
Associative Memory in Iterated Overparameterized Sigmoid Autoencoders

Yibo Jiang¹ Cengiz Pehlevan²

Abstract

Recent work showed that overparameterized autoencoders can be trained to implement associative memory via iterative maps, when the trained input-output Jacobian of the network has all of its eigenvalue norms strictly below one. Here, we theoretically analyze this phenomenon for sigmoid networks by leveraging recent developments in deep learning theory, especially the correspondence between training neural networks in the infinite-width limit and performing kernel regression with the Neural Tangent Kernel (NTK). We find that overparameterized sigmoid autoencoders can have attractors in the NTK limit for both training with a single example and multiple examples under certain conditions. In particular, for multiple training examples, we find that the norm of the largest Jacobian eigenvalue drops below one with increasing input norm, leading to associative memory.

1. Introduction

The mechanisms behind memory have been a long interest of neuroscientists. Hopfield’s seminal work proposed that associative memory can be implemented by attractor neural dynamics (Hopfield, 1982), and has been the dominant model that shapes thinking in this domain (Hertz, 2018). Recently, Radhakrishnan et al. (2019; 2018) proposed an alternative mechanism and showed that overparameterized autoencoders trained with gradient descent could also implement associative memory in an iterative fashion. These networks are reported to be easy to train and to not suffer from spurious attractors, unlike Hopfield networks (Amit & Treves, 1989; Hertz, 2018). The potential benefits of this approach make a theoretical account of it necessary, which we aim to provide here.

¹John A. Paulson School of Engineering and Applied Sciences, Harvard University, Cambridge, MA, USA ²Center for Brain Science, Harvard University, Cambridge, MA, USA. Correspondence to: Cengiz Pehlevan <cpehlevan@seas.harvard.edu>.

We study auto-encoders in a limit of neural networks that makes theoretical analysis possible. Specifically, as the width of the hidden layers of a feedforward neural network is taken to infinity with a particular initialization scheme, its training dynamics simplifies and can be described by ridgeless kernel interpolation with a kernel called the Neural Tangent Kernel (NTK) (Jacot et al., 2018). Working in the NTK limit, we examine the input-output Jacobian matrices of the trained networks as they control the stability of trained fixed points. We focus on networks with sigmoid activation functions, but our results can be extended to other sigmoidal functions such as erf and tanh using similar techniques. We make a distinction between the cases of a single training example and multiple training examples, which exhibit different memory behaviors.

Our main contributions and results are summarized below:

- First, we analyze autoencoders in the NTK limit trained on a single training example. We argue that the trained Jacobian will stay close to initialization under certain conditions. Therefore, if the initial Jacobian has all eigenvalue norms smaller than 1, this training example will be an **attractor**.
- Next, we specialize to 2-layer networks. We show that when the norms of training examples are small, attractor formation can **fail** due to the presence of eigenvalue 1 in the spectrum of trained Jacobian matrices.
- We show that a 2-layer network can have **attractors** when the norms of training examples are large as the induced NTK relies more on the non-linear, saturated region of the activation function. This suggests that the network transitions from a regime where attractor formation fails to a regime where it succeeds as the input norm grows.
- We verify the predictions of our theoretical results in simulations.

Many previous works on generalization (Allen-Zhu et al., 2019; Cao & Gu, 2019) and training (Jacot et al., 2018; Du et al., 2018; Allen-Zhu et al., 2018) in the NTK limit focus on input data on the unit sphere which is violated in practice (Krizhevsky et al., 2012). We highlight how the input norm can set trained neural networks in different learning regions with respect to the trained input-output Jacobian.

Associative memory behavior induced by training overparameterized autoencoders could provide insights into the

implicit bias and generalization of neural networks. Autoencoders are trained to learn identity maps, however the existence of attractors indicates failure to learn such maps and thus failure to generalize. Recently, [Zhang et al. \(2019\)](#) observed that fully connected networks tend to learn a constant function, a global attractor, when trained on a single example. Our single training example results explain this behavior. However, we also show that attractor formation is dependent on input norm when multiple examples are present, demonstrating that training with a single example is not sufficient to explain the implicit bias of neural networks.

2. Related Work

Our results use ideas related to NTK and signal propagation in deep networks with random weights. Attractor behavior can also be associated with the implicit bias of deep learning. We review relevant literature from these domains.

Neural Tangent Kernel: We first review literature on neural networks optimization, especially the NTK theory which we will use extensively through this paper. Training of neural networks poses a challenging non-convex optimization problem. Analysis simplifies if focused on the linearized training dynamics of gradient flow using the NTK theory ([Jacot et al., 2018](#)). The basic idea is that, if initialized properly, in the infinite width limit, parameters of the network stay close to initialization ([Chizat et al., 2019](#)). Thus, NTK stays relatively constant throughout training. Because NTK governs the training dynamics, positive-definite kernel ensures global convergence of optimization. Subsequent papers expanded this idea to finite network widths, gradient descent or stochastic gradient descent as opposed to gradient flow, and different loss functions for regression and classification ([Du et al., 2018](#); [Allen-Zhu et al., 2018](#); [Zou & Gu, 2019](#); [Lee et al., 2019](#)). Also related is research pointing that neural networks at initialization in the infinite width limit behave as Gaussian Processes ([Lee et al., 2017](#); [Matthews et al., 2018](#)).

Signal Propagation in Deep Networks with Random Weights: In a set of ideas that we will make use of later, [Poole et al. \(2016\)](#) developed a mean-field formalism to study layer-to-layer propagation of activation variances and covariances in deep networks with random weights. This line of work ([Poole et al., 2016](#); [Schoenholz et al., 2016](#)) identifies a phase transition between ordered and chaotic regime, where nearby input points converge or diverge as they propagate through the layers, induced by different variances of weight and bias initializations. Using random matrix theories developed for Gram matrices of neural networks ([Pennington & Worah, 2017](#)), [Pennington et al. \(2017; 2018\)](#) calculate singular value spectra for input-output Jacobians at initialization and identify its relation to ordered/chaotic training regime. These results cannot be

applied to our problem, as they are in a different setting assuming a large depth limit such that the variances for all layers are at the fixed points of the layer-to-layer iterative maps.

Generalization and Implicit Bias: Associative memory behavior of neural networks can provide insight into generalization of deep learning and implicit bias of gradient descent. Here, we review some recent works in this area, focusing on overparameterized networks and NTK. A pair of recent papers ([Hayou et al., 2019](#); [Xiao et al., 2019](#)) demonstrate how to use the theories developed for signal propagation to understand NTK regression better and thus trainability and generalization of neural networks. Methods derived from statistical physics also give insight ([Cohen et al., 2019](#); [Bordelon et al., 2020](#)). However, there is a gap in terms of generalization between NTK regression and training neural networks ([Allen-Zhu & Li, 2019](#); [Arora et al., 2019](#)) (see however ([Lee et al., 2019](#))), which prompts research on generalization in deep learning beyond NTK ([Allen-Zhu et al., 2019](#); [Bai & Lee, 2019](#)). Another line of research on generalization looks at the implicit bias of gradient descent. Gradient descent on logistic regression can lead to the max-margin solution ([Soudry et al., 2018](#); [Ji & Telgarsky, 2018](#)), while optimization on mean squared regression has a shortest path solution ([Oymak & Soltanolkotabi, 2018](#)).

3. Preliminaries

In this section, we set up our notation and review background material.

3.1. Neural Networks

The output of a neural network is defined by $f(\mathbf{x}) = \tilde{\alpha}^{(L)}(\mathbf{x})$, where the functions $\tilde{\alpha}^{(\ell)}(\cdot) : \mathbb{R}^{n_0} \rightarrow \mathbb{R}^{n_\ell}$ (*pre-activations*) and $\alpha^{(\ell)}(\cdot) : \mathbb{R}^{n_0} \rightarrow \mathbb{R}^{n_\ell}$ (*activations*) follow the recursive relation:

$$\alpha^{(0)}(\mathbf{x}) \equiv \mathbf{x}$$

$$\tilde{\alpha}^{(\ell+1)}(\mathbf{x}) \equiv \frac{1}{\sqrt{n_\ell}} \mathbf{W}^{(\ell)} \alpha^{(\ell)}(\mathbf{x}), \quad \alpha^{(\ell)}(\mathbf{x}) \equiv \sigma(\tilde{\alpha}^{(\ell)}(\mathbf{x})),$$

where σ is an element-wise activation function and weights are initialized by sampling from an i.i.d. standard Gaussian. We are mostly interested in

$$\text{sigmoid} = \frac{1}{1 + e^{-x}} \quad (1)$$

as the activation function. We drop the bias term for simplicity. We expect our results to be qualitatively the same with a bias term for sigmoid as the behavior is governed by the activation's shape. Throughout the paper, we will use $f_0(\mathbf{x})$ and $f_\infty(\mathbf{x})$ to denote neural networks at initialization and training to zero loss respectively. As shown later, the

attractor behavior is mostly governed by the shape of the sigmoidal activations.

Inputs: Given n training points $\{\mathbf{x}_i\}_1^n$, we define the following two matrices.

$$\hat{\mathbf{X}} = \begin{pmatrix} | & \cdots & | \\ \mathbf{x}_1 & \cdots & \mathbf{x}_n \\ | & \cdots & | \end{pmatrix}, f(\hat{\mathbf{X}}) = \begin{pmatrix} | & \cdots & | \\ f(\mathbf{x}_1) & \cdots & f(\mathbf{x}_n) \\ | & \cdots & | \end{pmatrix},$$

where $\hat{\mathbf{X}} \in \mathbb{R}^{n_0 \times n}$ is the data matrix and each column of $f(\hat{\mathbf{X}}) \in \mathbb{R}^{n_0 \times n}$ is the output of the corresponding training example. We further assume that all input examples share the same norm (i.e. $\forall i \|\mathbf{x}_i\|_2 = r$).

Jacobian Matrix: Given the network setup, the input-output Jacobian matrix can be computed as:

$$\mathbf{J}(\mathbf{x}) = \frac{1}{\sqrt{n_L}} \mathbf{W}^{(L)} \prod_{k=1}^L \left(\mathbf{D}^{(k)} \frac{1}{\sqrt{n_{k-1}}} \mathbf{W}^{(k-1)} \right)$$

where

$$\mathbf{D}^{(k)} = \text{diag}(\dot{\sigma}(\tilde{\alpha}^{(k)}(\mathbf{x}))).$$

Here $\dot{\cdot}$ denotes first derivative, and diag takes in a vector and outputs a diagonal matrix with the vector at the diagonal. We will use $J_0(\mathbf{x})$ and $J_\infty(\mathbf{x})$ to denote Jacobian at initialization and training to zero loss respectively.

Autoencoder: An autoencoder network is trained via gradient flow to optimize the following loss function:

$$\arg \min_f \frac{1}{2n} \sum_{i=1}^n \|f(\mathbf{x}_i) - \mathbf{x}_i\|_2^2,$$

where f is the network defined above.

3.2. Neural Tangent Kernel

In the large-width regime, the neural network f can be approximated by a linearization with respect to its parameters θ (Lee et al., 2019):

$$f(\mathbf{x}; \theta_t) \approx f_0(\mathbf{x}) + \partial_\theta f(\mathbf{x})|_{\theta=\theta_0} (\theta_t - \theta_0),$$

where θ_0 and θ_t are vectors of the network parameters at initialization and time t . The first term remains unchanged throughout training. Moreover, we can view $\partial_\theta f(\mathbf{x})|_{\theta=\theta_0}$ as a feature map in Hilbert space and derive the following matrix kernel,

$$(\Theta_0^{(L)}(\hat{\mathbf{x}}, \mathbf{x}))_{dd'} = \left\langle \partial_\theta f_d(\hat{\mathbf{x}})|_{\theta=\theta_0}, \partial_\theta f_{d'}(\mathbf{x})|_{\theta=\theta_0} \right\rangle,$$

where $\langle \cdot, \cdot \rangle$ denotes inner product. Indeed, in the infinite width limit (Jacot et al., 2018), $\Theta_0^{(L)}$ converges in probability (stochasticity induced by random initialization) to a

deterministic limiting kernel, $\Theta_0^{(L)}(\hat{\mathbf{x}}, \mathbf{x}) \rightarrow \Theta_\infty^{(L)}(\hat{\mathbf{x}}, \mathbf{x}) \mathbf{I}_{n_L}$, where $\Theta_\infty^{(L)}$ is a scalar kernel and the training dynamics is entirely governed by it.

On the other hand, it can also be shown that in the NTK limit ($n_1, \dots, n_L \rightarrow \infty$ sequentially) and when the weights are initialized by sampling from an i.i.d. standard Gaussian distribution, output functions at each layer tend to be i.i.d. centered Gaussian Processes at initialization (Lee et al., 2017), and the covariance matrix of layer ℓ can be defined recursively by

$$\Sigma^{(1)}(\hat{\mathbf{x}}, \mathbf{x}) = \frac{1}{n_0} \hat{\mathbf{x}}^T \mathbf{x},$$

$$\Sigma^{(\ell+1)}(\hat{\mathbf{x}}, \mathbf{x}) = \mathbb{E}_{g \sim \mathcal{N}(0, \Sigma^{(\ell)})} [\sigma(g(\mathbf{x})) \sigma(g(\hat{\mathbf{x}}))].$$

Under the same limit, $\Theta_\infty^{(L)}$ can also be recursively defined (Jacot et al., 2018):

$$\Theta_\infty^{(1)}(\hat{\mathbf{x}}, \mathbf{x}) = \Sigma^{(1)}(\hat{\mathbf{x}}, \mathbf{x}),$$

$$\Theta_\infty^{(\ell+1)}(\hat{\mathbf{x}}, \mathbf{x}) = \Theta_\infty^{(\ell)}(\hat{\mathbf{x}}, \mathbf{x}) \dot{\Sigma}^{(\ell+1)}(\hat{\mathbf{x}}, \mathbf{x}) + \Sigma^{(\ell+1)}(\hat{\mathbf{x}}, \mathbf{x}),$$

where

$$\dot{\Sigma}^{(\ell+1)}(\hat{\mathbf{x}}, \mathbf{x}) = \mathbb{E}_{g \sim \mathcal{N}(0, \Sigma^{(\ell)})} [\dot{\sigma}(g(\mathbf{x})) \dot{\sigma}(g(\hat{\mathbf{x}}))].$$

For gradient flow training with a least squares loss to zero training error, there is a closed form solution for $f_\infty(\mathbf{x})$ using NTK in the infinite width limit (Jacot et al., 2018). In the case of autoencoder, we get that,

$$f_\infty(\mathbf{x}) = \left(\hat{\mathbf{X}} - f_0(\hat{\mathbf{X}}) \right) \tilde{\mathbf{K}}^{-1} \mathbf{k}_x + f_0(\mathbf{x}), \quad (2)$$

and

$$\mathbf{J}_\infty(\mathbf{x}) = \left(\hat{\mathbf{X}} - f_0(\hat{\mathbf{X}}) \right) \tilde{\mathbf{K}}^{-1} \frac{\partial \mathbf{k}_x}{\partial \mathbf{x}} + \mathbf{J}_0(\mathbf{x}), \quad (3)$$

where

$$\tilde{\mathbf{K}}_{ij} = \Theta_\infty^L(\mathbf{x}_i, \mathbf{x}_j), \quad (\mathbf{k}_x)_i = \Theta_\infty^L(\mathbf{x}_i, \mathbf{x}).$$

3.3. Iterative Maps, Attractors, Associative Memory and Jacobian

We define attractors with respect to an iterative map. Notice that an autoencoder f is a map from \mathbb{R}^{n_0} to \mathbb{R}^{n_0} . Therefore, we can apply f iteratively to input \mathbf{x} . Formally, we define this sequence for any input \mathbf{x} , $\{f^k(\mathbf{x})\}_{k \in \mathbb{N}}$ where $f^k = \underbrace{f(\dots f(\mathbf{x}))}_k$.

Definition 1. A fixed point \mathbf{x}^* of map f ($f(\mathbf{x}^*) = \mathbf{x}^*$) is an *attractor* if there exists an open neighborhood of \mathbf{x}^* such that for any \mathbf{x} in this neighborhood, $\{f^k(\mathbf{x})\}_{k \in \mathbb{N}}$ converges to \mathbf{x}^* as $k \rightarrow \infty$. The set of all such points is called *basin of attraction* of \mathbf{x}^* .

Fixed points attractors can be used to implement associative memory (Hopfield, 1982). A memory clue sets the initial condition of the network dynamics, positioning the network state in a basin of attraction, and the actual memory is recapitulated by the attractor dynamics converging to the corresponding fixed point.

There is a well-know condition for a fixed point to be an attractor (Rudin et al., 1964).

