

Understanding Learning Invariance in Deep Linear Networks

Hao Duan¹ Guido Montúfar^{1,2}

Abstract

Equivariant and invariant machine learning models exploit symmetries and structural patterns in data to improve sample efficiency. While empirical studies suggest that data-driven methods such as regularization and data augmentation can perform comparably to explicitly invariant models, theoretical insights remain scarce. In this paper, we provide a theoretical comparison of three approaches for achieving invariance: data augmentation, regularization, and hard-wiring. We focus on mean squared error regression with deep linear networks, which parametrize rank-bounded linear maps and can be hard-wired to be invariant to specific group actions. We show that the critical points of the optimization problems for hard-wiring and data augmentation are identical, consisting solely of saddles and the global optimum. By contrast, regularization introduces additional critical points, though they remain saddles except for the global optimum. Moreover, we demonstrate that the regularization path is continuous and converges to the hard-wired solution.

1. Introduction

Equivariant and invariant models are a class of machine learning models designed to incorporate specific symmetries or invariances that are known to exist in the data. An equivariant model ensures that when the input undergoes a certain transformation, the model’s output transforms in a predictable way. Many powerful hard-wired equivariant and invariant structures have been proposed over the recent years (see, e.g., [Cohen & Welling, 2016](#); [Zaheer et al., 2017](#); [Geiger & Smidt, 2022](#); [Liao et al., 2024](#)). Such models are widely employed and have achieved state-of-the-art level performance across various scientific fields, including condensed-matter physics ([Fang et al., 2023](#)), catalyst de-

sign ([Zitnick et al., 2020](#)), drug discovery ([Igashov et al., 2024](#)), as well as several others.

Given an explicit description of the desired invariance and equivariance structures, a direct way to implement them is by hard-wiring a neural network in a way that constraints the types of functions that it can represent so that they are contained within the desired class. Another intuitive method to approximately enforce invariance and equivariance is data augmentation, where one instead supplies additional data in order to guide the network towards selecting functions from the desired class. Both approaches have shown to be viable for obtaining invariant or equivariant solutions (see, e.g., [Gerken & Kessel, 2024](#); [Moskalev et al., 2023](#)). However, it is not entirely clear how the learning processes and in particular the optimization problems compare. An obvious drawback of data augmentation is that the number of model parameters as well as the number of training data points may be large. On the other hand, it is known that constrained models ([Finzi et al., 2021](#)), or underparameterized models, can have a more complex optimization landscape, but the specific interplay between the amount of data and the structure of the data is not well understood. We are interested in the following question: how do invariance, regularization, and data augmentation influence the optimization process and the resulting solutions of learning? To start developing an understanding, we investigate the simplified setting of invariant linear networks, for which we study the static loss landscape of the three respective optimization problems.

The loss landscapes of neural networks are among the most intriguing and actively studied topics in theoretical deep learning. In particular, a series of works has documented the benefits of overparameterization in making the optimization landscape more benevolent (see, e.g., [Poston et al., 1991](#); [Gori & Tesi, 1992](#); [Soltanolkotabi et al., 2019](#); [Simsek et al., 2021](#); [2023](#); [Karhadkar et al., 2024](#)). This stands at odds with the success of data augmentation, since when using data augmentation as done in practice, even enormous models may no longer be overparameterized and may have fewer parameters than the number of training data points (see, e.g., [Garg et al., 2022](#); [Belkin et al., 2019](#)). Beyond overparameterization, the effects of different architecture choices on the loss landscape are of interest (see, e.g., [Li et al., 2018](#)). As mentioned above, equivariant and invariant architectures are of particular interest, as they could potentially help dra-

¹Department of Statistics & Data Science, University of California, Los Angeles, CA ²Department of Mathematics, University of California, Los Angeles, CA. Correspondence to: Hao Duan <hduan7@ucla.edu>.

matically reduce the sample complexity of learning within a clearly defined framework. This has been documented theoretically in a recent stream of works (see, e.g., [Mei et al., 2021](#); [Tahmasebi & Jegelka, 2023](#)). However, the impact of these architecture choices on the optimization landscape is still underexplored. Equivariant linear networks have received interest as simplified models to obtain concrete and actionable insights for more complex neural networks (see, e.g., [Chen & Zhu, 2023](#); [Kohn et al., 2022](#); [Zhao et al., 2023](#); [Nordenfors et al., 2024](#)).

Our work advances this line of investigation by considering the optimization problems arising from data augmentation, regularization and hard-wiring. We consider linear networks whose end-to-end functions are rank-constrained and thus cannot be simply re-parameterized as linear models. The non-convexity of the function space makes it nontrivial to draw conclusions about the impact of the constraints imposed by invariances. These models are a natural point of departure to study other networks with nonlinear function space, such as networks with nonlinear activation functions. We observe in particular that rank constraints are common in practice. For example, in generative models such as variational autoencoders (VAEs) ([Kingma & Welling, 2022](#)), the hidden layer is usually narrower than the input and output layers with the purpose of capturing a low-dimensional latent representation of the data. In large language models, low-rank adaptation (LoRA) ([Hu et al., 2021](#)) is also used to reduce the number of trainable parameters for downstream tasks. In these and other cases where the architecture has a narrow intermediate linear layer, rank constraints arise.

1.1. Contributions

In this work, we study the impact of invariance in learning by considering and comparing the optimization problems that arise in linear invariant neural networks with a non-linear function space.

- We consider three optimization problems: data augmentation, constrained model, and regularization. We show that these problems are equivalent in terms of their global optima, in the limit of strong regularization and full data augmentation.
- We study the regularization path and show that it continuously connects the global optima of the regularized problem and the global optima of the constrained invariant model.
- We are able to characterize all the critical points in function space for all three problems. In fact, the critical points for data augmentation and the constrained model are the same. There are more critical points for the unconstrained model with regularization.

1.2. Related Work

Loss Landscapes The static optimization landscape of linear networks has been studied in numerous works, whereby most works consider fully-connected networks. In particular, the seminal work of [Baldi & Hornik \(1989\)](#) showed for a two-layer linear network that the square loss has a single minimum up to trivial symmetry and all other critical points are saddles. [Kawaguchi \(2016\)](#) considered the deep case and showed the existence of bad saddles in parameter space for networks with three or more layers. [Laurent & Brecht \(2018\)](#) showed that for deep linear networks with no bottlenecks, all local minima are global for arbitrary convex differentiable losses, and [Zhou & Liang \(2018\)](#) offered a full characterization of the critical points for the square loss. The more recent work of [Trager et al. \(2020\)](#) found that for deep linear networks with bottlenecks, the non-existence of non-global local minima is very particular to the square loss. Several works have also considered more specialized linear network architectures, such as symmetric parametrization ([Tarmoun et al., 2021](#)) and deep linear convolutional networks ([Kohn et al., 2022; 2024a](#)). These and the recent work of ([Shahverdi, 2024](#)) also discuss the critical points in parameter and in function space. In this context we may also highlight the work of [Levin et al. \(2024\)](#), which studies the effect of parametrization on an optimization landscape. In contrast to these works, we focus on deep linear networks with bottlenecks that are invariant to a given group action. The corresponding functions are rank-constrained and thus cannot be simply re-parameterized as linear models.

Training Dynamics Although analyzing the training dynamics is not the main focus of our work, we would like to briefly highlight several related works in this direction. Many works have studied the convergence of parameter optimization in deep linear networks, which remains an interesting topic even in the case of fully-connected layers ([Arora et al., 2018; 2019a;b](#); [Xu et al., 2023](#); [Bah et al., 2021](#); [Br chet et al., 2023](#); [Saxe et al., 2013](#)). For certain types of linear convolutional networks, [Gunasekar et al. \(2018\)](#) studied the implicit bias of parameter optimization. In the context of equivariant models, ([Chen & Zhu, 2023](#)) discuss the implicit bias of gradient flow on linear equivariant steerable networks in group-invariant binary classification.

Invariance, regularization, and data augmentation A few works try to understand the difference between the various aforementioned methods to achieve invariance. [Geiping et al. \(2023\)](#) seek to disentangle the mechanisms through which data augmentation operates and suggest that data augmentation that promotes invariances may provide greater value than enforcing invariance alone, particularly when working with small to medium-sized datasets. Beside data augmentation, [Botev et al. \(2022\)](#) claims that explicit regu-

larization can improve generalization and outperform models that achieve invariance by averaging predictions of non-invariant models. Moskalev et al. (2023) empirically show that the invariance learned by data augmentation deteriorates rapidly, while models with regularization maintain low invariance error even under substantial distribution drift. Our work is inspired by their experiments, and we seek to theoretically study whether data augmentation can learn genuine invariance. A recent work by Kohn et al. (2024b) investigates linear neural networks through the lens of algebraic geometry and computes the dimension, singular points, and the Euclidean distance degree, which serves as an upper bound on the complexity of the optimization problem. We also consider the number of critical points but are primarily interested in the comparison of the loss landscapes arising from different methods. The work of (Gideoni, 2023) investigates the training dynamics of linear regression with data augmentation. In contrast, we consider regression with rank-bounded linear maps and also discuss the effect of regularization. Nordenfors et al. (2024) investigate the optimization dynamics of a neural network with data augmentation and compare it to an invariance hard-wired model. The authors show that the data augmented model and the hard-wired model have the same stationary points within the set of representable equivariant maps \mathcal{E} , but does not offer conclusions about stationary points that are not in \mathcal{E} . In contrast, we obtain a result that describes all critical points in a non-linear function space of rank-constrained linear maps and show that all of them are indeed invariant.

2. Preliminaries

We use $[n]$ to denote the set $\{1, 2, \dots, n\}$. \mathbf{I}_d represents a d by d identity matrix. For any square matrix $U \in \mathbb{C}^{n \times n}$, we use $U_r \in \mathbb{C}^{n \times r}$ to denote the truncation of U to its first r columns. In a slight abuse of notation, for any non-square matrix $\Sigma \in \mathbb{C}^{n \times m}$, we use $\Sigma_r \in \mathbb{C}^{r \times r}$ to denote the truncation of Σ to its first r columns and r rows. For any matrix $M \in \mathbb{C}^{n \times m}$, we denote the Hermitian as M^\dagger , the Moore-Penrose pseudoinverse as M^+ , and the transpose as M^T . We use $\|M\|_2$ and $\|M\|_F$ to denote the operator norm and the Frobenius norm of M , respectively. For a matrix $M \in \mathbb{R}^{n \times m}$, we use $\text{vec}(M)$ to denote the column by column vectorization of M in \mathbb{R}^{nm} . Given any two vector spaces V and W , we use $V \otimes W$ to denote the tensor (Kronecker) product of V and W .

2.1. Equivariance and Invariance

To set up our problem, we need to borrow some concepts from representation theory.

Definition 2.1. A representation of group \mathcal{G} on vector space \mathcal{X} is a homomorphism $\rho: \mathcal{G} \rightarrow GL(\mathcal{X})$, where $GL(\mathcal{X})$ is the group of invertible linear transformations on \mathcal{X} .

Definition 2.2. Let \mathcal{X} and \mathcal{Y} be two vector spaces with representations $\rho_{\mathcal{X}}$ and $\rho_{\mathcal{Y}}$ of the same group \mathcal{G} , respectively. A function $f: \mathcal{X} \rightarrow \mathcal{Y}$ is said to be *equivariant* with respect to $\rho_{\mathcal{X}}$ and $\rho_{\mathcal{Y}}$ if

$$f \circ \rho_{\mathcal{X}}(g) = \rho_{\mathcal{Y}}(g) \circ f, \quad \forall g \in \mathcal{G}. \quad (1)$$

If f is a linear function, we say f is a \mathcal{G} -linear map or a \mathcal{G} -intertwiner. For simplicity of notation, we write $f(gx) = gf(x)$ when $\rho_{\mathcal{X}}$ and $\rho_{\mathcal{Y}}$ are clear.

If $\rho_{\mathcal{Y}}$ is the trivial representation, i.e., $\rho_{\mathcal{Y}}(g)$ is the identity map for all $g \in \mathcal{G}$, then f is said to be *invariant* with respect to $\rho_{\mathcal{X}}$. We then write $f(gx) = f(x)$ when $\rho_{\mathcal{X}}$ is clear.

For a finite cyclic group \mathcal{G} there is a generator $g \in \mathcal{G}$ such that $\mathcal{G} = \{e, g, g^2, \dots, g^{n-1}\}$, where e is the identity element, n is the order of the group, and $g^i = g^j$ whenever $i \equiv j \pmod{n}$.

Example 2.3. For example, the rotational symmetries of a polygon with n sides in \mathbb{R}^2 form a group. The group is a cyclic group of order n , i.e., $\mathcal{G} = C_n$ with generator g , and the representation is generated by $\rho(g) = \begin{bmatrix} \cos \frac{\pi}{n} & -\sin \frac{\pi}{n} \\ \sin \frac{\pi}{n} & \cos \frac{\pi}{n} \end{bmatrix}$.

2.2. Deep Linear Neural Networks

A linear neural network $\Phi(\theta, \mathbf{x})$ with L layers of widths d_1, \dots, d_L is a model of linear functions

$$\Phi(\theta, \mathbf{x}) : \mathbb{R}^{d_\theta} \times \mathbb{R}^{d_x} \rightarrow \mathbb{R}^{d_y}; \quad \mathbf{x} \mapsto W_L \cdots W_1 \mathbf{x}, \quad (2)$$

parameterized by weight matrices $W_j \in \mathbb{R}^{d_j \times d_{j-1}}, \forall j \in [L]$. We write $\theta = (W_L, \dots, W_1) \in \Theta \subseteq \mathbb{R}^{d_\theta}$ for the tuple of weight matrices. The dimension of the *parameter space* Θ is $d_\theta = \sum_{j \in [L]} d_j d_{j-1}$, where $d_0 := d_x, d_L := d_y$ are the input and output dimensions, respectively. For simplicity of the notation, we will write $W := W_L \cdots W_1$ for the end-to-end matrix, and write $W_{j:i} := W_j \cdots W_i$ for the matrix product of layer i up to j for $1 \leq i \leq j \leq L$. We denote the network's parameterization map by

$$\begin{aligned} \mu : \Theta &\rightarrow \mathbb{R}^{d_L \times d_0}; \\ \theta &= (W_1, \dots, W_L) \mapsto W = W_L \cdots W_1. \end{aligned} \quad (3)$$

The network's *function space* is the image of the parametrization map μ , which is the set of linear functions it can represent, i.e., the set of $d_L \times d_0$ matrices of rank at most $r := \min\{d_0, \dots, d_L\}$. We denote the function space by $\mathcal{M}_r \subseteq \mathbb{R}^{d_L \times d_0}$. When $r = \min\{d_0, d_L\}$, the function space is a vector space which can represent any linear function mapping from \mathbb{R}^{d_0} to \mathbb{R}^{d_L} . On the other hand, when $r < \min\{d_0, d_L\}$, it is a non-convex subset of $\mathbb{R}^{d_L \times d_0}$, known as a *determinantal variety* (see Harris, 1992, Chapter 9), which is determined by polynomial constraints, namely

the vanishing of the $(r + 1) \times (r + 1)$ minors. We adopt the following terminology from [Trager et al. \(2020\)](#). The parametrization map μ is *filling* if $r = \min\{d_0, d_L\}$. If $r < \min\{d_0, d_L\}$, then μ is *non-filling*. In the filling case, $\mathcal{M}_r = \mathbb{R}^{d_L \times d_0}$, which is convex. In the non-filling case, $\mathcal{M}_r \subsetneq \mathbb{R}^{d_L \times d_0}$ is a determinantal variety, which is non-convex. We focus primarily on the non-filling case since a convex function space makes the problem less interesting.

Given a group \mathcal{G} , a representation $\rho_{\mathcal{X}}$ on the input space \mathcal{X} and a representation $\rho_{\mathcal{Y}}$ on the output space, an *equivariant linear network* is a linear neural network $\Phi(\theta, \mathbf{x})$ that is equivariant with respect to ρ , i.e., $W_L \cdots W_1 \rho_{\mathcal{X}}(g)x = \rho_{\mathcal{Y}}(g)W_L \cdots W_1 x$ for all $g \in \mathcal{G}$ and $x \in \mathcal{X}$. When $\rho_{\mathcal{Y}}$ is trivial, the network is called an *invariant linear network*. Though we focus on invariant linear networks, it is easy to extend all the results to equivariant linear networks by constructing a new representation taking the tensor product of $\rho_{\mathcal{X}}$ and $\rho_{\mathcal{Y}}$ (see Appendix A.2). In section 4 we will discuss how to define a deep linear network that is hardwired to be invariant to a given group.

2.3. Low Rank Approximation

For a linear network with $r = \min\{d_0, \dots, d_L\}$, the function space consists of $d_L \times d_0$ matrices of rank at most r . The optimization problem in such models is closely related to the problem of approximating a given matrix by a rank bounded matrix.

When the approximation error is measured in Frobenius norm, [Eckart & Young \(1936a\)](#) show that the optimal bounded-rank approximation of a matrix is given in terms of the top components in its singular value decomposition (see, e.g., [Strang, 2019, I.9](#)): If $A = U\Sigma V^T = \sigma_1 u_1 v_1^T + \dots + \sigma_n u_n v_n^T$ and B is any matrix of rank r , then $\|A - B\|_F \geq \|A - A_r\|_F$, where $A_r = \sigma_1 u_1 v_1^T + \dots + \sigma_r u_r v_r^T$. [Mirsky \(1960\)](#) showed that this result in fact holds for any matrix norm that depends only on the singular values. There are several generalizations of this result, for instance to bounded-rank approximation with some fixed entries ([Golub et al., 1987](#)), weighted least squares ([Ruben & Zamir, 1979](#); [Dutta & Li, 2017](#)), and approximation of symmetric matrices by rank-bounded symmetric positive semidefinite matrices ([Dax, 2014](#)). However, for general norms or general matrix constraints, the problem is known to be hard ([Song et al., 2017](#); [Gillis & Shitov, 2019](#)). We are interested in the problem of approximating a given matrix with a rank-bounded matrix that is constrained to within the set of matrices that represent linear maps that are invariant to given group actions.

3. Main Results

3.1. Global Optimum in Constrained Function Space

As we want our function space to contain only the \mathcal{G} -intertwiners, we need to constrain it accordingly. Due to the linearity of the representation $\rho_{\mathcal{X}}$, the constraints are also linear in $\mathbb{R}^{d_L \times d_0}$. Prior research has investigated the constraints for different groups (see, e.g., [Maron et al., 2019](#); [Puny et al., 2023](#); [Finzi et al., 2021](#)). We have the following proposition to explicitly characterize the constraints, proved in Appendix A.1. We will focus on the case where the group \mathcal{G} is finite and cyclic, the representation $\rho_{\mathcal{X}}$ is given and nontrivial, and the representation $\rho_{\mathcal{Y}}$ is trivial.

Proposition 3.1. *Given a cyclic group \mathcal{G} and a representation $\rho_{\mathcal{X}}$ of \mathcal{G} on vector space $\mathcal{X} = \mathbb{R}^{d_0}$, a linear function W mapping from \mathcal{X} to $\mathcal{Y} = \mathbb{R}^{d_L}$ is invariant with respect to $\rho_{\mathcal{X}}$ if and only if $WG = 0$, where $G = \mathbf{I}_{d_0} - \rho_{\mathcal{X}}(g)$, and g is the generator of \mathcal{G} .*

Remark 3.2. Though we assume that \mathcal{G} is cyclic, the above proposition can be generalized to any finitely generated group \mathcal{G} by replacing the single generator g with a set of generators $\{g_1, \dots, g_M\}$. For that, define $G_m = \mathbf{I}_{d_0} - \rho_{\mathcal{X}}(g_m)$ for all $m \in [M]$, and set $G = [G_1, \dots, G_M]$ a $d_0 \times (Md_0)$ matrix. In fact, we can even extend this proposition to continuous groups such as Lie groups. As discussed by [Finzi et al. \(2021\)](#), for any Lie group \mathcal{G} of dimension M with its corresponding Lie algebra \mathfrak{g} , we are able to find a basis $\{A_1, \dots, A_M\}$ for \mathfrak{g} . If the exponential map is surjective in \mathcal{G} , we can then use it to parameterize all elements in \mathcal{G} , i.e., for any $g \in \mathcal{G}$, we can find weights $\{\alpha_m \in \mathbb{R}\}_{m \in [M]}$ such that $g = \exp(\sum_{m=1}^M \alpha_m A_m)$. Therefore, $G_m = d\rho_{\mathcal{X}}(A_m)$ and $G = [G_1, \dots, G_M]$, where $d\rho$ is the Lie algebra representation. See Appendix A.1 for more details.

