

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**

$$\mathcal{F}_{\text{NN},N} = \left\{ f(\mathbf{x}) = \sum_{i=1}^N a_i \sigma(\langle \mathbf{w}_i, \mathbf{x} \rangle) : a_i \in \mathbb{R}, \mathbf{w}_i \in \mathbb{R}^d \right\}.$$

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

$$f_{\text{NN}}(\mathbf{x}; \theta) \approx f_{\text{NN}}(\mathbf{x}; \theta^0) + \langle \theta - \theta^0, \nabla_{\theta} f_{\text{NN}}(\mathbf{x}; \theta^0) \rangle$$
- Lazy training [1]: under certain initialization and for a large number of parameters N , the parameters θ learned by SGD stay close to the initialization θ^0 and the above approximation is accurate [2].
- In this regime, learning the neural network is essentially the same as learning the linearized part:

$$f_{\text{NN}}(\mathbf{x}; \theta) \approx 0 + \underbrace{\sum_{i=1}^N \Delta a_i \sigma(\langle \mathbf{w}_i^0, \mathbf{x} \rangle)}_{\text{Second layer linearization}} + \underbrace{\sum_{i=1}^N a_i^0 \sigma'(\langle \mathbf{w}_i^0, \mathbf{x} \rangle) \langle \Delta \mathbf{w}_i, \mathbf{x} \rangle}_{\text{First layer linearization}}$$

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 $\mathbf{w}_i \stackrel{i.i.d.}{\sim} \mathcal{N}(0, \mathbf{I}_d)$,

$$\mathcal{F}_{\text{RF},N}(\mathbf{W}) = \left\{ f_N(\mathbf{x}) = \sum_{i=1}^N a_i \sigma(\langle \mathbf{w}_i, \mathbf{x} \rangle) : a_i \in \mathbb{R} \right\},$$

$$\mathcal{F}_{\text{NT},N}(\mathbf{W}) = \left\{ f_N(\mathbf{x}) = \sum_{i=1}^N \langle a_i, \mathbf{x} \rangle \sigma'(\langle \mathbf{w}_i, \mathbf{x} \rangle) : a_i \in \mathbb{R}^d \right\}.$$

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 $\mathbf{M} \in \{\text{RF}, \text{NT}, \text{NN}\}$ is defined by

$$R_{\mathbf{M},N}(f_*) = \min_{\hat{f} \in \mathcal{F}_{\mathbf{M},N}(\mathbf{W})} \mathbb{E} \left\{ (f_*(\mathbf{x}) - \hat{f}(\mathbf{x}))^2 \right\},$$

$$R_{\text{NN},N}(f_*; \ell, \varepsilon) = \mathbb{E} \left\{ (f_*(\mathbf{x}) - \hat{f}(\mathbf{x}; \ell, \varepsilon))^2 \right\},$$

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

Quadratic Functions (QF)

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

$$y_i = f_*(\mathbf{x}_i) \equiv b_0 + \langle \mathbf{x}_i, \mathbf{B} \mathbf{x}_i \rangle, \text{ with } \mathbf{B} \succeq 0.$$

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

Results [5]:

