
On the Training Instability of Shuffling SGD with Batch Normalization

David X. Wu¹ Chulhee Yun² Suvrit Sra³

Abstract

We uncover how SGD interacts with batch normalization and can exhibit undesirable training dynamics such as divergence. More precisely, we study how Single Shuffle (SS) and Random Reshuffle (RR)—two widely used variants of SGD—interact surprisingly differently in the presence of batch normalization: *RR leads to much more stable evolution of training loss than SS*. As a concrete example, for regression using a linear network with batch normalized inputs, we prove that SS and RR converge to distinct global optima that are “distorted” away from gradient descent. Thereafter, for classification we characterize conditions under which training divergence for SS and RR can, and cannot occur. We present explicit constructions to show how SS leads to distorted optima in regression and divergence for classification, whereas RR avoids both distortion and divergence. We validate our results empirically in realistic settings, and conclude that the separation between SS and RR used with batch normalization is relevant in practice.

1. Introduction

Recent work in deep learning theory attempts to uncover how the choice of optimization algorithm and architecture influence training stability and efficiency. On the optimization front, stochastic gradient descent (SGD) is the *de facto* workhorse, and its importance has correspondingly led to the development of many different variants that aim to increase the ease and speed of training, such as AdaGrad (Duchi et al., 2011) and Adam (Kingma & Ba, 2014).

In reality, practitioners often do not use with-replacement sampling of gradients as required by SGD. Instead they use

¹Department of EECS, UC Berkeley, Berkeley, CA, USA ²Kim Jaechul Graduate School of AI, KAIST, Seoul, Korea ³Department of EECS, LIDS, MIT, Cambridge, MA, USA. Correspondence to: David X. Wu <david_wu@berkeley.edu>.

Proceedings of the 40th International Conference on Machine Learning, Honolulu, Hawaii, USA. PMLR 202, 2023. Copyright 2023 by the author(s).

without-replacement sampling, leading to two main variants of SGD: single-shuffle (SS) and random-reshuffle. SS randomly samples and fixes a permutation at the beginning of training, while RR randomly resamples permutations at each epoch. These shuffling algorithms are often more practical and can have improved convergence rates (Haochen & Sra, 2019; Safran & Shamir, 2020; Yun et al., 2021b; 2022; Cho & Yun, 2023; Cha et al., 2023).

Architecture design offers another avenue for practitioners to train networks more efficiently and encode salient inductive biases. Normalizing layers such as BatchNorm (BN) (Ioffe & Szegedy, 2015), LayerNorm (Ba et al., 2016), or InstanceNorm (Ulyanov et al., 2016) are often used with SGD to accelerate convergence and stabilize training. Recent work studies how these scale-invariant layers affect training through the effective learning rate (Li & Arora, 2019; Li et al., 2020; Wan et al., 2021; Lyu et al., 2022).

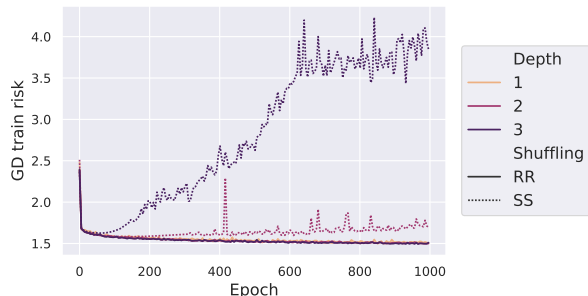
Motivated by these practical choices, we study how SS and RR interact with batch normalization at *training time*. Our experiments (Fig. 1) suggest that combining SS and BN can lead to surprising and undesirable training phenomena:

- (i) The training risk often diverges when using SS+BN to train linear networks (i.e. without nonlinear activations) on real datasets (see Figure 1a), while using SS without BN does not cause divergence (see Figure 10).
- (ii) Divergence persists after tuning the learning rate and other hyperparameters (Section 4.3) and also manifests more quickly in deeper linear networks (Figure 1a).
- (iii) SS+BN usually converges slower than RR+BN in nonlinear architectures such as ResNet18 (see Figure 1b).

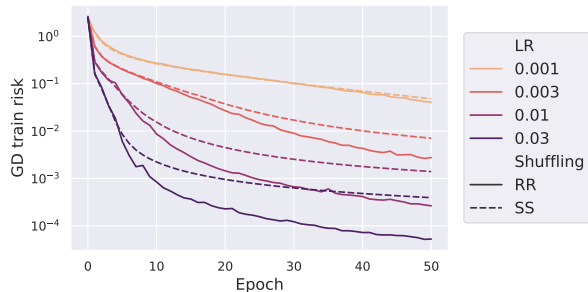
1.1. Summary of our contributions

In light of these experimental findings, we seek to develop a theoretical and experimental understanding of how shuffling SGD and BN collude to create divergence and other undesirable training behavior. Since these phenomena manifest themselves on the training risk, our results are not strictly coupled with generalization.

Put simply, the aberrant training dynamics stem from BN *not* being permutation invariant across epochs. This simple property interacts with SS undesirably, although *a priori* it is not obvious whether it should. More concretely, one expects SGD+BN to optimize the gradient descent (GD)



(a) Depths 1, 2, and 3 linear+BN networks.



(b) Finetuned ResNet18.

Figure 1. Surprising training phenomena using SS/RR+BN.

risk in expectation. However, due to BN’s sensitivity to permutations, both SS+BN and RR+BN implicitly train induced risks different from GD, and also from each other.

- In Section 3.2, we prove that the network $f(\mathbf{X}; \Theta) = \mathbf{W}\Gamma\text{BN}(\mathbf{X})$ which batch normalizes the input features converges to the optimum for the distorted risk induced by SS and RR (Theorems 3.2.2 and 3.2.3); the diagonal matrix Γ denotes the *trainable* scale parameters in the BN layer. Our proof requires a delicate analysis of the evolution of gradients, the noise arising from SS, and the two-layer architecture. Due to the presence of Γ , our results do not assume a fully-connected linear network, which distinguishes them from prior convergence results. In Section 3.3, we present a toy dataset for which SS is distorted away from GD with constant probability while RR averages out the distortion to align with GD. We validate our theoretical findings on synthetic data in Section 3.4.
- In Section 4.1, we connect properties of the distorted risks to divergence of training risk, gaining insights into which regimes can lead to divergence (Theorems 4.1.3 and 4.1.4). We show that in certain regimes, SS+BN can diverge, whereas RR+BN provably avoids divergence. These results motivate us to construct a toy dataset where SS diverges with constant probability, while RR avoids divergence (Section 4.2). In Section 4.3, we empirically validate our results on deeper linear+BN networks on a variety of datasets and hyperparameters. Our experiments also demonstrate that SS trains more slowly than RR in

more realistic nonlinear settings, including ReLU+BN networks and ResNet18. In doing so, we extend the relevance of our theoretical results to more complex and realistic settings.

It is worth noting that to obtain our results, our analysis had to overcome complications due to the non-i.i.d. stochastic gradients, the non-i.i.d. data \mathbf{X} (Assumption 2), and the intricacies introduced by BN’s permutation sensitivity.

1.2. Related work

Theoretical understanding of BN. Since the introduction of BN by Ioffe & Szegedy (2015), there has been a long line of work investigating the theoretical properties of BN; see e.g. (Bjorck et al., 2018; Kohler et al., 2018; Arora et al., 2018; Li & Arora, 2019; Kohler et al., 2019; Daneshmand et al., 2020; Li et al., 2020; Lobacheva et al., 2021). Much attention has been devoted to studying how BN can benefit optimization (Arora et al., 2018; Santurkar et al., 2018; Kohler et al., 2018), for example by implicitly tuning the learning rate or smoothing the loss function. The effect of BN on the intermediate representations of random networks, such as orthogonality or rank collapse, has also been studied (Daneshmand et al., 2020; 2021). We study the general setting with nonrandom linear activations. The scale invariance induced by BN also interacts with other optimization choices such as weight decay, which can lead to instability phenomena (Lobacheva et al., 2021; Wan et al., 2021; Lyu et al., 2022). However, these phenomena have a different origin than the distorted risks studied in this paper.

Interplay between BN and SGD. Prior theoretical work primarily studied how BN interacts with GD or with-replacement SGD (Arora et al., 2018; Santurkar et al., 2018; Li & Arora, 2019; Cai et al., 2019; Wan et al., 2021; Lyu et al., 2022). Arora et al. (2018); Wan et al. (2021) assumed global bounds on the smoothness with respect to network parameters and the SGD noise to analyze convergence to stationary points. We instead prove convergence to the global minimum of the SS distorted risk \mathcal{L}_π with *no* such assumptions (Theorem 3.2.2). Li & Arora (2019) assumed the batch size is large enough to ignore SGD noise, whereas we explicitly exhibit and study the separation between shuffling SGD and GD. For fully scale-invariant networks trained with GD, Lyu et al. (2022) identified an oscillatory edge of stability behavior around a manifold of minimizers. Our BN network has trainable scale-variant parameters \mathbf{W} and Γ , and we train with shuffling SGD instead of GD. Hence, the noise that leads to distorted risks is fundamentally different.