Consider a data set $\mathcal{D} = \{(\mathbf{x}_i, \mathbf{y}_i)\}_{i=1}^n$, a cyclic group \mathcal{G} , and a representation $\rho_{\mathcal{X}}$ of \mathcal{G} on vector space $\mathcal{X} = \mathbb{R}^{d_0}$. Let $X = \{\mathbf{x}_1, \dots, \mathbf{x}_n\} \in \mathbb{R}^{d_0 \times n}$, $Y = \{\mathbf{y}_1, \dots, \mathbf{y}_n\} \in \mathbb{R}^{d_L \times n}$. Given a positive integer $r < \min\{d_0, d_L\}$, we want to find an invariant linear and rank-bounded function that minimizes the empirical risk, i.e., we want to solve the following optimization problem:

$$\begin{aligned} \widehat{W} &= \arg \min_{W \in \mathbb{R}^{d_L \times d_0}} \frac{1}{n} \|WX - Y\|_F^2, \\ \text{s.t. } & WG = 0, \text{ rank}(W) \leq r. \end{aligned} \quad (4)$$

We assume XX^T has full rank d_0 such that we can use its positive definite square root $P = (XX^T)^{1/2} \in \mathbb{R}^{d_0 \times d_0}$ to derive:

$$\begin{aligned} & \|WX - Y\|_F^2 \\ &= \langle WX, WX \rangle_F - 2\langle WX, Y \rangle_F + \langle Y, Y \rangle_F \\ &= \langle WP, WP \rangle_F - 2\langle WP, YX^T P^{-1} \rangle_F + \langle Y, Y \rangle_F \end{aligned}$$

$$\begin{aligned}
 &= \|WP - YX^T P^{-1}\|_F^2 + \text{const} \\
 &= \|\widetilde{W} - YX^T P^{-1}\|_F^2 + \text{const},
 \end{aligned} \tag{5}$$

where $\widetilde{W} := WP$. We can see that the above optimization problem (4) is equivalent to the following low-rank approximation problem:

$$\begin{aligned}
 \widehat{W} &= \arg \min_{\widetilde{W} \in \mathbb{R}^{d_L \times d_0}} \frac{1}{n} \|\widetilde{W} - Z\|_F^2, \\
 \text{s.t. } \quad &\widetilde{W}\widetilde{G} = 0, \text{ rank}(\widetilde{W}) \leq r,
 \end{aligned} \tag{6}$$

where $Z = YX^T P^{-1}$ and $\widetilde{G} = P^{-1}G$. If we get the solution \widehat{W} , then we can recover the solution to (4) by $\widehat{W} = \widehat{W}P^{-1}$. Since $\widetilde{W}\widetilde{G} = 0$, we know that the rows of \widetilde{W} are in the left null space of \widetilde{G} . Then $\text{rank}(\widetilde{W}) \leq \text{nullity}(\widetilde{G}) = d_0 - \text{rank}(\widetilde{G})$. In order to make this low rank constraints nontrivial, we suppose $r < d := \text{nullity}(\widetilde{G})$. In the case where $r \geq d$, the projection of the unique least square estimator onto the left null space already satisfies the rank constraint, making the rank constraint meaningless. The following theorem characterizes the solution to the above optimization problem, proved in Appendix A.3.

Theorem 3.3. Denote $\bar{Z}^{inv} := Z(\mathbf{I}_{d_0} - \widetilde{G}\widetilde{G}^+)$. We assume $\text{rank}(\bar{Z}^{inv}) > r$. Let $\bar{Z}^{inv} = \bar{U}^{inv} \bar{\Sigma}^{inv} \bar{V}^{invT}$ be the SVD of \bar{Z}^{inv} . Then the solution to (4) is $\widehat{W}^{inv} = \bar{U}_r^{inv} \bar{\Sigma}_r^{inv} \bar{V}_r^{invT} P^{-1}$.

Remark 3.4. The assumption that $\text{rank}(\bar{Z}^{inv}) > r$ is mild. Fix any full row rank data matrix X and suppose $Y = WX + E$, where $E \in \mathbb{R}^{d_L \times n}$ is a random noise matrix. If each column of E is drawn independently from any continuous distribution with full support on \mathbb{R}^{d_L} , then with probability 1, $\text{rank}(\bar{Z}^{inv}) = \min\{d, d_L, d_0\} > r$. In Appendix A.10 we verified this on the MNIST data set.

The key observation is that if the target matrix lives in the invariant linear subspace, then the low-rank approximator of that matrix also lives in the invariant linear subspace. Theorem 3.3 shows how to find the global optima in the optimization problem of constrained space. Indeed, we can project the target matrix to the left null space of \widetilde{G} and find its low-rank approximator.

3.2. Global Optimum in Function Space with Regularization

Instead of imposing constraints on the function space, we can also regularize the optimization problem. We consider the following optimization problem:

$$\widehat{W} = \arg \min_{W \in \mathbb{R}^{d_L \times d_0}} \frac{1}{n} \|WX - Y\|_F^2 + \lambda \|WG\|_F^2,$$

$$\text{s.t. } \text{rank}(W) \leq r. \tag{7}$$

Similarly to optimization problem (4), we can rewrite problem (7) in the following form:

$$\begin{aligned}
 \widehat{W} &= \arg \min_{\widetilde{W} \in \mathbb{R}^{d_L \times d_0}} \frac{1}{n} \|\widetilde{W} - Z\|_F^2 + \lambda \|\widetilde{W}\widetilde{G}\|_F^2, \\
 \text{s.t. } \quad &\text{rank}(\widetilde{W}) \leq r.
 \end{aligned} \tag{8}$$

The optimization problem (8) is referred to as *manifold regularization* (Zhang & Zhao, 2013). In the context of manifold regularization, the input data points are assumed to lie on a low-dimensional manifold embedded in a high-dimensional space. The following proposition, characterizing the solution to the above optimization problem, can be established directly by following the manifold regularization result of Zhang & Zhao (2013, Theorem 1).

Proposition 3.5. Denote $B(\lambda)$ the square root of the symmetric positive definite matrix $\mathbf{I}_{d_0} + n\lambda\widetilde{G}\widetilde{G}^T$, i.e., $B(\lambda)^2 = \mathbf{I}_{d_0} + n\lambda\widetilde{G}\widetilde{G}^T$. Denote $\bar{Z}(\lambda)^{reg} = ZB(\lambda)^{-1}$, and $\bar{Z}(\lambda)^{reg} = \bar{U}(\lambda)^{reg} \bar{\Sigma}(\lambda)^{reg} \bar{V}(\lambda)^{regT}$ as the SVD of $\bar{Z}(\lambda)^{reg}$. Then the solution to problem 7 is $\widehat{W}(\lambda) = \bar{Z}_r(\lambda)^{reg} B(\lambda)^{-1} P^{-1} = \bar{U}_r(\lambda)^{reg} \bar{\Sigma}_r(\lambda)^{reg} \bar{V}_r(\lambda)^{regT} B(\lambda)^{-1} P^{-1}$.

Beside characterizing the global optimum of problem (7), we can also study the regularization path and relate it with the global optimum in the constrained function space. The following theorem states that the regularization path is continuous, and it connects the global optimum in the constrained function space and the global optimum without constraints or regularization. Although the regularization path for ℓ_2 regularization is usually continuous in a vector space, in the case of rank constraints that we consider here the theorem is not trivial.

Theorem 3.6. Assume $\bar{Z}(\lambda)^{reg} = ZB(\lambda)^{-1}$ is full rank for all $\lambda \geq 0$. Then, the regularization path of $\widehat{W}(\lambda)^{reg}$ is continuous on $(0, \infty)$. Moreover, we have $\lim_{\lambda \rightarrow \infty} \widehat{W}^{reg}(\lambda) = \widehat{W}^{inv}$.

Remark 3.7. Similar to Remark 3.4, the assumption that $\bar{Z}(\lambda)^{reg}$ is full rank for all $\lambda \geq 0$ is mild. If we fix any full row rank data matrix X , then $B(\lambda)$ is full rank for all $\lambda \geq 0$. Then, with probability 1, $\bar{Z}(\lambda)^{reg} = ZB(\lambda)^{-1}$ is full rank for all $\lambda \geq 0$.

3.3. Global Optimum in Function Space with Data Augmentation

Data augmentation is another, data-driven, method to achieve invariance. As an informed regularization strategy, it increases the sample size by applying all possible group

actions to the original data. The corresponding optimization problem is then given as follows:

$$\begin{aligned} \widehat{W} &= \arg \min_{W \in \mathbb{R}^{d_L \times d_0}} \frac{1}{n|\mathcal{G}|} \sum_{g \in \mathcal{G}} \|W \rho_{\mathcal{X}}(g)X - Y\|_F^2, \\ \text{s.t. } \text{rank}(W) &\leq r. \end{aligned} \quad (9)$$

We can rewrite the above optimization problem in the following form:

$$\begin{aligned} \widehat{W} &= \arg \min_{\widetilde{W} \in \mathbb{R}^{d_L \times d_0}} \frac{1}{n|\mathcal{G}|} \|\widetilde{W} - |\mathcal{G}|YX^T\bar{G}^TQ^{-1}\|_F^2, \\ \text{s.t. } \text{rank}(\widetilde{W}) &\leq r, \end{aligned} \quad (10)$$

where $\widetilde{W} := WQ$, $\bar{G} = \frac{1}{|\mathcal{G}|} \sum_{g \in \mathcal{G}} \rho_{\mathcal{X}}(g)$, and Q is the square root of the symmetric positive definite matrix $\sum_{g \in \mathcal{G}} \rho_{\mathcal{X}}(g)XX^T\rho_{\mathcal{X}}(g)^T$, i.e., $Q^2 = \sum_{g \in \mathcal{G}} \rho_{\mathcal{X}}(g)XX^T\rho_{\mathcal{X}}(g)^T$. The following proposition characterizes the solution to the above optimization problem (9), proved in Appendix A.6.

Proposition 3.8. Denote $\bar{Z}^{da} = |\mathcal{G}|YX^T\bar{G}^TQ^{-1}$, and $\bar{Z}^{da} = \bar{U}^{da}\bar{\Sigma}^{da}\bar{V}^{daT}$ as the SVD of \bar{Z}^{da} . Then the solution to the above optimization problem (9) is $\widehat{W}^{da} = \bar{Z}_r^{da}Q^{-1} = \bar{U}_r^{da}\bar{\Sigma}_r^{da}\bar{V}_r^{daT}Q^{-1}$. Moreover, if $\rho_{\mathcal{X}}$ is unitary, then \widehat{W}^{da} is an invariant linear map, i.e., $\widehat{W}^{da}G = 0$.

All together, we arrive at the following statement.

Theorem 3.9. Assume $\rho_{\mathcal{X}}$ is unitary. Then the global optima in the function space with data augmentation and the global optima in the constrained function space are the same, i.e., $\widehat{W}^{da} = \widehat{W}^{inv}$.

This theorem tells us that data augmentation and constrained model have the same global optima, which is also the limit of the global optima in the optimization problem with regularization. Beside the global optima, we are also interested in comparing the critical points of all three optimization problems. The following section discusses this in detail.

3.4. Critical Points in the Function Space

We consider a fixed matrix $Z \in \mathbb{R}^{d_L \times d_0}$ with SVD $Z = U\Sigma V^T$. Let $m = \min\{d_0, d_L\}$ and denote by $[m]_r$ the set of all subsets of $[m]$ of cardinality r . For $\mathcal{I} \in [m]_r$, we define $\Sigma_{\mathcal{I}} \in \mathbb{R}^{d_L \times d_0}$ to be the diagonal matrix with entries $\sigma_{\mathcal{I},1}, \sigma_{\mathcal{I},2}, \dots, \sigma_{\mathcal{I},m}$, where $\sigma_{\mathcal{I},i} = \sigma_i$ if $i \in \mathcal{I}$ and $\sigma_{\mathcal{I},i} = 0$ otherwise. Define $\ell_Z(W) := \|Z - W\|_F^2$ as the loss function in the function space \mathcal{M}_r . The function space \mathcal{M}_r is a manifold with singularities. A point $P \in \mathcal{M}_r$ is a critical point of ℓ_Z if and only if $Z - P \in N_P\mathcal{M}_r$. Following Trager et al. (2020, Theorem 28) we can characterize the critical points of the loss function ℓ_Z in the function space \mathcal{M}_r as follows (see Appendix A.8).

Proposition 3.10. Assume all non-zero singular values of \bar{Z}^{inv} , \bar{Z}^{da} , $\bar{Z}(\lambda)^{reg}$ are pairwise distinct.

1. (Constrained Space) The number of critical points in the optimization problem (4) is $\binom{d}{r}$. They are all in the form of $\bar{U}^{inv}\bar{\Sigma}_{\mathcal{I}}^{inv}\bar{V}^{invT}P^{-1}$, where $\mathcal{I} \in [d]_r$. The unique global minimum is $\bar{U}^{inv}\bar{\Sigma}_{[r]}^{inv}\bar{V}^{invT}P^{-1}$, which is also the unique local minimum.
2. (Data Augmentation) The number of critical points in the optimization problem (9) is $\binom{d}{r}$. They are all in the form of $\bar{U}^{da}\bar{\Sigma}_{\mathcal{I}}^{da}\bar{V}^{daT}Q^{-1}$, where $\mathcal{I} \in [d]_r$. These critical points are the same as the critical points in the constrained function space. The unique global minimum is $\bar{U}^{da}\bar{\Sigma}_{[r]}^{da}\bar{V}^{daT}Q^{-1}$, which is also the unique local minimum.
3. (Regularization) The number of critical points in the optimization problem (7) is $\binom{m}{r}$. They are all in the form of $\bar{U}^{reg}\bar{\Sigma}_{\mathcal{I}}^{reg}\bar{V}^{regT}B(\lambda)^{-1}P^{-1}$, where $\mathcal{I} \in [m]_r$. The unique global minimum is $\bar{U}^{reg}\bar{\Sigma}_{[r]}^{reg}\bar{V}^{regT}B(\lambda)^{-1}P^{-1}$, which is also the unique local minimum.

According to this result, we can say that the critical points in the constrained function space are the same as the critical points in the function space with data augmentation. Furthermore, the number of critical points in function space for a model trained with regularization is larger than the number of critical points in the other two cases. We observe that fully-connected linear networks have no spurious local minima, meaning that each local minimum in parameter space corresponds to a local minimum in function space (Trager et al., 2020). This is a consequence of the geometry of determinantal varieties that also holds in our cases, suggesting that also for our three optimization problems there are no spurious local minima.

4. Experiments

4.1. Convergence to an Invariant Critical Point via Data Augmentation

The following experiment demonstrates that *gradient descent* on the optimization problem (9) converges to a critical point that parameterizes an invariant function. The training data, consisting of 1000 samples before data augmentation, is a subset of the MNIST dataset. For computational efficiency, the images are downsampled to 14×14 pixels, resulting in a vectorized representation of dimension 196 for each image. The classification task involves 9 classes, and we aim to train a linear model mapping from \mathbb{R}^{196} to

\mathbb{R}^9 that is invariant under 90-degree rotations. Since digits 6 and 9 are rotationally equivalent, we exclude digit 9 from the dataset. The group associated with this invariance is the cyclic group of order 4, denoted as $\mathcal{G} = C_4$, where the representation $\rho_{\mathcal{X}}$ of \mathcal{G} on \mathbb{R}^{196} is the rotation operator. We employ a data augmentation technique that applies all possible group actions to the original data, yielding a total of 4000 training samples. We also conducted experiments with the same setup but with cross-entropy loss, and the results are similar. Details can be found in Appendix A.11.

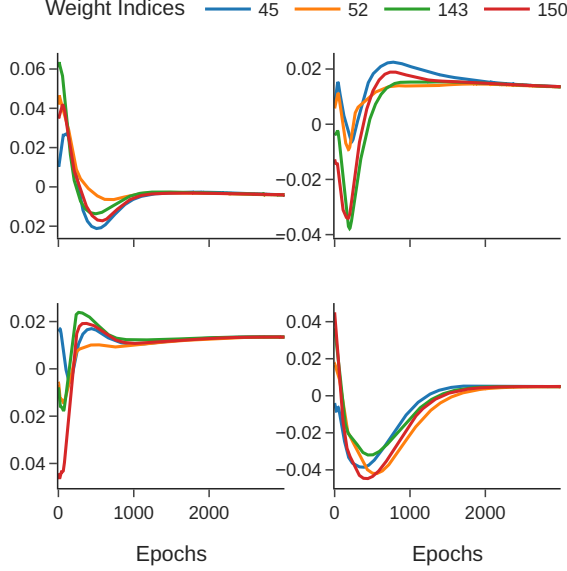


Figure 1. Weights in a two-layer linear neural network trained using data augmentation with mean squared error loss.

We train a two-layer linear neural network with 5 hidden units using the *mean squared error* (MSE) loss function, which parameterizes all $\mathbb{R}^{9 \times 196}$ matrices with rank at most 5. The targets are the one-hot encoded labels. The model is trained using the *Adam* optimizer (Kingma & Ba, 2015) with a learning rate of 0.001 and Adam parameter $\beta = (0.9, 0.999)$, which is the default value in PyTorch (Paszke et al., 2019). The following Figure 1 depicts the evolution of certain entries in the end-to-end matrix W . In our setup, the learned linear map is invariant if and only if specific columns are identical. For example, according to the linear constraints in W (see Proposition 3.1), columns 45, 52, 143, and 150 of W should be exactly the same to achieve invariance. From Figure 1, we observe that the entries in W converge to approximately the same values, indicating that the learned map is nearly invariant.

4.2. Training Curves of Three Approaches

In the same setup as the previous experiment, we compare the performance of the model trained with all three approaches: data augmentation, hard-wiring, and regular-

ization with different choices of the penalty parameter λ . In practice, we parameterize the model in the constrained function space by multiplying a basis matrix B to the weight matrix of the linear model, i.e., $f(x) = W_2 W_1 B x$, where $W_2 \in \mathbb{R}^{9 \times 5}$ and $W_1 \in \mathbb{R}^{5 \times 49}$ are the learnable weight matrices of the linear model, and the basis matrix $B \in \mathbb{R}^{49 \times 196}$ is a matrix that satisfies $BG = 0$. It is worth noting that it is equivalent to perform feature-averaging before feeding the data to the model if we parameterize the invariant function space in this way. Regarding the regularization method, $\lambda \in \{0.001, 0.01, 0.1\}$ when using MSE as the loss.

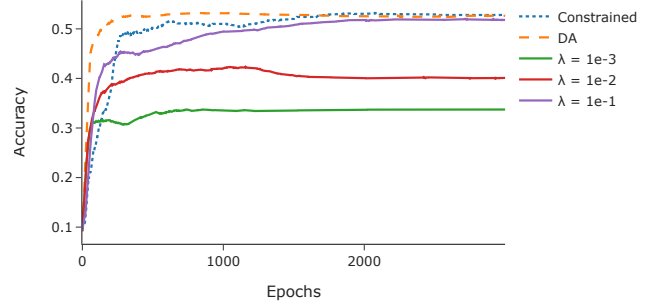


Figure 2. Training curves for data augmentation (DA), regularization (λ), and constrained model under mean squared error loss.

Figure 2 shows the training curves of all three methods under MSE loss. In terms of regularization, though the models are trained without data augmentation, the curves we show here are accuracy for the augmented dataset to make a fair comparison. We can see that data augmentation and hard-wiring have similar performance in the late stage of training. When λ is suitable, regularization can also achieve very similar performance to the previous two methods. All three methods converge to the critical point at a similar rate (around 500 epochs). In fact, when trained with MSE, the hard-wired model converges to the same global optimum as the model trained with data augmentation. This result is consistent with the theoretical analysis in Theorem 3.9. Regarding the amount of time required for training, training with data augmentation is computationally much more expensive than hard-wiring. This is because the model trained with data augmentation requires more samples (4 times more in this case) and more parameters (about 4 times more in this case) than the hard-wired model. Regularization is in between the other two methods since it only requires more parameters but not more samples.

4.3. Comparison between Data Augmentation and Regularization

In this section, we empirically study the training dynamics in both data augmentation and regularization. Using the

same setup as the above experiments, we are showing the evolution of the non-invariant part of the learned end-to-end matrix \widehat{W} . For any \widehat{W} , we can decompose it into two parts, an invariant part and a non-invariant part, i.e., $\widehat{W} = (\widehat{W} - \widehat{W}_\perp) + \widehat{W}_\perp$. A similar decomposition has been used by Gideoni (2023). In Figure 3, we track the evolution of \widehat{W}_\perp by computing $\|\widehat{W}_\perp\|_F$ and $\|\widehat{W} - \widehat{W}_\perp\|_F^2 / \|\widehat{W}\|_F^2$ after each training epoch. When \widehat{W} is very close to an invariant function, $\|\widehat{W}_\perp\|_F$ should be close to 0 and $\|\widehat{W} - \widehat{W}_\perp\|_F^2 / \|\widehat{W}\|_F^2$ should be close to 1 (see Figure 4). Figure 3 shows that with data augmentation $\|\widehat{W}_\perp\|_F$ first increases and then tends to decrease to zero. For regularization, since the penalty coefficient λ is finite, the critical points are actually not invariant. Therefore, in this case we can see that $\|\widehat{W}_\perp\|_F$ does not converge to zero.