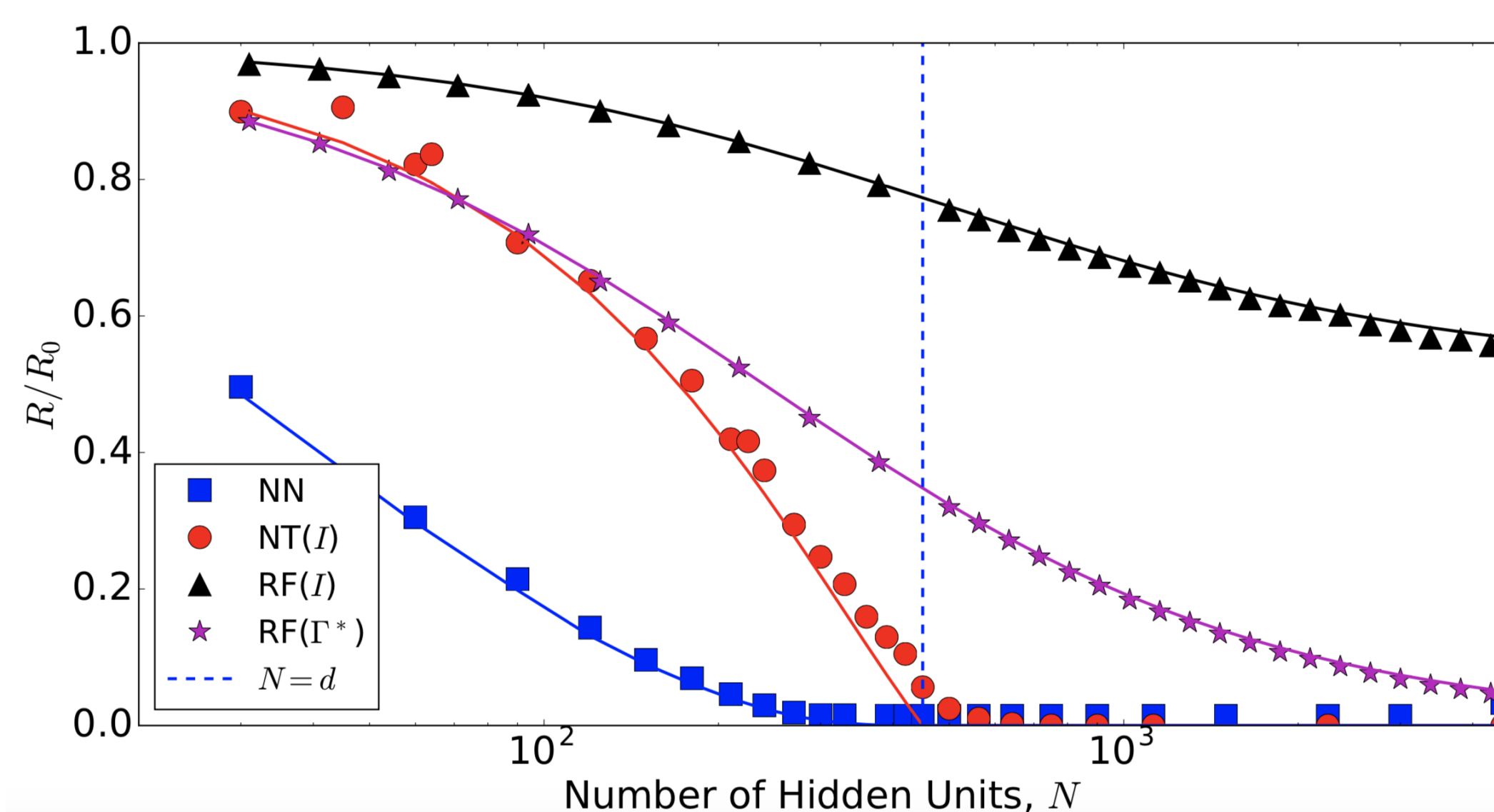


Figure 1: Prediction (test) error in fitting a quadratic function in $d = 450$ dimensions, as a function of the number of neurons N . Lines are analytical predictions obtained in this paper [5], and dots are empirical results.

- 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 $\mathbf{x}_i | y_i = +1 \sim \mathcal{N}(0, \mathbf{I}_d + \Delta)$, $\mathbf{x}_i | y_i = -1 \sim \mathcal{N}(0, \mathbf{I}_d - \Delta)$. Take $\sigma(u) = u^2 + c_0$ and $\mathbf{w}_i \sim \mathcal{N}(0, \mathbf{I}_d/d)$.