Proposition 1. *A fixed point \mathbf{x}^* is an attractor of a differentiable map f if all eigenvalues of the Jacobian of f at \mathbf{x}^* are strictly less than 1 in absolute value.*

Therefore, for a point \mathbf{x} to be an attractor, we need two conditions: (1) $f(\mathbf{x}) = \mathbf{x}$ and (2) all the eigenvalues of $J(\mathbf{x})$ have norm strictly smaller than 1. Condition (1) can be justified in the NTK limit (Jacot et al., 2018; Du et al., 2018; Allen-Zhu et al., 2018) as long as $\tilde{\mathbf{K}}$ is positive definite. This is theoretically true if all data points live on a sphere, the network has non-polynomial Lipschitz activation function and $L \geq 2$ (c.f. Proposition 2 in (Jacot et al., 2018)). In practice, it is easy to achieve $f(\mathbf{x}) \approx \mathbf{x}$ in overparameterized networks. Therefore, our focus is on the second condition.

4. Theoretical Results

In this section, we present our theoretical results about the attractor behavior of iterated overparametrized autoencoders. We focus on two distinct settings:

1. In the first setting, we consider a neural network of any depth in the NTK limit with a single training example.
2. In the second setting, we focus on a two-layer network in the NTK limit and multiple training examples.

We first give two key results about the Jacobian norm at network initialization that will be used later. We use the operator norm (induced l_2 -norm) of Jacobian as a proxy to control the eigenvalue, because the norm is given by the largest singular value, which upper bounds the norm of the largest eigenvalue.

4.1. A Bound on the Norm of the Jacobian at Network Initialization

In this section, we first show that with high probability, the operator norm of the initial Jacobian for sigmoid networks drops with increasing depth. To prove this, we use a mathematical technique called the ϵ -net argument (Tao, 2012) (Theorem 1) valid for any activation function. We then argue that for sigmoid networks the upper bound of the initial Jacobian norm is concentrated around 0.5 for large n_0 . This explains the formation of attractors when the trained Jacobian norm stays close to initialization (c.f. Section 4.2).

We first prove a proposition to be used in the ϵ -net argument.

Given a unit vector $\mathbf{z}^{(0)}$ and input $\hat{\mathbf{x}}$, we recursively define $\tilde{\mathbf{z}}^{(L)} \equiv \mathbf{J}(\hat{\mathbf{x}})\mathbf{z}^{(0)}$:

$$\tilde{\mathbf{z}}^{(\ell)} \equiv \frac{1}{\sqrt{n^{(\ell-1)}}} \mathbf{W}^{(\ell-1)} \mathbf{z}^{(\ell-1)}, \quad \mathbf{z}^{(\ell)} \equiv \mathbf{D}^{(\ell)} \tilde{\mathbf{z}}^{(\ell)}.$$

The following proposition provides the distribution of $\mathbf{z}^{(\ell)}$.

Proposition 2. *For a fixed unit vector $\mathbf{z}^{(0)}$, fixed input data $\hat{\mathbf{x}}$ and a network of depth L at random initialization, with a Lipschitz nonlinearity σ , and in the limit $n_1, \dots, n_{L-1} \rightarrow \infty$, $\mathbf{J}(\hat{\mathbf{x}})\mathbf{z}^{(0)}$ has the following recursion with $\mathbf{z}_i^{(\ell)} = \hat{z}^{(\ell)}$:*

$$\hat{z}^{(1)} = \sigma'(a) b \quad (a, b) \sim \mathcal{N} \left(\mathbf{0}, \begin{bmatrix} \frac{\|\hat{\mathbf{x}}\|_2^2}{n_0}, & \frac{\hat{\mathbf{x}}^T \mathbf{z}^{(0)}}{n_0} \\ \frac{\hat{\mathbf{x}}^T \mathbf{z}^{(0)}}{n_0}, & \frac{\|\mathbf{z}^{(0)}\|_2^2}{n_0} \end{bmatrix} \right),$$

$$\hat{z}^{(\ell+1)} = \sigma'(a) b$$

$$(a, b) \sim \mathcal{N} \left(\mathbf{0}, \begin{bmatrix} \mathbb{E}[(\hat{\alpha}^{(\ell)})^2], & \mathbb{E}[\hat{\alpha}^{(\ell)} \hat{z}^{(\ell)}] \\ \mathbb{E}[\hat{\alpha}^{(\ell)} \hat{z}^{(\ell)}], & \mathbb{E}[(\hat{z}^{(\ell)})^2] \end{bmatrix} \right),$$

$$\tilde{\mathbf{z}}_i^{(L)} = \hat{z}^{(L)} \sim \mathcal{N} \left(0, \mathbb{E}[(\hat{z}^{(L-1)})^2] \right),$$

where

$$\hat{\alpha}^{(1)} = \sigma(a) \quad a \sim \mathcal{N} \left(0, \frac{\|\hat{\mathbf{x}}\|_2^2}{n_0} \right),$$

$$\hat{\alpha}^{(\ell+1)} = \sigma(a) \quad a \sim \mathcal{N} \left(0, \mathbb{E}[(\hat{\alpha}^{(\ell)})^2] \right).$$

Proof. See Appendix A. □

Using Proposition 2, we apply the ϵ -net argument (Tao, 2012) (Appendix A) to obtain a bound on the Jacobian operator norm.

Theorem 1. *For any data point \mathbf{x}_i , $i \in [1, \dots, n]$, with probability at least $1 - O(n)e^{-O(n_0)}$,*

$$\|\mathbf{J}(\mathbf{x}_i)\|_{op} \leq c\sqrt{n_0\tau}$$

where c is a constant and

$$\tau = \sup_{\mathbf{x}_i \in \tilde{\mathbf{X}}, \|\mathbf{z}^{(0)}\|_2=1} \mathbb{E}[(\hat{z}^{(L-1)})^2 | \mathbf{z}^{(0)}, \mathbf{x}_i]$$

Proof. See Appendix A. □

This bound on singular values of $\mathbf{J}(\mathbf{x})$ provides a way to understand how the Jacobian at network initialization changes with respect to activation functions and the number of layers. Specifically, for sigmoid activation function, we know that $\sigma'(x) \in (0, \frac{1}{4}]$ and for any arbitrary unit vector $\mathbf{z}^{(0)}$,

$$\mathbb{E}[(\hat{z}^{(L-1)})^2 | \mathbf{z}^{(0)}, \mathbf{x}] = \mathbb{E}[\sigma'(a)^2 b^2] \leq \frac{\mathbb{E}[(\hat{z}^{(L-2)})^2 | \mathbf{z}^{(0)}, \mathbf{x}]}{16}.$$

Thus,

$$\mathbb{E}[(\hat{z}^{(L-1)})^2 | \mathbf{z}^{(0)}, \mathbf{x}] \leq \frac{1}{n_0 16^{L-1}} \|\mathbf{z}^{(0)}\|_2^2 = \frac{1}{n_0 16^{L-1}},$$

and with probability at least $1 - O(n)e^{-O(n_0)}$,

$$\|\mathbf{J}(\mathbf{x}_i)\|_{op} \leq \frac{c}{4^{L-1}} \quad \forall i \in [n]. \quad (4)$$

Therefore, with high probability, the initial Jacobian norm decreases with increasing layers. Based on this result, we choose to study a 2-layer sigmoid network because it would give an upper bound on the initial Jacobian norm.

Next, we argue that the largest initial Jacobian norm for a two-layer sigmoid network is concentrated around $1/2$. For a given training point $\hat{\mathbf{x}}$,

$$\mathbf{J}(\hat{\mathbf{x}}) = \frac{1}{\sqrt{n_1}} \mathbf{W}^{(1)} \mathbf{D}^{(1)} \frac{1}{\sqrt{n_0}} \mathbf{W}^{(0)}.$$

From Proposition 2, we can reach maximum $\mathbb{E}[(\hat{z}^{L-1})^2 | \mathbf{z}^{(0)}, \hat{\mathbf{x}}]$ for any fixed $\mathbf{z}^{(0)}$ when $\hat{\mathbf{x}} = \mathbf{0}$ because every gradient of the hidden layer is at max $\frac{1}{4}$ and the covariance between \mathbf{x} and $\mathbf{z}^{(0)}$ is zero, which would informally suggest that the upper bound of initial Jacobian norm is likely achieved at $\hat{\mathbf{x}} = \mathbf{0}$ (Another way to view this comes from Proposition 2 where for any given unit vector $\mathbf{z}^{(0)}$, the norm of output vector $\hat{\mathbf{z}}_i^{(2)}$ has a higher mean when $\hat{\mathbf{x}} = \mathbf{0}$). This gives us the simplification:

$$\mathbf{J}(\hat{\mathbf{x}}) = \frac{1}{4} \frac{1}{\sqrt{n_0}} \frac{1}{\sqrt{n_1}} \mathbf{W}^{(1)} \mathbf{W}^{(0)}.$$

Observe that $\mathbf{W} = \frac{1}{\sqrt{n_1}} \mathbf{W}^{(1)} \mathbf{W}^{(0)}$ is a Gaussian random matrix with $n_1 \rightarrow \infty$ where each entry is i.i.d. and the largest singular value of $\frac{1}{\sqrt{n_0}} \mathbf{W}$ is concentrated at 2 for large n_0 (Vershynin, 2012). Consequently, the largest initial Jacobian norm for a sigmoid network is concentrated around $1/2$.

4.2. Training a Multilayer Network with a Single Example

In this section, we consider the special case when there is only one training example, \mathbf{x}_1 . We show that under certain conditions the Jacobian stays close to initialization, and, combined with the result from the previous section, the trained network can form attractors.

We start by analyzing the NTK solution. Note that in this case (2) can be simplified to:

$$f_\infty(\mathbf{x}) = \frac{\Theta_\infty^L(\mathbf{x}, \mathbf{x}_1)}{\Theta_\infty^L(\mathbf{x}_1, \mathbf{x}_1)} (\mathbf{x}_1 - f_0(\mathbf{x}_1)) + f_0(\mathbf{x}).$$

As we describe below, $\Theta_\infty^{(L)}$ tends to a constant kernel as $L \rightarrow \infty$ for some network initialization (Hayou et al., 2019) and therefore the trained Jacobian equals the initial one.

To see this, first, we pay close attention to the covariance matrix $\Sigma^{(L)}$, which is the building block of $\Theta_\infty^{(L)}$. We define

$q_{ab}^{(\ell)} \equiv \Sigma^\ell(\mathbf{x}_a, \mathbf{x}_b)$ and $c_{ab}^{(\ell)} \equiv q_{ab}^{(\ell)} / \sqrt{q_{aa}^{(\ell)} q_{bb}^{(\ell)}}$. It can be shown that $c^* = 1$ is a fixed point of $c_{ab}^{(\ell)}$ as $\ell \rightarrow \infty$ (Poole et al., 2016). For sigmoidal networks, the stability of c^* is governed by

$$\chi_1 \equiv \left. \frac{\partial c_{ab}^{(\ell)}}{\partial c_{ab}^{(\ell-1)}} \right|_{c=1} = \mathbb{E}[\dot{\sigma}(\sqrt{q^*} z)^2], \quad z \sim \mathcal{N}(0, 1).$$

where q^* is what $q_{aa}^{(\ell)}$ converges to. If $\chi_1 < 1$, c^* is a stable fixed point, suggesting that all points become equally similar as they progress through layers. A network under such initialization is said to be in the *ordered region* (Schoenholz et al., 2016). For such networks, $\Theta_\infty^{(\ell)}$ converges to a constant kernel with increasing layers (Xiao et al., 2019; Hayou et al., 2019). As $L \rightarrow \infty$, $\partial_{\mathbf{x}} \Theta_\infty^{(L)}(\mathbf{x}, \mathbf{x}_1)|_{\mathbf{x}=\mathbf{x}_1} \rightarrow 0$ and,

$$\mathbf{J}_\infty(\mathbf{x}_1) \stackrel{L \rightarrow \infty}{=} \mathbf{J}_0(\mathbf{x}_1),$$

as long as $\Theta_\infty^L(\mathbf{x}_1, \mathbf{x}_1)$ does not converge to 0 as $L \rightarrow \infty$. We note that this argument applies to other weight and bias variance scaling factors at initialization than the special case (1 and 0 respectively) we focus on.

Specializing to sigmoid networks, we first observe that they are in the ordered region because $\dot{\sigma}(x) \in (0, \frac{1}{4}]$, giving an upper bound on $\chi_1 \leq 1/16$. Further the lower bound of $\Theta_\infty^{(L)}(\mathbf{x}_1, \mathbf{x}_1)$ is $1/4$ (Lemma 5 in Appendix B). Therefore, we expect the Jacobian to be constant during training in the large depth limit. In practice, the trained Jacobian stays close to initialization for 2 - 3 layer sigmoid networks as shown in Section 5.2.

Our analysis can explain the empirical results of (Zhang et al., 2019) that fully connected networks trained with single example tend to learn constant functions, leading to memorization. To see this, remember that for a sigmoid network the norm of the initial Jacobian falls with increasing the number of layers, eq. (4). Other sigmoidal activation functions may show the same behavior. In the same limit, the Jacobian remains constant during training, implying $f_\infty(\mathbf{x})$ is approximately constant around \mathbf{x}_1 .

The large-depth analysis presented in this section does not carry over to multiple training examples. The limiting constant kernel is singular and fails in training the network to zero loss (Jacot et al., 2018). Instead, we will study 2-layer sigmoid networks for multiple training examples.

4.3. Training a 2-Layer Network with Multiple Examples: Linear Region

Next, we consider our second setting of a 2-layer network trained with multiple examples. In this section, we argue that for small input norm r , a network in the NTK regime behaves like a linear network resulting in a Jacobian with eigenvalue 1, and thus may not form associative memory.

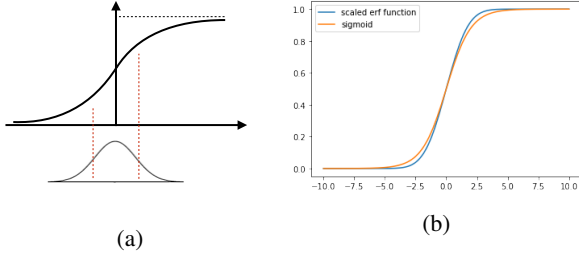


Figure 1: (a) Linear Region Illustration for Sigmoid: The top graph is a sigmoid function while the bottom one is a Gaussian distribution. The majority of Gasussian distribution falls in the linear region of sigmoid. (b) Sigmoid vs Rescaled Erf.

To see this, first note that for a 2-layer network and any two training points \mathbf{x}_i and \mathbf{x}_j , NTK can be written as

$$\Theta_{\infty}^2(\mathbf{x}_i, \mathbf{x}_j) = \Sigma^1(\mathbf{x}_i, \mathbf{x}_j) \dot{\Sigma}^2(\mathbf{x}_i, \mathbf{x}_j) + \Sigma^2(\mathbf{x}_i, \mathbf{x}_j),$$

where $\Sigma^1(\mathbf{x}_i, \mathbf{x}_j) = \frac{\mathbf{x}_i^T \mathbf{x}_j}{n_0}$ and both Σ^2 and $\dot{\Sigma}^2$ are expectations governed by Σ^1 . Because Σ^1 is a covariance matrix defined by inner products, small input norms means small variance for the Gaussian process and the expectation is mostly concentrated to the linear region of the activation. This justifies a linear approximation to activation function. Figure 1a shows the sigmoid, which can be approximated linearly around $x = 0$ as $\sigma(x) \approx \frac{1}{4}x + \frac{1}{2}$.

Most of activation functions have similar linear behavior though the range in which this approximation is accurate may differ. Thus, we focus on an arbitrary linear activation function $\alpha x + \beta$ which leads to an initial Jacobian $J_0(\mathbf{x}) = \alpha \frac{1}{\sqrt{n_1}} \frac{1}{\sqrt{n_0}} \mathbf{W}^{(1)} \mathbf{W}^{(0)}$ with initial weights $\mathbf{W}^{(1)}$ and $\mathbf{W}^{(0)}$. We will return to the discussion of sigmoid at the end of this section. For simplicity, we also assume that $\dot{\mathbf{X}}$ is full rank and $n_0 \geq n$ as memory tasks tend to have high dimension inputs like images and audios (Radhakrishnan et al., 2019).

We start examining the Jacobian with the easiest case, $\sigma(x) = \alpha x$ and $n = n_0$, which leads to $\mathbf{J}_{\infty}(\mathbf{x}) = \mathbf{I}_{n_0}$.

Lemma 2. *Suppose there is a 2-layer network. If the activation function is $\sigma(x) = \alpha x$, $n = n_0$ and the data matrix is full rank. Then at NTK limit, $\mathbf{J}_{\infty}(\mathbf{x}) = \mathbf{I}_{n_0}$.*

Proof. See Appendix C.1. \square

In fact, the multiplicity of eigenvalue 1 is directly related to n under certain conditions.