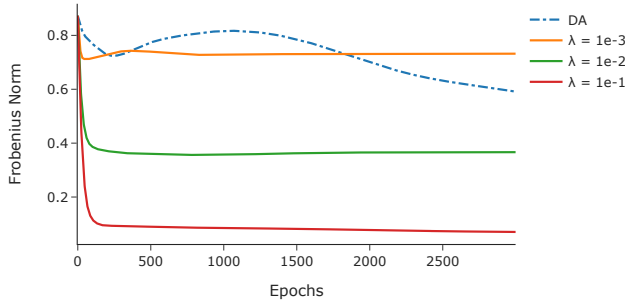


Figure 3. $\|\widehat{W}_\perp\|_F$, where W_\perp is the non-invariant part of W .

Interestingly, we can see that for both data augmentation and regularization, $\|\widehat{W}_\perp\|_F$ displays a “double descent”. Our conjecture is that the loss may also be decomposed into two parts, one controlling the error from invariance, and the other one controlling the error from the target. Therefore, the gradient of the weights during training can be decomposed into two directions as well, resulting in this phenomenon. This intuition could help us better understand the training dynamics and identify methods to accelerate training. Further research needs to be done to investigate this both theoretically and empirically.

5. Conclusion

This work explores learning with invariances from the perspective of the associated optimization problems. We investigate the loss landscape of linear invariant neural networks across the settings of data augmentation, constrained models, and explicit regularization, for which we characterized the form of the global optima (Proposition 3.8, Theorem 3.3, Proposition 3.5). We find that data augmentation and constrained models share the same global optima (Theorem 3.9), which also correspond to the limit of the global optima in the regularized problem (Theorem 3.6). Additionally, the

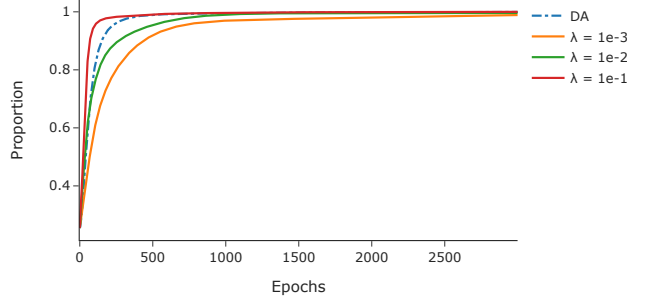


Figure 4. $\|\widehat{W} - \widehat{W}_\perp\|_F^2 / \|\widehat{W}\|_F^2$.

critical points in both data augmentation and constrained models are identical, while regularization generally introduces more critical points (Proposition 3.10). Though our theoretical results are for linear networks, since we consider a non-convex function space, it is natural to conjecture that some of the conclusions might carry over to other models with non-convex function space. Empirical results in Appendix A.12 indicate that data augmentation and a constrained model indeed achieve a similar loss in the late phase of training for two-layer neural networks with different activation functions. At the same time, we observe that for models with a higher expressive power it is more difficult to learn invariance from the data. Based on our theoretical results, we suggest that data augmentation may have similar performance to constrained models, but will incur higher data and computing costs. The regularized model does not require more data and should have a performance close to the constrained model, but it may induce more critical points. The constrained model should have the best performance though one might need to design the invariant architecture carefully before feeding the data to the model.

6. Limitations and future work

We are focusing on deep linear networks, which are a simplified model of neural networks. Nonetheless, we considered the interesting case of rank-bounded end-to-end maps, which is a non-convex function space. Owing to the geometry of this model and the mean squared error loss, the global optima in all three optimization problems are the same. However, this is generally not true when the function class is more complicated or the loss is not the MSE. Moskalev et al. (2023) empirically suggest that data-driven methods fail to learn genuine invariance as in weight-tying shallow ReLU networks for classification tasks with the cross-entropy loss. Experiments suggest that nonlinear networks may still learn invariance via data augmentation when trained with enough data. It would be interesting to investigate this phenomenon theoretically. Furthermore, as mentioned in Section 4.3, the training dynamics of our setup is also worth studying.

Impact Statement

Our work provides a theoretical foundation for incorporating invariances in neural networks, comparing hard-wired constraints, data augmentation, and regularization. We show that constrained models and data augmentation lead to the same global optima, while regularization introduces additional critical points. These findings inform network design by clarifying when to enforce symmetry explicitly versus learning it from data. This has practical implications for models in physics, molecular modeling, and vision, where leveraging symmetry improves efficiency and generalization.

Acknowledgement

This project has been supported by NSF DMS-2145630, NSF CCF-2212520, DFG SPP 2298 project 464109215, and BMBF in DAAD project 57616814.

References

- Arora, S., Cohen, N., and Hazan, E. On the Optimization of Deep Networks: Implicit Acceleration by Overparameterization. In *Proceedings of the 35th International Conference on Machine Learning*, pp. 244–253, 2018.
- Arora, S., Cohen, N., Golowich, N., and Hu, W. A Convergence Analysis of Gradient Descent for Deep Linear Neural Networks. In *International Conference on Learning Representations*, 2019a.
- Arora, S., Cohen, N., Hu, W., and Luo, Y. Implicit Regularization in Deep Matrix Factorization. In *Advances in Neural Information Processing Systems*, volume 32, 2019b.
- Bah, B., Rauhut, H., Terstiege, U., and Westdickenberg, M. Learning deep linear neural networks: Riemannian gradient flows and convergence to global minimizers. *Information and Inference: A Journal of the IMA*, 11(1): 307–353, 2021.
- Baldi, P. and Hornik, K. Neural networks and principal component analysis: Learning from examples without local minima. *Neural Networks*, 2(1):53–58, 1989.
- Belkin, M., Hsu, D., Ma, S., and Mandal, S. Reconciling modern machine-learning practice and the classical bias–variance trade-off. *Proceedings of the National Academy of Sciences*, 116(32):15849–15854, 2019. doi: 10.1073/pnas.1903070116. URL <https://www.pnas.org/doi/abs/10.1073/pnas.1903070116>.
- Botev, A., Bauer, M., and De, S. Regularising for invariance to data augmentation improves supervised learning, 2022. URL <https://arxiv.org/abs/2203.03304>.
- Bourbaki, N. *Integration II: Chapters 7–9*. Springer Berlin Heidelberg, Berlin, Heidelberg, 2004. ISBN 978-3-642-05821-9 978-3-662-07931-7. doi: 10.1007/978-3-662-07931-7. URL <https://link.springer.com/10.1007/978-3-662-07931-7>.
- Bréchet, P., Papagiannouli, K., An, J., and Montúfar, G. Critical points and convergence analysis of generative deep linear networks trained with Bures-Wasserstein loss. In *Proceedings of the 40th International Conference on Machine Learning*, volume 202 of *Proceedings of Machine Learning Research*, pp. 3106–3147. PMLR, 2023. URL <https://proceedings.mlr.press/v202/brechet23a.html>.
- Chen, Z. and Zhu, W. On the implicit bias of linear equivariant steerable networks. In *Advances in Neural Information Processing Systems*, volume 36, pp. 6132–6155. Curran Associates, Inc., 2023. URL <https://openreview.net/forum?id=DnVjDRLwVu>.
- Cohen, T. and Welling, M. Group equivariant convolutional networks. In Balcan, M. F. and Weinberger, K. Q. (eds.), *Proceedings of The 33rd International Conference on Machine Learning*, volume 48 of *Proceedings of Machine Learning Research*, pp. 2990–2999, New York, New York, USA, 20–22 Jun 2016. PMLR. URL <https://proceedings.mlr.press/v48/cohenc16.html>.
- Dax, A. Low-rank positive approximants of symmetric matrices. *Advances in Linear Algebra & Matrix Theory*, 4:172–185, 2014.
- Dieci, L., Gasparo, M. G., and Papini, A. Continuation of singular value decompositions. *Mediterranean Journal of Mathematics*, 2(2):179–203, 2005. ISSN 1660-5454. doi: 10.1007/s00009-005-0038-6. URL <https://doi.org/10.1007/s00009-005-0038-6>.
- Dutta, A. and Li, X. On a problem of weighted low-rank approximation of matrices. *SIAM Journal on Matrix Analysis and Applications*, 38(2):530–553, 2017. doi: 10.1137/15M1043145. URL <https://doi.org/10.1137/15M1043145>.
- Eckart, C. and Young, G. The approximation of one matrix by another of lower rank. *Psychometrika*, 1(3):211–218, 1936a.
- Eckart, C. and Young, G. The approximation of one matrix by another of lower rank. *Psychometrika*, 1:211–218, 1936b. doi: 10.1007/BF02288367. URL <https://doi.org/10.1007/BF02288367>.

- Fang, S., Geiger, M., Checkelsky, J., and Smidt, T. Phonon predictions with e(3)-equivariant graph neural networks. In *AI for Accelerated Materials Design - NeurIPS 2023 Workshop*, 2023. URL <https://openreview.net/forum?id=xxyHjer00Y>.
- Finzi, M., Welling, M., and Wilson, A. G. A practical method for constructing equivariant multilayer perceptrons for arbitrary matrix groups. In *International Conference on Machine Learning*, 2021. URL <https://api.semanticscholar.org/CorpusID:233296901>.
- Garg, S., Jha, S., Mahloujifar, S., Mahmood, M., and Wang, M. Overparameterization from computational constraints. In Oh, A. H., Agarwal, A., Belgrave, D., and Cho, K. (eds.), *Advances in Neural Information Processing Systems*, 2022. URL https://openreview.net/forum?id=7uIGl1AB_M_.
- Geiger, M. and Smidt, T. e3nn: Euclidean neural networks, 2022. URL <https://arxiv.org/abs/2207.09453>.
- Geiping, J., Goldblum, M., Somepalli, G., Shwartz-Ziv, R., Goldstein, T., and Wilson, A. G. How much data are augmentations worth? An investigation into scaling laws, invariance, and implicit regularization. In *The Eleventh International Conference on Learning Representations*, 2023. URL <https://openreview.net/forum?id=3aQs3MCSexD>.
- Gerken, J. E. and Kessel, P. Emergent equivariance in deep ensembles, 2024. URL <https://arxiv.org/abs/2403.03103>.
- Gideoni, Y. Implicitly learned invariance and equivariance in linear regression, 2023. URL <https://openreview.net/pdf?id=ZnxYNriPlg>.
- Gillis, N. and Shitov, Y. Low-rank matrix approximation in the infinity norm. *Linear Algebra and its Applications*, 581:367–382, 2019.
- Golikov, E., Pokonechnyy, E., and Korviakov, V. Neural tangent kernel: A survey, 2022. URL <https://arxiv.org/abs/2208.13614>.
- Golub, G., Hoffman, A., and Stewart, G. A generalization of the Eckart-Young-Mirsky matrix approximation theorem. *Linear Algebra and its Applications*, 88-89:317–327, 1987. URL <https://www.sciencedirect.com/science/article/pii/0024379587901145>.
- Gori, M. and Tesi, A. On the problem of local minima in backpropagation. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 14(1):76–86, 1992. URL <https://ieeexplore.ieee.org/document/107014>.
- Gunasekar, S., Lee, J. D., Soudry, D., and Srebro, N. Implicit bias of gradient descent on linear convolutional networks. In *Advances in Neural Information Processing Systems*, volume 31. Curran Associates, Inc., 2018. URL https://proceedings.neurips.cc/paper_files/paper/2018/file/0e98aeeb54acf612b9eb4e48a269814c-Paper.pdf.
- Harris, J. *Determinantal Varieties*, pp. 98–113. Springer New York, New York, NY, 1992. ISBN 978-1-4757-2189-8. doi: 10.1007/978-1-4757-2189-8_9. URL https://doi.org/10.1007/978-1-4757-2189-8_9.
- Horn, R. A. and Johnson, C. R. *Matrix analysis*. Cambridge University Press, second edition, corrected reprint edition, 2017. ISBN 978-0-521-54823-6 978-0-521-83940-2.
- Hu, E. J., Shen, Y., Wallis, P., Allen-Zhu, Z., Li, Y., Wang, S., Wang, L., and Chen, W. Lora: Low-rank adaptation of large language models, 2021. URL <https://arxiv.org/abs/2106.09685>.
- Igashov, I., Stärk, H., Vignac, C., Schneuing, A., Satorras, V. G., Frossard, P., Welling, M., Bronstein, M., and Correia, B. Equivariant 3d-conditional diffusion model for molecular linker design. *Nature Machine Intelligence*, 6 (4):417–427, 2024.
- Jacot, A., Gabriel, F., and Hongler, C. Neural tangent kernel: Convergence and generalization in neural networks. In Bengio, S., Wallach, H., Larochelle, H., Grauman, K., Cesa-Bianchi, N., and Garnett, R. (eds.), *Advances in Neural Information Processing Systems*, volume 31. Curran Associates, Inc., 2018. URL https://proceedings.neurips.cc/paper_files/paper/2018/file/5a4belfa34e62bb8a6ec6b91d2462f5a-Paper.pdf.
- Karhadkar, K., Murray, M., Tseran, H., and Montúfar, G. Mildly overparameterized ReLU networks have a favorable loss landscape. *Transactions on Machine Learning Research*, 2024. URL <https://openreview.net/forum?id=10WARaIwFn>.
- Kawaguchi, K. Deep learning without poor local minima. In *Advances in Neural Information Processing Systems*, volume 29. Curran Associates, Inc., 2016. URL <https://proceedings.neurips.cc/paper/2016/file/f2fc990265c712c49d51a18a32b39f0c-Paper.pdf>.
- Kingma, D. P. and Ba, J. Adam: A method for stochastic optimization. In *International Conference on Learning Representations*, 2015.

- Kingma, D. P. and Welling, M. Auto-encoding variational bayes, 2022. URL <https://arxiv.org/abs/1312.6114>.
- Kohn, K., Merkh, T., Montúfar, G., and Trager, M. Geometry of linear convolutional networks. *SIAM Journal on Applied Algebra and Geometry*, 6(3):368–406, 2022. doi: 10.1137/21M1441183. URL <https://doi.org/10.1137/21M1441183>.
- Kohn, K., Montúfar, G., Shahverdi, V., and Trager, M. Function space and critical points of linear convolutional networks. *SIAM Journal on Applied Algebra and Geometry*, 8(2):333–362, 2024a. doi: 10.1137/23M1565504. URL <https://doi.org/10.1137/23M1565504>.
- Kohn, K., Sattelberger, A.-L., and Shahverdi, V. Geometry of linear neural networks: Equivariance and invariance under permutation groups, 2024b. URL <https://arxiv.org/abs/2309.13736>.
- Laurent, T. and Brecht, J. Deep Linear Networks with Arbitrary Loss: All Local Minima Are Global. In *Proceedings of the 35th International Conference on Machine Learning*, pp. 2902–2907. PMLR, 2018. URL <https://proceedings.mlr.press/v80/laurent18a.html>.
- Levin, E., Kileel, J., and Boumal, N. The effect of smooth parametrizations on nonconvex optimization landscapes. *Math. Program.*, March 2024. doi: 10.1007/s10107-024-02058-3. URL <https://link.springer.com/10.1007/s10107-024-02058-3>.
- Li, H., Xu, Z., Taylor, G., Studer, C., and Goldstein, T. Visualizing the loss landscape of neural nets. In *Advances in Neural Information Processing Systems*, volume 31. Curran Associates, Inc., 2018. URL https://proceedings.neurips.cc/paper_files/paper/2018/file/a41b3bb3e6b050b6c9067c67f663b915-Paper.pdf.
- Li, Z., Wang, R., Yu, D., Du, S. S., Hu, W., Salakhutdinov, R., and Arora, S. Enhanced convolutional neural tangent kernels, 2019. URL <https://arxiv.org/abs/1911.00809>.
- Liao, Y.-L., Wood, B., Das, A., and Smidt, T. EquiformerV2: Improved Equivariant Transformer for Scaling to Higher-Degree Representations. In *International Conference on Learning Representations (ICLR)*, 2024. URL <https://openreview.net/forum?id=mCOBKZmrzD>.
- Maron, H., Ben-Hamu, H., Shamir, N., and Lipman, Y. Invariant and equivariant graph networks. In *International Conference on Learning Representations*, 2019. URL <https://openreview.net/forum?id=Syx72jC9tm>.
- Mei, S., Misiakiewicz, T., and Montanari, A. Learning with invariances in random features and kernel models. In *Proceedings of Thirty Fourth Conference on Learning Theory*, volume 134 of *Proceedings of Machine Learning Research*, pp. 3351–3418. PMLR, 15–19 Aug 2021. URL <https://proceedings.mlr.press/v134/mei21a.html>.
- Mirsky, L. Symmetric Gauge Functions and Unitary Invariant Norms. *The Quarterly Journal of Mathematics*, 11 (1):50–59, 1960.
- Moskalev, A., Sepiarskaia, A., Bekkers, E. J., and Smeulders, A. On genuine invariance learning without weight-tying. In *ICML workshop on Topology, Algebra, and Geometry in Machine Learning*, 2023.
- Nordenfors, O., Ohlsson, F., and Flinth, A. Optimization dynamics of equivariant and augmented neural networks, 2024. URL <https://arxiv.org/abs/2303.13458>.
- Paszke, A., Gross, S., Massa, F., Lerer, A., Bradbury, J., Chanan, G., Killeen, T., Lin, Z., Gimelshein, N., Antiga, L., Desmaison, A., Kopf, A., Yang, E., DeVito, Z., Raison, M., Tejani, A., Chilamkurthy, S., Steiner, B., Fang, L., Bai, J., and Chintala, S. Pytorch: An imperative style, high-performance deep learning library. In *Advances in Neural Information Processing Systems*, volume 32. Curran Associates, Inc., 2019. URL https://proceedings.neurips.cc/paper_files/paper/2019/file/bdbca288fee7f92f2bfa9f7012727740-Paper.pdf.
- Poston, T., Lee, C.-N., Choie, Y., and Kwon, Y. Local minima and back propagation. In *IJCNN-91-Seattle International Joint Conference on Neural Networks*, volume ii, pp. 173–176 vol.2, 1991. URL <https://ieeexplore.ieee.org/document/155333>.
- Puny, O., Lim, D., Kiani, B., Maron, H., and Lipman, Y. Equivariant polynomials for graph neural networks. In *Proceedings of the 40th International Conference on Machine Learning*, volume 202 of *Proceedings of Machine Learning Research*, pp. 28191–28222. PMLR, 23–29 Jul 2023. URL <https://proceedings.mlr.press/v202/puny23a.html>.
- Ruben, G. and Zamir, S. Lower rank approximation of matrices by least squares with any choice of weights. *Technometrics*, 21(4):489–498, 1979. URL <http://www.jstor.org/stable/1268288>.