BN’s effect on risk function. Previous work identified the distortion of risk function due to noisy batch statistics in BN. Yong et al. (2020) studied the asymptotic regularization effect of noisy batch statistics *in expectation* for with-replacement SGD. In contrast, we characterize this noise

nonasymptotically w.h.p. over π for SS and a.s. with respect to the data for RR. Wu & Johnson (2021) studied the difficulty of precisely estimating the population statistics at train time, especially when using an exponential moving average. We avoid these issues altogether by evaluating directly on the GD risk. Moreover, we prove concentration inequalities for without-replacement batch statistics (Proposition C.2.4).

Ghost batch normalization. In the presence of BN, it is common practice to use *ghost batch normalization*, a scheme which break up large batches into virtual “ghost” batches, as this tends to improve the generalization of the network (Hoffer et al., 2017; Shallue et al., 2019; Summers & Dinneen, 2020). Minibatch statistics are calculated with respect to the ghost batches, and each gradient step is computed by summing the gradient contributions from the ghost batches. This algorithm is closely related to our method of analysis for SS+BN/RR+BN. Indeed, in our setup we also break up the full batch into mini-batches, and our analysis reduces to showing that SS+BN and RR+BN trajectories track those obtained by following the aggregate gradient signal from summing over mini-batches. We comment more on the similarities between ghost BN and our setup in Section 3.1.

Shuffling and optimization. Outside SGD, the effect of random shuffling has also been studied for classical nonlinear optimization schemes such as coordinate gradient descent (CGD) and ADMM (see Sun et al. (2020); Gürbüzbalaban et al. (2020) and references therein). On convex quadratic optimization problems, they demonstrate separations in convergence rates between SS, RR, and with-replacement sampling. Our main focus is the optimum that the algorithms converge to rather than their convergence rates.

Implicit bias. Our work is also motivated by a burgeoning line of work which studies the *implicit bias* of different optimization algorithms (Soudry et al., 2018; Gunasekar et al., 2018; Ji & Telgarsky, 2018; 2019; 2020; Yun et al., 2021a; Jagadeesan et al., 2022). These results establish how optimization algorithms such as gradient flow (GF), gradient descent (GD) or even with-replacement SGD are biased towards certain optima. For example, in the interpolating regime, GD converges to the min-norm solution (Gunasekar et al., 2018; Woodworth et al., 2020) for linear regression and the max-margin classifier for classification (Soudry et al., 2018; Nacson et al., 2019b;a). While our work does not focus on generalization, it is connected in spirit to implicit bias. Indeed, our analysis centers the study of how the risk functions and optima are affected by choices of the optimizer (SS/RR) and the architecture (BN).

2. Problem setup

For $n \in \mathbb{Z}^+$ we use the notation $[n] \triangleq \{1, \dots, n\}$. We write π to denote a permutation of $[n]$, and \mathbb{S}_n is the symmetric

group of all such π . For any matrix $\mathbf{A} \in \mathbb{R}^{d \times n}$, $\pi \circ \mathbf{A} \in \mathbb{R}^{d \times n}$ is result of shuffling the columns of \mathbf{A} according to π . Also, $\|\mathbf{A}\|_2$ and $\|\mathbf{A}\|_F$ refer to the spectral norm and Frobenius norm, respectively. We write $\sigma_{\min}(\mathbf{A}) \triangleq \inf_{\|v\|=1} \|\mathbf{A}v\|$ to denote minimum singular value of \mathbf{A} . We use $\text{Span}(\mathbf{A})$ to denote the span of \mathbf{A} ’s columns. The (coordinatewise) sign function $\text{sgn}(\cdot) : \mathbb{R} \rightarrow \{-1, 0, 1\}$ is defined as $\text{sgn}(x) = x/|x|$ for $x \neq 0$ and $\text{sgn}(0) = 0$.

Data. Let $\mathbf{Z} = (\mathbf{X}, \mathbf{Y})$ be the given dataset, with $\mathbf{X} = [\mathbf{x}_1 \ \dots \ \mathbf{x}_n] \in \mathbb{R}^{d \times n}$ representing the feature matrix and corresponding labels $\mathbf{Y} = [\mathbf{y}_1 \ \dots \ \mathbf{y}_n] \in \mathbb{R}^{p \times n}$. In the classification setting we will assume $\mathbf{Y} \in \{\pm 1\}^{1 \times n}$.

Prediction model. A batch normalization (BN) layer can be separated into a normalizing component BN and a scaling component Γ ; we ignore the bias parameters for analysis. Given any matrix $\mathbf{B} = [\mathbf{x}_1 \ \dots \ \mathbf{x}_q] \in \mathbb{R}^{d \times q}$ (here, $q \geq 2$ is arbitrary), the normalizing transform $\text{BN}(\cdot)$ maps it to $\text{BN}(\mathbf{B}) \in \mathbb{R}^{d \times q}$ by operating coordinatewise on each \mathbf{x}_i in \mathbf{B} . In particular, for the k th coordinate of \mathbf{x}_i , denoted as $x_{i,k}$, the transform BN sends $x_{i,k} \mapsto \frac{x_{i,k} - \mu_k}{\sqrt{\sigma_k^2 + \epsilon}}$ where μ_k and σ_k^2 are the batch empirical mean and variance of the k th coordinate, respectively, and ϵ is an arbitrary positive constant used to avoid numerical instability. For technical reasons, we omit ϵ in our analysis. The scaling matrix $\Gamma \in \mathbb{R}^{d \times d}$ is a diagonal matrix which models the tunable coordinatewise scale parameters inside the BN layer.

Throughout the paper, we consider neural networks of the form $f(\cdot; \Theta) = \mathbf{W}\Gamma\text{BN}(\cdot)^1$. We use $\Theta = (\mathbf{W}, \Gamma)$ to denote the collection of all parameters in the network. With the presence of batch normalization layers, the output of f is a function of the input datapoint *as well as* the batch it belongs to. Even changing one point of a batch \mathbf{B} can affect the batch statistics (i.e., μ_k ’s and σ_k^2 ’s) and in turn change the outputs of f for the entire batch.

Loss functions. We study regression with squared loss $\ell(\hat{\mathbf{y}}, \mathbf{y}) \triangleq \|\hat{\mathbf{y}} - \mathbf{y}\|^2$ and binary classification with logistic loss $\ell(\hat{\mathbf{y}}, \mathbf{y}) \triangleq -\log(\rho(y\hat{\mathbf{y}}))$, where $\rho(t) = 1/(1 + e^{-t})$. Let $\hat{\mathbf{Y}}, \mathbf{Y} \in \mathbb{R}^{p \times q}$ denote network outputs and true labels for a mini-batch of q datapoints, respectively. Define the mini-batch risk as the columnwise sum $\mathcal{L}(\hat{\mathbf{Y}}, \mathbf{Y}) \triangleq \sum_{i=1}^q \ell(\hat{\mathbf{Y}}_{:,i}, \mathbf{Y}_{:,i})$, where $\mathbf{Y}_{:,i}$ denotes the i th column of \mathbf{Y} .

Optimization methods. We consider shuffling-based variants of SGD, namely single-shuffle (SS) and random-reshuffle (RR). These algorithms proceed in *epochs*, i.e., full passes through shuffled dataset. As the names suggest, SS randomly samples a permutation $\pi \in \mathbb{S}_n$ at the beginning of the first epoch and adheres to this permutation. RR

¹We can readily generalize to arbitrary learned (but frozen) feature mappings under suitable changes to the assumptions.

randomly resamples permutations $\pi_k \in \mathbb{S}_n$ at each epoch k .

Throughout, the (mini-)batch size will be denoted as B . For simplicity, we assume that the n datapoints can be divided into m batches of size B . With a permutation $\pi \in \mathbb{S}_n$, the dataset $\mathbf{Z} = (\mathbf{X}, \mathbf{Y})$ is thus perfectly partitioned into m batches $(\mathbf{X}_\pi^1, \mathbf{Y}_\pi^1), \dots, (\mathbf{X}_\pi^m, \mathbf{Y}_\pi^m)$, where $\mathbf{X}_\pi^j \in \mathbb{R}^{d \times B}$ and $\mathbf{Y}_\pi^j \in \mathbb{R}^{p \times B}$ consist of the $(j(B-1) + 1, \dots, jB)$ th columns of the shuffled $\pi \circ \mathbf{X}$ and $\pi \circ \mathbf{Y}$, respectively.