Lemma 3. *Suppose there is a 2-layer network with activation function $\sigma(x) = \alpha x$ and given initial weights $\mathbf{W}^{(1)} \in \mathbb{R}^{n_0 \times n_1}$, $\mathbf{W}^{(0)} \in \mathbb{R}^{n_1 \times n_0}$. If the data matrix is full rank with $n \leq n_0$, then, at the NTK limit ($n_1 \rightarrow \infty$),*

$\mathbf{J}_{\infty}(\mathbf{x})$ has eigenvalue 1 with multiplicity at least n . If at the NTK limit, α is chosen such that $\|\mathbf{J}_0(\mathbf{x})\|_{op} < 1$, then the multiplicity is exactly n and 1 is the largest eigenvalue norm.

Proof. See Appendix C.1. \square

Remark 1. Lemma 3 suggests that a network trained with a single example at convergence can have Jacobian eigenvalue 1 for $\sigma(x) = \alpha x$ regardless of initial Jacobians. This result may seem to contradict Section 4.2. However, in this case, the argument in Section 4.2 is violated because the diagonal value of $\Theta_{\infty}^{(2)}$ is also converging to zero as when $\alpha < 1$, the linear activation has a shrinking effect on the layer outputs.

The result can be naturally extended to $\sigma(x) = \alpha x + \beta$ with $\beta > 0$.

Lemma 4. *Suppose there is a 2-layer network with activation function $\sigma(x) = \alpha x + \beta$, given initial weights $\mathbf{W}^{(1)} \in \mathbb{R}^{n_0 \times n_1}$, $\mathbf{W}^{(0)} \in \mathbb{R}^{n_1 \times n_0}$ and every data point has the same norm r (i.e. $\forall i \in [n] \|\mathbf{x}\|_2 = r$). If the data matrix is full rank with $n \leq n_0$, then, at the NTK limit $n_1 \rightarrow \infty$, $\mathbf{J}_{\infty}(\mathbf{x})$ has eigenvalues 1 with multiplicity at least $n - 1$. If at the NTK limit, α and β are chosen such that*

$$\|\mathbf{J}_0(\mathbf{x})\|_{op} = 1 - \Delta, \quad \left\| \frac{1}{\sqrt{n_1}} \mathbf{W}^{(1)} \mathbf{1}_{n_1} \right\|_2 < \frac{\beta n_0 \Delta}{2r\alpha^2},$$

where $0 < \Delta \leq 1$, then the multiplicity is exactly $n - 1$ and 1 is the largest eigenvalue norm.

Proof. See Appendix C.1. \square

Collectively we proved some sufficiency conditions for the largest trained Jacobian eigenvalue to be 1. We emphasize that if the largest eigenvalue is 1, and not strictly below 1, the network can fail to form attractors.

Now, we return to sigmoid networks and discuss the implications of Lemma 4. In the linear region, we have $\alpha = \frac{1}{4}$ and $\beta = \frac{1}{2}$. As mentioned in Section 4.1, in the NTK limit, $\|\mathbf{J}_0(\mathbf{x})\|_2 \approx \frac{1}{2}$, implying $\Delta \approx \frac{1}{2}$. On the other hand, $\left\| \frac{1}{\sqrt{n_1}} \mathbf{W}^{(1)} \mathbf{1}_{n_1} \right\|_2$ is the norm of a standard Gaussian vector which follows the chi distribution, and it is concentrated around $\sqrt{n_0}$. Then, the conditions in Lemma 4 hold with high probability if $r < 2\sqrt{n_0}$. Attractor formation can fail in sigmoid networks for small input norms, which we observe in simulations (Section 5).

4.4. Training a 2-Layer Network with Multiple Examples: Beyond Linear Region and a Transition to Attractor Formation

Section 4.3 suggests that small r leads to linear behavior and networks may fail to form associative memory. In this

section, we explore larger norm inputs where a network in the NTK regime utilizes the non-linear region of the sigmoid. We will argue that in the large norm limit all Jacobian eigenvalues' norms are below 1. Taking into account our result that attractor formation may fail for small values of r , we identify a transition with increasing r from a regime where memory formation does not occur to a regime where it occurs. Our results can be adapted for other sigmoidal functions.

For simplicity, we further assume that there are no parallel inputs in the training data (no \mathbf{x} and $-\mathbf{x}$ at the same time). This is not a hard constraint and an analysis with parallel inputs is given in Appendix D.3.

Our strategy is to calculate the Jacobian in the large norm limit using Equation (3). We first note for large r , $\mathbf{X}_{ij} \gg f_0(\mathbf{X})_{ij}$ because all the hidden units in the network is between 0 and 1. Therefore Equation 3 can be approximated by

$$\mathbf{J}_\infty(\mathbf{x}) \approx \hat{\mathbf{X}} \tilde{\mathbf{K}}^{-1} \frac{\partial \mathbf{k}_x}{\partial \mathbf{x}} + \mathbf{J}_0(\mathbf{x}) \quad (5)$$

The items of interest here are $\tilde{\mathbf{K}}$ and $\frac{\partial \mathbf{k}_x}{\partial \mathbf{x}}$, for which we need to calculate the NTK and its gradient. We first define an approximation of the NTK and then use it to estimate the Jacobian.

Approximation of the NTK: Approximating sigmoid function (σ_s) by erf function there $\sigma_s(x) \approx \frac{1}{2} \text{erf}(\frac{1}{2}x) + \frac{1}{2}$ (shown in Figure 1b) allows us to use a known closed form solution for NTK (Lee et al., 2019; Williams, 1997). With this approximation, 2-layer NTK can be written as follows (full derivation can be found in Appendix D.1):

$$\begin{aligned} \Theta_\infty^{(2)}(\hat{\mathbf{x}}, \mathbf{x}) &\approx \frac{1}{2\pi} \frac{\hat{\mathbf{x}}^T \mathbf{x}}{\sqrt{(2n_0 + \mathbf{x}^T \mathbf{x})(2n_0 + \hat{\mathbf{x}}^T \hat{\mathbf{x}}) - (\hat{\mathbf{x}}^T \mathbf{x})^2}} \\ &+ \frac{1}{2\pi} \arcsin\left(\frac{\hat{\mathbf{x}}^T \mathbf{x}}{\sqrt{(\mathbf{x}^T \mathbf{x} + 2n_0)(\hat{\mathbf{x}}^T \hat{\mathbf{x}} + 2n_0)}}\right) + \frac{1}{4}. \end{aligned} \quad (6)$$

In order to calculate the gradient of NTK, we define:

$$\mathcal{T}(\Sigma, \sigma_1, \sigma_2)(\mathbf{x}_i, \mathbf{x}_j) \equiv \mathbb{E}_{f \sim \mathcal{N}(0, \Sigma)}[\sigma_1(f(\mathbf{x}_i))\sigma_2(f(\mathbf{x}_j))]$$

Using this definition, the NTK's gradient for a 2-layer network with arbitrary activation function σ is given by,

$$\begin{aligned} \frac{\partial \Theta_\infty^{(2)}(\hat{\mathbf{x}}, \mathbf{x})}{\partial \mathbf{x}} &= \frac{1}{n_0^2} \hat{\mathbf{x}}^T \mathbf{x} \mathcal{T}(\Sigma^1, \ddot{\sigma}, \ddot{\sigma}) \hat{\mathbf{x}} + \frac{2}{n_0} \mathcal{T}(\Sigma^1, \dot{\sigma}, \dot{\sigma}) \hat{\mathbf{x}} \\ &+ \frac{1}{n_0^2} \hat{\mathbf{x}}^T \mathbf{x} \mathcal{T}(\Sigma^1, \dot{\sigma}, \ddot{\sigma}) \mathbf{x} + \frac{1}{n_0} \mathcal{T}(\Sigma^1, \sigma, \ddot{\sigma}) \mathbf{x}. \end{aligned}$$

Now we can use these expressions to calculate $\tilde{\mathbf{K}}$ and $\frac{\partial \mathbf{k}_x}{\partial \mathbf{x}}$ in Equation 5.

Calculation of $\tilde{\mathbf{K}}$ and $\frac{\partial \mathbf{k}_x}{\partial \mathbf{x}}$: We first look at Equation (6), using the fact that because all the data points have the same norm $\mathbf{x}_i^T \mathbf{x}_j = r^2 \rho_{i,j}$ with $|\rho_{i,j}| < 1$ (since there are no parallel inputs):

$$\begin{aligned} \tilde{\mathbf{K}}_{ij} &= \frac{1}{4} + \frac{1}{2\pi} \arcsin\left(\frac{r^2 \rho_{i,j}}{r^2 + 2n_0}\right) \\ &+ \frac{1}{2\pi} \frac{r^2 \rho_{i,j}}{\sqrt{(2n_0 + r^2)^2 - r^4 \rho_{i,j}^2}}. \end{aligned}$$

To get insight into the behavior of $\tilde{\mathbf{K}}_{ij}$, let's consider diagonal and off-diagonal entries separately in the large- r limit. First the off-diagonals:

$$\tilde{\mathbf{K}}_{ij} \approx \frac{1}{4} + \frac{1}{2\pi} \arcsin(\rho_{i,j}) + \frac{1}{2\pi} \frac{\rho_{i,j}}{\sqrt{1 - \rho_{i,j}^2}}.$$

Next we look at the diagonals:

$$\tilde{\mathbf{K}}_{ii} \approx \frac{r}{4\pi\sqrt{n_0}}.$$

The diagonal grows linearly with r and dominates over the off-diagonals. The kernel tends to a scaled identity matrix in the $r \rightarrow \infty$ limit. Thus,

$$\left\| \tilde{\mathbf{K}}^{-1} \right\|_{op} \approx \frac{4\pi\sqrt{n_0}}{r}. \quad (7)$$

Next, we look at $\frac{\partial \mathbf{k}_x}{\partial \mathbf{x}}$. As we are interested in trained Jacobian at training examples, without loss of generality, we focus on $\mathbf{J}_\infty(\mathbf{x}_1)$. We already know that $\tilde{\mathbf{K}}$ tends to converge to a diagonal matrix with large r . One could use this result to suggest that $\frac{\partial \Theta_\infty^{(2)}(\mathbf{x}_i, \mathbf{x})}{\partial \mathbf{x}}|_{\mathbf{x}_j} \approx 0$ if $i \neq j$. In fact, with large r , it can be shown that (Appendix D.2):

$$\frac{\partial \mathbf{k}_x}{\partial \mathbf{x}} \Big|_{\mathbf{x}_1} \approx \left[\frac{\partial \Theta_\infty^{(2)}(\mathbf{x}_1, \mathbf{x})}{\partial \mathbf{x}} \Big|_{\mathbf{x}_1}, \mathbf{0}, \dots, \mathbf{0} \right]^T \quad (8)$$

and

$$\left\| \frac{\partial \mathbf{k}_x}{\partial \mathbf{x}} \right\|_{op} \approx \frac{1}{8\pi\sqrt{n_0}}. \quad (9)$$

Trained Jacobian and attractor formation: Finally, we can use our results in Equations (7) and (9) to calculate the trained Jacobian using Equation (5) with large r :

$$\begin{aligned} \|\mathbf{J}_\infty(\mathbf{x})\|_{op} &\approx \left\| \hat{\mathbf{X}} \tilde{\mathbf{K}}^{-1} \frac{\partial \mathbf{k}_x}{\partial \mathbf{x}} + \mathbf{J}_0(\mathbf{x}) \right\|_{op} \\ &\leq \left\| \hat{\mathbf{X}} \tilde{\mathbf{K}}^{-1} \frac{\partial \mathbf{k}_x}{\partial \mathbf{x}} \right\|_{op} + \|\mathbf{J}_0(\mathbf{x})\|_{op} \\ &\approx r \frac{4\pi\sqrt{n_0}}{r} \frac{1}{8\pi\sqrt{n_0}} + \|\mathbf{J}_0(\mathbf{x})\|_{op} \\ &= \frac{1}{2} + \|\mathbf{J}_0(\mathbf{x})\|_{op}. \end{aligned}$$

To go from the second to the third row of the equation above, we used the fact that each column of \mathbf{X} has norm r , and in the limit $\tilde{\mathbf{K}}^{-1}$ is a scaled identity matrix and $\frac{\partial \mathbf{k}_x}{\partial \mathbf{x}}$ has only one non-zero row and that row is a scaled \mathbf{x}_1 (Appendix D.2).

To see the implications of this result, we observe that as $r \rightarrow \infty$, $\mathbf{J}_0(\mathbf{x}) \approx \mathbf{0}$ because the pre-activation of the two layer network will be infinitely big and the units saturate leading to zero gradients in the Jacobian matrix, we can conclude in this limit

$$\|\mathbf{J}_\infty(\mathbf{x})\|_{op} \leq 1/2.$$

Combined with our previous result that attractor formation may fail for small values of r , this suggests that there is a transition from a region where associative memory formation fails to region where it succeeds with increasing r . We show this transition in simulations in Section 5 and verify that the largest eigenvalue of the Jacobian falls towards $1/2$ asymptotically.

4.5. Beyond NTK

In practice, large values of r often require a larger width to stay in NTK. If the large width condition is violated, trained weights will not stay close to initialization. However, in this case, a significant deviation of weights from initialization caused by gradient descent may saturate hidden units, leading to zero gradients in the Jacobian matrix. As shown in Figure 3, for fixed hidden size, larger input radius leads to near zero eigenvalues when the NTK conditions no longer hold. Therefore, the network can behave as if only the last layer is trained with saturated hidden features and close to zero Jacobian norm since all the $\mathbf{D}^{(\ell)}$ matrices have mostly zero diagonal entries. Note that, in overparameterized networks, this type of optimization can often result in zero training loss as well due to the over-parameterization of the last layer.

5. Simulations

5.1. Experiment Setup

Training and iterative convergence criteria: Training is stopped when the training loss of the auto-encoder drops below a threshold, which we chose to be 10^{-7} . Iterative convergence happens when a non-fixed point converges to a fixed point measured by mean-squared-error after passing through the trained autoencoder iteratively. This threshold was 10^{-2} .

Implementation details: We used vanilla gradient descent with learning rate 1, similar to (Jacot et al., 2018). The code is implemented with Pytorch. For each setting, we ran experiments 100 times to get 100 sets of samples. For all experiments except experiments on MNIST, the samples are

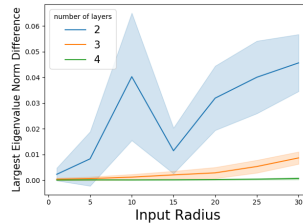


Figure 2: Difference of the largest eigenvalue norms at initialization and after training, i.e. $\|\lambda_1(\mathbf{J}_0)\| - \|\lambda_1(\mathbf{J}_\infty)\|$.

randomly generated.

5.2. Single Training Example

We first present experiments for a single training example. The sigmoid network is chosen to have hidden dimension 1000 with input dimension 32. Figure 2 shows that the difference between the largest eigenvalue norms at initialization and after training decreases with more number of layers.

5.3. Multiple Training Examples

Linear Region: In this section, we first illustrate the eigenvalue distribution in the linear region by sampling unit vectors as data ($r = 1$). Here, we trained 2 layer sigmoid networks with input dimension 10 and hidden size 1000 for 2, 5 and 8 training points. As suggested by Lemma 4, there should be $n - 1$ eigenvalues with norm around 1. This is supported by Figure E.1 in the Appendix, where 10%, 40% and 70% of the eigenvalues are near 1. Here, the presence of eigenvalue 1 indicates that the network operates in the linear region.

Beyond Linear Region: We demonstrate how the largest eigenvalue norm varies with the input radius. We test with various hidden dimensions to achieve the NTK limit on input dimension 32 with the number of training data 5, 20, and 40. Only experiments that can be trained to have loss below 10^{-7} or fit into single Titan V GPU are included.

The general trend, as shown in Figure 3, is that as we move away from the linear region, the largest eigenvalue norm will drop, and as we keep increasing the hidden layer size, it will move close to the $1/2$ limit as suggested by our analysis. It is also worth noting that the input radius needs to be increased with large number of training examples to get training loss below 10^{-7} under reasonable iterations, implying the capacity of the networks is controlled by input radius and agreeing with the intuition and theoretical results from (Allen-Zhu et al., 2018), that there needs to be some level of separation between data points to train the network.