- Saxe, A. M., McClelland, J. L., and Ganguli, S. Exact solutions to the nonlinear dynamics of learning in deep linear neural networks. In *Second International Conference on Learning Representations*, 2013. URL https://openreview.net/forum?id=_wzZwKpTDF_9C.
- Shahverdi, V. Algebraic complexity and neurovariety of linear convolutional networks, 2024. URL <https://arxiv.org/abs/2401.16613>.
- Simsek, B., Ged, F., Jacot, A., Spadaro, F., Hongler, C., Gerstner, W., and Brea, J. Geometry of the loss landscape in overparameterized neural networks: Symmetries and invariances. In Meila, M. and Zhang, T. (eds.), *Proceedings of the 38th International Conference on Machine Learning*, volume 139 of *Proceedings of Machine Learning Research*, pp. 9722–9732. PMLR, 18–24 Jul 2021.
- Simsek, B., Bendjeddou, A., Gerstner, W., and Brea, J. Should under-parameterized student networks copy or average teacher weights? In *Thirty-seventh Conference on Neural Information Processing Systems*, 2023. URL <https://openreview.net/forum?id=MG0mYskXN2>.
- Soltanolkotabi, M., Javanmard, A., and Lee, J. D. Theoretical insights into the optimization landscape of over-parameterized shallow neural networks. *IEEE Transactions on Information Theory*, 65(2):742–769, 2019. URL <https://ieeexplore.ieee.org/document/8409482>.
- Song, Z., Woodruff, D. P., and Zhong, P. Low Rank Approximation with Entrywise L1-Norm Error. In *Proceedings of the 49th Annual ACM SIGACT Symposium on Theory of Computing*, STOC 2017, pp. 688–701, New York, NY, USA, 2017. Association for Computing Machinery.
- Strang, G. *Linear Algebra and Learning from Data*. Wellesley-Cambridge Press, Philadelphia, PA, 2019. doi: 10.1137/1.9780692196380. URL <https://epubs.siam.org/doi/abs/10.1137/1.9780692196380>.
- Tahmasebi, B. and Jegelka, S. The exact sample complexity gain from invariances for kernel regression. In *Thirty-seventh Conference on Neural Information Processing Systems*, 2023. URL <https://openreview.net/forum?id=6iouUxI45W>.
- Tarmoun, S., Franca, G., Haeffele, B. D., and Vidal, R. Understanding the Dynamics of Gradient Flow in Overparameterized Linear models. In *Proceedings of the 38th International Conference on Machine Learning*, pp. 10153–10161. PMLR, 2021. URL <https://proceedings.mlr.press/v139/tarmoun21a.html>.
- Trager, M., Kohn, K., and Bruna, J. Pure and spurious critical points: a geometric study of linear networks. In *International Conference on Learning Representations*, 2020. URL <https://openreview.net/forum?id=rkgOlCVYvB>.
- Xu, Z., Min, H., Tarmoun, S., Mallada, E., and Vidal, R. Linear convergence of gradient descent for finite width over-parametrized linear networks with general initialization. In *Proceedings of The 26th International Conference on Artificial Intelligence and Statistics*, volume 206 of *Proceedings of Machine Learning Research*, pp. 2262–2284. PMLR, 2023. URL <https://proceedings.mlr.press/v206/xu23c.html>.
- Zaheer, M., Kottur, S., Ravanbakhsh, S., Poczos, B., Salakhutdinov, R. R., and Smola, A. J. Deep sets. In *Advances in Neural Information Processing Systems*, volume 30. Curran Associates, Inc., 2017. URL https://proceedings.neurips.cc/paper_files/paper/2017/file/f22e4747da1aa27e363d86d40ff442fe-Paper.pdf.
- Zhang, Z. and Zhao, K. Low-rank matrix approximation with manifold regularization. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 35(7):1717–1729, 2013. doi: 10.1109/TPAMI.2012.274.
- Zhao, B., Ganey, I., Walters, R., Yu, R., and Dehmamy, N. Symmetries, flat minima, and the conserved quantities of gradient flow. In *The Eleventh International Conference on Learning Representations*, 2023. URL <https://openreview.net/forum?id=9ZpciCounFb>.
- Zhou, Y. and Liang, Y. Critical points of linear neural networks: Analytical forms and landscape properties. In *International Conference on Learning Representations*, 2018. URL <https://openreview.net/forum?id=SysEexbRb>.
- Zitnick, C. L., Chanussot, L., Das, A., Goyal, S., Heras-Domingo, J., Ho, C., Hu, W., Lavril, T., Palizhati, A., Riviere, M., Shuaibi, M., Sriram, A., Tran, K., Wood, B., Yoon, J., Parikh, D., and Ulissi, Z. An introduction to electrocatalyst design using machine learning for renewable energy storage, 2020. URL <https://arxiv.org/abs/2010.09435>.

A. Appendix

A.1. Proof of Proposition 3.1

Proposition 3.1. *Given a cyclic group \mathcal{G} and a representation $\rho_{\mathcal{X}}$ of \mathcal{G} on vector space $\mathcal{X} = \mathbb{R}^{d_0}$, a linear function W mapping from \mathcal{X} to $\mathcal{Y} = \mathbb{R}^{d_L}$ is invariant with respect to $\rho_{\mathcal{X}}$ if and only if $WG = 0$, where $G = \mathbf{I}_{d_0} - \rho_{\mathcal{X}}(g)$, and g is the generator of \mathcal{G} .*

Proof. Suppose \mathcal{G} is a cyclic group of order k with generator g , i.e., $\mathcal{G} = \langle g \rangle$, $g^k = e$. If W is invariant with respect to $\rho_{\mathcal{X}}$, then $W\rho_{\mathcal{X}}(h) = W$ for all $h \in \mathcal{G}$. Then we have $W(\mathbf{I}_{d_0} - \rho_{\mathcal{X}}(g)) = 0$ for the generator g .

Conversely, if $W(\mathbf{I}_{d_0} - \rho_{\mathcal{X}}(g)) = 0$ for the generator g , then we have $W\rho_{\mathcal{X}}(g) = W$. Multiplying both sides by $\rho_{\mathcal{X}}(g)$, we have $W\rho_{\mathcal{X}}(g^2) = W\rho_{\mathcal{X}}^2(g) = W\rho_{\mathcal{X}}(g) = W$. By induction, we can see that $W\rho_{\mathcal{X}}(g^j) = W$ for all $j \in [k]$. \square

The following proposition extends the above proposition to cases when the group is continuous. The key point is that we can parameterize any element in the continuous group in terms of basis in its corresponding Lie algebra, along with a discrete set of generators.

Proposition A.1. *[Theorem 1 in Finzi et al. (2021)] Let \mathcal{G} be a real connected Lie group of dimension M with finitely many connected components. Given a representation ρ on vector space V of dimension D , the constraint equations*

$$\rho(g)v = v, \forall v \in V, g \in \mathcal{G} \quad (11)$$

holds if and only if

$$d\rho(A_m)v = 0, \quad \forall m \in [M], \quad (12)$$

$$(\rho(h_p) - \mathbf{I}_D)v = 0, \quad \forall p \in [P], \quad (13)$$

where $\{A_m\}_{m=1}^M$ are M basis vectors for the M dimensional Lie Algebra \mathfrak{g} with induced representation $d\rho$, and for some finite collection $\{h_p\}_{p=1}^P$ of discrete generators.

A.2. Extension from invariance to equivariance

Extension from invariance to equivariance is straightforward due to the fact that the constraints are still linear in the vector space of linear maps from \mathcal{X} to \mathcal{Y} . The following proposition shows how to find the linear constraints.

Proposition A.2. *Given a group \mathcal{G} , an input vector space \mathcal{X} with representation $\rho_{\mathcal{X}}$ of \mathcal{G} and an output space \mathcal{Y} with representation $\rho_{\mathcal{Y}}$ of \mathcal{G} , a linear function $f : \mathcal{X} \rightarrow \mathcal{Y}, x \mapsto Wx$ is equivariant with respect to $\rho_{\mathcal{X}}$ and $\rho_{\mathcal{Y}}$ if and only if $\text{vec}(W) \in \bigcap_{g \in \mathcal{G}} \ker \left(\rho_{\mathcal{X}}(g) \otimes \rho_{\mathcal{Y}}(g^{-1})^T - \mathbf{I}_{d_{\mathcal{X}}d_{\mathcal{Y}}} \right)$, where $d_{\mathcal{X}}$ is the dimension of \mathcal{X} and $d_{\mathcal{Y}}$ is the dimension of \mathcal{Y} .*

Proof. By definition, f is equivariant if and only if $W\rho_{\mathcal{X}}(g) = \rho_{\mathcal{Y}}(g)W$ for all $g \in \mathcal{G}$. We can then get $\rho_{\mathcal{Y}}(g^{-1})W\rho_{\mathcal{X}}(g) = W$. By vectorizing both sides, we can see that

$$\text{vec}(\rho_{\mathcal{Y}}(g^{-1})W\rho_{\mathcal{X}}(g)) = \left(\rho_{\mathcal{X}}(g)^T \otimes \rho_{\mathcal{Y}}(g^{-1}) \right) \text{vec}(W) = \text{vec}(W),$$

implying that $\text{vec}(W) \in \bigcap_{g \in \mathcal{G}} \ker \left(\rho_{\mathcal{X}}(g) \otimes \rho_{\mathcal{Y}}(g^{-1})^T - \mathbf{I}_{d_{\mathcal{X}}d_{\mathcal{Y}}} \right)$. \square

A.3. Proof of Theorem 3.3

The following lemma proves a key observation that if a matrix lives in a left null space of another matrix, then the low rank approximator remains in the left null space of the other matrix.

Lemma A.3. *Given a matrix $A \in \mathbb{R}^{n \times m}$ and a matrix $B \in \mathbb{R}^{m \times p}$, $AB = 0$, where $d = \text{nullity}(B)$. Let $A = U\Sigma V^T$ be the SVD of A , where $U \in \mathbb{R}^{n \times n}$, $\Sigma \in \mathbb{R}^{n \times m}$, and $V \in \mathbb{R}^{m \times m}$. Then for any $r \leq \text{rank}(A) \leq d$, V_r^T lives in the left null space of B , namely, $V_r^T B = 0$, and $A_r B = 0$.*

Proof.

$$\begin{aligned}
 A &= U\Sigma V^T, \quad AB = 0 \\
 \Rightarrow U\Sigma V^T B &= 0 \\
 \Rightarrow \Sigma V^T B &= 0 \\
 \Rightarrow \Sigma_d V_d^T B &= 0, \quad d = \text{nullity}(B).
 \end{aligned}$$

Since Σ_d is a diagonal matrix, and the diagonal entries are non-zero, we have that $V_d^T B = 0$. And $V_d = [V_r \quad V_{d-r}]$, we have $V_r^T B = 0$. We can now see that $A_r B = U_r \Sigma_r V_r^T B = 0$. \square

Theorem 3.3. Denote $\bar{Z}^{inv} := Z(\mathbf{I}_{d_0} - \tilde{G}\tilde{G}^+)$. We assume $\text{rank}(\bar{Z}^{inv}) > r$. Let $\bar{Z}^{inv} = \bar{U}^{inv} \bar{\Sigma}^{inv} \bar{V}^{invT}$ be the SVD of \bar{Z}^{inv} . Then the solution to (4) is $\widehat{W}^{inv} = \bar{U}_r^{inv} \bar{\Sigma}_r^{inv} \bar{V}_r^{invT} P^{-1}$.

Proof. As stated in the main text, we can rewrite the optimization problem 4 as the following form:

$$\widehat{W} = \arg \min_{\widetilde{W} \in \mathbb{R}^{d_L \times d_0}} \frac{1}{n} \|\widetilde{W} - Z\|_F^2, \quad \text{s.t.} \quad \widetilde{W}\tilde{G} = 0, \quad \text{rank}(\widetilde{W}) \leq r, \quad (14)$$

where $Z = YX^T P^{-1}$, and $\tilde{G} = P^{-1}G$. There are two cases to consider.

Case 1: $Z\tilde{G} = 0$. We assume Z has rank d . Then we can perform SVD on $Z = U\Sigma V^T = U_d \Sigma_d V_d^T$. Eckart & Young (1936b) have shown that the best rank- r approximation of Z is given by $Z_r = U_r \Sigma_r V_r^T$. According to Lemma A.3, we can see that $Z_r \tilde{G} = 0$. Therefore, the solution to the above optimization problem is $\widehat{W} = Z_r$.

Case 2: $Z\tilde{G} \neq 0$. We can then decompose $Z = \bar{Z} + Z_\perp$, where $\bar{Z}\tilde{G} = 0$, $\langle \bar{Z}, Z_\perp \rangle_F = 0$. Therefore, we can see that

$$\begin{aligned}
 \|\widetilde{W} - Z\|_F^2 &= \|(\widetilde{W} - \bar{Z}) - Z_\perp\|_F^2 \\
 &= \|\widetilde{W} - \bar{Z}\|_F^2 + \|Z_\perp\|_F^2 - 2\langle \widetilde{W} - \bar{Z}, Z_\perp \rangle_F \\
 &= \|\widetilde{W} - \bar{Z}\|_F^2 + \|Z_\perp\|_F^2
 \end{aligned} \quad (15)$$

Thus, the solution to the above optimization problem is

$$\widehat{W} = \arg \min_{\widetilde{W} \in \mathbb{R}^{d_L \times d_0}} \|\widetilde{W} - Z\|_F^2 = \arg \min_{\widetilde{W} \in \mathbb{R}^{d_L \times d_0}} \|\widetilde{W} - \bar{Z}\|_F^2.$$

This is then reduced to the low-rank approximation problem of \bar{Z} , which is the same as in **Case 1**. Let $\bar{Z} = \bar{U}\bar{\Sigma}\bar{V}^T$ be the SVD of \bar{Z} . Then the solution is $\widehat{W} = \bar{U}_r \bar{\Sigma}_r \bar{V}_r^T$.

Note that \bar{Z} can be found by projecting Z onto the left null space of \tilde{G} . An easy construction is $\bar{Z} = Z(\mathbf{I}_{d_0} - \tilde{G}\tilde{G}^+)$. To see this, we can check that $\bar{Z}\tilde{G} = 0$ and $\langle \bar{Z}, Z_\perp \rangle_F = 0$. We have

$$\bar{Z}\tilde{G} = Z(\mathbf{I}_{d_0} - \tilde{G}\tilde{G}^+)\tilde{G} = Z\tilde{G} - Z\tilde{G}\tilde{G}^+\tilde{G} = Z\tilde{G} - Z\tilde{G} = 0. \quad (16)$$

To check $\langle \bar{Z}, Z_\perp \rangle_F = 0$, we have

$$\begin{aligned}
 \langle \bar{Z}, Z_\perp \rangle_F &= \text{tr} [\bar{Z}^T Z_\perp] = \text{tr} [\bar{Z}^T (Z - \bar{Z})] \\
 &= \text{tr} [(\mathbf{I}_{d_0} - \tilde{G}\tilde{G}^+)^T Z^T Z \tilde{G}\tilde{G}^+] \\
 &= \text{tr} [Z^T Z \tilde{G}\tilde{G}^+] - \text{tr} [\tilde{G}\tilde{G}^+ Z^T Z \tilde{G}\tilde{G}^+] \\
 &= \text{tr} [Z^T Z \tilde{G}\tilde{G}^+] - \text{tr} [Z^T Z \tilde{G}\tilde{G}^+ \tilde{G}\tilde{G}^+] \\
 &= \text{tr} [Z^T Z \tilde{G}\tilde{G}^+] - \text{tr} [Z^T Z \tilde{G}\tilde{G}^+] = 0.
 \end{aligned} \quad (17)$$

\square

A.4. Proof of Proposition 3.5

Proposition 3.5. Denote $B(\lambda)$ the square root of the symmetric positive definite matrix $\mathbf{I}_{d_0} + n\lambda\tilde{G}\tilde{G}^T$, i.e., $B(\lambda)^2 = \mathbf{I}_{d_0} + n\lambda\tilde{G}\tilde{G}^T$. Denote $\overline{Z(\lambda)}^{reg} = ZB(\lambda)^{-1}$, and $\overline{Z(\lambda)}^{reg} = \overline{U(\lambda)}^{reg}\overline{\Sigma(\lambda)}^{reg}\overline{V(\lambda)}^{regT}$ as the SVD of $\overline{Z(\lambda)}^{reg}$. Then the solution to problem 7 is $\widehat{W(\lambda)}^{reg} = \overline{Z_r(\lambda)}^{reg}B(\lambda)^{-1}P^{-1} = \overline{U_r(\lambda)}^{reg}\overline{\Sigma_r(\lambda)}^{reg}\overline{V_r(\lambda)}^{regT}B(\lambda)^{-1}P^{-1}$.

Proof. The loss function is defined as:

$$\mathcal{L}(\widetilde{W}) = \frac{1}{n} \|\widetilde{W} - Z\|_F^2 + \lambda \|\widetilde{W}\tilde{G}\|_F^2 \quad (18)$$

$$= \frac{1}{n} \text{tr}[(\widetilde{W} - Z)^T(\widetilde{W} - Z)] + \lambda \text{tr}[(\widetilde{W}\tilde{G})^T(\widetilde{W}\tilde{G})]$$

$$= \frac{1}{n} \text{tr}[\widetilde{W}\widetilde{W}^T - 2\widetilde{W}^TZ + Z^TZ] + \lambda \text{tr}[\widetilde{W}(\tilde{G}\tilde{G}^T)\widetilde{W}^T]$$

$$= \frac{1}{n} \text{tr}[\widetilde{W}(\mathbf{I}_{d_0} + n\lambda\tilde{G}\tilde{G}^T)\widetilde{W}^T - 2\widetilde{W}^TZ + Z^TZ]$$

$$= \frac{1}{n} \text{tr}[\widetilde{W}B(\lambda)B(\lambda)^T\widetilde{W}^T - 2B(\lambda)^T\widetilde{W}^TZB(\lambda)^{-1} + Z^TZ]$$

$$= \frac{1}{n} \|\widetilde{W}B(\lambda) - ZB(\lambda)^{-1}\|_F^2 + \text{const.} \quad (19)$$

Therefore, the optimization problem is equivalent to the following low rank approximation problem:

$$\widehat{W(\lambda)} := \arg \min_{\widetilde{W} \in \mathbb{R}^{d_L \times d_0}} \frac{1}{n} \|\widetilde{W}B(\lambda) - ZB(\lambda)^{-1}\|_F^2, \quad \text{rank}(\widetilde{W}) \leq r \quad (20)$$

$$= \overline{Z_r(\lambda)}^{reg}B(\lambda)^{-1} \quad (21)$$

$$= \overline{U_r(\lambda)}^{reg}\overline{\Sigma_r(\lambda)}^{reg}\overline{V_r(\lambda)}^{regT}B(\lambda)^{-1} \quad (22)$$

Since $\widehat{W} = \widehat{W}P^{-1}$, we have $\widehat{W(\lambda)}^{reg} = \overline{U_r(\lambda)}^{reg}\overline{\Sigma_r(\lambda)}^{reg}\overline{V_r(\lambda)}^{regT}B(\lambda)^{-1}P^{-1}$. \square

A.5. Proof of Theorem 3.6

To prove the theorem, we need the following lemma:

Lemma A.4 (Theorem 2.1 in Dieci et al. (2005)). Let A be a \mathcal{C}^s , $s \geq 1$, matrix valued function, $t \in [0, 1] \rightarrow A(t) \in \mathbb{R}^{m \times n}$, $m \geq n$, of rank n , having p disjoint groups of singular values ($p \leq n$) that vary continuously for all $t : \Sigma_1, \dots, \Sigma_p$. Let $z = m - n$. Consider the function $M \in \mathcal{C}^s([0, 1], \mathbb{R}^{(m+n) \times (m+n)})$ given by

$$M(t) = \begin{bmatrix} 0 & A(t) \\ A^T(t) & 0 \end{bmatrix}. \quad (23)$$

Then, there exists orthogonal $Q \in \mathcal{C}^s([0, 1], \mathbb{R}^{(m+n) \times (m+n)})$ of the form

$$Q(t) = \begin{bmatrix} U_2(t) & U_1(t)/\sqrt{2} & U_1(t)/\sqrt{2} \\ 0 & V(t)/\sqrt{2} & -V(t)/\sqrt{2} \end{bmatrix}, \quad (24)$$

such that

$$Q^T(t)M(t)Q(t) = \begin{bmatrix} 0 & 0 & 0 \\ 0 & S(t) & 0 \\ 0 & 0 & -S(t) \end{bmatrix}, \quad (25)$$

where S is $S = \text{diag}(S_i, i = 1, \dots, p)$, and each S_i is symmetric positive definite, and its eigenvalues coincide with the $\Sigma_i, i = 1, \dots, p$. We have $U_2 \in \mathcal{C}^s([0, 1], \mathbb{R}^{m \times z})$, $U_1 \in \mathcal{C}^s([0, 1], \mathbb{R}^{m \times n})$, and $V \in \mathcal{C}^s([0, 1], \mathbb{R}^{n \times n})$. Equivalently, if we let $U = \begin{bmatrix} U_1 & U_2 \end{bmatrix}$, then

$$U^T(t)A(t)V(t) = \begin{bmatrix} S(t) \\ 0 \end{bmatrix},$$

with the previous form of S .

Theorem 3.6. Assume $\overline{Z(\lambda)}^{reg} = ZB(\lambda)^{-1}$ is full rank for all $\lambda \geq 0$. Then, the regularization path of $\widehat{W(\lambda)}^{reg}$ is continuous on $(0, \infty)$. Moreover, we have $\lim_{\lambda \rightarrow \infty} \widehat{W}^{reg}(\lambda) = \widehat{W}^{inv}$.