For a parameter Θ optimized with SS or RR, we denote the j th iterate on the k th epoch by Θ_j^k . The starting iterate of the k th epoch is Θ_0^k which is equal to the last iterate of the previous epoch Θ_m^{k-1} . For each $j \in [m]$, SS and RR perform a mini-batch SGD update with stepsize $\eta_k > 0$:

$$\Theta_j^k \leftarrow \Theta_{j-1}^k - \eta_k \nabla_{\Theta} \mathcal{L}(f(\mathbf{X}_{\pi_k}^j; \Theta_{j-1}^k), \mathbf{Y}_{\pi_k}^j).$$

3. Main regression results: convergence to optima of distorted risks

In this section, we introduce the framework of distorted risks to elucidate the distinction between SS+BN and RR+BN. For ease of theoretical analysis, we consider the simplified setup where BN is applied only to the input features, although we note that the framework can be readily generalized to any learned but frozen features. This framework also applies to classification; we continue to study it in Section 4. We then present our global convergence results (Theorems 3.2.2 and 3.2.3) for the distorted risks induced by SS and RR for squared loss regression. In the one-dimensional case, we uncover an averaging relationship between the SS and RR optima (Proposition 3.3.1) which can help RR reduce distortion. We exemplify this averaging relationship with a simple example and extend it to higher dimensions with experiments on synthetic data.

3.1. Framework: the idea of distorted risks

We now formally introduce the notion of a *distorted risk*. Distorted risks are crucial to our analysis, as they encode the interaction between shuffling SGD and BN. We show that these distorted risks \mathcal{L}_π and \mathcal{L}_{RR} are respectively induced by certain batch normalized datasets $\overline{\mathbf{X}}_\pi$ and $\overline{\mathbf{X}}_{\text{RR}}$ obtained by batch normalizing the input features.

The *undistorted* risk we actually want to minimize is the risk that corresponds to full-batch GD. Define the GD features $\overline{\mathbf{X}}_{\text{GD}} \triangleq \text{BN}(\mathbf{X})$, which induces the *GD risk*:

$$\mathcal{L}_{\text{GD}}(\Theta) \triangleq \mathcal{L}(f(\mathbf{X}; \Theta), \mathbf{Y}) = \mathcal{L}(\mathbf{W}\mathbf{T}\overline{\mathbf{X}}_{\text{GD}}, \mathbf{Y}).$$

However, during epoch k , SS or RR optimize a distorted risk dependent on π_k . To see why, define the SS dataset

$$\begin{aligned} \overline{\mathbf{X}}_\pi &\triangleq \text{BN}_\pi(\mathbf{X}) \triangleq [\text{BN}(\mathbf{X}_\pi^1) \quad \dots \quad \text{BN}(\mathbf{X}_\pi^m)] \\ \mathbf{Y}_\pi &\triangleq [\mathbf{Y}_\pi^1 \quad \dots \quad \mathbf{Y}_\pi^m], \end{aligned}$$

for every permutation $\pi \in \mathbb{S}_n$. Similarly, form the RR dataset $(\overline{\mathbf{X}}_{\text{RR}}, \mathbf{Y}_{\text{RR}}) \in \mathbb{R}^{d \times (n \cdot n!)} \times \mathbb{R}^{p \times (n \cdot n!)}$ by concatenating the SS datasets $(\overline{\mathbf{X}}_\pi, \mathbf{Y}_\pi)$ across all π .

Crucially, the SS data $\overline{\mathbf{X}}_\pi$ encodes the distortion due to the interaction between SS with permutation π and BN; the RR data $\overline{\mathbf{X}}_{\text{RR}}$ does the same for RR and BN. Indeed, since SS uses a fixed π , it implicitly optimizes the SS distorted risk

$$\mathcal{L}_\pi(\Theta) \triangleq \sum_{j=1}^m \mathcal{L}(f(\mathbf{X}_\pi^j; \Theta), \mathbf{Y}_\pi^j) = \mathcal{L}(\mathbf{W}\mathbf{T}\overline{\mathbf{X}}_\pi, \mathbf{Y}_\pi).$$

Likewise, by collapsing the epoch update into a noisy ‘‘SGD’’ update, we observe that RR over epochs implicitly optimizes the RR distorted risk

$$\mathcal{L}_{\text{RR}}(\Theta) \triangleq \frac{1}{n!} \sum_{\pi \in \mathbb{S}_n} \mathcal{L}_\pi(\Theta) = \frac{1}{n!} \mathcal{L}(\mathbf{W}\mathbf{T}\overline{\mathbf{X}}_{\text{RR}}, \mathbf{Y}_{\text{RR}}).$$

We reiterate that SS and RR distortions originate from using *both* shuffling and batch normalization: shuffling alters the batch-dependent affine transforms that BN applies. With this notation, the connection between SS+BN/RR+BN and ghost BN becomes more evident: one can view the full batch as the batch in ghost BN and the mini-batches as the virtual ghost batches. Moreover, the proofs of Theorems 3.2.2 and 3.2.3 demonstrate that ghost BN would witness the same type of distortion as SS+BN/RR+BN.

To aid clarity, we adopt the convention that overlines connote batch normalization with *some* batching, and vice versa. For example, the SS dataset $\overline{\mathbf{X}}_\pi \triangleq \text{BN}_\pi(\mathbf{X})$ is normalized, while the shuffled dataset $\mathbf{X}_\pi = \pi \circ \mathbf{X}$ is not.

3.2. Convergence results for regression

We now present our main regression results: SS+BN and RR+BN converge to the global optima of their respective distorted risks encoded by the SS dataset $\overline{\mathbf{X}}_\pi$ and the RR dataset $\overline{\mathbf{X}}_{\text{RR}}$. We require the following rank assumptions.

Assumption 1 (Full rank assumption).

- (a) $\overline{\mathbf{X}}_\pi \in \mathbb{R}^{d \times n}$ satisfies $\text{rank}(\overline{\mathbf{X}}_\pi) \geq d$. In particular, $\sigma_{\min}(\overline{\mathbf{X}}_\pi \overline{\mathbf{X}}_\pi^\top) > 0$.
- (b) $\overline{\mathbf{X}}_{\text{RR}} \in \mathbb{R}^{d \times (n \cdot n!)}$ satisfies $\text{rank}(\overline{\mathbf{X}}_{\text{RR}}) \geq d$. In particular, $\sigma_{\min}(\overline{\mathbf{X}}_{\text{RR}} \overline{\mathbf{X}}_{\text{RR}}^\top) > 0$.

It is natural to ask when Assumption 1 holds. We demonstrate that the following mild assumption implies it; the assumption states that the feature matrix \mathbf{X} is drawn from a joint density on matrices in a *potentially non-i.i.d. fashion*.

Assumption 2. \mathbf{X} is drawn from a density with respect to the Lebesgue measure on $\mathbb{R}^{d \times n}$.

Since BN centers the mini-batch features, we have $\text{rank}(\overline{\mathbf{X}}_\pi) \leq \min\{d, (B-1)\frac{n}{B}\}$ and $\text{rank}(\overline{\mathbf{X}}_{\text{RR}}) \leq$

$\min \{d, (B-1)\binom{n}{B}\}^2$. We now show that if $B > 2$ these upper bounds are achieved almost surely. Thus, we identify reasonable conditions under which Assumption 1 holds almost surely over the draw of data, irrespective of shuffling.

Proposition 3.2.1. *Assume Assumption 2 and $B > 2$. Then we have $\text{rank}(\overline{\mathbf{X}}_\pi) = \min \{d, (B-1)\frac{n}{B}\}$ and $\text{rank}(\overline{\mathbf{X}}_{\text{RR}}) = \min \{d, (B-1)\binom{n}{B}\}$ a.s.. Consequently, if $(B-1)\frac{n}{B} \geq d$, Assumption 1(a) holds a.s. for $\overline{\mathbf{X}}_\pi$, and if $(B-1)\binom{n}{B} \geq d$, Assumption 1(b) holds a.s. for $\overline{\mathbf{X}}_{\text{RR}}$.*

Although we could have just assumed Assumption 1, the nonlinearity introduced by BN makes it nontrivial to identify mild sufficient conditions on the original features to control the rank of SS and RR datasets. Furthermore, controlling the rank of these datasets is crucial to our analysis of GD risk divergence in the classification setting (see Section 4).

Next, we present our main SS convergence result: SS converges for appropriate stepsizes. We defer the proof and explicit convergence rates to Appendix A.1.

Theorem 3.2.2 (Convergence of SS). *Let $f(\cdot; \Theta) = \mathbf{W}\Gamma\text{BN}(\cdot)$ be a linear+BN network initialized at $\Theta_0^1 = (\mathbf{W}_0^1, \Gamma_0^1) = (\mathbf{0}, \mathbf{I})$. We train f using SS with permutation π and suppose that Assumption 1(a) holds for this π . SS uses the following decreasing stepsize, which is well-defined:*

$$\eta_k = \frac{1}{k^\beta} \cdot \min \left\{ O \left(\frac{1}{\sigma_{\min}(\overline{\mathbf{X}}_\pi \overline{\mathbf{X}}_\pi^\top)} \right), \frac{\sqrt{2\beta-1} \text{poly}(\sigma_{\min}(\overline{\mathbf{X}}_\pi^\top))}{\text{poly}(n, d, \|\mathbf{Y}\|_F)} \right\},$$

where $1/2 < \beta < 1$. Then the risk $\mathcal{L}_\pi(\Theta_0^k)$ converges to the global minimum \mathcal{L}_π^* as $k \rightarrow \infty$.