$$R_{\mathbf{M},N}(\mathbb{P}_{\mathbf{I},\Delta}) \approx \begin{cases} \frac{1}{1 + \frac{\rho}{1+2\rho} \cdot \frac{\tilde{r}(\Delta)^2}{2d}} & \text{for } \mathbf{M} = \text{RF}, \\ \frac{1}{1 + \kappa(\rho, \Delta) \|\Delta\|_F^2 / 2} & \text{for } \mathbf{M} = \text{NT}, \\ \frac{1}{1 + \sum_{i=1}^{N \wedge d} \lambda_i(\Delta)^2 / 2} & \text{for } \mathbf{M} = \text{NN}. \end{cases}$$

- 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 $\sigma(x) = x^2 - 1$, $\mathbf{w}_i \sim \mathcal{N}(0, \Gamma)$. Then, as $N, d \rightarrow \infty$ with $N/d \rightarrow \rho$

$$R_{\text{RF},N}(f_*) = \|f_*\|_{L_2}^2 \left(1 - \frac{\rho d \langle \mathbf{B}, \Gamma \rangle^2}{\|\mathbf{B}\|_F^2 (1 + \rho d \|\Gamma\|_F^2)} + o_{d,\mathbb{P}}(1) \right).$$

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

$$\lim_{\rho \rightarrow \infty} \lim_{d \rightarrow \infty, N/d \rightarrow \rho} \frac{R_{\text{RF},N}(f_*)}{\|f_*\|_{L_2}^2} = \lim_{d \rightarrow \infty} \left(1 - \frac{\langle \mathbf{B}, \Gamma \rangle^2}{\|\mathbf{B}\|_F^2 \|\Gamma\|_F^2} \right).$$

The risk vanishes only if $\Gamma \propto \mathbf{B}$ is chosen perfectly and $\rho \rightarrow \infty$. The asymptotic risk is independent of the non-linearity!

Neural Tangent model

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

$$\frac{\mathbb{E}_{\mathbf{W}}[R_{\text{NT},N}(f_*)]}{\|f_*\|_{L_2}^2} = \left\{ (1 - \rho)_+^2 + \rho(1 - \rho) + \frac{\text{Tr}(\mathbf{B})^2}{d \|\mathbf{B}\|_F^2} + o_d(1) \right\}.$$

- 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 $\sigma(x) = x^2$. Then, as $N, d \rightarrow \infty$ with $N/d \rightarrow \rho$

$$\lim_{t \rightarrow \infty} \lim_{\varepsilon \rightarrow 0} \mathbb{P} \left(\left| R_{\text{NN},N}(f_*; \ell = t/\varepsilon, \varepsilon) - R_{\text{NN},N}(f_*) \right| \geq \delta \right) = 0,$$

$$R_{\text{NN},N}(f_*) = 2 \sum_{i=N+1}^d \lambda_i(\mathbf{B})^2, \text{ with } \lambda_1(\mathbf{B}) \geq \dots \geq \lambda_d(\mathbf{B}) \text{ ordered eigenvalues of } \mathbf{B}.$$

- Here, we studied a one-pass version of SGD. The probability is over the random initialization \mathbf{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 $\rho \in (0, \infty)$ (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 $N < d$.