Proof. Let $U^{\tilde{G}} \Sigma^{\tilde{G}} V^{\tilde{G}^T}$ be the SVD of \tilde{G} . Since $\text{nullity}(\tilde{G}) = d$, then $\text{rank}(\tilde{G}) = d_0 - d$, suggesting that only the first $d_0 - d$ elements of $\Sigma^{\tilde{G}}$ are non-zero. Denote $\Sigma^{\tilde{G}} = \text{diag}(\sigma_1^{\tilde{G}}, \dots, \sigma_{d_0-d}^{\tilde{G}}, 0, \dots, 0)$, then we have $\tilde{G}^+ = V^{\tilde{G}} \text{diag}(1/\sigma_1^{\tilde{G}}, \dots, 1/\sigma_{d_0-d}^{\tilde{G}}, 0, \dots, 0) U^{\tilde{G}^T}$ according to the property of Moore-Penrose pseudoinverse. Therefore, we have

$$\begin{aligned} \mathbf{I}_{d_0} + n\lambda \tilde{G} \tilde{G}^T &= \mathbf{I}_{d_0} + n\lambda U^{\tilde{G}} \Sigma^{\tilde{G}^2} U^{\tilde{G}^T} = U^{\tilde{G}} \left(\mathbf{I}_{d_0} + n\lambda \Sigma^{\tilde{G}^2} \right) U^{\tilde{G}^T} \\ &= U^{\tilde{G}} \text{diag}(1 + n\lambda \sigma_1^{\tilde{G}^2}, \dots, 1 + n\lambda \sigma_{d_0-d}^{\tilde{G}^2}, 1, \dots, 1) U^{\tilde{G}^T}, \end{aligned} \quad (26)$$

$$\begin{aligned} B(\lambda) &:= (\mathbf{I}_{d_0} + n\lambda \tilde{G} \tilde{G}^T)^{\frac{1}{2}} \\ &= U^{\tilde{G}} \text{diag}(\sqrt{1 + n\lambda \sigma_1^{\tilde{G}^2}}, \dots, \sqrt{1 + n\lambda \sigma_{d_0-d}^{\tilde{G}^2}}, 1, \dots, 1) U^{\tilde{G}^T}, \end{aligned} \quad (27)$$

and

$$\begin{aligned} \lim_{\lambda \rightarrow \infty} B(\lambda)^{-1} &= \lim_{\lambda \rightarrow \infty} (\mathbf{I}_{d_0} + n\lambda \tilde{G} \tilde{G}^T)^{-\frac{1}{2}} \\ &= \lim_{\lambda \rightarrow \infty} U^{\tilde{G}} \text{diag}(1/\sqrt{1 + n\lambda \sigma_1^{\tilde{G}^2}}, \dots, 1/\sqrt{1 + n\lambda \sigma_{d_0-d}^{\tilde{G}^2}}, 1, \dots, 1) U^{\tilde{G}^T}, \\ &= U^{\tilde{G}} \text{diag}(0, \dots, 0, 1, \dots, 1) U^{\tilde{G}^T} \end{aligned} \quad (28)$$

On the other hand, we have

$$\begin{aligned} \mathbf{I}_{d_0} - \tilde{G} \tilde{G}^+ &= \mathbf{I}_{d_0} - U^{\tilde{G}} \text{diag}(1, \dots, 1, 0, \dots, 0) U^{\tilde{G}^T} \\ &= U^{\tilde{G}} \text{diag}(0, \dots, 0, 1, \dots, 1) U^{\tilde{G}^T}. \end{aligned}$$

Thus, we can see that $\lim_{\lambda \rightarrow \infty} B(\lambda)^{-1} = \mathbf{I}_{d_0} - \tilde{G} \tilde{G}^+$.

Recall that $\widehat{W(\lambda)}^{reg} = \overline{Z_r(\lambda)}^{reg} B(\lambda)^{-1} P^{-1} = \overline{U_r(\lambda)}^{reg} \overline{\Sigma_r(\lambda)}^{reg} \overline{V_r(\lambda)}^{reg^T} B(\lambda)^{-1} P^{-1}$ and $\widehat{W}^{inv} = \overline{U}_r^{inv} \overline{\Sigma}_r^{inv} \overline{V}_r^{inv^T}$.

First, we want to show that the regularization path is continuous on $(0, \infty)$. According to Weyl's inequality for singular values, we have the following inequalities:

$$|\sigma_k(\overline{Z(\lambda + \delta)}^{reg}) - \sigma_k(\overline{Z(\lambda)}^{reg})| \leq \|\overline{Z(\lambda + \delta)}^{reg} - \overline{Z(\lambda)}^{reg}\|_2, \quad \forall k \in [\min\{d_0, d_L\}]. \quad (29)$$

On the other hand, we have,

$$\|\overline{Z(\lambda + \delta)}^{reg} - \overline{Z(\lambda)}^{reg}\|_2 \quad (30)$$

$$= \|ZB(\lambda + \delta)^{-1} - ZB(\lambda)^{-1}\|_2 \quad (31)$$

$$= \|ZU^{\tilde{G}} \text{diag}\left(\frac{1}{\sqrt{1 + n\lambda \sigma_1^{\tilde{G}^2}}} - \frac{1}{\sqrt{1 + n(\lambda + \delta) \sigma_1^{\tilde{G}^2}}}, \dots, \right. \quad (32)$$

$$\left. \frac{1}{\sqrt{1 + n\lambda \sigma_{d_0-d}^{\tilde{G}^2}}} - \frac{1}{\sqrt{1 + n(\lambda + \delta) \sigma_{d_0-d}^{\tilde{G}^2}}}, 0, \dots, 0\right) U^{\tilde{G}^T}\|_2 \quad (33)$$

$$\leq \|Z\|_2 \max_{i \in [d_0-d]} \left| \frac{1}{\sqrt{1 + n\lambda \sigma_i^{\tilde{G}^2}}} - \frac{1}{\sqrt{1 + n(\lambda + \delta) \sigma_i^{\tilde{G}^2}}} \right| \rightarrow 0, \quad \text{as } \delta \rightarrow 0. \quad (34)$$

Therefore, the singular values of $\overline{Z(\lambda)}^{reg}$ are continuous with respect to λ on $(0, \infty)$. It is also easy to check that the function $f(\lambda) = \frac{1}{\sqrt{1+c\lambda}}$ is smooth on $[0, \infty)$ for any constant $c > 0$. Applying [Lemma A.4](#) to $\overline{Z(\lambda)}^{reg}$, we find that there exist smooth $\overline{U(\lambda)}^{reg}$ and $\overline{V(\lambda)}^{reg}$ such that $\overline{Z(\lambda)}^{reg} = \overline{U(\lambda)}^{reg} \overline{\Sigma(\lambda)}^{reg} \overline{V(\lambda)}^{regT}$. Thus, by truncating $\overline{U(\lambda)}^{reg}$ and $\overline{V(\lambda)}^{reg}$, $\overline{U_r(\lambda)}^{reg}$ and $\overline{V_r(\lambda)}^{reg}$ are also smooth functions of λ on $(0, \infty)$. Since the singular values are continuous with respect to λ , we have that $\overline{\Sigma_r(\lambda)}^{reg}$ is also continuous on $(0, \infty)$. Then $B(\lambda)$ is continuous on $(0, \infty)$. Since the product of continuous functions is continuous, the regularization path is continuous on $(0, \infty)$.

Finally, we want to show that $\lim_{\lambda \rightarrow \infty} \widehat{W(\lambda)}^{reg} = \widehat{W}^{inv}$ holds. Taking the limit of λ to infinity, we have that $\lim_{\lambda \rightarrow \infty} \overline{Z_r(\lambda)}^{reg} = \lim_{\lambda \rightarrow \infty} ZB(\lambda)^{-1} = Z(\mathbf{I}_{d_0} - \tilde{G}\tilde{G}^+) = \tilde{Z}^{inv}$. According to the continuity of the regularization path, we get $\lim_{\lambda \rightarrow \infty} \overline{U_r(\lambda)}^{reg} \overline{\Sigma_r(\lambda)}^{reg} \overline{V_r(\lambda)}^{regT} = \overline{U}^{inv} \overline{\Sigma}^{inv} \overline{V}^{invT}$.

Due to the fact that $\lim_{\lambda \rightarrow \infty} ZB(\lambda)^{-1}$ lives in the left null space of \tilde{G} , [Lemma A.3](#) tells us that the limit $\lim_{\lambda \rightarrow \infty} \overline{U_r(\lambda)}^{reg} \overline{\Sigma_r(\lambda)}^{reg} \overline{V_r(\lambda)}^{regT}$ also lives in the left null space of \tilde{G} . Thus, we have that

$$\lim_{\lambda \rightarrow \infty} \overline{U_r(\lambda)}^{reg} \overline{\Sigma_r(\lambda)}^{reg} \overline{V_r(\lambda)}^{regT} B(\lambda)^{-1} = \lim_{\lambda \rightarrow \infty} \overline{U_r(\lambda)}^{reg} \overline{\Sigma_r(\lambda)}^{reg} \overline{V_r(\lambda)}^{regT}. \quad (35)$$

The proof is complete. \square

A.6. Proof of Proposition 3.8

To prove the proposition, we need the following lemma:

Lemma A.5. *M and G are both real d by d matrices. G is diagonalizable, and M is positive definite. If $MG = GM$, then $M^{\frac{1}{2}}G = GM^{\frac{1}{2}}$, where $M^{\frac{1}{2}}$ is the positive definite square root of M.*

Proof. Let $M = P\Lambda P^T$ be the eigen decomposition of M . Since M is positive definite, we have that P is orthogonal, and Λ is a diagonal matrix with positive entries. According to theorem 1.3.12 in [Horn & Johnson \(2017\)](#), we know that PGP^T is also diagonal since M and G commute. Write $G = PDP^T$, then $GM^{\frac{1}{2}} = P^T D P P^T \Lambda P = P^T D \Lambda P = P^T \Lambda P P^T D P = M^{\frac{1}{2}}G$. \square

Lemma A.6. *Let $(\mathcal{G}, \mathcal{A}, \lambda)$ be a measure space. Consider a nontrivial representation $\rho_{\mathcal{X}}$ of a compact group \mathcal{G} , let λ be the normalized Haar measure on \mathcal{G} . The existence of the Haar measure is guaranteed by the compactness of \mathcal{G} ([Bourbaki, 2004](#)). Define $\bar{G} := \int_{\mathcal{G}} \rho_{\mathcal{X}}(g) d\lambda(g)$. Then we have the following properties:*

1. $\bar{G}\rho_{\mathcal{X}}(h) = \bar{G}$ for all $h \in \mathcal{G}$.
2. \bar{G} is idempotent, i.e., $\bar{G}^2 = \bar{G}$. That is to say, \bar{G} is a projection operator from \mathcal{X} to the subspace all \mathcal{G} -fixed points.
3. If $\rho_{\mathcal{X}}$ is unitary, i.e., $\rho_{\mathcal{X}}(h)^{\dagger} \rho_{\mathcal{X}}(h) = \mathbf{I}_d$ for all $h \in \mathcal{G}$, then \bar{G} is Hermitian.

Proof.

1. Here, we need to use the fact that the Haar measure is left-invariant, i.e., $\lambda(gA) = \lambda(A)$ for all $g \in \mathcal{G}$ and $A \in \mathcal{A}$. We have

$$\bar{G}\rho_{\mathcal{X}}(h) = \int_{\mathcal{G}} \rho_{\mathcal{X}}(g) d\lambda(g) \rho_{\mathcal{X}}(h) = \int_{\mathcal{G}} \rho_{\mathcal{X}}(gh) d\lambda(g) = \int_{\mathcal{G}} \rho_{\mathcal{X}}(gh) d\lambda(gh) = \bar{G}. \quad (36)$$

2. To show that \bar{G} is idempotent, we have

$$\begin{aligned} \bar{G}^2 &= \left(\int_{\mathcal{G}} \rho_{\mathcal{X}}(g) d\lambda(g) \right) \left(\int_{\mathcal{G}} \rho_{\mathcal{X}}(h) d\lambda(h) \right) = \int_{\mathcal{G}} \int_{\mathcal{G}} \rho_{\mathcal{X}}(g) \rho_{\mathcal{X}}(h) d\lambda(g) d\lambda(h) \\ &= \int_{\mathcal{G}} \int_{\mathcal{G}} \rho_{\mathcal{X}}(gh) d\lambda(g) d\lambda(h) = \int_{\mathcal{G}} \int_{\mathcal{G}} \rho_{\mathcal{X}}(gh) d\lambda(gh) d\lambda(h) = \int_{\mathcal{G}} \bar{G} d\lambda(h) = \bar{G}. \end{aligned} \quad (37)$$

3. To see the last property, we have

$$\bar{G}^\dagger = \int_{\mathcal{G}} \rho_{\mathcal{X}}(g)^\dagger d\lambda(g) = \int_{\mathcal{G}} \rho_{\mathcal{X}}(g)^{-1} d\lambda(g) = \int_{\mathcal{G}} \rho_{\mathcal{X}}(g) d\lambda(g) = \bar{G}. \quad (38)$$

□

Lemma A.7. *Given a finite group \mathcal{G} with order n and a representation ρ of \mathcal{G} on vector space V over field \mathbb{C} , then for every $g \in \mathcal{G}$, there exists a basis P_g in which the matrix of $\rho(g)$ is diagonal for all $g \in \mathcal{G}$, with n -th roots of unity on the diagonal.*

Proof. Since \mathcal{G} is finite with order n , let g be the generator of \mathcal{G} , i.e., $g^n = e$ and $\rho(g)^n = \rho(g^n) = \rho(e) = \mathbf{I}$. We can write $\rho(g)$ in the form of Jordan canonical form, i.e., $\rho(g) = P_g^{-1} J P_g$, where $P_g \in GL(V)$, J is a block diagonal matrix in the following form

$$J = \begin{bmatrix} J_1 & & \\ & \ddots & \\ & & J_p \end{bmatrix},$$

and each block J_i is a square matrix of the form

$$J_i = \begin{bmatrix} \lambda_i & 1 & & \\ & \lambda_i & \ddots & \\ & & \ddots & 1 \\ & & & \lambda_i \end{bmatrix}.$$

We know that $\rho(g)^n = \mathbf{I}$, then $J^n = \mathbf{I}$, which implies that $J_i^n = \mathbf{I}$ for all $i \in [p]$. Let N_i be the Jordan block matrix with $\lambda_i = 0$. Then

$$J_i^n = (\lambda_i \mathbf{I} + N_i)^n = \sum_{k=0}^n \binom{n}{k} \lambda_i^{n-k} N_i^k = \mathbf{I}.$$

Notice that N_i^q is the matrix with zeros and ones only, with the ones in index position (a, b) with $a = b + q$. Therefore, the sum can be \mathbf{I} if and only if $\lambda_i^n = 1$ and $N_i = \mathbf{0}$ for all $i \in [p]$. Therefore, λ_i is an n -th root of unity for all $i \in [p]$, and J_i is diagonal with n -th roots of unity on the diagonal. Let $m \in [n]$, then $\rho(g^m) = \rho(g)^m = P_g^{-1} J^m P_g$. Clearly, J^m is also a diagonal matrix with n -th roots of unity on the diagonal. Therefore, the basis P_g is the same for all $\rho(g^m)$. □

Proposition 3.8. *Denote $\bar{Z}^{da} = |\mathcal{G}| Y X^T \bar{G}^T Q^{-1}$, and $\bar{Z}^{da} = \bar{U}^{da} \bar{\Sigma}^{da} \bar{V}^{daT}$ as the SVD of \bar{Z}^{da} . Then the solution to the above optimization problem (9) is $\widehat{W}^{da} = \bar{Z}_r^{da} Q^{-1} = \bar{U}_r^{da} \bar{\Sigma}_r^{da} \bar{V}_r^{daT} Q^{-1}$. Moreover, if $\rho_{\mathcal{X}}$ is unitary, then \widehat{W}^{da} is an invariant linear map, i.e., $\widehat{W}^{da} G = 0$.*

Proof. It is easy to see that $\widehat{W}^{da} = \bar{U}_r^{da} \bar{\Sigma}_r^{da} \bar{V}_r^{daT} Q^{-1}$ is the solution to the optimization problem 9 since it is in the exact form of a low-rank approximation, and we can apply the Eckart-Young-Mirsky theorem [Eckart & Young \(1936b\)](#) to get the solution directly. We still need to check that \widehat{W}^{da} is an invariant linear map, i.e., $\widehat{W}^{da} G = 0$. We have First, we observe that $\left(\sum_{g \in \mathcal{G}} \rho_{\mathcal{X}}(g) X X^T \rho_{\mathcal{X}}(g)^T \right)^{-1} \rho_{\mathcal{X}}(h) = \rho_{\mathcal{X}}(h) \left(\sum_{g \in \mathcal{G}} \rho_{\mathcal{X}}(g) X X^T \rho_{\mathcal{X}}(g)^T \right)^{-1}$ for all $h \in \mathcal{G}$. To see this, we have

$$\begin{aligned} & \left(\sum_{g \in \mathcal{G}} \rho_{\mathcal{X}}(g) X X^T \rho_{\mathcal{X}}(g)^T \right)^{-1} \rho_{\mathcal{X}}(h) \\ &= \left(\sum_{g \in \mathcal{G}} \rho_{\mathcal{X}}(h^{-1}) \rho_{\mathcal{X}}(g) X X^T \rho_{\mathcal{X}}(g)^T \right)^{-1} \\ &= \rho_{\mathcal{X}}(h^{-1})^T \left(\sum_{g \in \mathcal{G}} \rho_{\mathcal{X}}(h^{-1}) \rho_{\mathcal{X}}(g) X X^T \rho_{\mathcal{X}}(g)^T \rho_{\mathcal{X}}(h^{-1})^T \right)^{-1} \quad \text{unitarity of } \rho_{\mathcal{X}} \end{aligned}$$

$$= \rho_{\mathcal{X}}(h) \left(\sum_{g \in \mathcal{G}} \rho_{\mathcal{X}}(h^{-1}g) X X^T \rho_{\mathcal{X}}(h^{-1}g)^T \right)^{-1}. \quad (39)$$

Then by Lemma A.5, we have $Q^{-1}\rho_{\mathcal{X}}(h) = \rho_{\mathcal{X}}(h)Q^{-1}$. And, we have $\bar{G} = \bar{G}\rho_{\mathcal{X}}(h)$ for all $h \in \mathcal{G}$ by Lemma A.6. Therefore, we have

$$\begin{aligned} \bar{Z}^{da} \rho_{\mathcal{X}}(h) &= |\mathcal{G}| Y X^T \bar{G}^T Q^{-1} \rho_{\mathcal{X}}(h) \\ &= |\mathcal{G}| Y X^T \bar{G}^T \rho_{\mathcal{X}}(h) Q^{-1} \\ &= |\mathcal{G}| Y X^T \bar{G}^T Q^{-1} = \bar{Z}^{da}. \end{aligned} \quad (40)$$

Thus, we can say that $\bar{Z}^{da} G = 0$. Based on Lemma A.3, we can get that $\bar{Z}_r^{da} G = \bar{U}_r^{da} \bar{\Sigma}_r^{da} \bar{V}_r^{da^T} G = 0$. Therefore,

$$\begin{aligned} \widehat{W}^{da} \rho_{\mathcal{X}}(h) &= \bar{U}_r^{da} \bar{\Sigma}_r^{da} \bar{V}_r^{da^T} Q^{-1} \rho_{\mathcal{X}}(h) \\ &= \bar{U}_r^{da} \bar{\Sigma}_r^{da} \bar{V}_r^{da^T} \rho_{\mathcal{X}}(h) Q^{-1} \\ &= \bar{U}_r^{da} \bar{\Sigma}_r^{da} \bar{V}_r^{da^T} Q^{-1} = \widehat{W}^{da}. \end{aligned} \quad (41)$$

□

A.7. Proof of Theorem 3.9

To prove the theorem, we need the following lemma:

Lemma A.8. Let $A = \begin{bmatrix} A_{11} & A_{21}^\dagger \\ A_{21} & A_{22} \end{bmatrix} \in \text{GL}(n+m, \mathbb{C})$ be Hermitian and positive definite and $B \in \text{GL}(n, \mathbb{C})$, where $A_{11} \in \text{GL}(n, \mathbb{C})$ and $A_{22} \in \text{GL}(m, \mathbb{C})$ are both Hermitian and positive definite. Define $E = A \times \begin{bmatrix} B & 0_{n,m} \\ 0_{m,n} & 0_{m,m} \end{bmatrix} = \begin{bmatrix} A_{11}B & 0_{n,m} \\ A_{21}B & 0_{m,m} \end{bmatrix}$. Then $E^+ = \begin{bmatrix} E_{11} & E_{12} \\ 0_{m,n} & 0_{m,m} \end{bmatrix}$, where $E_{11} = B^{-1} (A_{11}^2 + A_{21}^\dagger A_{21})^{-1} A_{11}$, and $E_{12} = B^{-1} (A_{11}^2 + A_{21}^\dagger A_{21})^{-1} A_{21}^\dagger$.