Theorem 3.2.2 shows that using both SS and BN induces the network to converge to the global optimum of the SS distorted risk instead of the usual GD risk. The proof proceeds by aggregating the epoch-wise gradient updates on the collapsed matrix $\mathbf{W}\Gamma$. The main difficulty lies in carefully bounding the accumulation of various types of noise.

We now turn to RR convergence. For the sake of analysis, we make the following compact iterates assumption which is common in the RR literature (Haochen & Sra, 2019; Nagaraj et al., 2019; Ahn et al., 2020; Rajput et al., 2020).

Assumption 3. *For all (i, k) , the iterates $\Theta_i^k = (\mathbf{W}_i^k, \Gamma_i^k)$ satisfy $\|\mathbf{W}_i^k \Gamma_i^k\|_2 \leq A_{\text{RR}}$ for some absolute constant A_{RR} .*

Finally, we can show that RR converges in expectation to the global optimum of the RR distorted risk \mathcal{L}_{RR} . We defer the proof and explicit convergence rates to Appendix A.2.

Theorem 3.2.3 (Convergence of RR). *Assume Assumption 1(b) and Assumption 3. Using the same f and initialization as in Theorem 3.2.2, we train training f using RR with*

²Note that $\overline{\mathbf{X}}_{\text{RR}}$ contains many duplicate batches; only $\binom{n}{B}$ of them are unique, up to permutations of B columns inside a batch.

the following decreasing stepsize, which is well-defined:

$$\eta_k = \frac{1}{k^\beta} \cdot \min \left\{ O \left(\frac{1}{\sigma_{\min}(\overline{\mathbf{X}}_{\text{RR}} \overline{\mathbf{X}}_{\text{RR}}^\top)} \right), \frac{\sqrt{2\beta-1}}{\text{poly}(n, d, \|\mathbf{Y}\|_F, A_{\text{RR}})} \right\},$$

where $1/2 < \beta < 1$. Then the risk $\mathcal{L}_{\text{RR}}(\Theta_0^k)$ converges in expectation to the global minimum $\mathcal{L}_{\text{RR}}^*$ as $k \rightarrow \infty$.

The proof of Theorem 3.2.3 is similar to the SS case; the main subtlety is using Assumption 3 to handle expectations.

The main takeaway of Theorems 3.2.2 and 3.2.3 is that SS+BN and RR+BN converge to the optima of the SS and RR distorted risks, respectively. These distorted optima may differ from optimum of the GD risk. Moreover, the required stepsize for convergence is usually smaller for SS (where the requirement depends on π) compared to RR.

3.3. RR averages out SS distortion

Having shown that the two different algorithms drive the network parameters to global optima of two different distorted risks, it behooves us to study these optima. By collapsing the final layers \mathbf{W} and Γ into a single matrix $\mathbf{M} = \mathbf{W}\Gamma \in \mathbb{R}^{p \times d}$, we can study the global optima M_π^* and M_{RR}^* on the normalized datasets $\overline{\mathbf{X}}_\pi$ and $\overline{\mathbf{X}}_{\text{RR}}$. These global optima naturally correspond to the global optima of \mathcal{L}_π and \mathcal{L}_{RR} . In this section we illustrate how RR can average out SS distortion in the one-dimensional case.

We first relate the SS optima M_π^* to the RR optimum M_{RR}^* . A simple gradient calculation reveals $M_{\text{RR}}^* = \sum_\pi \mathbf{Y}_\pi \overline{\mathbf{X}}_\pi^\top (\sum_\pi \overline{\mathbf{X}}_\pi \overline{\mathbf{X}}_\pi^\top)^{-1}$. Since BN enforces the unit variance constraint, $\overline{\mathbf{X}}_\pi \overline{\mathbf{X}}_\pi^\top = n$ if $d = 1$. Simple algebraic manipulation then implies the following proposition.

Proposition 3.3.1. *If $d = 1$, $M_{\text{RR}}^* = \frac{1}{n!} \sum_{\pi \in \mathbb{S}_n} M_\pi^*$.*

Proposition 3.3.1 identifies an explicit averaging relationship between RR and SS in the one-dimensional case. This motivates the following simple construction where RR's averaging behavior removes SS distortion.

Dataset: SS distorted with constant probability, RR averages out distortion.

We visualize our toy dataset with $16n$ datapoints where $d = p = 1$, $B = 2$, and $n = 3$ in Figure 2a, along with the possible SS optima M_π^* . The dataset is comprised of four clusters of $4n$ points in the square $[-1, 1]^2$. By vertical symmetry of the clusters and Proposition 3.3.1, the RR and GD optima coincide at zero. However, SS is distorted away from GD. An anticoncentration calculation shows $M_\pi^* \neq 0$ with probability $1 - O(\frac{1}{\sqrt{n}})$ and $|M_\pi^*| = \Omega(\frac{1}{\sqrt{n}})$ with constant probability. The key insight is linking SS distortion to breaking symmetry in the SS dataset (see Proposition E.1.1 for details).

3.4. Regression experiments

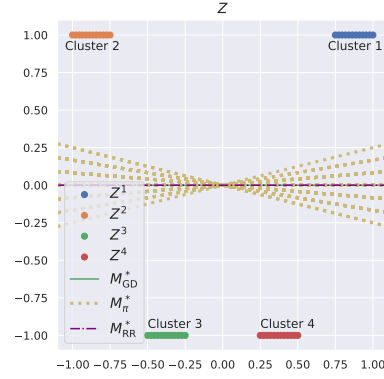
For our regression experiments, we used synthetic data with $n = 100$, $B = 10$, and $d = 10$. For $i \in [n]$, we sampled $\mathbf{x}_i \sim N(\mathbf{0}, \mathbf{I}_d)$ and generated $y_i = \mathbf{M}_{\text{true}} \mathbf{x}_i + \epsilon_i \in \mathbb{R}$ with $\mathbf{M}_{\text{true}} \sim U[-1, 1]^d$ and $\epsilon_i \sim N(0, 1)$. We trained the network $\mathbf{WTBN}(\mathbf{X})$ using SS and RR with an inverse learning rate schedule initialized at $\eta = 0.01$. We observed convergence to near optimal values on the SS and RR risks (Figure 7), which supports the convergence results (Theorems 3.2.2 and 3.2.3).

We also extended the toy dataset to the synthetic setup described above. As Figure 2b makes apparent, SS is consistently distorted away from the GD optimum, whereas RR averages out this distortion effect. We generated 500 datasets and evaluated the distortion for each one with the normalized distance $d(\mathbf{M}) \triangleq \frac{\|\mathbf{M} - \mathbf{M}_{\text{GD}}^*\|}{\|\mathbf{M}_{\text{GD}}^*\|}$. For SS, we computed the mean $d(\mathbf{M}_{\pi}^*)$ for 1000 random draws of π . For RR, we approximated $d(\mathbf{M}_{\text{RR}}^*)$ as follows. We sampled 1000 fresh random permutations to approximate the RR dataset $\overline{\mathbf{X}}_{\text{RR}}$, which we then used to approximate \mathbf{M}_{RR}^* (since it is intractable to average over all $n!$ permutations). We see that $d(\mathbf{M}_{\pi}^*) > 1$ for all of the SS experiments while $d(\mathbf{M}_{\pi}^*) \approx 0.1$ for all of the RR experiments.

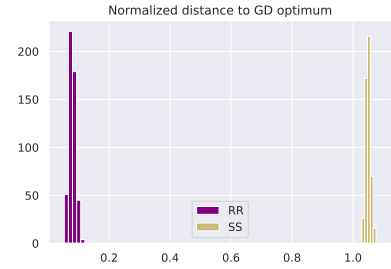
4. Main classification results: divergence regimes based on distorted risks

We now turn to analyzing linear+BN binary classifiers $f(\mathbf{X}; \Theta) = \text{sgn}(\mathbf{WTBN}(\mathbf{X}))$ trained with the logistic risk. To characterize divergence, we identify salient properties of the distorted risks first introduced in Section 3.1. These properties identify regimes where the SS+BN classifier can diverge on the GD risk (Theorem 4.1.3) yet the RR+BN classifier does not diverge (Theorem 4.1.4). This motivates the construction of a toy dataset (Section 4.2) where the optimal SS classifier diverges on the GD risk with constant probability. In Section 4.3 we extend our results to more realistic networks and datasets, demonstrating that these phenomena are not an artifact of our theoretical setup. Our theoretical results offer some justification for the empirical phenomenon of divergence when SS SGD is combined with BN for classification.

We briefly remark on why we analyze divergence conditions instead of directional convergence. The main difficulty lies in analyzing SGD instead of GD. One could hope to extend the techniques for directional convergence for homogeneous networks in Lyu & Li (2019) to the stochastic setting, but this is outside the scope of our paper. Furthermore, the analyses for deep linear networks such as Ji & Telgarsky (2020) rely on invariants which do not hold for us due to the diagonal $\mathbf{\Gamma}$ and the BN layers for deeper networks.