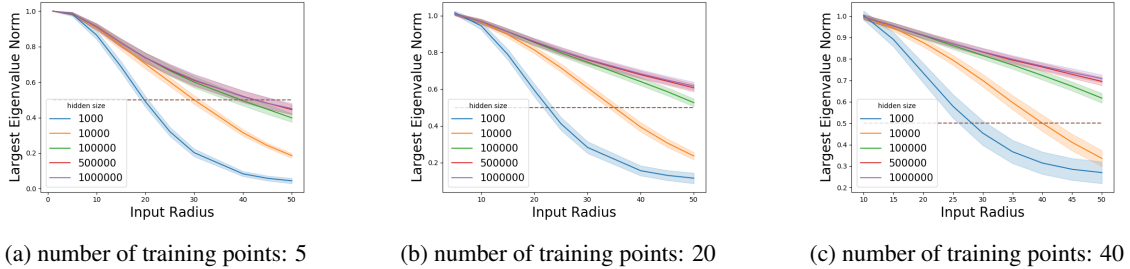


Figure 3: Largest eigenvalue norm vs input norm: input dimension 32.

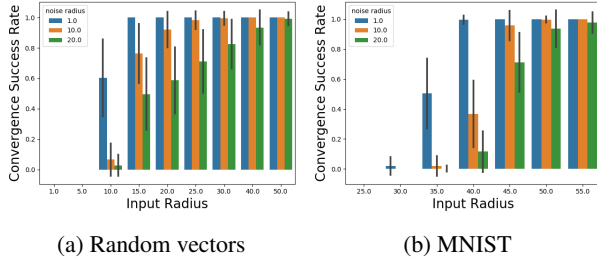


Figure 4: Convergence Success Rate vs Input Radius: 5 training examples

5.4. Basin of Attraction

We test basin of attraction by adding Gaussian noise to training examples and check if the modified examples can converge to the original ones via iterative maps under 50 iterations. And the convergence rate is the number of samples that could be successfully recovered. The standard deviation of the Gaussian noise is called the noise radius. The network has two layers with hidden size 10000 and input dimension 32. Not surprisingly, Figure 4a shows that larger input norm gives greater basin of attraction. More experiments can be found in Appendix E.2.

5.5. MNIST Data

We also test basin of attraction experiments on MNIST dataset to check if we can recover real training examples. The images are preprocessed by subtracting means and rescaled to have different input norms for testing. Similar to the setting before, Figure 4b also shows that larger input norm gives greater basin of attraction. Notice that because MNIST images have large input dimension, they need larger radius to move out of the linear region. More experiments can be found in Appendix E.3.

5.6. Sigmoidal Activation

Finally, we show that our results can be extended to different sigmoidal activation functions as well. We chose two-layer network with hidden size 10000 and input dimension 32

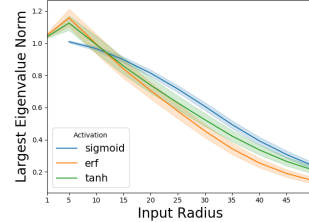


Figure 5: Input Radius and Eigenvalue Norm Curve for Different Activation Functions

and 20 training examples. As before, only settings that led to a training loss below 10^{-7} are included. Figure 5 clearly shows that all the activation functions share similar curves. Notice that both tanh and erf have large eigenvalue when r is small. This is not a contradiction to our Lemma 3 as their $\alpha = \dot{\sigma}(0)$ is too large to satisfy the conditions in Lemma 3. To further verify this result, we also include how the histogram of eigenvalue norm changes for those activations in the Appendix E.4.

6. Conclusion

In this paper, we theoretically and empirically show that training overparameterized sigmoid autoencoders can lead to attractors for a single training example and multiple training examples in the non-linear region, with the help of theories developed in the NTK limit. We identified a behavior change governed by the input radius. Some future directions include generalizing our results to other activations, identifying other factors that can determine whether autoencoders can learn to have attractors or not, and how the formations of attractors are related to generalization in deep learning.

Acknowledgements

C. Pehlevan thanks the Harvard Data Science Initiative, Google and Intel for support.

References

- Allen-Zhu, Z. and Li, Y. What can resnet learn efficiently, going beyond kernels? In *Advances in Neural Information Processing Systems*, pp. 9015–9025, 2019.
- Allen-Zhu, Z., Li, Y., and Song, Z. A convergence theory for deep learning via over-parameterization. *arXiv preprint arXiv:1811.03962*, 2018.
- Allen-Zhu, Z., Li, Y., and Liang, Y. Learning and generalization in overparameterized neural networks, going beyond two layers. In *Advances in neural information processing systems*, pp. 6155–6166, 2019.
- Amit, D. J. and Treves, A. Associative memory neural network with low temporal spiking rates. *Proceedings of the National Academy of Sciences*, 86(20):7871–7875, 1989.
- Arora, S., Du, S. S., Hu, W., Li, Z., Salakhutdinov, R. R., and Wang, R. On exact computation with an infinitely wide neural net. In *Advances in Neural Information Processing Systems*, pp. 8139–8148, 2019.
- Bai, Y. and Lee, J. D. Beyond linearization: On quadratic and higher-order approximation of wide neural networks. *arXiv preprint arXiv:1910.01619*, 2019.
- Bordelon, B., Canatar, A., and Pehlevan, C. Spectrum dependent learning curves in kernel regression and wide neural networks. *arXiv preprint arXiv:2002.02561*, 2020.
- Cao, Y. and Gu, Q. Generalization bounds of stochastic gradient descent for wide and deep neural networks. In *Advances in Neural Information Processing Systems*, pp. 10835–10845, 2019.
- Chizat, L., Oyallon, E., and Bach, F. On lazy training in differentiable programming. In *Advances in Neural Information Processing Systems*, pp. 2933–2943, 2019.
- Cohen, O., Malka, O., and Ringel, Z. Learning curves for deep neural networks: a gaussian field theory perspective. *arXiv preprint arXiv:1906.05301*, 2019.
- Du, S. S., Lee, J. D., Li, H., Wang, L., and Zhai, X. Gradient descent finds global minima of deep neural networks. *arXiv preprint arXiv:1811.03804*, 2018.
- Hayou, S., Doucet, A., and Rousseau, J. Mean-field behaviour of neural tangent kernel for deep neural networks. *arXiv preprint arXiv:1905.13654*, 2019.
- Hertz, J. A. *Introduction to the theory of neural computation*. CRC Press, 2018.
- Hopfield, J. J. Neural networks and physical systems with emergent collective computational abilities. *Proceedings of the national academy of sciences*, 79(8):2554–2558, 1982.
- Jacot, A., Gabriel, F., and Hongler, C. Neural tangent kernel: Convergence and generalization in neural networks. In *Advances in neural information processing systems*, pp. 8571–8580, 2018.
- Ji, Z. and Telgarsky, M. Risk and parameter convergence of logistic regression. *arXiv preprint arXiv:1803.07300*, 2018.
- Kolar, M. and Liu, H. Marginal regression for multitask learning. In *Artificial Intelligence and Statistics*, pp. 647–655, 2012.
- Krizhevsky, A., Sutskever, I., and Hinton, G. E. Imagenet classification with deep convolutional neural networks. In *Advances in neural information processing systems*, pp. 1097–1105, 2012.
- Lee, J., Bahri, Y., Novak, R., Schoenholz, S. S., Pennington, J., and Sohl-Dickstein, J. Deep neural networks as gaussian processes. *arXiv preprint arXiv:1711.00165*, 2017.
- Lee, J., Xiao, L., Schoenholz, S., Bahri, Y., Novak, R., Sohl-Dickstein, J., and Pennington, J. Wide neural networks of any depth evolve as linear models under gradient descent. In *Advances in neural information processing systems*, pp. 8570–8581, 2019.
- Matthews, A. G. d. G., Rowland, M., Hron, J., Turner, R. E., and Ghahramani, Z. Gaussian process behaviour in wide deep neural networks. *arXiv preprint arXiv:1804.11271*, 2018.
- Miller, K. S. On the inverse of the sum of matrices. *Mathematics magazine*, 54(2):67–72, 1981.
- Oymak, S. and Soltanolkotabi, M. Overparameterized nonlinear learning: Gradient descent takes the shortest path? *arXiv preprint arXiv:1812.10004*, 2018.
- Pennington, J. and Worah, P. Nonlinear random matrix theory for deep learning. In *Advances in Neural Information Processing Systems*, pp. 2637–2646, 2017.
- Pennington, J., Schoenholz, S., and Ganguli, S. Resurrecting the sigmoid in deep learning through dynamical isometry: theory and practice. In *Advances in neural information processing systems*, pp. 4785–4795, 2017.
- Pennington, J., Schoenholz, S. S., and Ganguli, S. The emergence of spectral universality in deep networks. *arXiv preprint arXiv:1802.09979*, 2018.

- Poole, B., Lahiri, S., Raghu, M., Sohl-Dickstein, J., and Ganguli, S. Exponential expressivity in deep neural networks through transient chaos. In *Advances in neural information processing systems*, pp. 3360–3368, 2016.
- Radhakrishnan, A., Yang, K., Belkin, M., and Uhler, C. Memorization in overparameterized autoencoders. *arXiv preprint arXiv:1810.10333*, 2018.
- Radhakrishnan, A., Belkin, M., and Uhler, C. Overparameterized neural networks can implement associative memory. *arXiv preprint arXiv:1909.12362*, 2019.
- Rudin, W. et al. *Principles of mathematical analysis*, volume 3. McGraw-hill New York, 1964.
- Schoenholz, S. S., Gilmer, J., Ganguli, S., and Sohl-Dickstein, J. Deep information propagation. *arXiv preprint arXiv:1611.01232*, 2016.
- Soudry, D., Hoffer, E., Nacson, M. S., Gunasekar, S., and Srebro, N. The implicit bias of gradient descent on separable data. *The Journal of Machine Learning Research*, 19(1):2822–2878, 2018.
- Tao, T. *Topics in random matrix theory*, volume 132. American Mathematical Soc., 2012.
- Vershynin, R. Introduction to the non-asymptotic analysis of random matrices. *Compressed Sensing*, pp. 210–268, 2012.
- Williams, C. K. Computing with infinite networks. In *Advances in neural information processing systems*, pp. 295–301, 1997.
- Xiao, L., Pennington, J., and Schoenholz, S. S. Disentangling trainability and generalization in deep learning. *arXiv preprint arXiv:1912.13053*, 2019.
- Zhang, C., Bengio, S., Hardt, M., and Singer, Y. Identity crisis: Memorization and generalization under extreme overparameterization. *arXiv preprint arXiv:1902.04698*, 2019.
- Zou, D. and Gu, Q. An improved analysis of training overparameterized deep neural networks. In *Advances in Neural Information Processing Systems*, pp. 2053–2062, 2019.

A. Proofs for Sec 4.1

Proposition 2. For a fixed unit vector $\mathbf{z}^{(0)}$, fixed input data $\hat{\mathbf{x}}$ and a network of depth L at random initialization, with a Lipschitz nonlinearity σ , and in the limit $n_1, \dots, n_{L-1} \rightarrow \infty$, $\mathbf{J}(\hat{\mathbf{x}})\mathbf{z}^{(0)}$ has the following recursion with $\mathbf{z}_i^{(\ell)} = \hat{z}^{(\ell)}$:

$$\begin{aligned}\hat{z}^{(1)} &= \sigma'(a)b \quad (a, b) \sim \mathcal{N}\left(\mathbf{0}, \begin{bmatrix} \frac{\|\hat{\mathbf{x}}\|_2^2}{n_0}, & \frac{\hat{\mathbf{x}}^T \mathbf{z}^{(0)}}{n_0} \\ \frac{\hat{\mathbf{x}}^T \mathbf{z}^{(0)}}{n_0}, & \frac{\|\mathbf{z}^{(0)}\|_2^2}{n_0} \end{bmatrix}\right), \\ \hat{z}^{(\ell+1)} &= \sigma'(a)b \\ (a, b) &\sim \mathcal{N}\left(\mathbf{0}, \begin{bmatrix} \mathbb{E}[(\hat{\alpha}^{(\ell)})^2], & \mathbb{E}[\hat{\alpha}^{(\ell)} \hat{z}^{(\ell)}] \\ \mathbb{E}[\hat{\alpha}^{(\ell)} \hat{z}^{(\ell)}], & \mathbb{E}[(\hat{z}^{(\ell)})^2] \end{bmatrix}\right), \\ \hat{z}_i^{(L)} &= \hat{z}^{(L)} \sim \mathcal{N}\left(0, \mathbb{E}[(\hat{z}^{(L-1)})^2]\right),\end{aligned}$$

where

$$\begin{aligned}\hat{\alpha}^{(1)} &= \sigma(a) \quad a \sim \mathcal{N}\left(0, \frac{\|\hat{\mathbf{x}}\|_2^2}{n_0}\right), \\ \hat{\alpha}^{(\ell+1)} &= \sigma(a) \quad a \sim \mathcal{N}\left(0, \mathbb{E}[(\hat{\alpha}^{(\ell)})^2]\right).\end{aligned}$$

Proof. We will prove this by induction for $\ell = 1, \dots, L-1$.

Basic Step

$$\mathbf{z}_i^{(1)} = \sigma' \left(\frac{1}{\sqrt{n_0}} (\mathbf{W}_i^{(0)})^T \hat{\mathbf{x}} \right) \frac{1}{\sqrt{n_0}} (\mathbf{W}_i^{(0)})^T \mathbf{z}^{(0)}$$

Notice that $\mathbf{W}_i^{(0)} \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_{n_0})$. Thus, we have the following:

$$\begin{aligned}a &= \frac{1}{\sqrt{n_0}} (\mathbf{W}_i^{(0)})^T \hat{\mathbf{x}} \sim \mathcal{N}\left(0, \frac{\|\hat{\mathbf{x}}\|_2^2}{n_0}\right) \\ b &= \frac{1}{\sqrt{n_0}} (\mathbf{W}_i^{(0)})^T \mathbf{z}^{(0)} \sim \mathcal{N}\left(0, \frac{\|\mathbf{z}^{(0)}\|_2^2}{n_0}\right)\end{aligned}$$

a and b are not independent:

$$\mathbb{E}[ab] = \mathbb{E}\left[\left(\frac{1}{\sqrt{n_0}} (\mathbf{W}_i^{(0)})^T \hat{\mathbf{x}}\right) \left(\frac{1}{\sqrt{n_0}} (\mathbf{W}_i^{(0)})^T \mathbf{z}^{(0)}\right)\right] = \frac{1}{n_0} \hat{\mathbf{x}}^T \mathbb{E}[(\mathbf{W}_i^{(0)})(\mathbf{W}_i^{(0)})^T] \mathbf{z}^{(0)} = \frac{1}{n_0} \hat{\mathbf{x}}^T \mathbf{I}_{n_0} \mathbf{z}^{(0)} = \frac{\hat{\mathbf{x}}^T \mathbf{z}^{(0)}}{n_0}$$

Note that the result is independent of the index i , we can define $\hat{z}^{(1)} = \mathbf{z}_i^{(1)}$. Therefore, the base step has been proven.