Proof. We need to verify that our solution satisfies the properties of the Moore-Penrose pseudoinverse. Notice the following property:

$$E_{11}A_{11} + E_{12}A_{21} = B^{-1} \quad (42)$$

First, we need to show that $EE^+E = E$ and $E^+EE^+ = E^+$. We have

$$\begin{aligned} EE^+E &= \begin{bmatrix} A_{11}B & 0_{n,m} \\ A_{21}B & 0_{m,m} \end{bmatrix} \begin{bmatrix} E_{11} & E_{12} \\ 0_{m,n} & 0_{m,m} \end{bmatrix} \begin{bmatrix} A_{11}B & 0_{n,m} \\ A_{21}B & 0_{m,m} \end{bmatrix} \\ &= \begin{bmatrix} A_{11}BE_{11} & A_{11}BE_{12} \\ A_{21}BE_{11} & A_{21}BE_{12} \end{bmatrix} \begin{bmatrix} A_{11}B & 0_{n,m} \\ A_{21}B & 0_{m,m} \end{bmatrix} \\ &= \begin{bmatrix} A_{11}B(E_{11}A_{11} + E_{12}A_{21})B & 0_{n,m} \\ A_{21}B(E_{11}A_{11} + E_{12}A_{21})B & 0_{m,m} \end{bmatrix} \\ &= \begin{bmatrix} A_{11}B & 0_{n,m} \\ A_{21}B & 0_{m,m} \end{bmatrix} = E. \end{aligned} \quad (43)$$

Similarly, we want to show that $E^+EE^+ = E^+$. We have

$$\begin{aligned} E^+EE^+ &= \begin{bmatrix} E_{11} & E_{12} \\ 0_{m,n} & 0_{m,m} \end{bmatrix} \begin{bmatrix} A_{11}B & 0_{n,m} \\ A_{21}B & 0_{m,m} \end{bmatrix} \begin{bmatrix} E_{11} & E_{12} \\ 0_{m,n} & 0_{m,m} \end{bmatrix} \\ &= \begin{bmatrix} (E_{11}A_{11} + E_{12}A_{21})B & 0_{n,m} \\ 0_{m,n} & 0_{m,m} \end{bmatrix} \begin{bmatrix} E_{11} & E_{12} \\ 0_{m,n} & 0_{m,m} \end{bmatrix} \end{aligned}$$

$$\begin{aligned}
 &= \begin{bmatrix} \mathbf{I}_n & 0_{n,m} \\ 0_{m,n} & 0_{m,m} \end{bmatrix} \begin{bmatrix} E_{11} & E_{12} \\ 0_{m,n} & 0_{m,m} \end{bmatrix} \\
 &= \begin{bmatrix} E_{11} & E_{12} \\ 0_{m,n} & 0_{m,m} \end{bmatrix} = E^+.
 \end{aligned} \tag{44}$$

We also need to verify that EE^+ and E^+E are Hermitian. We have

$$EE^+ = \begin{bmatrix} A_{11} \left(A_{11}^2 + A_{21}^\dagger A_{21} \right)^{-1} A_{11} & A_{11} \left(A_{11}^2 + A_{21}^\dagger A_{21} \right)^{-1} A_{21}^\dagger \\ A_{21} \left(A_{11}^2 + A_{21}^\dagger A_{21} \right)^{-1} A_{11} & A_{21} \left(A_{11}^2 + A_{21}^\dagger A_{21} \right)^{-1} A_{21}^\dagger \end{bmatrix}, \tag{45}$$

and

$$E^+E = \begin{bmatrix} \mathbf{I}_n & 0_{n,m} \\ 0_{m,n} & 0_{m,m} \end{bmatrix}. \tag{46}$$

It is clear that both EE^+ and E^+E are Hermitian. Therefore, we have shown that E^+ is indeed the Moore-Penrose pseudoinverse of E . \square

Lemma A.9. *Let $Z \in \mathbb{C}^{m \times n}$ be a full-rank matrix. $Q \in \mathbb{C}^{n \times n}$ is Hermitian and positive semi-definite, and $P \in \mathbb{C}^{n \times n}$ satisfying $Q^2 = PP^\dagger$. Given $r < \text{rank}(Q)$, let Z_1 and Z_2 be the best rank- r approximation of ZQ and ZP with respect to the Frobenius norm, respectively, then $Z_1Q = Z_2P^\dagger$.*

Proof. Let $P = USV^\dagger$ be the SVD of P , then we have $Q = USU^\dagger$. Since $ZQ^2 = ZPP^\dagger$, we can see that $ZQUSU^\dagger = ZPV SU^\dagger$. Therefore, we have $ZQ = ZP(VU^\dagger)$. VU^\dagger is a unitary matrix, and according to the rotational invariance of SVD, we can say that $Z_1 = Z_2(VU^\dagger)$, i.e., if $ZP = \tilde{U}\tilde{S}\tilde{V}^\dagger$, then $ZQ = \tilde{U}\tilde{S}(UV^\dagger\tilde{V})^\dagger$, $Z_2 = \tilde{U}_r\tilde{S}_r\tilde{V}_r^\dagger$, and $Z_1 = \tilde{U}_r\tilde{S}_r(UV^\dagger\tilde{V})^\dagger = \tilde{U}_r\tilde{S}_r\tilde{V}_r^\dagger(UV^\dagger)^\dagger$. It is easy to check that $Z_1Q = Z_2P^\dagger$. \square

Theorem 3.9. *Assume ρ_X is unitary. Then the global optima in the function space with data augmentation and the global optima in the constrained function space are the same, i.e., $\bar{W}^{da} = \bar{W}^{inv}$.*

Proof. First, we want to prove that

$$|\mathcal{G}|\bar{G} \left(\sum_{g \in \mathcal{G}} \rho_X(g) X X^\top \rho_X(g)^\top \right)^{-1} = P^{-1} \left(\mathbf{I}_{d_0} - (P^{-1}G)(P^{-1}G)^+ \right) P^{-1} \tag{47}$$

Similar to the proof of [Proposition 3.8](#), we know that $\left(\sum_{g \in \mathcal{G}} \rho_X(g) X X^\top \rho_X(g)^\top \right)^{-1}$ commutes with $\rho_X(g)$ for all $g \in \mathcal{G}$. Then, $\left(\sum_{g \in \mathcal{G}} \rho_X(g) X X^\top \rho_X(g)^\top \right)^{-1}$ commutes with \bar{G} as well. According to [Lemma A.5](#), $Q^{-1} = \left(\sum_{g \in \mathcal{G}} \rho_X(g) X X^\top \rho_X(g)^\top \right)^{-\frac{1}{2}}$ commutes with \bar{G} . We also know that $|\mathcal{G}|\bar{G} \left(\sum_{g \in \mathcal{G}} \rho_X(g) X X^\top \rho_X(g)^\top \right)^{-1}$ is a \mathcal{G} -fixed point. Therefore, we have

$$\begin{aligned}
 &|\mathcal{G}|\bar{G} \left(\sum_{g \in \mathcal{G}} \rho_X(g) X X^\top \rho_X(g)^\top \right)^{-1} = |\mathcal{G}|\bar{G} \left(\sum_{g \in \mathcal{G}} \rho_X(g) X X^\top \rho_X(g)^\top \right)^{-1} \bar{G} \\
 &= |\mathcal{G}|\bar{G}Q^{-1}Q^{-1}\bar{G} = |\mathcal{G}|\bar{G}Q^{-1}\bar{G}Q^{-1} = (|\mathcal{G}|^{\frac{1}{2}}\bar{G}Q^{-1})^2.
 \end{aligned}$$

On the other hand, $\mathbf{I}_{d_0} - (P^{-1}G)(P^{-1}G)^+$ is an idempotent projection matrix. Therefore, we have

$$\begin{aligned}
 &P^{-1} \left(\mathbf{I}_{d_0} - (P^{-1}G)(P^{-1}G)^+ \right) P^{-1} \\
 &= P^{-1} \left(\mathbf{I}_{d_0} - (P^{-1}G)(P^{-1}G)^+ \right) P^{-1} = P^{-1} \left(\mathbf{I}_{d_0} - (P^{-1}G)(P^{-1}G)^+ \right)^2 P^{-1} \\
 &= P^{-1} \left(\mathbf{I}_{d_0} - (P^{-1}G)(P^{-1}G)^+ \right) (P^{-1} \left(\mathbf{I}_{d_0} - (P^{-1}G)(P^{-1}G)^+ \right))^\dagger
 \end{aligned}$$

If Equation 47 holds, then we can apply Lemma A.9 directly to get the result. Therefore, we only need to prove Equation 47.

Let $\rho_{\mathcal{X}}(g) = V\Lambda_g V^{-1}$ be the eigen-decomposition of $\rho_{\mathcal{X}}(g)$, where g is the generator of \mathcal{G} and Λ_g is a diagonal matrix with the eigenvalues of $\rho_{\mathcal{X}}(g)$ on the diagonal. This can be done according to Lemma A.7. Furthermore, under the assumption that $\rho_{\mathcal{X}}$ is unitary, we have $V^{-1} = V^\dagger$. It is worth noting that Λ_g is a diagonal matrix with $|\mathcal{G}|$ -th roots of unity on the diagonal, and among the $|\mathcal{G}|$ -th roots of unity, d of them are 1. Without loss of generality, we assume that the first d eigenvalues are 1. Define $\tilde{X} = V^{-1}X$, and let $\tilde{X}_{1:d}$ be the first d rows of \tilde{X} , and $\tilde{X}_{(d+1):d_0}$ be the last $d_0 - d$ rows of \tilde{X} . Now, let's simplify the LHS of Equation 47:

$$\begin{aligned}\bar{G} &= \frac{1}{|\mathcal{G}|} \sum_{h \in \mathcal{G}} \rho_{\mathcal{X}}(h) = \frac{1}{|\mathcal{G}|} V \left(\sum_{h \in \mathcal{G}} \Lambda_h \right) V^{-1} \\ &= V \left(\frac{1}{|\mathcal{G}|} \sum_{i \in [\mathcal{G}]} \Lambda_g^i \right) V^{-1} = V \begin{bmatrix} \mathbf{I}_d & 0_{d, d_0-d} \\ 0_{d_0-d, d} & 0_{d_0-d, d_0-d} \end{bmatrix} V^{-1},\end{aligned}\quad (48)$$

The last equality in Equation 48 holds because the partial geometric series to order $|\mathcal{G}|$ is 0 for any root of unity other than 1, i.e., $\sum_{j=1}^{|\mathcal{G}|} (e^{\frac{2\pi k i}{|\mathcal{G}|}})^j = 0$ for any $k \neq 0$. On the other hand,

$$\begin{aligned}& \sum_{g \in \mathcal{G}} \frac{1}{|\mathcal{G}|} \rho_{\mathcal{X}}(g) X X^T \rho_{\mathcal{X}}(g)^T \\ &= \sum_{g \in \mathcal{G}} \frac{1}{|\mathcal{G}|} \rho_{\mathcal{X}}(g) X X^\dagger \rho_{\mathcal{X}}(g)^\dagger \\ &= V \left(\sum_{g \in \mathcal{G}} \frac{1}{|\mathcal{G}|} \Lambda_g \tilde{X} \tilde{X}^\dagger \Lambda_g^\dagger \right) V^{-1} \\ &= V \left(\left(\sum_{g \in \mathcal{G}} \frac{1}{|\mathcal{G}|} \text{diag}(\Lambda_g) \text{diag}(\Lambda_g)^\dagger \right) \odot \tilde{X} \tilde{X}^\dagger \right) V^{-1} \\ &= V \left(\begin{bmatrix} \mathbf{I}_d & 0_{d, d_0-d} \\ 0_{d_0-d, d} & \dots \end{bmatrix} \odot \tilde{X} \tilde{X}^\dagger \right) V^{-1} \\ &= V \begin{bmatrix} \tilde{X}_{1:d} \tilde{X}_{1:d}^\dagger & 0_{d, d_0-d} \\ 0_{d_0-d, d} & \dots \end{bmatrix} V^{-1}.\end{aligned}\quad (49)$$

Therefore, the LHS of Equation 47 is

$$\begin{aligned}& \bar{G} \left(\sum_{g \in \mathcal{G}} \frac{1}{|\mathcal{G}|} \rho_{\mathcal{X}}(g) X X^T \rho_{\mathcal{X}}(g)^T \right)^{-1} \\ &= V \begin{bmatrix} \mathbf{I}_d & 0_{d, d_0-d} \\ 0_{d_0-d, d} & 0_{d_0-d, d_0-d} \end{bmatrix} \begin{bmatrix} \tilde{X}_{1:d} \tilde{X}_{1:d}^\dagger & 0_{d, d_0-d} \\ 0_{d_0-d, d} & \dots \end{bmatrix}^{-1} V^{-1} \\ &= V \begin{bmatrix} \left(\tilde{X}_{1:d} \tilde{X}_{1:d}^\dagger \right)^{-1} & 0_{d, d_0-d} \\ 0_{d_0-d, d} & 0_{d_0-d, d_0-d} \end{bmatrix} V^{-1}.\end{aligned}\quad (50)$$

The RHS of Equation 47 is

$$\begin{aligned}& P^{-1} (\mathbf{I}_{d_0} - (P^{-1}G)(P^{-1}G)^+) P^{-1} \\ &= V \tilde{P}^{-1} V^{-1} \left(\mathbf{I}_{d_0} - \left(V \tilde{P}^{-1} (\Lambda_g - \mathbf{I}_{d_0}) V^{-1} \right) \left(V \tilde{P}^{-1} (\Lambda_g - \mathbf{I}_{d_0}) V^{-1} \right)^+ \right) V \tilde{P}^{-1} V^{-1} \\ &= V \tilde{P}^{-1} \left(\mathbf{I}_{d_0} - \left(\tilde{P}^{-1} (\Lambda_g - \mathbf{I}_{d_0}) \right) \left(\tilde{P}^{-1} (\Lambda_g - \mathbf{I}_{d_0}) \right)^+ \right) \tilde{P}^{-1} V^{-1},\end{aligned}\quad (51)$$

where $\tilde{P}^2 = \tilde{X}\tilde{X}^\dagger$.

To prove that the LHS equals the RHS, we need to show that

$$\begin{bmatrix} \left(\tilde{X}_{1:d}\tilde{X}_{1:d}^\dagger\right)^{-1} & 0_{d,d_0-d} \\ 0_{d_0-d,d} & 0_{d_0-d,d_0-d} \end{bmatrix} = \tilde{P}^{-1} \left(\mathbf{I}_{d_0} - \left(\tilde{P}^{-1}(\Lambda_g - \mathbf{I}_{d_0})\right) \left(\tilde{P}^{-1}(\Lambda_g - \mathbf{I}_{d_0})\right)^+ \right) \tilde{P}^{-1}. \quad (54)$$

We can see that

$$\tilde{P}^2 \left(\tilde{P}^{-2} - \begin{bmatrix} \left(\tilde{X}_{1:d}\tilde{X}_{1:d}^\dagger\right)^{-1} & 0_{d,d_0-d} \\ 0_{d_0-d,d} & 0_{d_0-d,d_0-d} \end{bmatrix} \right) \quad (55)$$

$$\begin{aligned} &= \mathbf{I}_{d_0} - \tilde{P}^2 \begin{bmatrix} \left(\tilde{X}_{1:d}\tilde{X}_{1:d}^\dagger\right)^{-1} & 0_{d,d_0-d} \\ 0_{d_0-d,d} & 0_{d_0-d,d_0-d} \end{bmatrix} \\ &= \mathbf{I}_{d_0} - \begin{bmatrix} \tilde{X}_{1:d}\tilde{X}_{1:d}^\dagger & \tilde{X}_{1:d}\tilde{X}_{(d+1):d_0}^\dagger \\ \tilde{X}_{(d+1):d_0}\tilde{X}_{1:d}^\dagger & \tilde{X}_{(d+1):d_0}\tilde{X}_{(d+1):d_0}^\dagger \end{bmatrix} \begin{bmatrix} \left(\tilde{X}_{1:d}\tilde{X}_{1:d}^\dagger\right)^{-1} & 0_{d,d_0-d} \\ 0_{d_0-d,d} & 0_{d_0-d,d_0-d} \end{bmatrix} \\ &= \mathbf{I}_{d_0} - \begin{bmatrix} \mathbf{I}_d & 0_{d,d_0-d} \\ \tilde{X}_{(d+1):d_0}\tilde{X}_{1:d}^\dagger \left(\tilde{X}_{1:d}\tilde{X}_{1:d}^\dagger\right)^{-1} & 0_{d_0-d,d_0-d} \end{bmatrix} \\ &= \begin{bmatrix} 0_{d,d} & 0_{d,d_0-d} \\ -\tilde{X}_{(d+1):d_0}\tilde{X}_{1:d}^\dagger \left(\tilde{X}_{1:d}\tilde{X}_{1:d}^\dagger\right)^{-1} & \mathbf{I}_{d_0-d} \end{bmatrix}. \end{aligned} \quad (56)$$

On the other hand, we rewrite \tilde{P}^{-1} block-wisely, i.e., $\tilde{P}^{-1} = \begin{bmatrix} \tilde{P}_{11} & \tilde{P}_{12} \\ \tilde{P}_{12}^\dagger & \tilde{P}_{22} \end{bmatrix}$. By Lemma A.8, we have

$$(\Lambda_g - \mathbf{I}_{d_0}) \left(\tilde{P}^{-1}(\Lambda_g - \mathbf{I}_{d_0}) \right)^+ \tilde{P}^{-1} \quad (57)$$

$$\begin{aligned} &= \begin{bmatrix} 0_{d,d} & 0_{d,d_0-d} \\ (\tilde{P}_{22}^2 + \tilde{P}_{12}^\dagger \tilde{P}_{12})^{-1} \tilde{P}_{12}^\dagger & (\tilde{P}_{22}^2 + \tilde{P}_{12}^\dagger \tilde{P}_{12})^{-1} \tilde{P}_{22} \end{bmatrix} \begin{bmatrix} \tilde{P}_{11} & \tilde{P}_{12} \\ \tilde{P}_{12}^\dagger & \tilde{P}_{22} \end{bmatrix} \\ &= \begin{bmatrix} 0_{d,d} & 0_{d,d_0-d} \\ (\tilde{P}_{22}^2 + \tilde{P}_{12}^\dagger \tilde{P}_{12})^{-1} (\tilde{P}_{12}^\dagger \tilde{P}_{11} + \tilde{P}_{22} \tilde{P}_{12}^\dagger) & \mathbf{I}_{d_0-d} \end{bmatrix} \end{aligned} \quad (58)$$

By definition, we know that $\tilde{P}^{-2} \tilde{X}\tilde{X}^\dagger = \mathbf{I}_{d_0}$. Therefore,

$$\begin{bmatrix} \tilde{P}_{11} & \tilde{P}_{12} \\ \tilde{P}_{12}^\dagger & \tilde{P}_{22} \end{bmatrix}^2 \begin{bmatrix} \tilde{X}_{1:d}\tilde{X}_{1:d}^\dagger & \tilde{X}_{1:d}\tilde{X}_{(d+1):d_0}^\dagger \\ \tilde{X}_{(d+1):d_0}\tilde{X}_{1:d}^\dagger & \tilde{X}_{(d+1):d_0}\tilde{X}_{(d+1):d_0}^\dagger \end{bmatrix} = \mathbf{I}_{d_0}, \quad (59)$$

$$\begin{bmatrix} \tilde{P}_{11}^2 + \tilde{P}_{12}\tilde{P}_{12}^\dagger & \tilde{P}_{11}\tilde{P}_{12} + \tilde{P}_{12}\tilde{P}_{22} \\ \tilde{P}_{12}^\dagger \tilde{P}_{11} + \tilde{P}_{22}\tilde{P}_{12}^\dagger & \tilde{P}_{22}^2 + \tilde{P}_{12}^\dagger \tilde{P}_{12} \end{bmatrix} \begin{bmatrix} \tilde{X}_{1:d}\tilde{X}_{1:d}^\dagger & \tilde{X}_{1:d}\tilde{X}_{(d+1):d_0}^\dagger \\ \tilde{X}_{(d+1):d_0}\tilde{X}_{1:d}^\dagger & \tilde{X}_{(d+1):d_0}\tilde{X}_{(d+1):d_0}^\dagger \end{bmatrix} = \mathbf{I}_{d_0}. \quad (60)$$

By equating the LHS and RHS of the above equation, we can get that

$$\begin{aligned} &(\tilde{P}_{12}^\dagger \tilde{P}_{11} + \tilde{P}_{22} \tilde{P}_{12}^\dagger) \tilde{X}_{1:d}\tilde{X}_{1:d}^\dagger + (\tilde{P}_{22}^2 + \tilde{P}_{12}^\dagger \tilde{P}_{12}) \tilde{X}_{(d+1):d_0}\tilde{X}_{1:d}^\dagger = 0_{d_0-d,d}, \\ &-\tilde{X}_{(d+1):d_0}\tilde{X}_{1:d}^\dagger \left(\tilde{X}_{1:d}\tilde{X}_{1:d}^\dagger\right)^{-1} = (\tilde{P}_{22}^2 + \tilde{P}_{12}^\dagger \tilde{P}_{12})^{-1} (\tilde{P}_{12}^\dagger \tilde{P}_{11} + \tilde{P}_{22} \tilde{P}_{12}^\dagger) \end{aligned} \quad (61)$$

We have shown that the LHS equals the RHS in Equation 47. The theorem is proved. \square

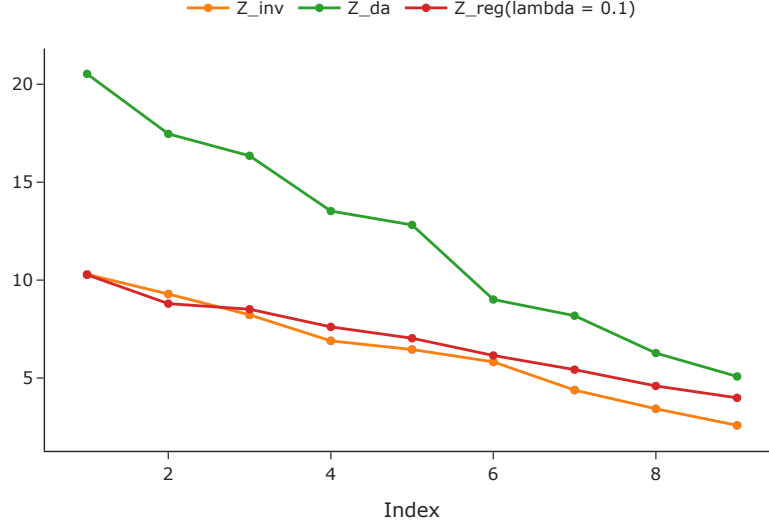


Figure 5. The spectrum of target matrices in the MNIST dataset.