(a) Dataset with 48 datapoints demonstrating distortion of SS optima \mathbf{M}_{π}^* .



(b) Normalized distance to GD optimum $d(\mathbf{M}) = \|\mathbf{M} - \mathbf{M}_{\text{GD}}^*\| / \|\mathbf{M}_{\text{GD}}^*\|$.

Figure 2. Top: toy dataset for regression, showing how RR can average out the distortion of SS. Bottom: histogram of distortion of SS and RR optima on synthetic data for $d = 10$. The SS optima significantly deviate from the GD optima, whereas the RR optima are relatively close. This supports the intuition that RR can nontrivially smooth out the bias of SS in higher dimensions.

Throughout, we use $\mathbf{v} = (\mathbf{WT})^\top \in \mathbb{R}^d$ to refer to the vector that determines the decision boundary of our classifier f . We remind the reader of the datasets which induce the different distorted risks (Section 3.1). Given dataset $\mathbf{Z} = (\mathbf{X}, \mathbf{Y})$, the GD dataset is $\overline{\mathbf{Z}}_{\text{GD}} \triangleq (\overline{\mathbf{X}}_{\text{GD}}, \mathbf{Y}_{\text{GD}}) = (\text{BN}(\mathbf{X}), \mathbf{Y})$. Similarly define the SS dataset $\overline{\mathbf{Z}}_{\pi} \triangleq (\overline{\mathbf{X}}_{\pi}, \mathbf{Y}_{\pi}) = (\text{BN}_{\pi}(\mathbf{X}), \pi \circ \mathbf{Y})$ and the RR dataset $\overline{\mathbf{Z}}_{\text{RR}} \triangleq (\overline{\mathbf{X}}_{\text{RR}}, \mathbf{Y}_{\text{RR}})$ by concatenating $\overline{\mathbf{Z}}_{\pi}$ over all permutations π . If the labels are clear from context, we occasionally abuse terminology and refer to the features as the dataset.

4.1. Analysis of problem structure for classification

To analyze the optima of the distorted risks, we introduce relevant concepts from Ji & Telgarsky (2019). Given a dataset $\mathbf{Z} = (\mathbf{X}, \mathbf{Y}) = \{(\mathbf{x}_i, y_i)\}_{i=1}^n$, with labels $y_i \in \{\pm 1\}$, greedily define a *maximal linearly separable subset* $\mathbf{S}^{\text{LS}} \triangleq (\mathbf{X}^{\text{LS}}, \mathbf{Y}^{\text{LS}})$ as follows. Include (\mathbf{x}_i, y_i) in \mathbf{S}^{LS} if there exists a classifier $\mathbf{u}_i \in \mathbb{R}^d$ with $y_i \mathbf{u}_i^\top \mathbf{x}_i > 0$ and $y_j \mathbf{u}_i^\top \mathbf{x}_j \geq 0$ for all j . For reasons that will be clear shortly, denote the complement of \mathbf{S}^{LS} in \mathbf{Z} by $\mathbf{S}^{\text{SC}} \triangleq (\mathbf{X}^{\text{SC}}, \mathbf{Y}^{\text{SC}})$.

In particular, there exists a classifier \mathbf{u} such that: (1) \mathcal{S}^{LS} is perfectly separated by \mathbf{u} (2) the datapoints \mathcal{X}^{SC} in \mathcal{S}^{SC} are orthogonal to \mathbf{u} , so they are on the decision boundary. We can choose \mathbf{u} to be the max-margin classifier \mathbf{u}^{MM} on \mathcal{S}^{LS} . The notation \mathcal{S}^{SC} is chosen because the logistic risk is strongly convex when restricted to bounded subsets of $\text{Span}(\mathcal{X}^{\text{SC}})$, meaning there is a unique finite minimizer \mathbf{v}^{SC} in this subspace. Ji & Telgarsky (2019) show that linear classifiers trained on the logistic risk with GD are implicitly biased towards the ray $\mathbf{v}^{\text{SC}} + t \cdot \mathbf{u}^{\text{MM}}$ for $t > 0$.

We now identify a salient property of the distorted risks.

Definition 1 (Separability decomposition). *The separability decomposition of dataset \mathcal{Z} refers to $\mathcal{Z} = \mathcal{S}^{\text{LS}} \sqcup \mathcal{S}^{\text{SC}}$.*

If $\mathcal{S}^{\text{LS}} = \mathcal{Z}$, we say \mathcal{Z} is linearly separable (LS). If both \mathcal{S}^{LS} and \mathcal{S}^{SC} are nonempty, we say \mathcal{Z} is *partially linearly separable* (PLS). Finally, if $\mathcal{S}^{\text{SC}} = \mathcal{Z}$, we slightly abuse terminology and say \mathcal{Z} is strongly convex (SC).³

Because the logistic loss does not always have finite infima, we now introduce the notion of an optimal direction.

Definition 2 (Optimal direction). *Given dataset $\mathcal{Z} = (\mathbf{X}, \mathbf{Y})$, we say a sequence of iterates $\mathbf{v}(t)$ infimizes \mathcal{L} if $\mathcal{L}(\mathbf{v}(t)^\top \mathbf{X}, \mathbf{Y}) \rightarrow \inf_{\mathbf{w} \in \mathbb{R}^d} \mathcal{L}(\mathbf{w}^\top \mathbf{X}, \mathbf{Y})$. We call $\mathbf{v} \in \mathbb{R}^d$ an optimal direction if there exists $\mathbf{u} \in \mathbb{R}^d$ such that $\{\mathbf{u} + t\mathbf{v}\}_{t \geq 1}$ infimizes \mathcal{L} .⁴*

Definition 1 is motivated by the following results which identify how the separability decomposition affects optimal directions. Their proofs are deferred to Appendix B.4.

Lemma 4.1.1. *Let $\mathcal{Z} = \mathcal{S}^{\text{LS}} \sqcup \mathcal{S}^{\text{SC}}$. If \mathbf{v} is an optimal direction for \mathcal{L} , then $\mathbf{v}^\top \mathbf{x} = 0$ for all $\mathbf{x} \in \text{Span}(\mathcal{X}^{\text{SC}})$ and $y_i \mathbf{v}^\top \mathbf{x}_i > 0$ for every $(\mathbf{x}_i, y_i) \in \mathcal{S}^{\text{LS}}$.*

Combining the above lemma and the definitions yields the following proposition, which characterizes SS and RR divergence using the separability decomposition.

Proposition 4.1.2. *Suppose Assumption 1(a) holds, the iterates $\mathbf{v}_\pi(t)$ infimize \mathcal{L}_π , and their projections onto $\text{Span}(\overline{\mathcal{X}}_\pi^{\text{SC}})^\perp$ converge in direction to some optimal direction \mathbf{v}_π^* for \mathcal{L}_π . Then the GD risk \mathcal{L}_{GD} diverges if and only if $\overline{\mathcal{Z}}_\pi$ is PLS or LS and there exists some $(\mathbf{x}_i, y_i) \in \overline{\mathcal{Z}}_{\text{GD}}$ such that $y_i \mathbf{v}_\pi^{*\top} \mathbf{x}_i < 0$. The analogous statement holds true for $\overline{\mathcal{Z}}_{\text{RR}}$ under Assumption 1(b). Furthermore, the “if” part holds true for SS and RR without Assumption 1.*

In particular, Proposition 4.1.2 implies that if the RR dataset is SC and rank d , the GD risk does not diverge. Moreover,

³Here, PLS refers to the “general case” discussed in Ji & Telgarsky (2019), but we chose to use this alternative terminology because we found the term “general” can lead to confusion.

⁴This definition is catered towards the SC+full rank \mathbf{X} or PLS/LS case. However, since Proposition 3.2.1 provides sufficient conditions for full-rank data, this subtlety is unimportant.

it naturally leads to the question of understanding ranks and separability decompositions of the SS and RR datasets; the former question is already answered by Proposition 3.2.1.

To analyze the separability decomposition with high probability or almost surely, we assume the labels are balanced.

Assumption 4 (Balanced classes). *The data \mathcal{Z} either has*

- an equal number of positive and negative examples; or*
- at least B positive and B negative examples.*

Finally, we informally state our main classification result: SS+BN can diverge in some regimes (see Theorem B.2.1 for details).

Theorem 4.1.3 (SS+BN can diverge (informal)). *Assume Assumption 2, Assumption 4(a), and $B > 2$. If $d \leq (B - 1) \frac{n}{B}$, SS can diverge if $B = \Omega(\log n)$ and $\overline{\mathcal{Z}}_{\text{GD}}$ ’s separability decomposition can change with small perturbations. Otherwise, SS can diverge regardless of the batch size and the separability decomposition of $\overline{\mathcal{Z}}_{\text{GD}}$.*

Whereas Theorem 4.1.3 establishes regimes where SS+BN can diverge, we can show that RR+BN prevents divergence in a much larger regime (see Theorem B.3.1 for details).