Inductive Step

$$\mathbf{z}_i^{(\ell+1)} = \sigma' \left(\frac{1}{\sqrt{n_\ell}} (\mathbf{W}_i^{(\ell)})^T \tilde{\alpha}^{(\ell)}(\hat{\mathbf{x}}) \right) \frac{1}{\sqrt{n_\ell}} (\mathbf{W}_i^{(\ell)})^T \mathbf{z}^{(\ell)}$$

Then,

$$a = \frac{1}{\sqrt{n_\ell}} (\mathbf{W}_i^{(\ell)})^T \tilde{\alpha}^{(\ell)}(\hat{\mathbf{x}}) \sim \mathcal{N}\left(\mathbf{0}, \frac{1}{n_\ell} \sum_{i=0}^{n_\ell} ((\tilde{\alpha}^{(\ell)}(\hat{\mathbf{x}})_i)^2)\right)$$

With $n_1, \dots, n_\ell \rightarrow \infty$, $\text{Var}(a) = \mathbb{E}[(\hat{\alpha}^{(\ell)})^2]$. Similarly,

$$\begin{aligned}b &= \frac{1}{\sqrt{n_\ell}} (\mathbf{W}_i^{(\ell)})^T \mathbf{z}^{(\ell)} \\ b &\sim \mathcal{N}(0, \mathbb{E}[(\hat{z}^{(\ell)})^2]) \quad \text{if } n_1, \dots, n_\ell \rightarrow \infty\end{aligned}$$

On the other hand,

$$\begin{aligned}\mathbb{E}[ab] &= \mathbb{E}\left[\left(\frac{1}{\sqrt{n_\ell}}(\mathbf{W}_i^{(\ell)})^T \tilde{\alpha}^{(\ell)}(\mathbf{x})\right)\left(\frac{1}{\sqrt{n_\ell}}(\mathbf{W}_i^{(\ell)})^T \mathbf{z}^{(\ell)}\right)\right] = \frac{1}{n_\ell}(\tilde{\alpha}^{(\ell)}(\mathbf{x}))^T \mathbb{E}[(\mathbf{W}_i^{(\ell)})(\mathbf{W}_i^{(\ell)})^T] \mathbf{z}^{(\ell)} = \frac{1}{n_\ell}(\tilde{\alpha}^{(\ell)}(\mathbf{x}))^T \mathbf{z}^{(\ell)} \\ &= \mathbb{E}[\hat{\alpha}^{(\ell)} \hat{z}^{(\ell)}] \quad \text{if } n_1, \dots, n_\ell \rightarrow \infty\end{aligned}$$

The recursive definition is now proven up to layer $\ell - 1$. Now let's look at the last layer.

$$\tilde{\mathbf{z}}_i^{(L)} = \frac{1}{\sqrt{n_{L-1}}} \mathbf{W}_i^{(L-1)} \mathbf{z}^{(L-1)}$$

By similar arguments as before, it is easy to show that with $n_1, \dots, n_{L-1} \rightarrow \infty$, $\tilde{\mathbf{z}}_i^{(L)} \sim \mathcal{N}(0, \mathbb{E}[(\hat{z}^{(L-1)})^2])$. This concludes the proof. \square

Theorem 1. For any data point \mathbf{x}_i , $i \in [1, \dots, n]$, with probability at least $1 - O(n)e^{-O(n_0)}$,

$$\|\mathbf{J}(\mathbf{x}_i)\|_{op} \leq c\sqrt{n_0\tau}$$

where c is a constant and

$$\tau = \sup_{\mathbf{x}_i \in \hat{\mathbf{X}}, \|\mathbf{z}^{(0)}\|_2=1} \mathbb{E}[(\hat{z}^{(L-1)})^2 | \mathbf{z}^{(0)}, \mathbf{x}_i]$$

Proof. For a fixed unit vector $\mathbf{z}^{(0)}$ and fixed input $\hat{\mathbf{x}}$, we know that based on Proposition 2, $\tilde{\mathbf{z}}_i^{(L)} \sim \mathcal{N}(0, \mathbb{E}[(\hat{z}^{(L-1)})^2 | \mathbf{z}^{(0)}, \hat{\mathbf{x}}])$. Define z as

$$z = \frac{1}{\mathbb{E}[(\hat{z}^{(L-1)})^2 | \mathbf{z}^{(0)}, \hat{\mathbf{x}}]} \|\tilde{\mathbf{z}}^{(L)}\|_2^2 = \chi_{n_0}^2$$

First, notice that we can have the following tail bound for chi-square distribution (for instance, (Kolar & Liu, 2012))

$$\Pr[|z/n_0 - 1| \geq \epsilon] \leq \exp\left(-\frac{3}{16}n_0\epsilon^2\right)$$

when $\epsilon \in [0, 1/2)$. In this case, let $\epsilon = \frac{1}{3}$. Consider a subset of coordinates M with cardinality $|M| \leq O(n_0)$ (Allen-Zhu et al., 2018). Taking the ϵ ball \mathcal{B} of this subspace with $\epsilon = 1/3$, we know what

$$|\mathcal{B}| \leq 7^{|M|} = e^{|M|\ln 7} = e^{O(n_0)}$$

Then, taking the union bound for all unit vectors in \mathcal{B} , we know that

$$\begin{aligned}\forall \mathbf{z}_0 \in \mathcal{B} \quad & \bigcup_{z_0} \Pr[|z/n_0 - 1| \geq \frac{1}{3}] \\ & \leq \exp\left(-\frac{1}{48}n_0\right) \exp(O(n_0)) \leq \exp(-O(n_0))\end{aligned}$$

Therefore, by the ϵ -net argument (Tao, 2012), for any unit vector \mathbf{u} with only non-zero entries in M , we have with probability $1 - \exp(-O(n_0))$,

$$\|\mathbf{J}(\hat{\mathbf{x}})\mathbf{u}\|_2^2 \leq 2n_0\tau \|\mathbf{u}\|_2^2 = C^2 \|\mathbf{u}\|_2^2$$

For any arbitrary vector \mathbf{v} , we can decompose it in the following way: $\mathbf{v} = \mathbf{u}_1 + \mathbf{u}_2 + \dots + \mathbf{u}_K$ with $K = O(1)$ where each \mathbf{u}_i comes from a different non-overlapping coordinate set M .

$$\begin{aligned}\|\mathbf{J}(\hat{\mathbf{x}})\mathbf{v}\|_2 &\leq C \sum_{i=1}^K \|\mathbf{u}_i\|_2 \leq C\sqrt{K} \left(\sum_{i=1}^K \|\mathbf{u}_i\|_2^2\right)^{1/2} \\ &\leq O(1)C \|\mathbf{v}\|.\end{aligned}$$

Thus, with probability at least $1 - O(1)\exp(-O(n_0))$,

$$\begin{aligned}\|\mathbf{J}(\hat{\mathbf{x}})\|_{op} &\leq O(1)C = O(1)\sqrt{2n_0\tau} \\ &= c\sqrt{n_0\tau},\end{aligned}$$

where c is a constant. Taking the union bound over all the data points concludes the proof. \square

B. Proofs for Sec 4.2

Lemma 5. *Under the setting in Section 3.1 with sigmoid as the activation function,*

$$\Theta_{\infty}^{(L)}(\mathbf{x}, \mathbf{x}) \geq \frac{1}{4}$$

Proof. For any ℓ , we have

$$\begin{aligned} \Theta_{\infty}^{(\ell+1)}(\mathbf{x}, \mathbf{x}) &= \Theta_{\infty}^{(\ell)}(\mathbf{x}, \mathbf{x}) \dot{\Sigma}^{(\ell+1)}(\mathbf{x}, \mathbf{x}) + \Sigma^{(\ell+1)}(\mathbf{x}, \mathbf{x}) \\ &\geq \Sigma^{(\ell+1)}(\mathbf{x}, \mathbf{x}) \\ &= \mathbb{E}_{g \sim \mathcal{N}(0, \Sigma^{(\ell)})} [\sigma(g(\mathbf{x}))^2] \\ &= \mathbb{E}_{g \sim \mathcal{N}(0, \Sigma^{(\ell)})} \left[\left(\sigma(g(\mathbf{x})) - \frac{1}{2} \right)^2 \right] + \frac{1}{4} \\ &\geq \frac{1}{4} \end{aligned}$$

where $\sigma(f(\mathbf{x})) - \frac{1}{2}$ moves sigmoid function to the origin such that it is an odd function. \square

C. Proofs for Sec 4.3

C.1. Main Lemmas

Lemma 2. *Suppose there is a 2-layer network. If the activation function is $\sigma(x) = \alpha x$, $n = n_0$ and the data matrix is full rank. Then at NTK limit, $\mathbf{J}_{\infty}(\mathbf{x}) = \mathbf{I}_{n_0}$.*

Proof.

$$\mathbf{J}_{\infty}(\mathbf{x}) = \frac{2\alpha^2}{n_0} (\hat{\mathbf{X}} - f_0(\hat{\mathbf{X}})) \left(\frac{2\alpha^2}{n_0} \hat{\mathbf{X}}^T \hat{\mathbf{X}} \right)^{-1} \hat{\mathbf{X}}^T + \mathbf{J}_0(\mathbf{x})$$

Notice that $\mathbf{J}_0(\mathbf{x}) = \alpha \frac{1}{\sqrt{n_1}} \frac{1}{\sqrt{n_0}} \mathbf{W}^{(1)} \mathbf{W}^{(0)}$ and $f_0(\hat{\mathbf{X}}) = \alpha \frac{1}{\sqrt{n_1}} \frac{1}{\sqrt{n_0}} \mathbf{W}^{(1)} \mathbf{W}^{(0)} \hat{\mathbf{X}} = \mathbf{J}_0(\mathbf{x}) \hat{\mathbf{X}}$.

$$\begin{aligned} \mathbf{J}_{\infty}(\mathbf{x}) &= \mathbf{J}_0(\mathbf{x}) - f_0(\hat{\mathbf{X}}) (\hat{\mathbf{X}}^T \hat{\mathbf{X}})^{-1} \hat{\mathbf{X}}^T + \hat{\mathbf{X}} (\hat{\mathbf{X}}^T \hat{\mathbf{X}})^{-1} \hat{\mathbf{X}}^T \\ &= \alpha \frac{1}{\sqrt{n_1}} \frac{1}{\sqrt{n_0}} \mathbf{W}^{(1)} \mathbf{W}^{(0)} - \alpha \frac{1}{\sqrt{n_1}} \frac{1}{\sqrt{n_0}} \mathbf{W}^{(1)} \mathbf{W}^{(0)} \hat{\mathbf{X}} (\hat{\mathbf{X}}^T \hat{\mathbf{X}})^{-1} \hat{\mathbf{X}}^T + \hat{\mathbf{X}} (\hat{\mathbf{X}}^T \hat{\mathbf{X}})^{-1} \hat{\mathbf{X}}^T \\ &= \alpha \frac{1}{\sqrt{n_1}} \frac{1}{\sqrt{n_0}} \mathbf{W}^{(1)} \mathbf{W}^{(0)} - \alpha \frac{1}{\sqrt{n_1}} \frac{1}{\sqrt{n_0}} \mathbf{W}^{(1)} \mathbf{W}^{(0)} \mathbf{I}_{n_0} + \mathbf{I}_{n_0} \\ &= \mathbf{I}_{n_0} \end{aligned}$$

\square

Lemma 3. *Suppose there is a 2-layer network with activation function $\sigma(x) = \alpha x$ and given initial weights $\mathbf{W}^{(1)} \in \mathbb{R}^{n_0 \times n_1}$, $\mathbf{W}^{(0)} \in \mathbb{R}^{n_1 \times n_0}$. If the data matrix is full rank with $n \leq n_0$, then, at the NTK limit ($n_1 \rightarrow \infty$), $\mathbf{J}_{\infty}(\mathbf{x})$ has eigenvalue 1 with multiplicity at least n . If at the NTK limit, α is chosen such that $\|\mathbf{J}_0(\mathbf{x})\|_{op} < 1$, then the multiplicity is exactly n and 1 is the largest eigenvalue norm.*

Proof. Based on the proof of last section, we know that

$$\begin{aligned} \mathbf{J}_{\infty}(\mathbf{x}) &= \mathbf{J}_0(\mathbf{x}) - f_0(\hat{\mathbf{X}}) (\hat{\mathbf{X}}^T \hat{\mathbf{X}})^{-1} \hat{\mathbf{X}}^T + \hat{\mathbf{X}} (\hat{\mathbf{X}}^T \hat{\mathbf{X}})^{-1} \hat{\mathbf{X}}^T \\ &= \mathbf{J}_0(\mathbf{x}) - \mathbf{J}_0(\mathbf{x}) \hat{\mathbf{X}} (\hat{\mathbf{X}}^T \hat{\mathbf{X}})^{-1} \hat{\mathbf{X}}^T + \hat{\mathbf{X}} (\hat{\mathbf{X}}^T \hat{\mathbf{X}})^{-1} \hat{\mathbf{X}}^T \end{aligned}$$

In this case,

$$\hat{\mathbf{X}} (\hat{\mathbf{X}}^T \hat{\mathbf{X}})^{-1} \hat{\mathbf{X}}^T = V \Sigma V^T$$

where V is an orthogonal matrix and

$$\Sigma = \begin{bmatrix} 1 & 0 & \dots & 0 \\ 0 & 1 & \dots & 0 \\ \dots & \dots & \dots & \dots \\ 0 & \dots & 1 & 0 \\ \dots & \dots & \dots & \dots \\ 0 & 0 & \dots & 0 \end{bmatrix}$$

So

$$\begin{aligned} \mathbf{J}_\infty(\mathbf{x}) &= \mathbf{J}_0(\mathbf{x})(\mathbf{I}_{n_0} - V\Sigma V^T) + V\Sigma V^T \\ &= \alpha \frac{1}{\sqrt{n_1}} \frac{1}{\sqrt{n_0}} \mathbf{W}^{(1)} \mathbf{W}^{(0)} (\mathbf{I}_{n_0} - V\Sigma V^T) + V\Sigma V^T \end{aligned}$$

Interestingly, $(I_d - V\Sigma V^T)$ and $V\Sigma V^T$ contain orthogonal eigenvectors. For convenience, let $\{v_i\}_{i=1}^n$ be the set of eigenvectors of $V\Sigma V^T$ with eigenvalue 1. Furthermore, let $V_{\parallel} = \text{span}(\{v_i\}_{i=1}^n)$ and $V_{\perp} = \text{span}(\{v_i\}_{i=1}^n)^{\perp}$. Because we are in the linear region, $\mathbf{J}_\infty(\mathbf{x})$ and $\mathbf{J}_0(\mathbf{x})$ do not depend on \mathbf{x} . We'll use \mathbf{J}_∞ to refer $\mathbf{J}_\infty(\mathbf{x})$ and \mathbf{J}_0 as $\mathbf{J}_0(\mathbf{x})$.

- For any vector $v^{\parallel} \in V_{\parallel}$,

$$\mathbf{J}_0(\mathbf{I}_{n_0} - V\Sigma V^T)v^{\parallel} = 0$$

and

$$\mathbf{J}_\infty v^{\parallel} = v^{\parallel}$$

Thus, all vectors in $\{v_i\}_{i=1}^n$ are eigenvectors of \mathbf{J}_∞ with eigenvalue 1 regardless of the choice of α .

- On the other hand, let v be any complex vector such that

$$v = \text{Re}(v) + i\text{Im}(v)$$

If v is an eigenvector of \mathbf{J}_∞ with eigenvalue $\lambda = a + ib$, then

$$\mathbf{J}_\infty \text{Re}(v) = a\text{Re}(v) - b\text{Im}(v)$$

$$\mathbf{J}_\infty \text{Im}(v) = b\text{Re}(v) + a\text{Im}(v)$$

Let's first decompose $\text{Re}(v)$ and $\text{Im}(v)$.

$$\text{Re}(v) = v_r^{\perp} + v_r^{\parallel}$$

$$\text{Im}(v) = v_i^{\perp} + v_i^{\parallel}$$

where $v_r^{\perp}, v_i^{\perp} \in V_{\perp}$ and $v_r^{\parallel}, v_i^{\parallel} \in V_{\parallel}$.

$$\mathbf{J}_\infty(v_r^{\perp} + v_r^{\parallel}) = J_0 v_r^{\perp} + v_r^{\parallel} = (a v_r^{\perp} - b v_i^{\perp}) + (a v_r^{\parallel} - b v_i^{\parallel})$$

$$\mathbf{J}_\infty(v_i^{\perp} + v_i^{\parallel}) = J_0 v_i^{\perp} + v_i^{\parallel} = (b v_r^{\perp} + a v_i^{\perp}) + (b v_r^{\parallel} + a v_i^{\parallel})$$

By adding and subtracting two equations,

$$\mathbf{J}_0(v_r^{\perp} + v_i^{\perp}) + v_r^{\parallel} + v_i^{\parallel} = \left[(a+b)v_r^{\perp} + (a-b)v_i^{\perp} \right] + \left[(a+b)v_r^{\parallel} + (a-b)v_i^{\parallel} \right]$$

$$\mathbf{J}_0(v_r^{\perp} - v_i^{\perp}) + v_r^{\parallel} - v_i^{\parallel} = \left[(a-b)v_r^{\perp} - (a+b)v_i^{\perp} \right] + \left[(a-b)v_r^{\parallel} - (a+b)v_i^{\parallel} \right]$$

When α is chosen such that $\|\mathbf{J}_0\| < 1$,

$$\|(a+b)v_r^{\perp} + (a-b)v_i^{\perp}\|_2 < \|v_r^{\perp} + v_i^{\perp}\|_2$$

$$\|(a-b)v_r^{\perp} - (a+b)v_i^{\perp}\|_2 < \|v_r^{\perp} - v_i^{\perp}\|_2$$

Then,

$$\begin{aligned} (a^2 + b^2)\|v_r^\perp\|_2^2 + (a^2 + b^2)\|v_i^\perp\|_2^2 &< \|v_r^\perp\|_2^2 + \|v_i^\perp\|_2^2 \\ |\lambda|^2 = a^2 + b^2 &< 1 \end{aligned}$$

This suggests that any complex eigenvector with components from V_\perp would have eigenvalue with norm smaller than 1.

□

Lemma 4. Suppose there is a 2-layer network with activation function $\sigma(x) = \alpha x + \beta$, given initial weights $\mathbf{W}^{(1)} \in \mathbb{R}^{n_0 \times n_1}$, $\mathbf{W}^{(0)} \in \mathbb{R}^{n_1 \times n_0}$ and every data point has the same norm r (i.e. $\forall i \in [n] \|\mathbf{x}\|_2 = r$). If the data matrix is full rank with $n \leq n_0$, then, at the NTK limit $n_1 \rightarrow \infty$, $\mathbf{J}_\infty(\mathbf{x})$ has eigenvalues 1 with multiplicity at least $n - 1$. If at the NTK limit, α and β are chosen such that

$$\|\mathbf{J}_0(\mathbf{x})\|_{op} = 1 - \Delta, \quad \left\| \frac{1}{\sqrt{n_1}} \mathbf{W}^{(1)} \mathbf{1}_{n_1} \right\|_2 < \frac{\beta n_0 \Delta}{2r\alpha^2},$$

where $0 < \Delta \leq 1$, then the multiplicity is exactly $n - 1$ and 1 is the largest eigenvalue norm.