A.8. Proof of Proposition 3.10

Proposition A.10. Suppose the target matrix $Z \in \mathbb{R}^{d_L \times d_0}$ has rank $m > d > r$. The critical points of ℓ_Z restricted to the function space \mathcal{M}_r are all matrices of the form $U\Sigma_{\mathcal{I}}V^T$ where $\mathcal{I} \in [d]_r$. If $0 < \sigma_{r+1} < \sigma_r$, then the local minimum is the critical point with $\mathcal{I} = [r]$. It is the global minimum.

The proof is adapted from the proof of (Trager et al., 2020, Theorem 28).

Proof. A matrix $P \in \mathcal{M}_r$ is a critical point if and only if $Z - P \in N_P\mathcal{M}_r = \text{Col}(P)^\perp \otimes \text{Row}(P)^\perp$, where $N_P\mathcal{M}_r$ denotes the normal space of \mathcal{M}_r at point P . If $P = \sum_{i=1}^r \sigma'_i (u'_i \otimes v'_i)$ and $Z - P = \sum_{j=1}^e \sigma''_j (u''_j \otimes v''_j)$ are SVD with $\sigma'_i \neq 0$ and $\sigma''_j \neq 0$, the column spaces of P and $Z - P$ are spanned by the u'_i and u''_j , respectively. Similarly, the row spaces of P and $Z - P$ are spanned by the v'_i and v''_j , respectively. So P is a critical point if and only if the vectors u'_i, u''_j and v'_i, v''_j are orthonormal, i.e., if

$$Z = P + (Z - P) = \sum_{i=1}^r \sigma'_i (u'_i \otimes v'_i) + \sum_{j=1}^e \sigma''_j (u''_j \otimes v''_j)$$

is a SVD of Z . This proves that the critical points are of the form $U\Sigma_{\mathcal{I}}V^T$ where $Z = U\Sigma V^T$ is a SVD and $\mathcal{I} \in [d]_r$. Since $\ell_Z(U\Sigma_{\mathcal{I}}V^T) = \|U\Sigma_{[d] \setminus \mathcal{I}}V^T\|^2 = \|\Sigma_{[d] \setminus \mathcal{I}}\|^2 = \sum_{i \notin \mathcal{I}} \sigma_i^2$, we see that the global minima are exactly the critical points selecting r of the largest singular values of Z , i.e., with $\mathcal{I} = [r]$. It is left to show that there are no other local minima. For this, we consider a critical point $P = U\Sigma_{\mathcal{I}}V^T$ such that at least one selected singular value σ_i for $i \in \mathcal{I}$ is strictly smaller than σ_r . This is possible since $0 < \sigma_{r+1} < \sigma_r$. To see that P cannot be a local minimum, one can follow the proofs in (Trager et al., 2020, Theorem 28). \square

Proposition 3.10. Assume all non-zero singular values of $\bar{Z}^{inv}, \bar{Z}^{da}, \bar{Z}(\lambda)^{reg}$ are pairwise distinct.

1. (Constrained Space) The number of critical points in the optimization problem (4) is $\binom{d}{r}$. They are all in the form of $\bar{U}^{inv} \bar{\Sigma}_{\mathcal{I}}^{inv} \bar{V}^{invT} P^{-1}$, where $\mathcal{I} \in [d]_r$. The unique global minimum is $\bar{U}^{inv} \bar{\Sigma}_{[r]}^{inv} \bar{V}^{invT} P^{-1}$, which is also the unique local minimum.
2. (Data Augmentation) The number of critical points in the optimization problem (9) is $\binom{d}{r}$. They are all in the form of $\bar{U}^{da} \bar{\Sigma}_{\mathcal{I}}^{da} \bar{V}^{daT} Q^{-1}$, where $\mathcal{I} \in [d]_r$. These critical points are the same as the critical points in the constrained function space. The unique global minimum is $\bar{U}^{da} \bar{\Sigma}_{[r]}^{da} \bar{V}^{daT} Q^{-1}$, which is also the unique local minimum.

3. (Regularization) The number of critical points in the optimization problem (7) is $\binom{m}{r}$. They are all in the form of $\bar{U}^{reg} \bar{\Sigma}_{\mathcal{I}}^{reg} \bar{V}^{regT} B(\lambda)^{-1} P^{-1}$, where $\mathcal{I} \in [m]_r$. The unique global minimum is $\bar{U}^{reg} \bar{\Sigma}_{[r]}^{reg} \bar{V}^{regT} B(\lambda)^{-1} P^{-1}$, which is also the unique local minimum.

Proof. This follows directly from Proposition A.10 and the fact that \bar{Z}^{da} and \bar{Z}^{inv} are both rank d matrices while \bar{Z}^{reg} has rank m . \square

A.9. Extension to Shallow Nonlinear Networks

Consider a two-layer homogeneous nonlinear network with d_0 input units, d_1 hidden units, and $d_L = 1$ output units for simplicity. The output of the network is given by

$$f(\mathbf{x}; \Theta) = \frac{1}{\sqrt{d_1}} \sum_{d=1}^{d_1} a_d \sigma(\mathbf{w}_d^T \mathbf{x}), \quad (62)$$

where the parameters are initialized as $a_d \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(0, 1)$, $\mathbf{w}_d \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(0, \mathbf{I}_{d_0})$, and $\sigma : \mathbb{R} \rightarrow \mathbb{R}$ is an element-wise nonlinear activation function. Denote the collection of all parameters as $\Theta = \{a_d, \mathbf{w}_d\}_{d=1}^{d_1}$.

The seminal work of Jacot et al. (2018) shows that the dynamics of the training process of an infinite-width two-layer neural network can be captured by the neural tangent kernel (NTK) $\mathcal{K}(\mathbf{x}, \mathbf{x}')$, which is defined as

$$\mathcal{K}(\mathbf{x}, \mathbf{x}') = \nabla_{\Theta} f(\mathbf{x}; \Theta)^T \nabla_{\Theta} f(\mathbf{x}'; \Theta), \quad (63)$$

In the case of a two-layer homogeneous nonlinear network Equation (62), the NTK can be expressed as

$$\mathcal{K}(\mathbf{x}, \mathbf{x}') = \frac{1}{d_1} \sum_{d=1}^{d_1} [a_d^2 \sigma'(\mathbf{w}_d^T \mathbf{x}) \sigma'(\mathbf{w}_d^T \mathbf{x}') \mathbf{x}^T \mathbf{x}' + \sigma(\mathbf{w}_d^T \mathbf{x}) \sigma(\mathbf{w}_d^T \mathbf{x}')] \quad (64)$$

According to the law of large numbers, as $d_1 \rightarrow \infty$, the NTK converges to a deterministic kernel $\mathcal{K}_{\infty}(\mathbf{x}, \mathbf{x}')$, also known as the limiting NTK, which can be expressed as,

$$\begin{aligned} \mathcal{K}_{\infty}(\mathbf{x}, \mathbf{x}') &= \mathbb{E}_{a \sim \mathcal{N}(0,1), \mathbf{w} \sim \mathcal{N}(0, \mathbf{I}_{d_0})} [a^2 \sigma'(\mathbf{w}^T \mathbf{x}) \sigma'(\mathbf{w}^T \mathbf{x}') \mathbf{x}^T \mathbf{x}' + \sigma(\mathbf{w}^T \mathbf{x}) \sigma(\mathbf{w}^T \mathbf{x}')] \\ &= \mathbb{E}_{\mathbf{w} \sim \mathcal{N}(0, \mathbf{I}_{d_0})} [\sigma'(\mathbf{w}^T \mathbf{x}) \sigma'(\mathbf{w}^T \mathbf{x}') \mathbf{x}^T \mathbf{x}' + \sigma(\mathbf{w}^T \mathbf{x}) \sigma(\mathbf{w}^T \mathbf{x}')] \end{aligned} \quad (65)$$

Proposition A.11. Let \mathcal{G} be a finite group, and $\rho_{\mathcal{X}}$ be a unitary representation of \mathcal{G} on \mathbb{R}^{d_0} . If the activation function σ and its derivative are both integrable with respect to Gaussian measure, i.e., $\sigma(x), \sigma'(x) \in \mathcal{L}^1(\mathbb{R}, \gamma)$, where γ is the standard Gaussian measure, then the limiting NTK $\mathcal{K}_{\infty}(\mathbf{x}, \mathbf{x}')$ is equivariant, i.e., $\mathcal{K}_{\infty}(\rho_{\mathcal{X}}(g)\mathbf{x}, \rho_{\mathcal{X}}(g)\mathbf{x}') = \mathcal{K}_{\infty}(\mathbf{x}, \mathbf{x}')$ for all $g \in \mathcal{G}$.

Proof. The proof follows from the fact that isotropic Gaussian random variables are invariant under orthogonal transformations. According to Cauchy's inequality, we have

$$\mathbb{E}_{\mathbf{w} \sim \mathcal{N}(0, \mathbf{I}_{d_0})} [\sigma(\mathbf{w}^T \mathbf{x}) \sigma(\mathbf{w}^T \mathbf{x}')] \leq \mathbb{E}_{\mathbf{w} \sim \mathcal{N}(0, \mathbf{I}_{d_0})} [\sigma(\mathbf{w}^T \mathbf{x})] \mathbb{E}_{\mathbf{w} \sim \mathcal{N}(0, \mathbf{I}_{d_0})} [\sigma(\mathbf{w}^T \mathbf{x}')] < \infty.$$

Therefore, the kernel $\mathcal{K}_{\infty}(\mathbf{x}, \mathbf{x}')$ is well-defined. The proof of the equivariance property is as follows:

$$\begin{aligned} \mathcal{K}_{\infty}(\rho_{\mathcal{X}}(g)\mathbf{x}, \rho_{\mathcal{X}}(g)\mathbf{x}') &= \mathbb{E}_{\mathbf{w} \sim \mathcal{N}(0, \mathbf{I}_{d_0})} [\sigma'(\mathbf{w}^T \rho_{\mathcal{X}}(g)\mathbf{x}) \sigma'(\mathbf{w}^T \rho_{\mathcal{X}}(g)\mathbf{x}') (\rho_{\mathcal{X}}(g)\mathbf{x})^T (\rho_{\mathcal{X}}(g)\mathbf{x}')] + \\ &\quad \mathbb{E}_{\mathbf{w} \sim \mathcal{N}(0, \mathbf{I}_{d_0})} [\sigma(\mathbf{w}^T \rho_{\mathcal{X}}(g)\mathbf{x}) \sigma(\mathbf{w}^T \rho_{\mathcal{X}}(g)\mathbf{x}')] \\ &= \mathbb{E}_{\mathbf{v} \sim \mathcal{N}(0, \mathbf{I}_{d_0})} [\sigma'(\mathbf{v}^T \mathbf{x}) \sigma'(\mathbf{v}^T \mathbf{x}') \mathbf{x}^T \mathbf{x}' + \sigma(\mathbf{v}^T \mathbf{x}) \sigma(\mathbf{v}^T \mathbf{x}')] \\ &= \mathcal{K}_{\infty}(\mathbf{x}, \mathbf{x}'), \end{aligned}$$

where $\mathbf{v} = \rho_{\mathcal{X}}(g)\mathbf{w}$ is also isotropic Gaussian, and the second equality follows from the fact that $\rho_{\mathcal{X}}(g)$ is an orthogonal transformation, i.e., $\rho_{\mathcal{X}}(g)^T = \rho_{\mathcal{X}}(g)^{-1}$. \square

Remark A.12. Most of the commonly used activation functions, such as ReLU, leaky ReLU, sigmoid, and tanh, satisfy the integrability condition. For example (Golikov et al., 2022), when the activation function is ReLU, the limiting NTK can be expressed as

$$\mathcal{K}_\infty(\mathbf{x}, \mathbf{x}') = \frac{\mathbf{x}^\top \mathbf{x}'}{2\pi} (\pi - \theta) + \frac{\|\mathbf{x}\| \|\mathbf{x}'\|}{2\pi} ((\pi - \theta) \cos \theta + \sin \theta), \quad (66)$$

where $\theta = \arccos\left(\frac{\mathbf{x}^\top \mathbf{x}'}{\|\mathbf{x}\| \|\mathbf{x}'\|}\right)$.

Proposition A.13 (Theorem 4.1 in Li et al. (2019)). *Let \mathcal{G} be a finite group, and $\rho_{\mathcal{X}}$ be a unitary representation of \mathcal{G} on \mathbb{R}^{d_0} . Let \mathcal{K} be an equivariant kernel. Define the augmented kernel $\mathcal{K}^{\mathcal{G}}(\mathbf{x}, \mathbf{x}') := \mathbb{E}_{g \in \mathcal{G}} [\mathcal{K}(\rho_{\mathcal{X}}(g)\mathbf{x}, \mathbf{x}')]$. Let $\mathcal{D} = \{\mathbf{x}_i, y_i\}_{i=1}^n$ be a dataset, and $\mathcal{D}^{\mathcal{G}} = \{\rho_{\mathcal{X}}(g)\mathbf{x}_i, y_i\}_{i \in [n], g \in \mathcal{G}}$ be the augmented dataset. We use X and \mathbf{y} to denote the data matrix and target vector of \mathcal{D} , and $X^{\mathcal{G}}$ and $\mathbf{y}^{\mathcal{G}}$ to denote the data matrix and target vector of $\mathcal{D}^{\mathcal{G}}$. We also use \mathcal{K}_X to denote the kernel matrix of \mathcal{K} on data matrix X . Then the prediction of $\mathcal{K}^{\mathcal{G}}$ on \mathcal{D} is equivalent to the prediction of \mathcal{K} on $\mathcal{D}^{\mathcal{G}}$, i.e.,*

$$\sum_{i=1}^n \alpha_i \mathcal{K}^{\mathcal{G}}(\mathbf{x}, \mathbf{x}_i) = \sum_{i \in [n], g \in \mathcal{G}} \beta_{i,g} \mathcal{K}(\mathbf{x}, \rho_{\mathcal{X}}(g)\mathbf{x}_i), \quad \forall \mathbf{x} \in \mathbb{R}^{d_0}$$

where $\boldsymbol{\alpha} = \{\alpha_i\}_{i=1}^n = \mathcal{K}_X^{\mathcal{G}}{}^{-1} \mathbf{y}$, and $\boldsymbol{\beta} = \{\beta_i\}_{i=1}^n = \mathcal{K}_{X^{\mathcal{G}}}^{-1} \mathbf{y}^{\mathcal{G}}$.

Corollary A.14. *Let \mathcal{G} be a finite group, and $\rho_{\mathcal{X}}$ be a unitary representation of \mathcal{G} on \mathbb{R}^{d_0} . Let $\mathcal{D} = \{\mathbf{x}_i, y_i\}_{i=1}^n$ be a dataset, and $\mathcal{D}^{\mathcal{G}} = \{\rho_{\mathcal{X}}(g)\mathbf{x}_i, y_i\}_{i \in [n], g \in \mathcal{G}}$ be its augmented dataset. For a univariate two-layer homogeneous ReLU network, the prediction of its limiting NTK \mathcal{K}_∞ on the augmented dataset $\mathcal{D}^{\mathcal{G}}$ is invariant.*

Proof. According to Proposition A.11, the limiting NTK \mathcal{K}_∞ of the two-layer homogeneous ReLU network is equivariant. Then we can apply Proposition A.13 to show that the prediction of \mathcal{K}_∞ on the augmented dataset $\mathcal{D}^{\mathcal{G}}$ is equivalent to the prediction of $\mathcal{K}_\infty^{\mathcal{G}}$ on the original dataset \mathcal{D} . We only need to check that the augmented kernel $\mathcal{K}_\infty^{\mathcal{G}}$ will give us an invariant predictor. To verify this, we have

$$\begin{aligned} \sum_{i=1}^n \alpha_i \mathcal{K}_\infty^{\mathcal{G}}(\rho_{\mathcal{X}}(h)\mathbf{x}, \mathbf{x}_i) &= \frac{1}{|\mathcal{G}|} \sum_{i=1}^n \alpha_i \sum_{g \in \mathcal{G}} \mathcal{K}_\infty(\rho_{\mathcal{X}}(g)\rho_{\mathcal{X}}(h)\mathbf{x}, \mathbf{x}_i) \\ &= \frac{1}{|\mathcal{G}|} \sum_{i=1}^n \alpha_i \sum_{g \in \mathcal{G}} \mathcal{K}_\infty(\rho_{\mathcal{X}}(gh)\mathbf{x}, \mathbf{x}_i) \\ &= \frac{1}{|\mathcal{G}|} \sum_{i=1}^n \alpha_i \sum_{g \in \mathcal{G}} \mathcal{K}_\infty^{\mathcal{G}}(\rho_{\mathcal{X}}(g)\mathbf{x}, \mathbf{x}_i) \\ &= \sum_{i=1}^n \alpha_i \mathcal{K}_\infty^{\mathcal{G}}(\mathbf{x}, \mathbf{x}_i) \end{aligned}$$

□

Now let's consider a group invariant convolutional network (Cohen & Welling, 2016) with the same number of parameters as in Equation (62). It is worth noting that this convolutional network only contains one group-lifting layer and a group-average pooling layer. Middle group-convolutional layers are not included in this model for the sake of simplicity. But it is enough to illustrate the main idea. The output of the network is given by

$$f^{\text{conv}}(\mathbf{x}; \Theta) = \frac{1}{\sqrt{d_1}} \sum_{d=1}^{d_1} a_d \left(\frac{1}{|\mathcal{G}|} \sum_{g \in \mathcal{G}} \sigma(\mathbf{w}_d^\top \rho_{\mathcal{X}}(g)\mathbf{x}) \right), \quad (67)$$

where the parameters are initialized as $a_d \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(0, 1)$, $\mathbf{w}_d \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(0, \mathbf{I}_{d_0})$, and $\sigma : \mathbb{R} \rightarrow \mathbb{R}$ is an element-wise nonlinear activation function. One can verify that the above model Equation (67) is invariant with respect to $\rho_{\mathcal{X}}$.

Proposition A.15. *Let \mathcal{G} be a finite group, and $\rho_{\mathcal{X}}$ be a unitary representation of \mathcal{G} on \mathbb{R}^{d_0} . Let f be a univariate two-layer homogeneous network with activation function σ , as defined in Equation (62), and f^{conv} be a univariate two-layer*

group-convolutional network with the same number of parameters, as f as defined in Equation (67). For any dataset $\mathcal{D} = \{\mathbf{x}_i, y_i\}_{i=1}^n$, let $\mathcal{D}^{\mathcal{G}} = \{\rho_{\mathcal{X}}(g)\mathbf{x}_i, y_i\}_{i \in [n], g \in \mathcal{G}}$ be the augmented dataset. We use X and \mathbf{y} to denote the data matrix and target vector of \mathcal{D} , and $X^{\mathcal{G}}$ and $\mathbf{y}^{\mathcal{G}}$ to denote the data matrix and target vector of $\mathcal{D}^{\mathcal{G}}$. We also use \mathcal{K}_X to denote the kernel matrix of \mathcal{K} on data matrix X . Then the NTK predictor of f on $\mathcal{D}^{\mathcal{G}}$ is equivalent to the NTK predictor of f^{conv} on \mathcal{D} , i.e.,

$$\sum_{i \in [n], g \in \mathcal{G}} \alpha_i \mathcal{K}_{\infty}(\mathbf{x}, \rho_{\mathcal{X}}(g)\mathbf{x}_i) = \sum_{i=1}^n \beta_i \mathcal{K}_{\infty}^{\text{conv}}(\mathbf{x}, \mathbf{x}_i), \quad \forall g \in \mathcal{G}, \forall \mathbf{x} \in \mathbb{R}^{d_0}$$

where $\boldsymbol{\alpha} = \{\alpha_i\}_{i=1}^n = \mathcal{K}_{X^{\mathcal{G}}}^{-1} \mathbf{y}^{\mathcal{G}}$, and $\boldsymbol{\beta} = \{\beta_i\}_{i=1}^n = \mathcal{K}_X^{-1} \mathbf{y}$.

