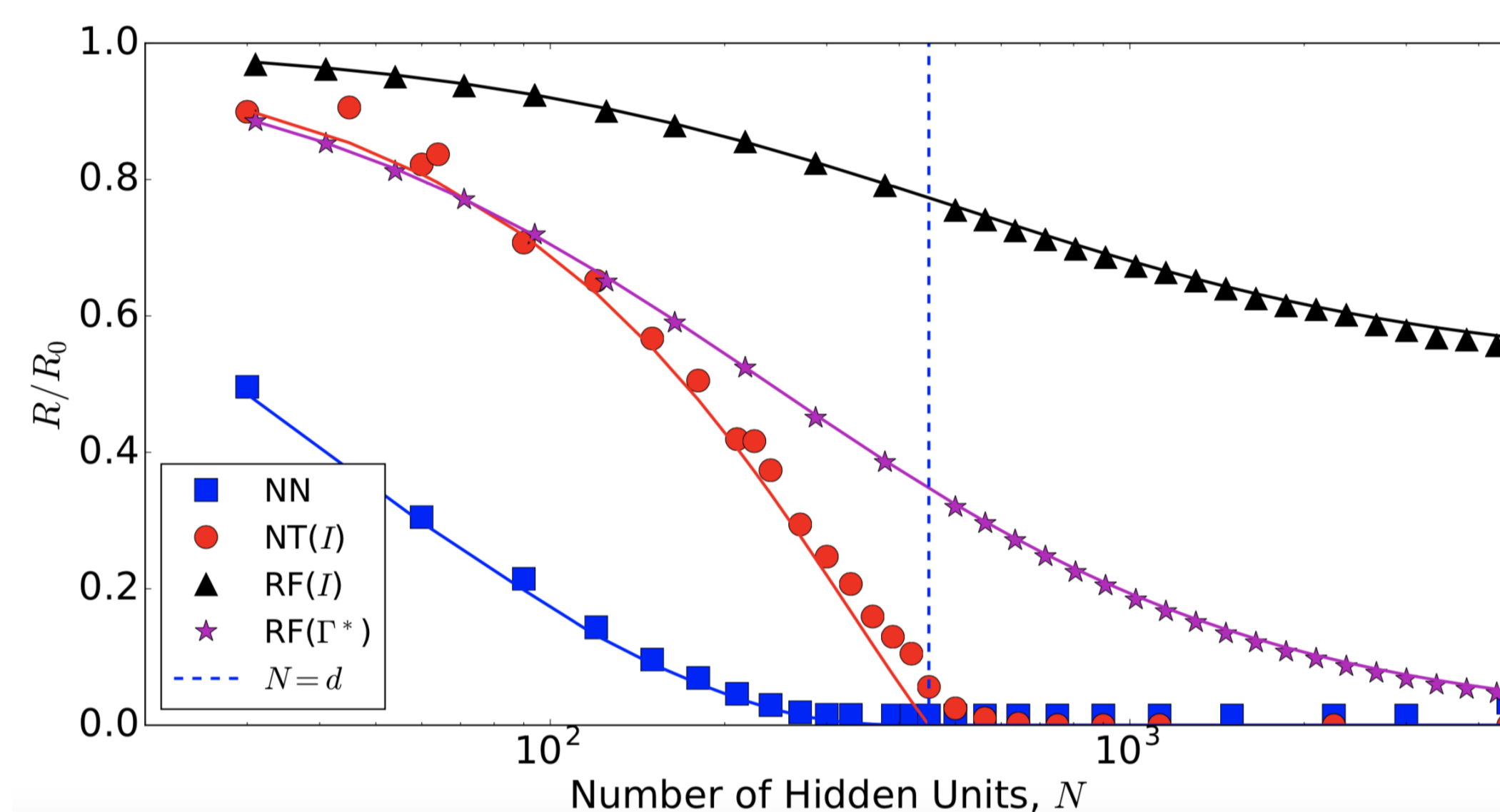


Figure 1: Prediction (test) error in fitting a quadratic function in $d = 450$ dimensions, as a function of the number of neurons N . Lines are analytical predictions obtained in this paper [5], and dots are empirical results.

- Naive RF/NT do not learn good representations of the data.
- SGD-trained NN learns the most important eigendirections of f and fits them, hence surpassing the NT model which remains confined to a random subspace spanned by \mathbf{w}_i .
- There exists an arbitrary large gap between the SGD-trained networks and the neural tangent model.

Neural networks are superior to linearized model such as RF and NT, because they can learn a good representation of the data.