Proof. First of all, let \mathbf{B} be an all-one matrix

$$\begin{aligned} \mathbf{J}_\infty(\mathbf{x}) &= \left(\hat{\mathbf{X}} - f_0(\hat{\mathbf{X}}) \right) \tilde{\mathbf{K}}^{-1} \frac{\partial k_x}{\partial \mathbf{x}} + \mathbf{J}_0(\mathbf{x}) \\ &= \left(\hat{\mathbf{X}} - f_0(\hat{\mathbf{X}}) \right) \left(\frac{2\alpha^2}{n_0} \hat{\mathbf{X}}^T \hat{\mathbf{X}} + \beta^2 \mathbf{B} \right)^{-1} \left(\frac{2\alpha^2}{n_0} \hat{\mathbf{X}}^T \right) + \mathbf{J}_0(\mathbf{x}) \\ &= \left(\hat{\mathbf{X}} - f_0(\hat{\mathbf{X}}) \right) \left(\hat{\mathbf{X}}^T \hat{\mathbf{X}} + \frac{n_0 \beta^2}{2\alpha^2} \mathbf{B} \right)^{-1} \hat{\mathbf{X}}^T + \mathbf{J}_0(\mathbf{x}) \\ &= \mathbf{J}_0(\mathbf{x}) + \hat{\mathbf{X}} \left(\hat{\mathbf{X}}^T \hat{\mathbf{X}} + \frac{n_0 \beta^2}{2\alpha^2} \mathbf{B} \right)^{-1} \hat{\mathbf{X}}^T - \left(\frac{\alpha}{\sqrt{n_1 n_0}} \mathbf{W}^{(1)} \mathbf{W}^{(0)} \hat{\mathbf{X}} + \beta \frac{1}{\sqrt{n_1}} \mathbf{W}^{(1)} \mathbf{1}_{n_1} \mathbf{1}_n^T \right) \left(\hat{\mathbf{X}}^T \hat{\mathbf{X}} + \frac{n_0 \beta^2}{2\alpha^2} \mathbf{B} \right)^{-1} \hat{\mathbf{X}}^T \\ &= \mathbf{J}_0(\mathbf{x}) + \hat{\mathbf{X}} \left(\hat{\mathbf{X}}^T \hat{\mathbf{X}} + \frac{n_0 \beta^2}{2\alpha^2} \mathbf{B} \right)^{-1} \hat{\mathbf{X}}^T - \left(\mathbf{J}_0(\mathbf{x}) \hat{\mathbf{X}} + \beta \frac{1}{\sqrt{n_1}} \mathbf{W}^{(1)} \mathbf{1}_{n_1} \mathbf{1}_n^T \right) \left(\hat{\mathbf{X}}^T \hat{\mathbf{X}} + \frac{n_0 \beta^2}{2\alpha^2} \mathbf{B} \right)^{-1} \hat{\mathbf{X}}^T \end{aligned}$$

Because in the linearized region, $\mathbf{J}_\infty(\mathbf{x})$ and $\mathbf{J}_0(\mathbf{x})$ do not depend on \mathbf{x} . We'll use \mathbf{J}_∞ to refer $\mathbf{J}_\infty(\mathbf{x})$ and \mathbf{J}_0 as $\mathbf{J}_0(\mathbf{x})$. For simplicity, we'll also use $c = \frac{n_0 \beta^2}{2\alpha^2}$.

Based on Lemma 6,

$$\hat{\mathbf{X}} \left(\hat{\mathbf{X}}^T \hat{\mathbf{X}} + c\mathbf{B} \right)^{-1} \hat{\mathbf{X}}^T = V\Lambda V^T$$

where $\Lambda = \text{diag}(\underbrace{1, \dots, 1}_{n-1}, \hat{\lambda}, \underbrace{0, \dots, 0}_{n_0-n})$ where $0 < \hat{\lambda} < 1$. Now,

$$\mathbf{J}_\infty = \mathbf{J}_0(I_{n_0} - V\Lambda V^T) + V\Lambda V^T - \beta \frac{1}{\sqrt{n_1}} \mathbf{W}^{(1)} \mathbf{1}_{n_1} \mathbf{1}_n^T \left(\hat{\mathbf{X}}^T \hat{\mathbf{X}} + c\mathbf{B} \right)^{-1} \hat{\mathbf{X}}^T$$

From Corollary 7, we know that the following two vectors are eigenvectors of $V\Lambda V^T$ with eigenvalue $\hat{\lambda}$,

$$\hat{\mathbf{X}} (\hat{\mathbf{X}}^T \hat{\mathbf{X}} + c\mathbf{B})^{-1} \mathbf{1}_n \quad \hat{\mathbf{X}} (\hat{\mathbf{X}}^T \hat{\mathbf{X}})^{-1} \mathbf{1}_n$$

Furthermore,

$$\hat{\mathbf{X}} (\hat{\mathbf{X}}^T \hat{\mathbf{X}} + c\mathbf{B})^{-1} \mathbf{1}_n = \hat{\lambda} \hat{\mathbf{X}} (\hat{\mathbf{X}}^T \hat{\mathbf{X}})^{-1} \mathbf{1}_n$$

And

$$\hat{\lambda} = \frac{1}{1 + cg}$$

where

$$g = \text{trace}(\mathbf{B}(\hat{\mathbf{X}}^T \hat{\mathbf{X}})^{-1}) = \|\hat{\mathbf{X}}(\hat{\mathbf{X}}^T \hat{\mathbf{X}})^{-1} \mathbf{1}_n\|_2^2$$

Let \hat{u} be a rescaled unit vector of $\hat{\mathbf{X}}(\hat{\mathbf{X}}^T \hat{\mathbf{X}})^{-1} \mathbf{1}_n$, then

$$\begin{aligned} \mathbf{J}_\infty \hat{u} &= \mathbf{J}_0(1 - \hat{\lambda})\hat{u} + \hat{\lambda}\hat{u} - \sqrt{g}\hat{\lambda}\beta \frac{1}{\sqrt{n_1}} \mathbf{W}^{(1)} \mathbf{1}_{n_1} \hat{u} \\ \|\mathbf{J}_\infty \hat{u}\|_2 &= \|\mathbf{J}_0(1 - \hat{\lambda})\hat{u} + \hat{\lambda}\hat{u} - \sqrt{g}\hat{\lambda}\beta \frac{1}{\sqrt{n_1}} \mathbf{W}^{(1)} \mathbf{1}_{n_1} \hat{u}\|_2 \\ &\leq \|\mathbf{J}_0\|_{op} \|(1 - \hat{\lambda})\hat{u}\|_2 + \|\hat{\lambda}\hat{u}\|_2 + \|\sqrt{g}\hat{\lambda}\beta \frac{1}{\sqrt{n_1}} \mathbf{W}^{(1)} \mathbf{1}_{n_1} \hat{u}\|_2 \\ &= (1 - \hat{\lambda})\|\mathbf{J}_0\|_{op} + \hat{\lambda} + \sqrt{g}\hat{\lambda}\|\beta \frac{1}{\sqrt{n_1}} \mathbf{W}^{(1)} \mathbf{1}_{n_1}\|_2 \\ &< (1 - \hat{\lambda})(1 - \Delta) + \hat{\lambda} + \sqrt{g}\hat{\lambda} \frac{\beta^2 n_0 \Delta}{2r\alpha^2} \\ &\leq (1 - \hat{\lambda})(1 - \Delta) + \hat{\lambda} + g\hat{\lambda} \frac{\beta^2 n_0 \Delta}{2\alpha^2} \quad (\text{Lemma 9}) \\ &= \frac{(1 - \Delta)cg + 1 + \frac{g\beta^2 n_0 \Delta}{2\alpha^2}}{1 + cg} = 1 \end{aligned}$$

Therefore, \mathbf{J}_∞ will shrink every vectors orthogonal to the eigenvectors in V with eigenvalue 1. By the same arguments in the proof of Lemma 3, we can conclude the proof. \square

C.2. Useful Lemmas

Lemma 6. Suppose $\mathbf{X} \in \mathbb{R}^{k \times m}$ is a full-rank matrix with $k \geq m$ and $m \geq 2$. Let c be an arbitrary positive constant and \mathbf{B} an all-one matrix. Consider the following real symmetric matrix,

$$\mathbf{X}(\mathbf{X}^T \mathbf{X} + c\mathbf{B})^{-1} \mathbf{X}^T$$

It can be characterized by having eigenvalue 1 with multiplicity $m - 1$, eigenvalue 0 with multiplicity $k - m$ and another eigenvalue λ such that $0 < \lambda < 1$.

Proof. By (Miller, 1981), if P and $P + Q$ are invertible, and Q has rank 1, then let $g' = \text{trace}(QP^{-1})$, we know that $g' \neq 1$, and

$$(P + Q)^{-1} = P^{-1} - \frac{1}{1 + g'} P^{-1} Q P^{-1}$$

First of all, it is easy to see that $(\mathbf{X}^T \mathbf{X} + c\mathbf{B})^{-1}$ is invertible. This is because $\mathbf{X}^T \mathbf{X}$ is positive definite and $c\mathbf{B}$ is positive semi-definite.

Since B is a rank one matrix,

$$(\mathbf{X}^T \mathbf{X} + c\mathbf{B})^{-1} = \underbrace{(\mathbf{X}^T \mathbf{X})^{-1}}_{I_1} - \frac{c}{1 + cg} \underbrace{(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{B} (\mathbf{X}^T \mathbf{X})^{-1}}_{I_2}$$

where $g = \text{trace}(\mathbf{B}(\mathbf{X}^T \mathbf{X})^{-1})$.

Let's consider the singular value decomposition of $\mathbf{X}^T = U\Sigma V^T$

- $\mathbf{X}^T \mathbf{X} = U\Sigma^2 U^T$ and $(\mathbf{X}^T \mathbf{X})^{-1} = U\Sigma^{-2} U^T$. So

$$\mathbf{X} I_1 \mathbf{X}^T = \mathbf{X} (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T = V \Sigma U^T U \Sigma^{-2} U^T U \Sigma V^T = V \Lambda_m V^T$$

where $\Lambda_m = \text{diag}(\underbrace{1, \dots, 1}_m, \underbrace{0, \dots, 0}_{k-m})$

•

$$\mathbf{X}^T I_2 \mathbf{X} = \frac{c}{1+cg} M = \frac{c}{1+cg} \mathbf{X}(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{B}(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T$$

The first thing to notice is that $\mathbf{B} = \mathbf{1}\mathbf{1}^T$ where $\mathbf{1}$ is a vector of ones. Therefore,

$$M = \mathbf{X}(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{B}(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T = \mathbf{X}(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{1}\mathbf{1}^T (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T = \mathbf{a}\mathbf{a}^T$$

where $\mathbf{a} = \mathbf{X}(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{1}$.

This implies that M is a rank one matrix with singular value $\|\mathbf{a}\|^2$. But we also know the following:

$$\begin{aligned} \|\mathbf{a}\|^2 &= \mathbf{a}^T \mathbf{a} = \text{trace}(\mathbf{a}\mathbf{a}^T) = \text{trace}(M) \\ &= \text{trace}(\mathbf{X}(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{B}(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T) = \text{trace}(\mathbf{X}^T \mathbf{X}(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{B}(\mathbf{X}^T \mathbf{X})^{-1}) \\ &= \text{trace}(\mathbf{B}(\mathbf{X}^T \mathbf{X})^{-1}) = g > 0 \end{aligned}$$

The last strict inequality comes from the fact that X is full rank so that $\mathbf{X}(\mathbf{X}^T \mathbf{X})^{-1}$ has no zero singular value. Furthermore,

$$\begin{aligned} \mathbf{X} I_1 \mathbf{X}^T \mathbf{a} &= \mathbf{X}(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{a} \\ &= \mathbf{X}(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{X}(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{1} = \mathbf{X}(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{1} \\ &= \mathbf{a} \end{aligned}$$

Because \mathbf{a} is not a zero vector, it is also one of the eigenvector of $\mathbf{X} I_1 X$ with eigenvalue 1.

And the eigenvalue of $X^T I_2 X$ is the following:

$$0 < \frac{cg}{1+cg} < 1$$

The inequalities comes from the fact that c is also non-negative. We'll denote $\sigma = \frac{cg}{1+cg}$. So

$$\mathbf{X} I_2 \mathbf{X}^T = \sigma \hat{\mathbf{a}} \hat{\mathbf{a}}^T$$

where $\hat{\mathbf{a}}$ is a rescaled to have unit length.

Now that we have examined two parts separately. Let's put them together. For convenience, we'll also denote $\mathbf{X}(\mathbf{X}^T \mathbf{X} + c\mathbf{B})^{-1} \mathbf{X}^T = \mathbf{X} I_1 \mathbf{X}^T - \mathbf{X} I_2 \mathbf{X}^T = \mathbf{M}_1 - \mathbf{M}_2$.

Based on the eigen decomposition of \mathbf{M}_1 ,

$$\mathbf{M}_1 = \sum_{k=1}^m \mathbf{u}_k \mathbf{u}_k^T$$

with lost of generality, let's also denote $\hat{\mathbf{a}} = \mathbf{u}_1$. Now,

$$\begin{aligned} \mathbf{M}_1 - \mathbf{M}_2 &= \sum_{k=1}^m \mathbf{u}_k \mathbf{u}_k^T - \sigma \mathbf{u}_1 \mathbf{u}_1^T \\ &= (1 - \sigma) \mathbf{u}_1 \mathbf{u}_1^T + \sum_{k=2}^m \mathbf{u}_k \mathbf{u}_k^T \end{aligned}$$

Because $0 < \sigma < 1$, $\mathbf{X}(\mathbf{X}^T \mathbf{X} + c\mathbf{B})^{-1} \mathbf{X}^T$ has eigenvalue 1 with multiplicity $m - 1$, eigenvalue 0 with multiplicity $k - m$ and another eigenvalue λ such that $0 < \lambda < 1$. \square

Corollary 7. Following the setup in Lemma 6, we could also know that $\mathbf{X}^T (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{1}$ is an eigenvector with with eigenvalue λ and

$$\left(\mathbf{X}(\mathbf{X}^T \mathbf{X} + c\mathbf{B})^{-1} \mathbf{X}^T \right) \mathbf{X}(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{1} = \mathbf{X}(\mathbf{X}^T \mathbf{X} + c\mathbf{B})^{-1} \mathbf{1} = \lambda \mathbf{X}(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{1}$$

Corollary 8. Suppose $\mathbf{X} \in \mathbb{R}^{k \times m}$ is a full-rank matrix with $k \geq m$ and $m \geq 2$. Let c be an arbitrary non-negative constant and \mathbf{B} an all-one matrix.

$$\|\mathbf{X}(\mathbf{X}^T \mathbf{X} + c\mathbf{B})^{-1} \mathbf{X}^T\|_{op} = 1$$

Remark 2. c can also take on negative values as long as cg is not close to -1 .