Theorem 4.1.4 (RR+BN does not diverge (informal)). *Assume Assumption 2, Assumption 4(b), and $B > 2$. If $d \leq (B - 1) \binom{n}{B}$, RR does not diverge almost surely.*

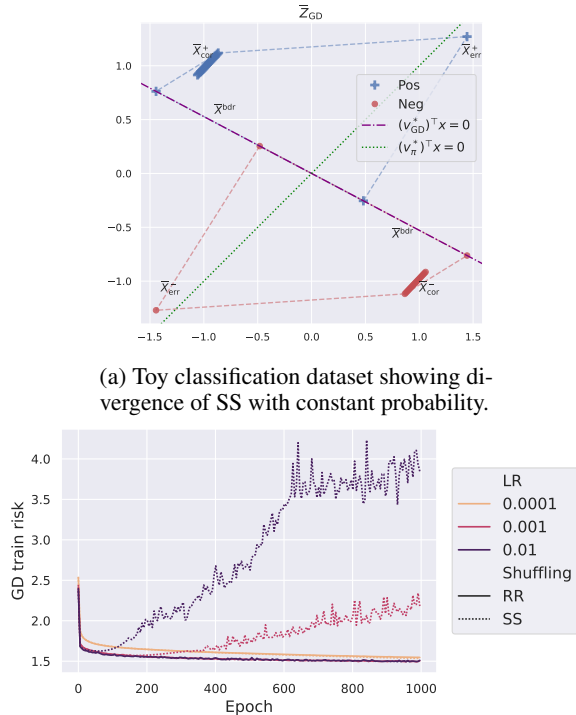
Theorem 4.1.3 implies that one cannot prevent SS divergence by simply increasing the batch size B ; it is also necessary for the GD dataset to be “robustly” LS or SC. Moreover, as soon as $d > (B - 1) \frac{n}{B}$, SS can diverge. In stark contrast, Theorem 4.1.4 establishes that even for small B , RR is *almost surely* robust to divergence as long as $d \leq (B - 1) \binom{n}{B}$. Although our theorems do not prove that SS+BN *necessarily* diverges, they offer some theoretical explanation for why SS+BN appears to be less stable than RR+BN for classification.

4.2. RR prevents divergence while SS diverges

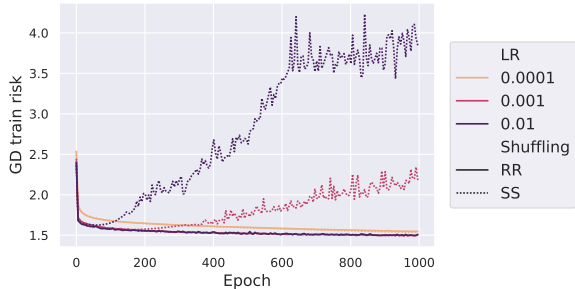
We present a toy dataset where SS drastically distorts the optimal direction, leading to divergence with constant probability. Meanwhile, RR does not diverge on this dataset. We use $d = B = 2$ to simplify the construction.⁵

Dataset: SS often diverges; RR does not. We describe our construction (Figure 3a) at a high level; see Proposition E.2.1 for details. The GD dataset is PLS with unique optimal direction \mathbf{v}_{GD}^* (its decision boundary is the purple dash-dotted line). Moreover, with constant probability the SS dataset is PLS with unique optimal direction \mathbf{v}_π^* (green dotted line) distorted away from \mathbf{v}_{GD}^* . Also, \mathbf{v}_π^* misclassifies points in the GD dataset ($\overline{\mathcal{X}}_{\text{err}}^+$ and $\overline{\mathcal{X}}_{\text{err}}^-$). Under the

⁵Since $B = 2$, there is no contradiction with Theorem 4.1.3.



(a) Toy classification dataset showing divergence of SS with constant probability.



(b) 3 layer linear+BN networks trained with varying stepsizes.

Figure 3. Top: Toy dataset demonstrating divergence of GD risk with constant probability. The dashed lines trace out the convex hulls of the positive and negative points. Bottom: Divergence of GD risk for a variety of stepsizes on CIFAR10. Note that there was eventually a separation for $\eta = 10^{-4}$ (see Figure 8).

additional assumptions in Proposition 4.1.2, the GD risk diverges. Finally, since the RR dataset is SC and rank d , RR does not diverge on the GD dataset.

4.3. Experiments on linear and nonlinear networks

We now verify our theoretical classification results on linear+BN and extend them to nonlinear networks on a variety of real-world datasets. This demonstrates that the separation between SS, RR, and GD is relevant in realistic settings and not merely an artifact of the linear setting. We refer to the linear+BN network $\mathcal{WTBN}(\mathbf{X})$ as 1-layer linear network, and also consider deeper linear networks with tunable parameters inside BN layers. We observe strikingly different training behaviors in the shallow and deep linear networks. The networks are formally defined in Appendix D; see <https://github.com/davidxwu/sgd-batchnorm-icml> for the experiment code.

As a motivating example, we ran an experiment on synthetic data (Figure 4) with the 2-layer linear network $f(\mathbf{X}) = \mathcal{WTBN}(\mathbf{A}\mathbf{X})$. Note that the tunable matrix \mathbf{A} acts before BN. Intriguingly, we observe that the SS dataset with features $\bar{\mathbf{X}}_\pi = \text{BN}_\pi(\mathbf{A}\mathbf{X})$ is SC at initialization, but up-

dating \mathbf{A} with SS makes it LS after training. Moreover, the batch size is large relative to n , so this dataset satisfies the necessary conditions for divergence in Proposition 4.1.2 and Theorem 4.1.3.

More specifically, Figures 4a and 4c plot the 2-dimensional GD and SS datasets, respectively, which are SC at initialization. However, after training with SS, we can see from Figures 4b and 4d that SS updates \mathbf{A} to make the SS dataset LS, whereas the GD dataset stays SC. Hence, by Proposition 4.1.2, the GD risk diverges. This example partially explains the discrepancy in training behavior between the 1-layer and deeper networks. Indeed, whereas the 1-layer architecture has static $\bar{\mathbf{Z}}_\pi$, the deeper networks have evolving weights inside BN which can push $\bar{\mathbf{Z}}_\pi$ to be LS/PLS.

To exhibit the above divergence on real data, we conducted experiments on the CIFAR10. Using SS and RR, we trained linear+BN networks of depths up to 3 for $T = 10^3$ epochs using stepsize $\eta = 10^{-2}$, batch size $B = 128$, and 512 hidden units per layer (see Appendix D for precise details).

As depicted in Figure 1a, we consistently observed SS divergence for the deeper networks (see Figure 9 for more evidence of divergence). As predicted by Theorem 4.1.4, RR did not exhibit divergence behavior. These phenomena persisted despite ablating the learning rate in $\{0.01, 0.001, 0.0001\}$, momentum in $\{0, 0.9, 0.99\}$, and batch size in $\{32, 64, 128\}$. The learning rate ablation is shown in Figure 3b; see Appendix D for the rest.

For the nonlinear experiments, we extended to the CIFAR10, MNIST, and CIFAR100 datasets. We used SS and RR to train 3-layer 512 hidden unit MLPs with BN and ReLU activation for $T = 10^3$ epochs, and also to finetune pretrained ResNet18 for $T = 50$ epochs. We consistently observed that in the final stages of training (i.e., relatively small training risk), SS trained slower than RR across all of the datasets, even after tuning the learning rate (see Figures 5 and 6).

5. Conclusion

This paper established that training BN networks with SS can lead to undesirable training behavior, including slower convergence or even divergence of the GD risk. However, RR provably mitigates this divergence behavior, and experimental evidence suggests that using RR usually converges faster than SS. This separation in training behavior between SS, RR, and GD is because data shuffling directly affects how BN operates on mini-batches. Our theoretical results establish a separation for the special case where BN is applied to the input features. The more general and realistic case where BN is applied to dynamically evolving layers is left as an important direction for future work. We also observed in preliminary experiments that a similar separation manifested for generalization, and we hope that adopting

Training Instability of Shuffling SGD with Batch Norm

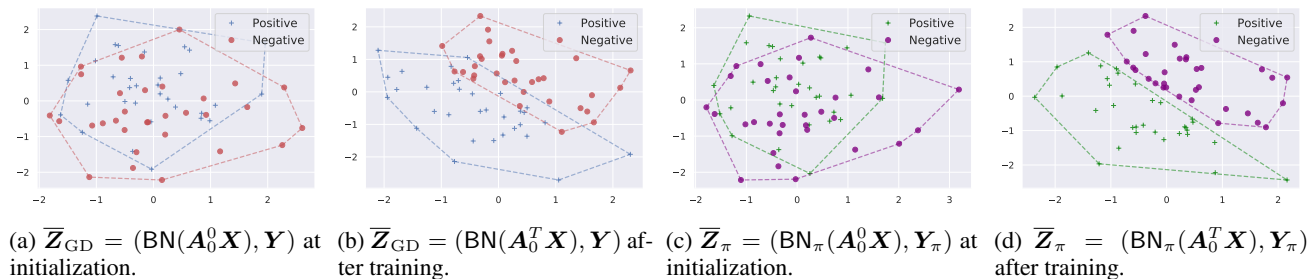


Figure 4. Snapshots of GD dataset $\bar{\mathcal{Z}}_{\text{GD}}$ and SS dataset $\bar{\mathcal{Z}}_{\pi}$ before and after running SS for $T = 10^4$ epochs with 32 positive and negative synthetic examples. While the GD dataset remains SC, the SS dataset become LS. Here $B = 16$, $\eta = 10^{-2}$, and $\epsilon = 10^{-5}$ for BN.