Mixture of Gaussians

Setting: $y_i = \pm 1$ with equal probability $1/2$, and $\mathbf{x}_i | y_i = +1 \sim \mathcal{N}(0, I_{d+})$, $\mathbf{x}_i | y_i = -1 \sim \mathcal{N}(0, I_{d-})$. Take $\phi(u) = u^2 + c_0$ and $\mathbf{w}_i \sim \mathcal{N}(0, I_{d/d})$.

$$R_{M,N}(P_{I, \gamma}) = \begin{cases} \frac{1}{1 + \frac{1}{1+\gamma} \cdot \frac{\bar{r}(\gamma)^2}{2d}} & \text{for } M = \text{RF}, \\ \frac{1}{1 + \left(\frac{\gamma}{1+\gamma}\right)^2 \frac{2}{F}} & \text{for } M = \text{NT}, \\ \frac{1}{1 + \mathbb{P}_{i=1}^N \frac{1}{d} \lambda_i(\gamma)^2 / 2} & \text{for } M = \text{NN}. \end{cases}$$

- See Figure 2 for analytical and empirical results.
- We recover a similar behavior as in the QF model.
- Note that the Bayes error is not achieved in this model.
- We do not show convergence of SGD in this setting but we expect a similar result to the QF model to hold.

Analytical Predictions for QF

Random features model

Theorem 1 ([5]) Take $\phi(x) = x^2 - 1$, $\mathbf{w}_i \sim \mathcal{N}(0, \gamma)$. Then, as $N, d \rightarrow \infty$ with $N/d \rightarrow \gamma$

$$R_{RF,N}(f) = \frac{f^2_{L_2} \mathbb{B}}{B^2_F} \left(1 - \frac{d}{B^2_F} \frac{2}{F} + o_{d,P}(1)\right)$$

- See [5] for the Theorem for general activation function ϕ .
- The risk highly depends on the weight distribution.
- In particular for any activation function,

$$\lim_{d \rightarrow \infty} \lim_{N/d \rightarrow \gamma} \frac{R_{RF,N}(f)}{f^2_{L_2}} = \lim_{d \rightarrow \infty} \left(1 - \frac{B^2}{B^2_F} \frac{2}{F}\right)$$

The risk vanishes only if \mathbf{B} is chosen perfectly and $\gamma = 1$. The asymptotic risk is independent of the non-linearity!

Neural Tangent model

Theorem 2 ([5]) Take $\phi(x) = x^2$, $\mathbf{w}_i \sim \mathcal{N}(0, I_{d/d})$. Then, as $N, d \rightarrow \infty$ with $N/d \rightarrow \gamma$

$$\frac{\mathbb{E} \mathbb{W}[R_{NT,N}(f)]}{f^2_{L_2}} = (1 - \gamma)_+^2 + (1 - \gamma)_+ \frac{\text{Tr}(\mathbf{B})^2}{d B^2_F} + o_d(1)$$

- For $N < d$, NT fits f along a random subspace determined by the weights \mathbf{w}_i (not the most important subspace).
- For $N \geq d$, weights span the whole space (vanishing risk).

Fully-trained NN model

Theorem 3 ([5]) Take $\phi(x) = x^2$. Then, as $N, d \rightarrow \infty$ with $N/d \rightarrow \gamma$

$$\lim_{d \rightarrow \infty} \lim_{N/d \rightarrow \gamma} \mathbb{P} \left(R_{NN,N}(f; \mathbf{w}^0) - R_{NN,N}(f) > \epsilon \right) = 0,$$

$$R_{NN,N}(f) = 2 \sum_{i=1}^d \lambda_i(\mathbf{B})^2$$

with $\lambda_1(\mathbf{B}) \geq \dots \geq \lambda_d(\mathbf{B})$ ordered eigenvalues of \mathbf{B} .

- Here, we studied a one-pass version of SGD. The probability is over the random initialization \mathcal{W}^0 and the samples.
- The global convergence is proved by showing convergence of SGD to the gradient flow in the population risk and then proving a strict saddle property for the population risk.
- SGD-learned NN fits f along the most important subspace (the N principal eigendirections of \mathbf{B}).

How General are these Phenomena?

- The separation between NN and NT is established only for $N < d$. We expect the separation to generalize to $N \geq d$ by considering higher order polynomials: for third- or higher-order polynomials, NT does not achieve vanishing risk at any $(0, \infty)$ (see [3]).
- While we are only able to provide theory for NN and NT for quadratic activation, we performed extensive experiments with other non-linearities. See Figure 3 for fitting a quadratic function with ReLU activation. In particular, the positive gap between NN and NT is still present for $N < d$.