Lemma 9. Suppose $\mathbf{X} \in \mathbb{R}^{k \times m}$ is a full-rank matrix with $k \geq m$ and \mathbf{B} an all-one matrix. If

$$\|\mathbf{X}_{\cdot, i}\|_2 = r \quad \forall i \in [m]$$

Then,

$$\text{trace}(\mathbf{B}(\mathbf{X}^T \mathbf{X})^{-1}) \geq \frac{1}{r^2}$$

Proof. First of all,

$$\begin{aligned} \text{trace}(\mathbf{B}(\mathbf{X}^T \mathbf{X})^{-1}) &\geq \text{trace}(\mathbf{1}^T (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{1}) \\ &\geq \|\mathbf{1}\|_2^2 \frac{1}{\|\mathbf{X}^T \mathbf{X}\|_{op}} = \frac{m}{\|\mathbf{X}^T \mathbf{X}\|_{op}} \end{aligned}$$

On the hand,

$$\|\mathbf{X}^T \mathbf{X}\|_{op} = \|\mathbf{X}^T\|_{op}^2 \leq \|\mathbf{X}^T\|_f^2 \leq \text{trace}(\mathbf{X}^T \mathbf{X}) \leq r^2 m$$

Therefore,

$$\text{trace}(\mathbf{B}(\mathbf{X}^T \mathbf{X})^{-1}) \geq \frac{1}{r^2}$$

□

D. Proofs for Sec 4.4

D.1. Derivation for the Approximated NTK

The closed form NTK of erf (Lee et al., 2019; Williams, 1997) can be written with the following two components:

$$\begin{aligned} \mathcal{T}(\Sigma, \text{erf}, \text{erf})(\mathbf{x}, \hat{\mathbf{x}}) &= \frac{2}{\pi} \arcsin \left(\frac{\Sigma(\mathbf{x}, \hat{\mathbf{x}})}{\sqrt{(\Sigma(\mathbf{x}, \mathbf{x}) + 0.5)(\Sigma(\hat{\mathbf{x}}, \hat{\mathbf{x}}) + 0.5)}} \right) \\ \mathcal{T}(\Sigma, \text{erf}, \text{erf})(\mathbf{x}, \hat{\mathbf{x}}) &= \frac{4}{\pi} \det(I + 2\Sigma)^{-\frac{1}{2}} = \frac{4}{\pi} \frac{1}{\sqrt{(1 + 2\Sigma(\mathbf{x}, \mathbf{x}))(1 + 2\Sigma(\hat{\mathbf{x}}, \hat{\mathbf{x}})) - 4\Sigma(\mathbf{x}, \hat{\mathbf{x}})^2}} \end{aligned}$$

Here, we can approximate sigmoid function σ_s by erf function:

$$\sigma_s(x) \approx \sigma_{\hat{s}}(x) = \frac{1}{2} \text{erf}\left(\frac{1}{2}x\right) + \frac{1}{2}$$

Then,

$$\begin{aligned} \mathcal{T}(\Sigma, \sigma_{\hat{s}}, \sigma_{\hat{s}})(\mathbf{x}, \hat{\mathbf{x}}) &= \mathbb{E}_{u, v \sim \mathcal{N}(0, \Sigma)}[\sigma_{\hat{s}}(u)\sigma_{\hat{s}}(v)] = \mathbb{E}\left[\frac{1}{4} \text{erf}\left(\frac{1}{2}u\right) \text{erf}\left(\frac{1}{2}v\right)\right] + \mathbb{E}\left[\frac{1}{4} \text{erf}\left(\frac{1}{2}u\right) + \frac{1}{4} \text{erf}\left(\frac{1}{2}v\right)\right] + \frac{1}{4} \\ &= \frac{1}{4} \mathbb{E}[\text{erf}\left(\frac{1}{2}u\right) \text{erf}\left(\frac{1}{2}v\right)] + \frac{1}{4} \\ &= \frac{1}{4} \mathcal{T}\left(\frac{1}{4}\Sigma, \text{erf}, \text{erf}\right)(\mathbf{x}, \hat{\mathbf{x}}) + \frac{1}{4} \end{aligned}$$

$$\boxed{\mathcal{T}(\Sigma, \sigma_{\hat{s}}, \sigma_{\hat{s}})(\mathbf{x}, \hat{\mathbf{x}}) = \frac{1}{4} + \frac{1}{2\pi} \arcsin \left(\frac{\Sigma(\mathbf{x}, \hat{\mathbf{x}})}{\sqrt{(\Sigma(\mathbf{x}, \mathbf{x}) + 2)(\Sigma(\hat{\mathbf{x}}, \hat{\mathbf{x}}) + 2)}} \right)}$$

and

$$\begin{aligned}\mathcal{T}(\Sigma, \dot{\sigma}_{\hat{s}}, \dot{\sigma}_s)(\mathbf{x}, \hat{\mathbf{x}}) &= \mathbb{E}_{u, v \sim \mathcal{N}(0, \Sigma)} [\dot{\sigma}_{\hat{s}}(u) \dot{\sigma}_s(v)] = \frac{1}{16} \mathbb{E}[\operatorname{erf}(\frac{1}{2}u) \operatorname{erf}(\frac{1}{2}v)] \\ &= \frac{1}{16} \mathcal{T}(\frac{1}{4}\Sigma, \operatorname{erf}, \operatorname{erf})(\mathbf{x}, \hat{\mathbf{x}})\end{aligned}$$

$$\boxed{\mathcal{T}(\Sigma, \dot{\sigma}_{\hat{s}}, \dot{\sigma}_s)(\mathbf{x}, \hat{\mathbf{x}}) = \frac{1}{2\pi} \frac{1}{\sqrt{(2 + \Sigma(\mathbf{x}, \mathbf{x}))(2 + \Sigma(\hat{\mathbf{x}}, \hat{\mathbf{x}})) - \Sigma(\mathbf{x}, \hat{\mathbf{x}})^2}}}$$

Based on the definition of NTK, we can derive the following for $\sigma_{\hat{s}}$

$$\begin{aligned}\Theta_{\infty}^1(\hat{\mathbf{x}}, \mathbf{x}) &= \Sigma^1(\hat{\mathbf{x}}, \mathbf{x}) = \frac{1}{n_0} \hat{\mathbf{x}}^T \mathbf{x} \\ \Theta_{\infty}^2(\hat{\mathbf{x}}, \mathbf{x}) &= \Theta_{\infty}^1(\hat{\mathbf{x}}, \mathbf{x}) \mathcal{T}(\Theta_{\infty}^1, \dot{\sigma}_{\hat{s}}, \dot{\sigma}_s)(\mathbf{x}, \hat{\mathbf{x}}) + \mathcal{T}(\Theta_{\infty}^1, \sigma_{\hat{s}}, \sigma_s)(\mathbf{x}, \hat{\mathbf{x}})\end{aligned}$$

Let's look at the first part

$$\begin{aligned}\Theta_{\infty}^1(\hat{\mathbf{x}}, \mathbf{x}) \mathcal{T}(\Theta_{\infty}^1, \dot{\sigma}_{\hat{s}}, \dot{\sigma}_s)(\mathbf{x}, \hat{\mathbf{x}}) &= \frac{1}{2\pi} \frac{1}{\sqrt{(2 + \frac{1}{n_0} \mathbf{x}^T \mathbf{x})(2 + \frac{1}{n_0} \hat{\mathbf{x}}^T \hat{\mathbf{x}}) - (\frac{1}{n_0} \hat{\mathbf{x}}^T \mathbf{x})^2}} \left[\frac{1}{n_0} \hat{\mathbf{x}}^T \mathbf{x} \right] \\ &= \frac{1}{2\pi} \frac{\hat{\mathbf{x}}^T \mathbf{x}}{\sqrt{(2n_0 + \mathbf{x}^T \mathbf{x})(2n_0 + \hat{\mathbf{x}}^T \hat{\mathbf{x}}) - (\hat{\mathbf{x}}^T \mathbf{x})^2}}\end{aligned}$$

and the second part

$$\begin{aligned}\mathcal{T}(\Theta_{\infty}^1, \sigma_{\hat{s}}, \sigma_s)(\mathbf{x}, \hat{\mathbf{x}}) &= \frac{1}{4} + \frac{1}{2\pi} \arcsin \left(\frac{\frac{1}{n_0} \hat{\mathbf{x}}^T \mathbf{x}}{\sqrt{(\frac{1}{n_0} \mathbf{x}^T \mathbf{x} + 2)(\frac{1}{n_0} \hat{\mathbf{x}}^T \hat{\mathbf{x}} + 2)}} \right) \\ &= \frac{1}{4} + \frac{1}{2\pi} \arcsin \left(\frac{\hat{\mathbf{x}}^T \mathbf{x}}{\sqrt{(\mathbf{x}^T \mathbf{x} + 2n_0)(\hat{\mathbf{x}}^T \hat{\mathbf{x}} + 2n_0)}} \right)\end{aligned}$$

D.2. Detailed Discussion of $\frac{\partial \mathbf{k}_x}{\partial \mathbf{x}}$

Without loss of generality, we will focus on $\frac{\partial \mathbf{k}_x}{\partial \mathbf{x}}|_{\mathbf{x}_1}$,

$$\frac{\partial \mathbf{k}_x}{\partial \mathbf{x}} = \begin{bmatrix} \frac{\Theta_{\infty}^L(\mathbf{x}_1, \mathbf{x})}{\partial \mathbf{x}} \\ \dots \\ \frac{\Theta_{\infty}^L(\mathbf{x}_n, \mathbf{x})}{\partial \mathbf{x}} \end{bmatrix}$$

where

$$\frac{\partial \Theta_{\infty}^L(\hat{\mathbf{x}}, \mathbf{x})}{\partial \mathbf{x}} = \underbrace{\frac{\partial \mathcal{T}(\Theta_{\infty}^1, \sigma_{\hat{s}}, \sigma_s)(\hat{\mathbf{x}}, \mathbf{x})}{\partial \mathbf{x}}}_{I_1^g(\hat{\mathbf{x}}, \mathbf{x})} + \underbrace{\frac{\partial \Theta_{\infty}^1(\hat{\mathbf{x}}, \mathbf{x}) \mathcal{T}(\Theta_{\infty}^1, \dot{\sigma}_{\hat{s}}, \dot{\sigma}_s)(\hat{\mathbf{x}}, \mathbf{x})}{\partial \mathbf{x}}}_{I_2^g(\hat{\mathbf{x}}, \mathbf{x})}$$

Let's look at each row separately, and break this down into two parts.

- $I_1^g(\hat{\mathbf{x}}, \mathbf{x})$

After deriving the derivative, we get this:

$$\begin{aligned}I_1^g(\hat{\mathbf{x}}, \mathbf{x}) &= \frac{1}{2\pi} \frac{1}{\sqrt{1 - A^2}} \frac{\hat{\mathbf{x}} \left[(\hat{\mathbf{x}}^T \hat{\mathbf{x}} + 2n_0)(\mathbf{x}^T \mathbf{x} + 2n_0) \right] - \mathbf{x} \left[(\hat{\mathbf{x}}^T \hat{\mathbf{x}} + 2n_0) \mathbf{x}^T \hat{\mathbf{x}} \right]}{\left[(\hat{\mathbf{x}}^T \hat{\mathbf{x}} + 2n_0)(\mathbf{x}^T \mathbf{x} + 2n_0) \right]^{\frac{3}{2}}} \\ A &= \frac{\hat{\mathbf{x}}^T \mathbf{x}}{\sqrt{(\mathbf{x}^T \mathbf{x} + 2n_0)(\hat{\mathbf{x}}^T \hat{\mathbf{x}} + 2n_0)}}\end{aligned}$$

Since we are only interested in $\mathbf{J}_\infty(\mathbf{x}_1)$ and each row of $\frac{\partial \mathbf{k}_x}{\partial \mathbf{x}}$, we'll examine $I_1^g(\mathbf{x}_i, \mathbf{x}_1)$.

$$\begin{aligned} I_1^g(\mathbf{x}_i, \mathbf{x}_1) &= \frac{1}{2\pi} \frac{r^2 + 2n_0}{\sqrt{(r^2 + 2n_0)^2 - (r^2 \rho_{i1})^2}} \frac{\mathbf{x}_i(r^2 + 2n_0)^2 - \mathbf{x}_1 \left[(r^2 + 2n_0)r^2 \rho_{i1} \right]}{(r^2 + 2n_0)^3} \\ &= \frac{1}{2\pi} \frac{1}{\sqrt{(r^2 + 2n_0)^2 - (r^2 \rho_{i1})^2}} \frac{\mathbf{x}_i(r^2 + 2n_0) - \mathbf{x}_1 r^2 \rho_{i1}}{r^2 + 2n_0} \end{aligned}$$

It is easy to see that $I_1^g(\mathbf{x}_i, \mathbf{x}_1) \rightarrow 0$ as r grows regardless of ρ_{i1} .

- $I_2^g(\hat{\mathbf{x}}, \mathbf{x})$

We know that

$$I_2^g(\hat{\mathbf{x}}, \mathbf{x}) = \frac{1}{2\pi} \frac{\hat{\mathbf{x}} \left[(\mathbf{x}^T \mathbf{x} + 2n_0)(\hat{\mathbf{x}}^T \hat{\mathbf{x}} + 2n_0) - (\hat{\mathbf{x}}^T \mathbf{x})^2 \right] - \hat{\mathbf{x}}^T \mathbf{x} \left[(2n_0 + \hat{\mathbf{x}}^T \hat{\mathbf{x}})\mathbf{x} - (\hat{\mathbf{x}}^T \mathbf{x})\hat{\mathbf{x}} \right]}{\left[(\mathbf{x}^T \mathbf{x} + 2n_0)(\hat{\mathbf{x}}^T \hat{\mathbf{x}} + 2n_0) - (\hat{\mathbf{x}}^T \mathbf{x})^2 \right]^{\frac{3}{2}}}$$

Again, let's examine $I_2^g(\mathbf{x}_i, \mathbf{x}_1)$.

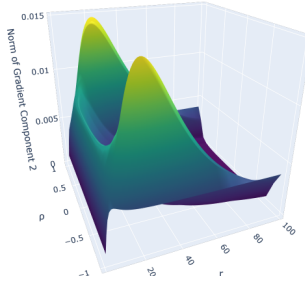
$$\begin{aligned} I_2^g(\mathbf{x}_i, \mathbf{x}_1) &= \frac{1}{2\pi} \frac{(r^2 + 2n_0)^2 \mathbf{x}_i - r^2 \rho_{i1} (2n_0 + r^2) \mathbf{x}_1}{\left[(r^2 + 2n_0)^2 - r^4 \rho_{i1}^2 \right]^{\frac{3}{2}}} \\ \|I_2^g(\mathbf{x}_i, \mathbf{x}_1)\|_2^2 &= \frac{1}{4\pi^2} \frac{r^2 \left[(r^2 + 2n_0)^4 + r^4 \rho_{i1}^2 (2n_0 + r^2)^2 - 2r^2 \rho_{i1}^2 (2n_0 + r^2)^3 \right]}{\left[(r^2 + 2n_0)^2 - r^4 \rho_{i1}^2 \right]^3} \\ &= \frac{1}{4\pi^2} \frac{16n_0^4 + r^2 \left[n_0^3 (32 - 16\rho_{i1}^2) + r^2 \left[n_0^2 (24 - 20\rho_{i1}^2) + r^2 [n_0(8 - 8\rho_{i1}^2) + r^2(1 - \rho_{i1}^2)] \right] \right]}{\left[r^4(1 - \rho_{i1}^2) + 4n_0 r^2 + 4n_0^2 \right]^3} \end{aligned}$$

Based on the equation for $\|I_2^g(\mathbf{x}_i, \mathbf{x}_1)\|_2^2$, we know that if $\rho_{i1}^2 \neq 1$, $\|I_2^g(\mathbf{x}_i, \mathbf{x}_1)\|_2^2$ eventually decays to zero with larger r . But $\|I_2^g(\mathbf{x}_i, \mathbf{x}_1)\|_2^2$ converges to a constant if $\rho_{i1}^2 = 1$. For simplicity, in this section, we do not assume there is any parallel input. Therefore, we can see that all the other terms will go to zero except $I_2^g(\mathbf{x}_1, \mathbf{x}_1)$. It is worth noting that if ρ_{i1} is close to one, the norm will see a spike before going down to zero. But in practice, the data is more than likely to be well separated with small $|\rho_{ij}|$. The discussion here is illustrated in Figure D.1.

Combining the above analysis on the two components of gradient, it is easy to see that with large r ,

$$\frac{\partial \mathbf{k}_x}{\partial \mathbf{x}} \Big|_{\mathbf{x}_1} \approx \begin{bmatrix} \frac{\Theta_\infty^L(\mathbf{x}_1, \mathbf{x})}{\partial \mathbf{x}} \Big|_{\mathbf{x}_1} \\ \mathbf{0} \\ \dots \\ \mathbf{0} \end{bmatrix}$$

$$\left\| \frac{\partial \mathbf{k}_x}{\partial \mathbf{x}} \Big|_{\mathbf{x}_1} \right\|_{op} \approx \left\| \frac{\Theta_\infty^L(\mathbf{x}_1, \mathbf{x})}{\partial \mathbf{x}} \Big|_{\mathbf{x}_1} \right\|_2 \approx \left\| \frac{1}{2\pi} \frac{2n_0(r^2 + 2n_0)}{(4n_0 r^2 + 4n_0^2)^{\frac{3}{2}}} \mathbf{x}_1 \right\|_2 = \frac{1}{2\pi} \frac{2n_0 r (r^2 + 2n_0)}{(4n_0 r^2 + 4n_0^2)^{\frac{3}{2}}} \approx \frac{1}{8\pi \sqrt{n_0}}$$


 Figure D.1: ρ, r vs Norm of Gradient Component 2

D.3. Parallel Inputs Analysis

In the previous section, we assume that there are no parallel inputs. But this assumption is not necessary. In fact, given training data $\{\mathbf{x}_i\}_1^n$, w.l.o.g, let's impose $\mathbf{x}_1 = -\mathbf{x}_2$. Based on the results we have in Section 4.4, we can still derive a similar approximation for the NTK regression solution.