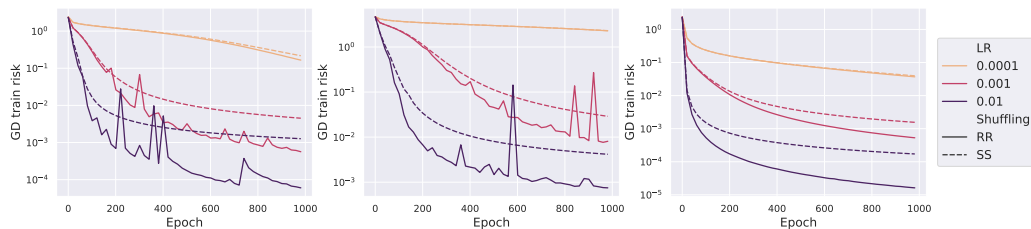


Figure 5. 3 layer ReLU+BN MLP on (left to right): CIFAR10, CIFAR100, and MNIST. Note the slower convergence for SS versus RR in the final stages of training for CIFAR10 and MNIST, and the early stages of training for CIFAR100.

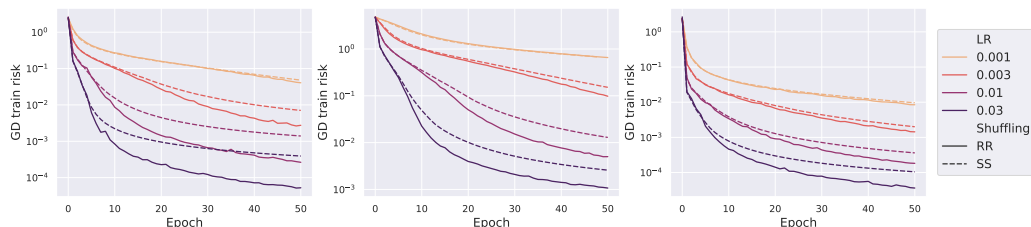


Figure 6. ResNet18 finetuned on (left to right): CIFAR10, CIFAR100, and MNIST. Note the slower convergence for SS versus RR across datasets in the final stages of training for CIFAR10 and MNIST, and the early stages of training for CIFAR100. For the smallest learning rate $\eta = 10^{-3}$, we observed a separation after 200 epochs.

a similar perspective will prove fruitful in pursuing this direction. We remark that similar surprising phenomena may arise when using other design choices that are implemented in a mini-batch fashion such as mixup (Zhang et al., 2017) and Sharpness-Aware Minimization (SAM) (Foret et al., 2020). For these reasons, we generally recommend that practitioners use RR instead of SS. Further future directions include establishing directional convergence for homogeneous classifiers trained with shuffling SGD and theoretically understanding conditions under which deeper networks diverge faster.

Acknowledgements

Part of the work was done while DW was an undergraduate at MIT. DW acknowledges support from NSF Graduate

Research Fellowship DGE-2146752. CY acknowledges support by Institute of Information & communications Technology Planning & evaluation (IITP) grant (No. 2019-0-00075, Artificial Intelligence Graduate School Program (KAIST)) funded by the Korea government (MSIT). CY is also supported by the National Research Foundation of Korea (NRF) grants (No. NRF-2019R1A5A1028324, RS-2023-00211352) funded by the Korea government (MSIT). CY acknowledges support from a grant funded by Samsung Electronics Co., Ltd. SS acknowledge support from an NSF CAREER grant (1846088), and NSF CCF-2112665 (TILOS AI Research Institute). DW appreciates helpful discussions with Xiang Cheng, Sidhanth Mohanty, Erik Jenner, Louis Golowich, Sam Gunn, and Thiago Bergamaschi.

References

- Agarwal, P. K., Guibas, L. J., Har-Peled, S., Rabinovitch, A., and Sharir, M. Penetration depth of two convex polytopes in 3d. *Nord. J. Comput.*, 7(3):227–240, 2000.
- Ahn, K., Yun, C., and Sra, S. Sgd with shuffling: optimal rates without component convexity and large epoch requirements. *arXiv preprint arXiv:2006.06946*, 2020.
- Arora, S., Li, Z., and Lyu, K. Theoretical analysis of auto rate-tuning by batch normalization. *arXiv preprint arXiv:1812.03981*, 2018.
- Ba, J. L., Kiros, J. R., and Hinton, G. E. Layer normalization. *arXiv preprint arXiv:1607.06450*, 2016.
- Bardenet, R. and Maillard, O.-A. Concentration inequalities for sampling without replacement. *Bernoulli*, 21:1361–1385, 2015.
- Bjorck, N., Gomes, C. P., Selman, B., and Weinberger, K. Q. Understanding batch normalization. *Advances in neural information processing systems*, 31, 2018.
- Boyd, S., Boyd, S. P., and Vandenberghe, L. *Convex optimization*. Cambridge university press, 2004.
- Cai, Y., Li, Q., and Shen, Z. A quantitative analysis of the effect of batch normalization on gradient descent. In *International Conference on Machine Learning*, pp. 882–890. PMLR, 2019.
- Cha, J., Lee, J., and Yun, C. Tighter lower bounds for shuffling SGD: Random permutations and beyond. *arXiv preprint arXiv:2303.07160*, 2023.
- Cho, H. and Yun, C. SGDA with shuffling: faster convergence for nonconvex-PL minimax optimization. In *International Conference on Learning Representations*, 2023.
- Daneshmand, H., Kohler, J., Bach, F., Hofmann, T., and Lucchi, A. Batch normalization provably avoids ranks collapse for randomly initialised deep networks. *Advances in Neural Information Processing Systems*, 33: 18387–18398, 2020.
- Daneshmand, H., Joudaki, A., and Bach, F. Batch normalization orthogonalizes representations in deep random networks. *Advances in Neural Information Processing Systems*, 34:4896–4906, 2021.
- Duchi, J., Hazan, E., and Singer, Y. Adaptive subgradient methods for online learning and stochastic optimization. *Journal of machine learning research*, 12(7), 2011.
- Foret, P., Kleiner, A., Mobahi, H., and Neyshabur, B. Sharpness-aware minimization for efficiently improving generalization. *arXiv preprint arXiv:2010.01412*, 2020.
- Gunasekar, S., Lee, J., Soudry, D., and Srebro, N. Characterizing implicit bias in terms of optimization geometry. In *International Conference on Machine Learning*, pp. 1832–1841. PMLR, 2018.
- Gürbüzbalaban, M., Ozdaglar, A., Vanli, N. D., and Wright, S. J. Randomness and permutations in coordinate descent methods. *Mathematical Programming*, 181:349–376, 2020.
- Haochen, J. and Sra, S. Random shuffling beats sgd after finite epochs. In *International Conference on Machine Learning*, pp. 2624–2633. PMLR, 2019.
- Hoffer, E., Hubara, I., and Soudry, D. Train longer, generalize better: closing the generalization gap in large batch training of neural networks. *Advances in Neural Information Processing Systems*, 30, 2017.
- Ioffe, S. and Szegedy, C. Batch normalization: Accelerating deep network training by reducing internal covariate shift. In *International conference on machine learning*, pp. 448–456. PMLR, 2015.
- Jagadeesan, M., Razenshteyn, I., and Gunasekar, S. Inductive bias of multi-channel linear convolutional networks with bounded weight norm. In *Conference on Learning Theory*, pp. 2276–2325. PMLR, 2022.
- Ji, Z. and Telgarsky, M. Gradient descent aligns the layers of deep linear networks. *arXiv preprint arXiv:1810.02032*, 2018.
- Ji, Z. and Telgarsky, M. The implicit bias of gradient descent on nonseparable data. In *Conference on Learning Theory*, pp. 1772–1798. PMLR, 2019.
- Ji, Z. and Telgarsky, M. Directional convergence and alignment in deep learning. *Advances in Neural Information Processing Systems*, 33:17176–17186, 2020.
- Johnson, C. R. *Matrix theory and applications*, volume 40. American Mathematical Soc., 1990.
- Kingma, D. P. and Ba, J. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014.
- Kohler, J., Daneshmand, H., Lucchi, A., Zhou, M., Neymeyr, K., and Hofmann, T. Towards a theoretical understanding of batch normalization. *stat*, 1050:27, 2018.
- Kohler, J., Daneshmand, H., Lucchi, A., Hofmann, T., Zhou, M., and Neymeyr, K. Exponential convergence rates for batch normalization: The power of length-direction decoupling in non-convex optimization. In *The 22nd International Conference on Artificial Intelligence and Statistics*, pp. 806–815. PMLR, 2019.