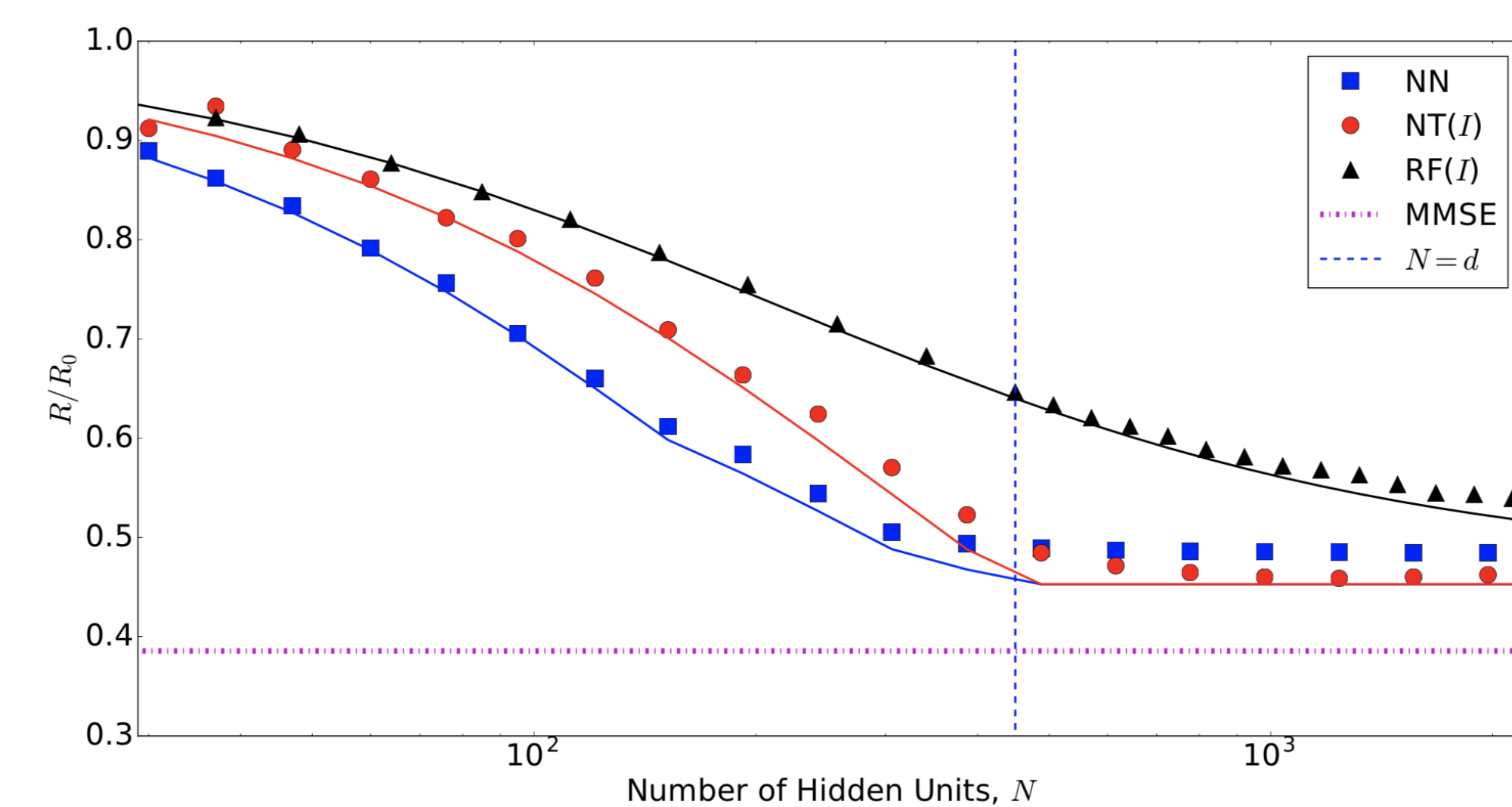


Figure 2: Prediction (test) error in fitting a mixture of Gaussians in $d = 450$ dimensions, as a function of N . Lines are analytical predictions obtained in this paper [5], and dots are empirical results. Dotted line is the Bayes error.