- $\tilde{\mathbf{K}}$

First of all,

$$\mathbf{K}_{ij}^1 = \mathcal{T}(\Theta_\infty^1, \sigma_{\hat{s}}, \sigma_{\hat{s}})(\mathbf{x}_i, \mathbf{x}_j) = \frac{1}{4} + \frac{1}{2\pi} \arcsin\left(\frac{r^2 \rho_{i,j}}{(r^2 + 2n_0)}\right)$$

If $\rho_{i,j} = 1$, then \mathbf{K}_{ij}^1 is going to converge to $\frac{1}{2}$ as r grows bigger. But if $\rho_{i,j} = -1$, this term is going to zero.

Therefore, $\tilde{\mathbf{K}}$ can be approximated by this block diagonal matrix.

$$\tilde{\mathbf{K}} \approx \begin{bmatrix} B_1 & \dots & 0 \\ 0 & B_2 & 0 \\ \dots & \dots & \dots \\ 0 & \dots & B_2 \end{bmatrix}$$

where

$$B_1 = \begin{bmatrix} I_k + \frac{1}{2} & -I_k \\ -I_k & I_k + \frac{1}{2} \end{bmatrix} \quad B_2 = I_k + \frac{1}{2} \quad I_k = \frac{1}{2\pi} \frac{r^2}{\sqrt{4n_0^2 + 4n_0 r^2}} \approx \frac{r}{4\pi\sqrt{n_0}}$$

The inverse of $\tilde{\mathbf{K}}$, is the following, as r grows large:

$$\tilde{\mathbf{K}}^{-1} \approx \begin{bmatrix} B_1^{-1} & \dots & 0 \\ 0 & \frac{1}{I_k + \frac{1}{2}} & 0 \\ \dots & \dots & \dots \\ 0 & \dots & \frac{1}{I_k + \frac{1}{2}} \end{bmatrix} \approx \begin{bmatrix} B_1^{-1} & \dots & 0 \\ 0 & \frac{4\pi\sqrt{n_0}}{r} & 0 \\ \dots & \dots & \dots \\ 0 & \dots & \frac{4\pi\sqrt{n_0}}{r} \end{bmatrix}$$

where

$$B_1^{-1} = \frac{1}{I_k + \frac{1}{4}} \begin{bmatrix} I_k + \frac{1}{2} & I_k \\ I_k & I_k + \frac{1}{2} \end{bmatrix}$$

- $\frac{\partial \mathbf{k}_x}{\partial \mathbf{x}}$

Based on the discussion from Section 4.4,

$$\frac{\partial \mathbf{k}_x}{\partial \mathbf{x}} \Big|_{\mathbf{x}_1} \approx \begin{bmatrix} \frac{\Theta_\infty^L(\mathbf{x}_1, \mathbf{x})}{\partial \mathbf{x}} \Big|_{\mathbf{x}_1} \\ -\frac{\Theta_\infty^L(\mathbf{x}_1, \mathbf{x})}{\partial \mathbf{x}} \Big|_{\mathbf{x}_1} \\ \dots \\ \mathbf{0} \end{bmatrix} = \begin{bmatrix} J_k \mathbf{x}_1 \\ -J_k \mathbf{x}_1 \\ \dots \\ \mathbf{0} \end{bmatrix}$$

where

$$J_k = \frac{1}{2\pi} \frac{2n_0(r^2 + 2n_0)}{(4n_0r^2 + 4n_0^2)^{\frac{3}{2}}} \approx \frac{1}{8\pi\sqrt{n_0}} \frac{1}{r}$$

Finally,

$$\begin{aligned} \left(\hat{\mathbf{X}} - f_0(\hat{\mathbf{X}}) \right) \tilde{\mathbf{K}}^{-1} \frac{\partial k_x}{\partial \mathbf{x}} &\approx \hat{\mathbf{X}} \tilde{\mathbf{K}}^{-1} \frac{\partial k_x}{\partial \mathbf{x}} \approx \hat{\mathbf{X}} \begin{bmatrix} B_1^{-1} & \cdots & 0 \\ 0 & \frac{1}{I_k + \frac{1}{2}} & 0 \\ \cdots & \cdots & \cdots \\ 0 & \cdots & \frac{1}{I_k + \frac{1}{2}} \end{bmatrix} \begin{bmatrix} J_k \mathbf{x}_1 \\ -J_k \mathbf{x}_1 \\ \mathbf{0} \\ \cdots \\ \mathbf{0} \end{bmatrix} \\ &= \hat{\mathbf{X}} \begin{bmatrix} \frac{\frac{1}{2} J_k}{I_k + \frac{1}{4}} \mathbf{x}_1 \\ -\frac{\frac{1}{2} J_k}{I_k + \frac{1}{4}} \mathbf{x}_1 \\ \mathbf{0} \\ \cdots \\ \mathbf{0} \end{bmatrix} = 2 \frac{\frac{1}{2} J_k}{I_k + \frac{1}{4}} \mathbf{x}_1 \mathbf{x}_1^T \end{aligned}$$

Thus,

$$\left\| \left(\hat{\mathbf{X}} - f_0(\hat{\mathbf{X}}) \right) \tilde{\mathbf{K}}^{-1} \frac{\partial \mathbf{k}_x}{\partial \mathbf{x}} \right\|_{op} \approx \frac{r^2 J_k}{I_k + \frac{1}{4}} = \frac{r^2 \frac{1}{8\pi\sqrt{n_0}} \frac{1}{r}}{\frac{r}{4\pi\sqrt{n_0}}} = \frac{1}{2}$$

By similar argument, as $r \rightarrow \infty$, we have

$$\|\mathbf{J}_\infty(\mathbf{x})\|_{op} \leq \frac{1}{2}$$

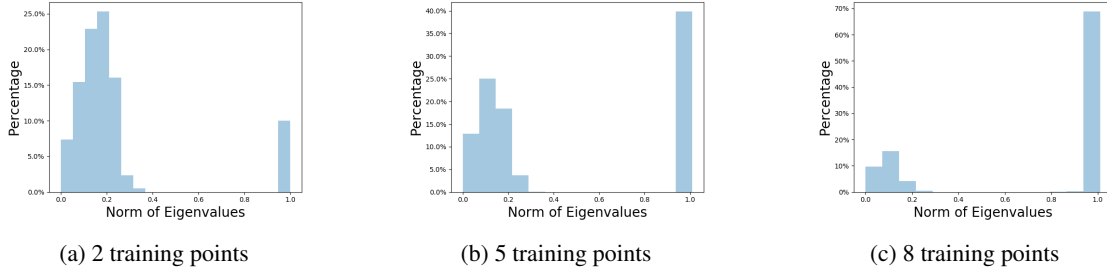


Figure E.1: Eigenvalue distribution of 2-layer sigmoid network trained with input dimension 10

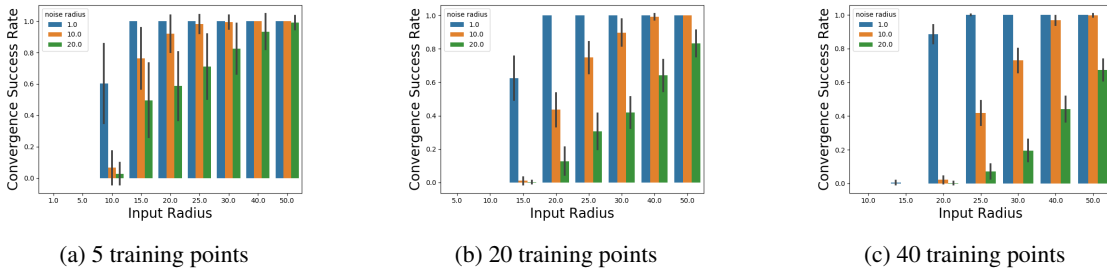


Figure E.2: Convergence success rate vs input norm: random data with input dimension 32

E. Additional Simulations

E.1. Multiple Points: Linear Region

In this section, we first illustrate the eigenvalue distribution in the linear region. Here, we trained 2 layer sigmoid networks with input dimension 10 and hidden size 1000 for 2, 5 and 8 training points. As suggested by Lemma 4, there should be $n - 1$ eigenvalues with norm around 1. This is supported by Figure E.1, as there are 10%, 40% and 70% eigenvalues around that region.

E.2. Basin of Attraction

We test basin of attraction by adding Gaussian noises to training examples and check if the modified examples can converge to the original ones via iterative maps under 50 iterations. The standard deviation of the Gaussian noise is called the noise radius. The network has 2 layers with hidden size 10000 and input dimension 32. Figure E.3 details experiments for 5, 20 and 40 examples. Not surprisingly, the basin of attraction is larger when there are fewer training examples and larger input norms since a level of separation between data is required.

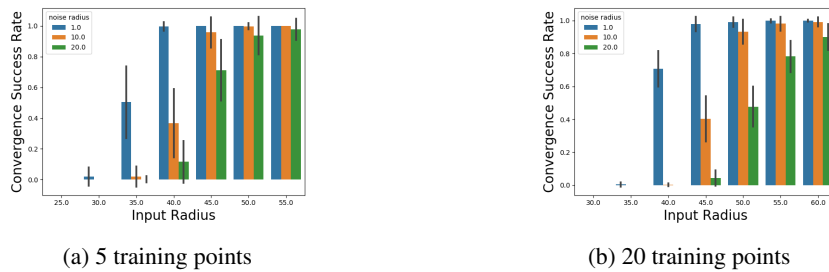


Figure E.3: Convergence success rate vs input norm: MNIST dataset

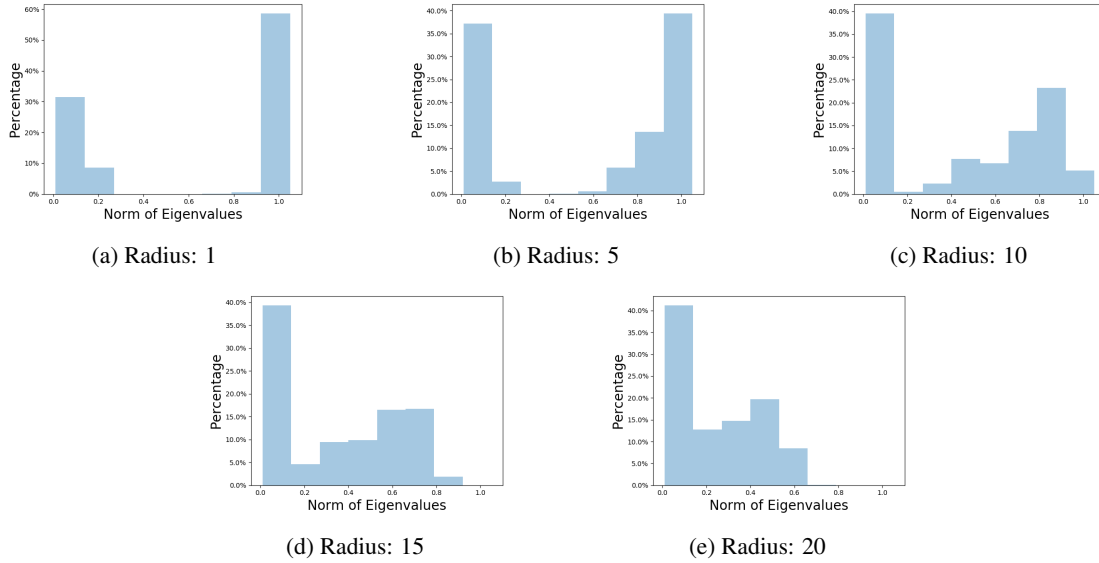


Figure E.4: Spectrum Change for Sigmoid

E.3. Basin of Attraction on Mnist

We also test basin of attraction experiments on MNIST dataset to check if we can recover real training examples. The images are preprocessed by subtracting means and rescaled to have different input norms for testing. Similar to the setting before, Figure 4b also shows that larger input norm gives greater basin of attraction for 5 and 20 examples. Notice that because MNIST images have large input dimension, they need larger radius to move out of the linear region.

E.4. Sigmoidal Activations

Finally, we show that our results can be extended to different sigmoidal activation functions as well. We chose 2 layer network with hidden size 10000, input dimension 32 and 20 training examples. As before, only settings that can let network converges to training loss below 10^{-7} are included. Figure 5 clearly suggests all the activation functions share similar curves. Notice that both tanh and erf have large eigenvalue when r is small. This is not a contradiction to our Lemma 3 as their $\alpha = \dot{\sigma}(0)$ is too large to satisfy the conditions in Lemma 3. The histogram of eigenvalue norm changes for those activation is shown in Figure E.4, Figure E.5, Figure E.6. It is clear that they all follow the same pattern.

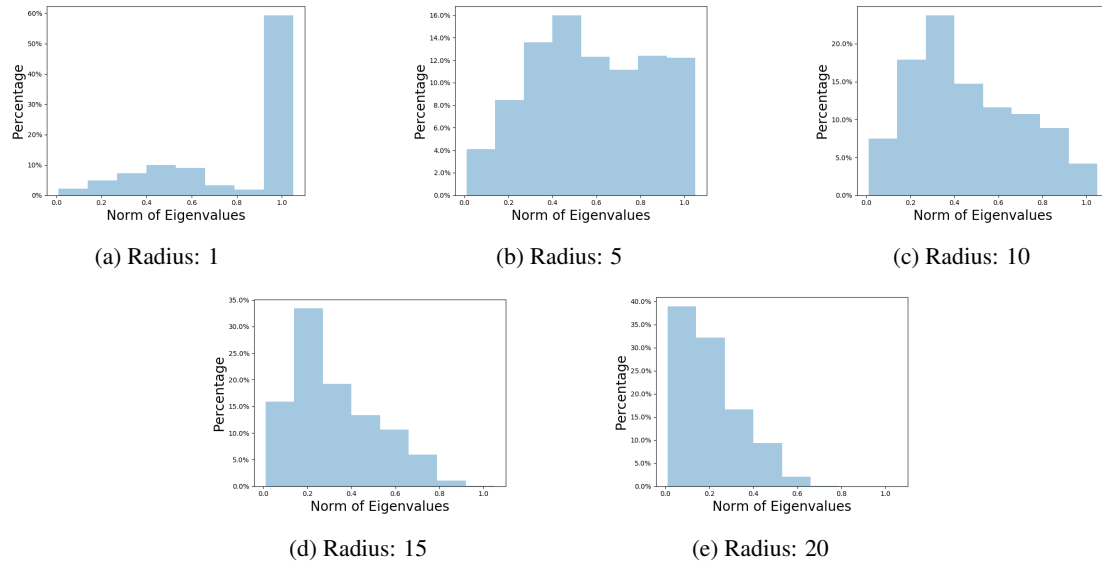


Figure E.5: Specturm Change for Erf

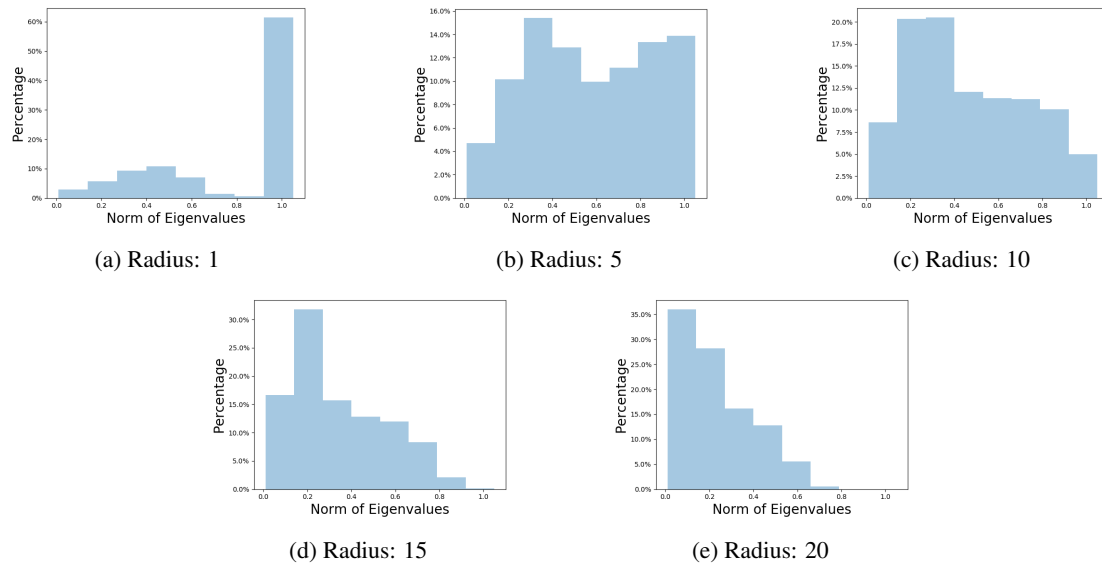


Figure E.6: Specturm Change for Sigmoid