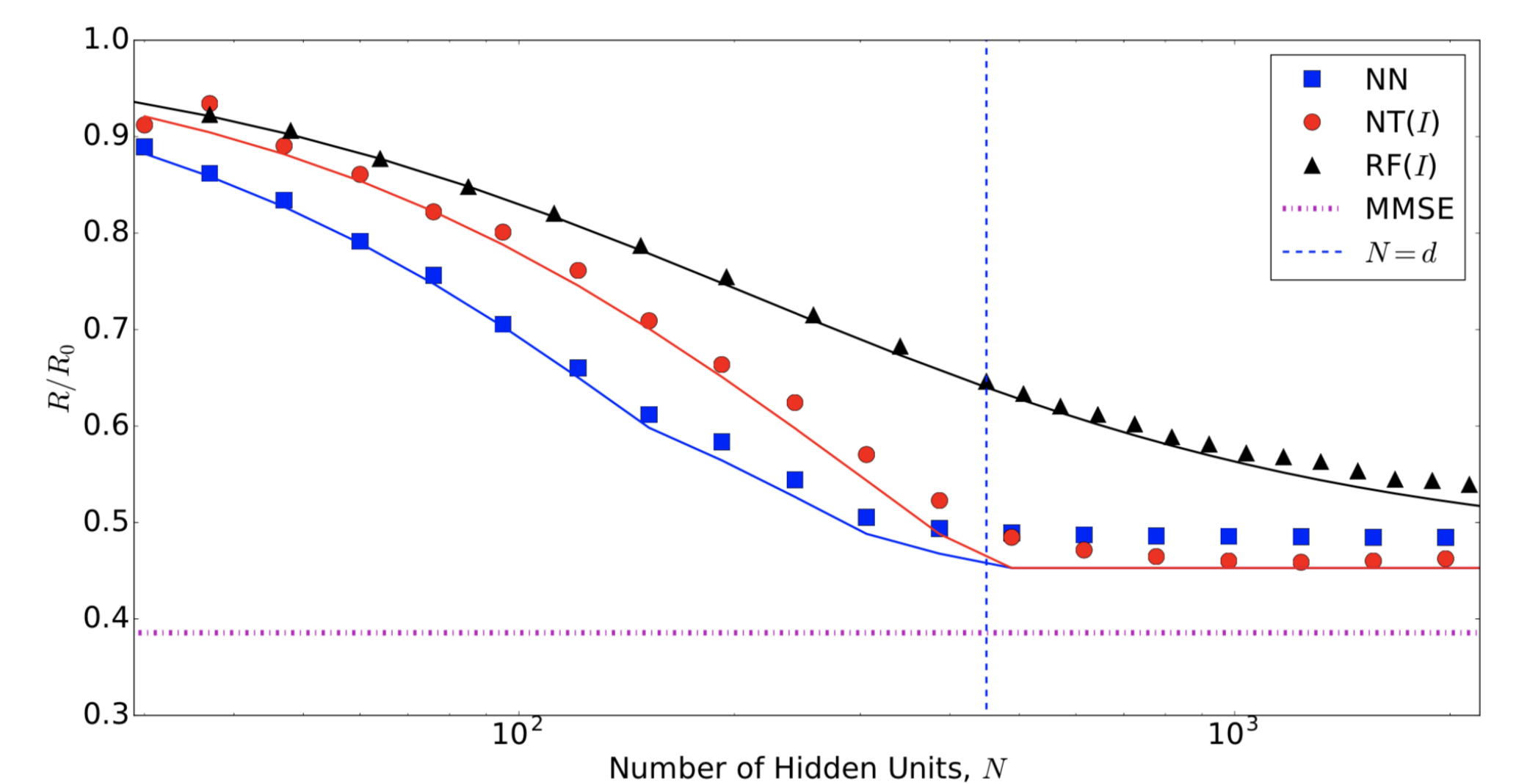


Figure 2: Prediction (test) error in fitting a mixture of Gaussians in $d = 450$ dimensions, as a function of N . Lines are analytical predictions obtained in this paper [5], and dots are empirical results. Dotted line is the Bayes error.