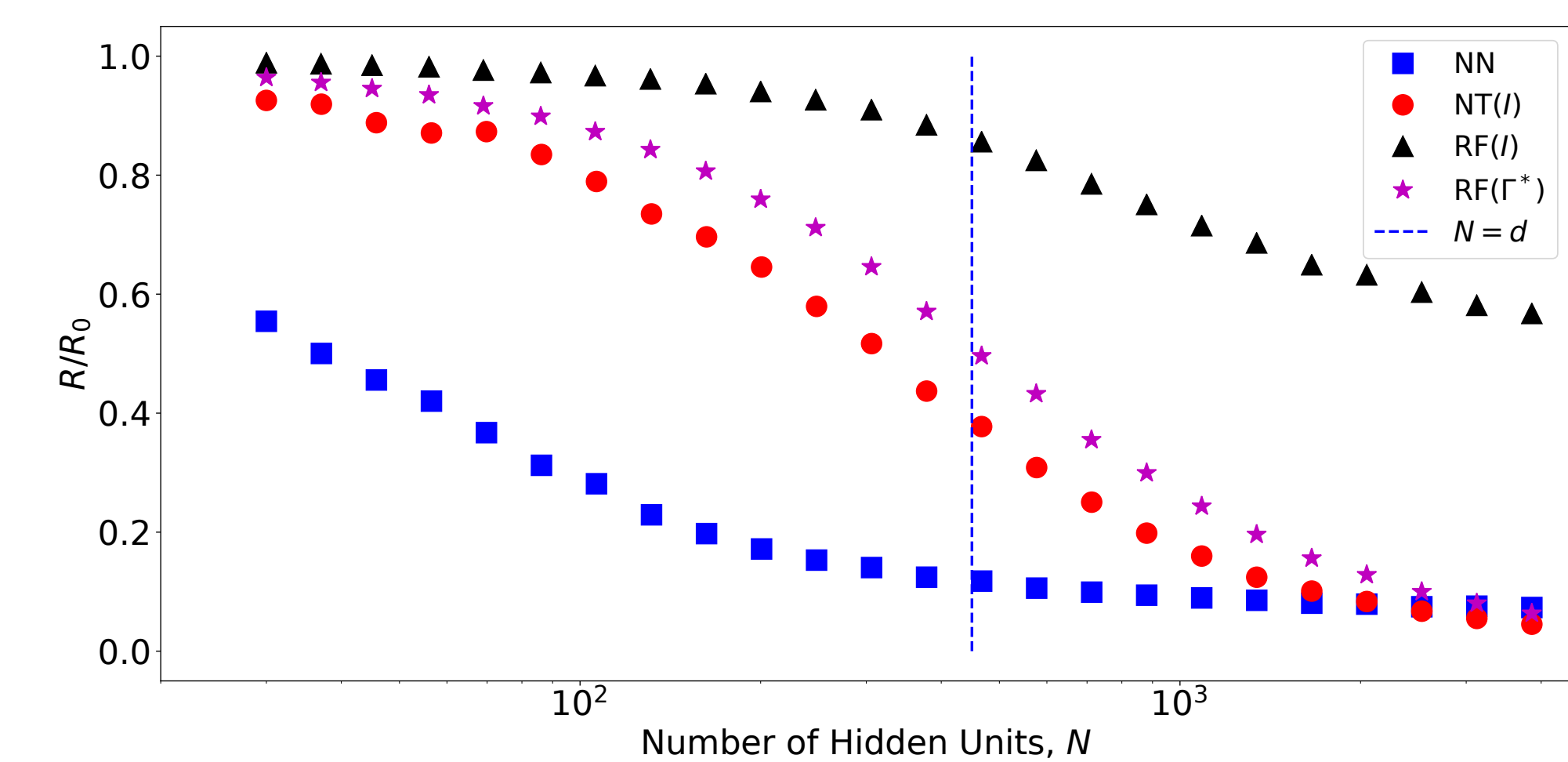


Figure 3: Empirical prediction (test) error in fitting a quadratic function in $d = 450$ dimensions with ReLu activation, as a function of N .

Bibliography

- [1] L. Chizat and F. Bach. A note on lazy training in supervised differentiable programming. *arXiv:1812.07956*, 2018.
- [2] 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.
- [3] B. Ghorbani, S. Mei, T. Misiakiewicz, and A. Montanari. Linearized two-layers neural networks in high dimension. *arXiv:1904.12191*, 2019.
- [4] 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.
- [5] 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.
- [6] A. Rahimi and B. Recht. Random features for large-scale kernel machines. In *Advances in neural information processing systems*, pages 1177–1184, 2008.