Proof. We need to show that the limiting NTK $\mathcal{K}_{\infty}^{\text{conv}}$ of the group-convolutional network f^{conv} is equivalent to the augmented limiting NTK $\mathcal{K}_{\infty}^{\mathcal{G}}$ of the original network f . The limiting NTK of f^{conv} can be expressed as

$$\begin{aligned} \mathcal{K}_{\infty}^{\text{conv}}(\mathbf{x}, \mathbf{x}') &= \mathbb{E}_{\mathbf{w} \sim \mathcal{N}(0, \mathbf{I}_{d_0})} \left[\left(\frac{1}{|\mathcal{G}|} \sum_{g \in \mathcal{G}} \sigma(\mathbf{w}^T \rho_{\mathcal{X}}(g)\mathbf{x}) \right) \left(\frac{1}{|\mathcal{G}|} \sum_{g' \in \mathcal{G}} \sigma(\mathbf{w}^T \rho_{\mathcal{X}}(g')\mathbf{x}') \right) \right] \\ &\quad + \mathbb{E}_{\mathbf{w} \sim \mathcal{N}(0, \mathbf{I}_{d_0})} \left[\left(\frac{1}{|\mathcal{G}|} \sum_{g \in \mathcal{G}} \sigma'(\mathbf{w}^T \rho_{\mathcal{X}}(g)\mathbf{x}) \mathbf{x}^T \rho_{\mathcal{X}}(g)^T \right) \left(\frac{1}{|\mathcal{G}|} \sum_{g' \in \mathcal{G}} \sigma'(\mathbf{w}^T \rho_{\mathcal{X}}(g')\mathbf{x}') \rho_{\mathcal{X}}(g')^T \mathbf{x}' \right) \right] \end{aligned} \quad (68)$$

The augmented NTK of f can be expressed as

$$\mathcal{K}_{\infty}^{\mathcal{G}}(\mathbf{x}, \mathbf{x}') = \mathbb{E}_{\mathbf{w} \sim \mathcal{N}(0, \mathbf{I}_{d_0})} \left[\frac{1}{|\mathcal{G}|} \sum_{g \in \mathcal{G}} \sigma(\mathbf{w}^T \rho_{\mathcal{X}}(g)\mathbf{x}) \sigma(\mathbf{w}^T \mathbf{x}') + \frac{1}{|\mathcal{G}|} \sum_{g \in \mathcal{G}} \sigma'(\mathbf{w}^T \rho_{\mathcal{X}}(g)\mathbf{x}) \sigma'(\mathbf{w}^T \mathbf{x}') \mathbf{x}^T \rho_{\mathcal{X}}(g)^T \mathbf{x}' \right] \quad (69)$$

For the first term, we have

$$\begin{aligned} &\mathbb{E}_{\mathbf{w} \sim \mathcal{N}(0, \mathbf{I}_{d_0})} \left[\left(\frac{1}{|\mathcal{G}|} \sum_{g \in \mathcal{G}} \sigma(\mathbf{w}^T \rho_{\mathcal{X}}(g)\mathbf{x}) \right) \left(\frac{1}{|\mathcal{G}|} \sum_{g' \in \mathcal{G}} \sigma(\mathbf{w}^T \rho_{\mathcal{X}}(g')\mathbf{x}') \right) \right] \\ &= \frac{1}{|\mathcal{G}|} \sum_{g' \in \mathcal{G}} \mathbb{E}_{\mathbf{w} \sim \mathcal{N}(0, \mathbf{I}_{d_0})} \left[\left(\frac{1}{|\mathcal{G}|} \sum_{g \in \mathcal{G}} \sigma(\mathbf{w}^T \rho_{\mathcal{X}}(g)\mathbf{x}) \right) \sigma(\mathbf{w}^T \rho_{\mathcal{X}}(g')\mathbf{x}') \right] \\ &= \frac{1}{|\mathcal{G}|} \sum_{g' \in \mathcal{G}} \mathbb{E}_{\mathbf{v} \sim \mathcal{N}(0, \mathbf{I}_{d_0})} \left[\left(\frac{1}{|\mathcal{G}|} \sum_{g \in \mathcal{G}} \sigma(\mathbf{v}^T \rho_{\mathcal{X}}(g'^{-1}g)\mathbf{x}) \right) \sigma(\mathbf{v}^T \mathbf{x}') \right] \\ &= \mathbb{E}_{\mathbf{v} \sim \mathcal{N}(0, \mathbf{I}_{d_0})} \left[\frac{1}{|\mathcal{G}|} \sum_{g \in \mathcal{G}} \sigma(\mathbf{v}^T \rho_{\mathcal{X}}(g)\mathbf{x}) \sigma(\mathbf{v}^T \mathbf{x}') \right], \end{aligned}$$

where the second equality follows from the fact that isotropic Gaussian random variables are invariant under orthogonal transformations, and the last equality applies the change of variable $g'^{-1}g \mapsto g$. Similarly, for the second term, we have

$$\begin{aligned} &\mathbb{E}_{\mathbf{w} \sim \mathcal{N}(0, \mathbf{I}_{d_0})} \left[\left(\frac{1}{|\mathcal{G}|} \sum_{g \in \mathcal{G}} \sigma'(\mathbf{w}^T \rho_{\mathcal{X}}(g)\mathbf{x}) \mathbf{x}^T \rho_{\mathcal{X}}(g)^T \right) \left(\frac{1}{|\mathcal{G}|} \sum_{g' \in \mathcal{G}} \sigma'(\mathbf{w}^T \rho_{\mathcal{X}}(g')\mathbf{x}') \rho_{\mathcal{X}}(g')^T \mathbf{x}' \right) \right] \\ &= \frac{1}{|\mathcal{G}|} \sum_{g' \in \mathcal{G}} \mathbb{E}_{\mathbf{w} \sim \mathcal{N}(0, \mathbf{I}_{d_0})} \left[\left(\frac{1}{|\mathcal{G}|} \sum_{g \in \mathcal{G}} \sigma'(\mathbf{w}^T \rho_{\mathcal{X}}(g)\mathbf{x}) \mathbf{x}^T \rho_{\mathcal{X}}(g)^T \right) \sigma'(\mathbf{w}^T \rho_{\mathcal{X}}(g')\mathbf{x}') \rho_{\mathcal{X}}(g')^T \mathbf{x}' \right] \\ &= \frac{1}{|\mathcal{G}|} \sum_{g' \in \mathcal{G}} \mathbb{E}_{\mathbf{v} \sim \mathcal{N}(0, \mathbf{I}_{d_0})} \left[\left(\frac{1}{|\mathcal{G}|} \sum_{g \in \mathcal{G}} \sigma'(\mathbf{v}^T \rho_{\mathcal{X}}(g'^{-1}g)\mathbf{x}) \mathbf{x}^T \rho_{\mathcal{X}}(g)^T \right) \sigma'(\mathbf{v}^T \rho_{\mathcal{X}}(g')\mathbf{x}') \rho_{\mathcal{X}}(g')^T \mathbf{x}' \right] \end{aligned}$$

$$= \mathbb{E}_{\mathbf{v} \sim \mathcal{N}(0, \mathbf{I}_{d_0})} \left[\frac{1}{|\mathcal{G}|} \sum_{g \in \mathcal{G}} \sigma'(\mathbf{v}^T \rho_{\mathcal{X}}(g) \mathbf{x}) \sigma'(\mathbf{v}^T \mathbf{x}') \mathbf{x}^T \rho_{\mathcal{X}}(g)^T \rho_{\mathcal{X}}(g) \mathbf{x}' \right].$$

Therefore, we have shown that $\mathcal{K}_{\infty}^{\text{conv}}(\mathbf{x}, \mathbf{x}') = \mathcal{K}_{\infty}^{\mathcal{G}}(\mathbf{x}, \mathbf{x}')$ for all $\mathbf{x}, \mathbf{x}' \in \mathbb{R}^{d_0}$. Finally, we can apply Proposition A.13 to show that the NTK predictor of f on $\mathcal{D}^{\mathcal{G}}$ is equivalent to the NTK predictor of f^{conv} on \mathcal{D} . \square

A.10. Empirical Spectrum of Target Matrices in MNIST Dataset

As discussed in Remark 3.4 and Proposition 3.10, we have assumptions about the rank and spectrum of the target matrices we are trying to approximate. As shown in Figure 5, we empirically computed the singular values of \bar{Z}^{da} , \bar{Z}^{inv} , $\bar{Z}(\lambda)^{reg}$ for MNIST dataset. We can see that all three target matrices have full rank. The singular values are pairwise different as well. Thus, the previous assumptions in Remark 3.4 and Proposition 3.10 are satisfied.

A.11. Experiments for Cross Entropy Loss

As mentioned in section 4, we have observed that even under cross-entropy loss, the specific entries in end-to-end matrix W converge to approximately the same values, indicating that the learned map is nearly invariant (see Figure 6). Furthermore, the training curves share similar patterns as those under MSE loss (see Figure 7 & Figure 2).

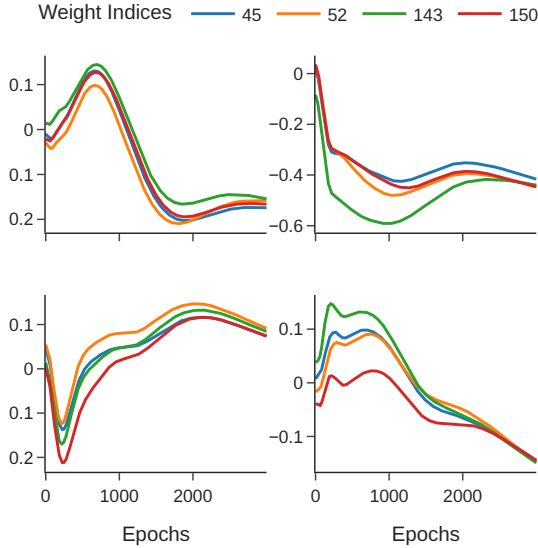


Figure 6. Weights in a two-layer linear neural network trained using data augmentation with cross-entropy loss.

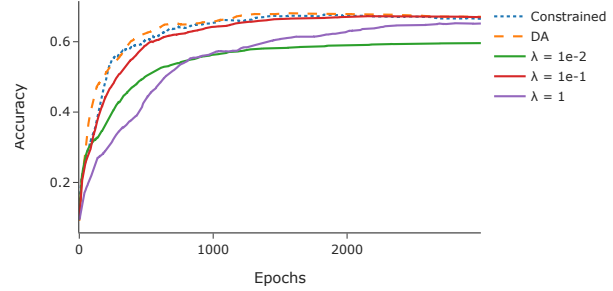


Figure 7. Training curves for data augmentation (DA), regularization (λ), and constrained model under cross-entropy loss.

In Figure 8 and Figure 9, we are still plotting $\|W^{\perp}\|_F$ for data augmentation and regularization trained on the same dataset, but with cross entropy loss. It is observed that, for larger λ , the dynamics of $\|W^{\perp}\|_F$ resemble those when trained with MSE (see Figure 3). On the other hand, for small λ , $\|W^{\perp}\|_F$ may increase at first, and then decrease. For data augmentation, if we allow more epochs, we can still observe that $\|W^{\perp}\|_F$ decreases after increasing. Our theoretical results only support the scenario for mean squared loss. Thus, when trained with cross entropy, we cannot say whether all the critical points are invariant or not. Future work can be done to investigate the critical points when trained with cross entropy loss.

A.12. Experiments for Two-layer Nonlinear Network

In Figure 10, we show the training curve for a two-layer neural network with different nonlinear activation functions trained with data augmentation and hard-wiring. The setup is the same as previous experiments in section 4. In this experiment, we used 5000 samples from the MNIST dataset for training with mean squared loss (MSE). Meanwhile, we also test the case when there is not a bottleneck middle layer. When the middle layer has a bottleneck, we set the number of hidden units as 7; otherwise, the number of hidden units is 15. We can see that both data augmentation and constrained model have similar

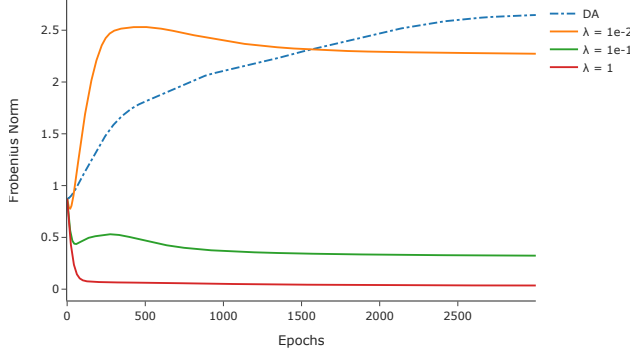


Figure 8. $\|W_{\perp}\|_F$ where W_{\perp} is the non-invariant part of W under cross-entropy loss.

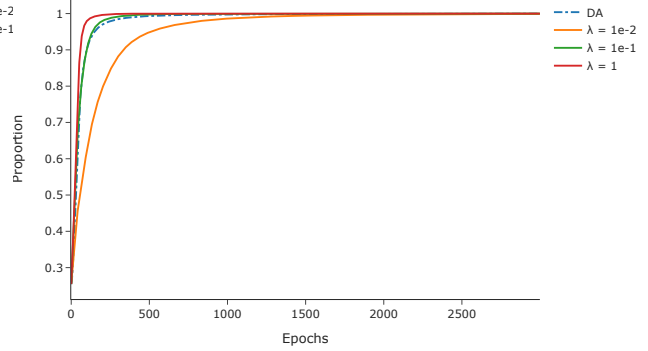


Figure 9. $\|W - W_{\perp}\|_F^2 / \|W\|_F^2$ under cross-entropy loss.

loss in the late phase of training for all four activation functions, especially when there is a bottleneck.

Besdies, [Moskalev et al. \(2023\)](#) suggests that invariance learned from data augmentation deteriorates under distribution shift in a classification setting. The architecture they choose is a 5-layer ReLU network. We would like to investigate this in a regression setting. The model we use is a 2-layer neural network with different activation functions trained with data augmentation and MSE. For any function f and an input point $x \in \mathcal{X}$, to measure the amount of invariance of f , we evaluate the scaled variance of outputs across the group orbit,

$$\epsilon_{inv}(f, x) := \mathbb{E}_{g \sim \lambda} \left(1 - \frac{f(gx)}{\bar{f}(x)} \right)^2,$$

where $\bar{f}(x) := \mathbb{E}_{g \sim \lambda}[f(gx)]$, and λ is the Haar measure on group \mathcal{G} . In [Figure 11](#), we train a 2-layer neural network with different activation functions on MNIST with data augmentation using training sample sizes $N = 1000$ and $N = 5000$. After training, we calculate $\epsilon_{inv}(f, x)$ for two different datasets: MNIST and Gaussian. For MNIST, we use 5000 samples from the original test set in MNIST. For Gaussian, we sample 5000 points from an isotropic Gaussian distribution in dimension 196. We are showing the median of $\{\epsilon_{inv}(f, x)\}_{x \in \mathcal{D}}$ in [Figure 11](#).

The observations can be summarized as follows:

1. **Effect of the size of the model:** Compared to the case without a bottleneck middle layer, $\epsilon_{inv}(f, x)$ is significantly smaller when there is a bottleneck. This suggests that it is more difficult to learn invariance from the data when the model has more parameters.
2. **Effect of the amount of training data:** We notice that $\epsilon_{inv}(f, x)$ is smaller when there are more training data. For underdetermined linear models, i.e., when the number of data points exceeds the input dimension, [Proposition 3.10](#) shows that all critical points are invariant. However, when the model is nonlinear, we need more data in order to learn the invariance via data augmentation.
3. **Robustness under distribution shift:** Though the model is trained on MNIST, $\epsilon_{inv}(f, x)$ does not increase significantly even when the model is tested on a completely different dataset. This suggests that the invariance learned from the data is fairly robust.

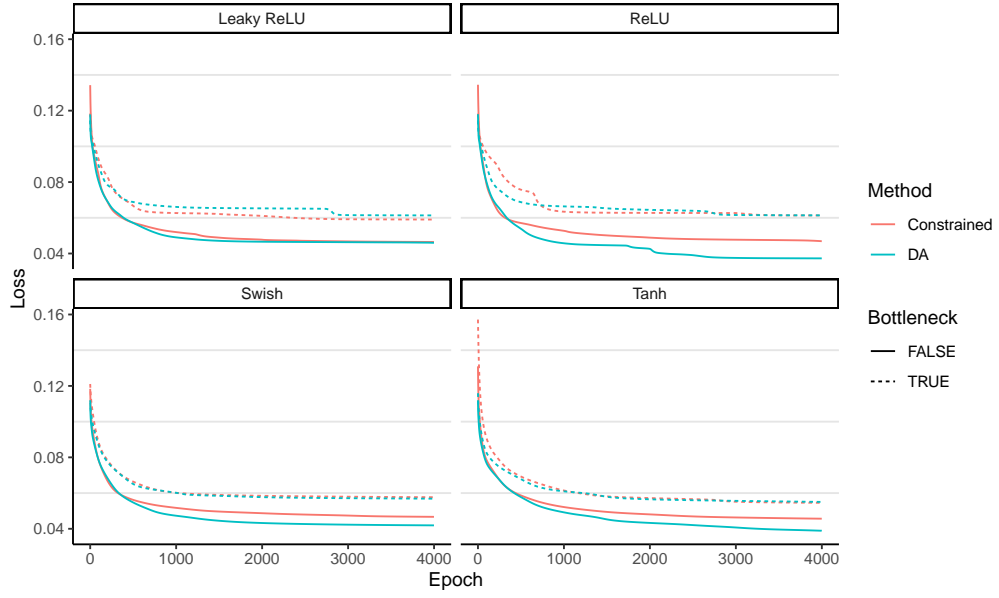


Figure 10. Training curves for a two-layer NN with different nonlinear activation functions via data augmentation and hard-wiring on MNIST.

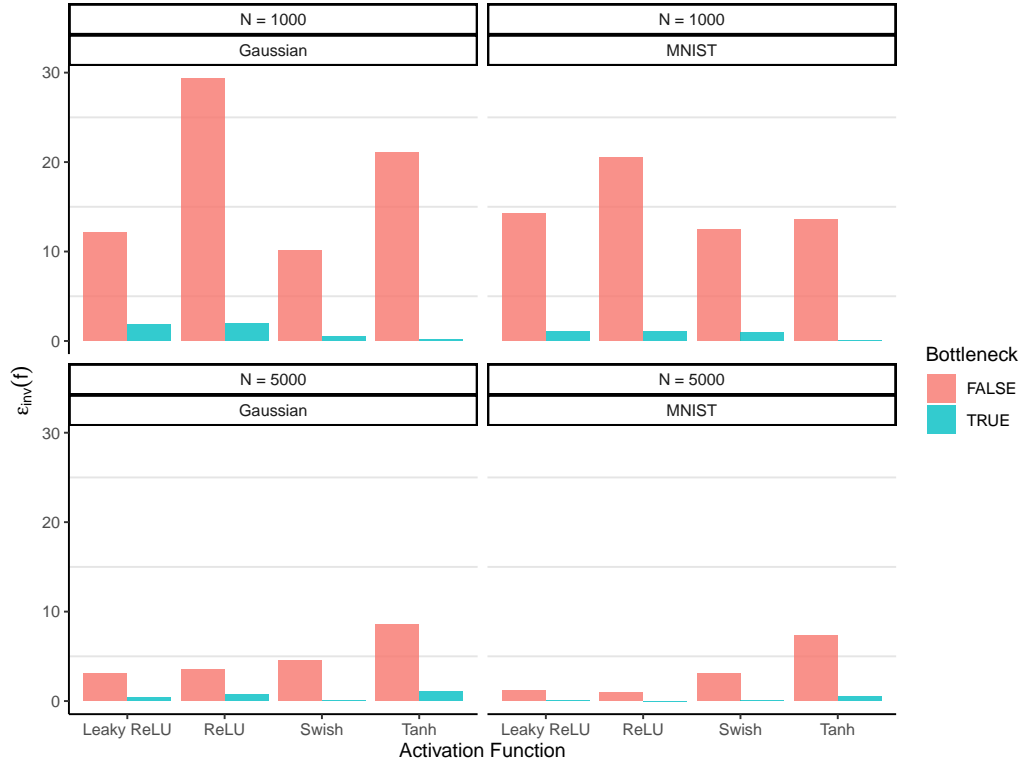


Figure 11. Median of the measure $\epsilon_{inv}(f)$ of discrepancy from invariance for 2-layer neural networks with different activation functions, trained on MNIST with data augmentation.