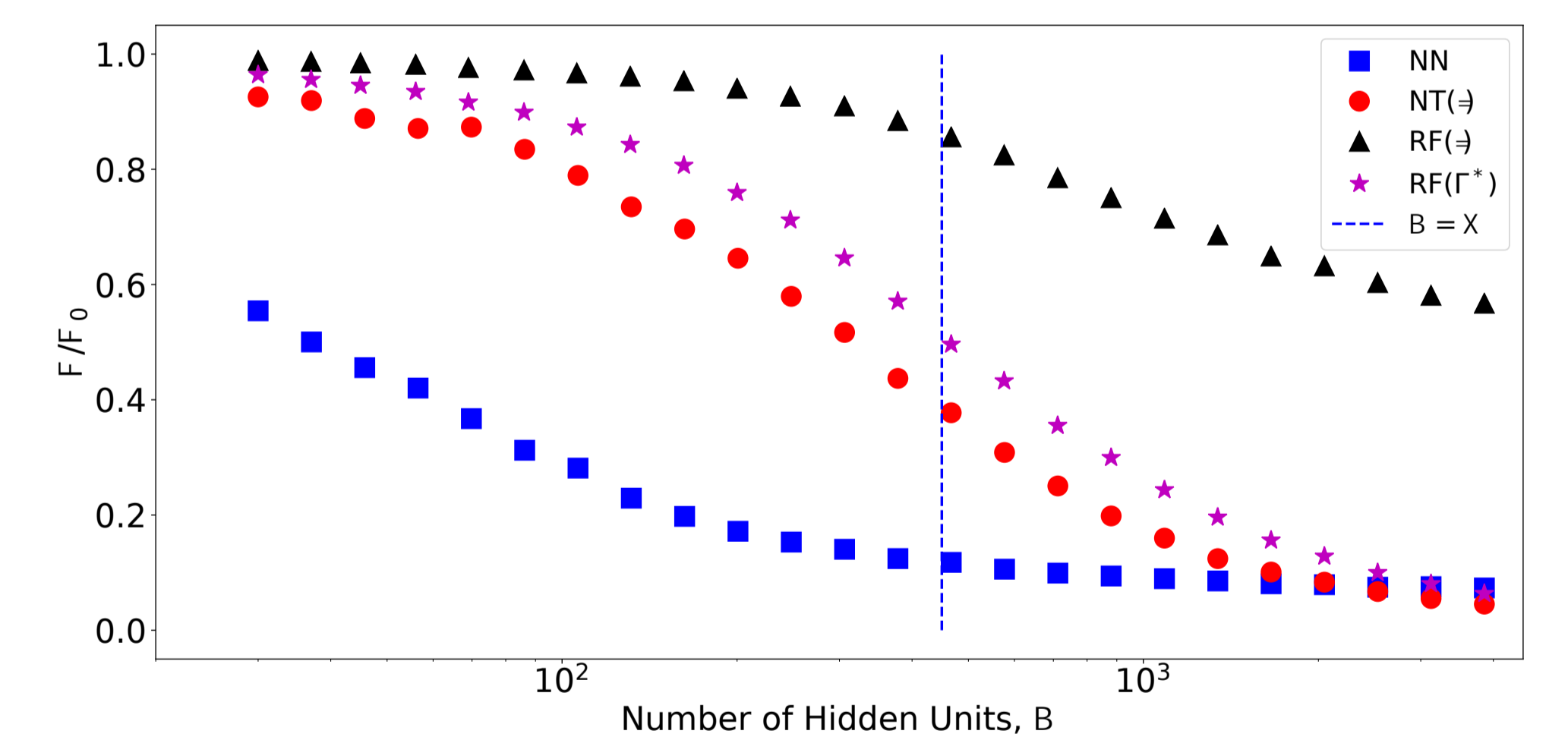


Figure 3: Empirical prediction (test) error in fitting a quadratic function in $d = 450$ dimensions with ReLU activation, as a function of N .

Bibliography

- L. Chizat and F. Bach. A note on lazy training in supervised differentiable programming. *arXiv:1812.07956*, 2018.
- S. S. Du, J. D. Lee, H. Li, L. Wang, and X. Zhai. Gradient descent finds global minima of deep neural networks. *arXiv:1811.03804*, 2018.
- B. Ghorbani, S. Mei, T. Misiakiewicz, and A. Montanari. Linearized two-layers neural networks in high dimension. *arXiv:1904.12191*, 2019.
- A. Jacot, F. Gabriel, and C. Hongler. Neural tangent kernel: Convergence and generalization in neural networks. In *Advances in neural information processing systems*, pages 8571–8580, 2018.
- S. Mei, T. Misiakiewicz, B. Ghorbani, and A. Montanari. Limitations of lazy training of two-layers neural network. In *Advances in Neural Information Processing Systems*, pages 9108–9118, 2019.
- A. Rahimi and B. Recht. Random features for large-scale kernel machines. In *Advances in neural information processing systems*, pages 1177–1184, 2008.