- Li, Z. and Arora, S. An exponential learning rate schedule for deep learning. *arXiv preprint arXiv:1910.07454*, 2019.
- Li, Z., Lyu, K., and Arora, S. Reconciling modern deep learning with traditional optimization analyses: The intrinsic learning rate. *Advances in Neural Information Processing Systems*, 33:14544–14555, 2020.
- Lobacheva, E., Kodryan, M., Chirkova, N., Malinin, A., and Vetrov, D. P. On the periodic behavior of neural network training with batch normalization and weight decay. *Advances in Neural Information Processing Systems*, 34: 21545–21556, 2021.
- Lyu, K. and Li, J. Gradient descent maximizes the margin of homogeneous neural networks. *arXiv preprint arXiv:1906.05890*, 2019.
- Lyu, K., Li, Z., and Arora, S. Understanding the generalization benefit of normalization layers: Sharpness reduction. *Advances in Neural Information Processing Systems*, 36, 2022.
- Maurer, A. Concentration inequalities for functions of independent variables. *Random Structures & Algorithms*, 29 (2):121–138, 2006.
- Mityagin, B. The zero set of a real analytic function. *arXiv preprint arXiv:1512.07276*, 2015.
- Nacson, M. S., Lee, J., Gunasekar, S., Savarese, P. H. P., Srebro, N., and Soudry, D. Convergence of gradient descent on separable data. In Chaudhuri, K. and Sugiyama, M. (eds.), *Proceedings of the Twenty-Second International Conference on Artificial Intelligence and Statistics*, volume 89 of *Proceedings of Machine Learning Research*, pp. 3420–3428. PMLR, 16–18 Apr 2019a. URL <https://proceedings.mlr.press/v89/nacson19b.html>.
- Nacson, M. S., Srebro, N., and Soudry, D. Stochastic gradient descent on separable data: Exact convergence with a fixed learning rate. In Chaudhuri, K. and Sugiyama, M. (eds.), *Proceedings of the Twenty-Second International Conference on Artificial Intelligence and Statistics*, volume 89 of *Proceedings of Machine Learning Research*, pp. 3051–3059. PMLR, 16–18 Apr 2019b. URL <https://proceedings.mlr.press/v89/nacson19a.html>.
- Nagaraj, D., Jain, P., and Netrapalli, P. Sgd without replacement: Sharper rates for general smooth convex functions. In *International Conference on Machine Learning*, pp. 4703–4711. PMLR, 2019.
- Nguyen, L. M., Tran-Dinh, Q., Phan, D. T., Nguyen, P. H., and van Dijk, M. A unified convergence analysis for shuffling-type gradient methods. *Journal of Machine Learning Research*, 22(207):1–44, 2021.
- Rajput, S., Gupta, A., and Papailiopoulos, D. Closing the convergence gap of sgd without replacement. In *International Conference on Machine Learning*, pp. 7964–7973. PMLR, 2020.
- Safran, I. and Shamir, O. How good is sgd with random shuffling? In *Conference on Learning Theory*, pp. 3250–3284. PMLR, 2020.
- Santurkar, S., Tsipras, D., Ilyas, A., and Madry, A. How does batch normalization help optimization? *Advances in neural information processing systems*, 31, 2018.
- Shallue, C. J., Lee, J., Antognini, J., Sohl-Dickstein, J., Frostig, R., and Dahl, G. E. Measuring the effects of data parallelism on neural network training. *Journal of Machine Learning Research*, 20(112):1–49, 2019.
- Soudry, D., Hoffer, E., Nacson, M. S., Gunasekar, S., and Srebro, N. The implicit bias of gradient descent on separable data. *The Journal of Machine Learning Research*, 19(1):2822–2878, 2018.
- Summers, C. and Dinneen, M. J. Four things everyone should know to improve batch normalization. In *International Conference on Learning Representations*, 2020. URL <https://openreview.net/forum?id=HJx8HANFDH>.
- Sun, R., Luo, Z.-Q., and Ye, Y. On the efficiency of random permutation for admm and coordinate descent. *Mathematics of Operations Research*, 45(1):233–271, 2020.
- Thomas, M. and Joy, A. T. *Elements of information theory*. Wiley-Interscience, 2006.
- Ulyanov, D., Vedaldi, A., and Lempitsky, V. Instance normalization: The missing ingredient for fast stylization. *arXiv preprint arXiv:1607.08022*, 2016.
- Wan, R., Zhu, Z., Zhang, X., and Sun, J. Spherical motion dynamics: Learning dynamics of normalized neural network using SGD and weight decay. In Beygelzimer, A., Dauphin, Y., Liang, P., and Vaughan, J. W. (eds.), *Advances in Neural Information Processing Systems*, 2021. URL <https://openreview.net/forum?id=RcbpHT7qjTq>.
- Woodworth, B., Gunasekar, S., Lee, J. D., Moroshko, E., Savarese, P., Golan, I., Soudry, D., and Srebro, N. Kernel and rich regimes in overparametrized models. In *Conference on Learning Theory*, pp. 3635–3673. PMLR, 2020.

- Wu, L., Wang, Q., and Ma, C. Global convergence of gradient descent for deep linear residual networks. *arXiv preprint arXiv:1911.00645*, 2019.
- Wu, Y. and Johnson, J. Rethinking” batch” in batchnorm. *arXiv preprint arXiv:2105.07576*, 2021.
- Yong, H., Huang, J., Meng, D., Hua, X., and Zhang, L. Momentum batch normalization for deep learning with small batch size. In *European Conference on Computer Vision*, pp. 224–240. Springer, 2020.
- Yun, C., Krishnan, S., and Mobahi, H. A unifying view on implicit bias in training linear neural networks. In *International Conference on Learning Representations*, 2021a.
- Yun, C., Sra, S., and Jadbabaie, A. Open problem: Can single-shuffle SGD be better than reshuffling SGD and GD? In *Conference on Learning Theory*, pp. 4653–4658. PMLR, 2021b.
- Yun, C., Rajput, S., and Sra, S. Minibatch vs local sgd with shuffling: Tight convergence bounds and beyond. In *International Conference on Learning Representations*, 2022.
- Zhang, H., Cisse, M., Dauphin, Y. N., and Lopez-Paz, D. mixup: Beyond empirical risk minimization. *arXiv preprint arXiv:1710.09412*, 2017.
- Zhou, Y. and Liang, Y. Characterization of gradient dominance and regularity conditions for neural networks. *arXiv preprint arXiv:1710.06910*, 2017.

A. Proofs for regression results

In this appendix, we provide the full details for the proof of convergence for SS and RR in the regression case.

Additional notation. We introduce some additional notation which we will use throughout the proof of Theorems 3.2.2 and 3.2.3. For a matrix \mathbf{A} , we use $\mathbf{A}_{i,:}$ and $\mathbf{A}_{:,j}$ to denote the i th row and j th column of \mathbf{A} , respectively. We also use $A_{i,j}$ to denote the (i, j) th entry of \mathbf{A} . The Hadamard product of two matrices $\mathbf{A}, \mathbf{B} \in \mathbb{R}^{m \times n}$ is denoted by $\mathbf{A} \odot \mathbf{B}$, with $(\mathbf{A} \odot \mathbf{B})_{i,j} = A_{i,j}B_{i,j}$. The diagonal operator $\text{diag} : \mathbb{R}^{m \times m} \rightarrow \mathbb{R}^{m \times m}$ is defined by $\text{diag}(\mathbf{A}) = \mathbf{I} \odot \mathbf{A}$. We denote the Frobenius inner product $\langle \mathbf{A}, \mathbf{B} \rangle_F = \sum_{i,j} A_{i,j}B_{i,j}$ and its induced norm by $\|\mathbf{A}\|_F$.

Also recall from Section 2 that when Θ is optimized with SS or RR, the i th iterate on the k th epoch is denoted by Θ_i^k . For simplicity, we will often say the (i, k) th iterate to refer to Θ_i^k . Denote the collapsed parameter matrix defined in Section 3 by $\mathbf{M} \triangleq \mathbf{W}\mathbf{\Gamma}$. We will abuse notation and sometimes denote the (i, k) th iterate by $\mathbf{M}_i^k \triangleq \mathbf{W}_i^k \mathbf{\Gamma}_i^k$.

Recall that the mini-batch risk used for updating the (i, k) th iterate of SS or RR is given by $\mathcal{L}(f(\mathbf{X}_\pi^{i+1}; \Theta_i^k), \mathbf{Y}_\pi^{i+1}) = \|\mathbf{Y}_\pi^{i+1} - \mathbf{W}_i^k \mathbf{\Gamma}_i^k \text{BN}(\mathbf{X}_\pi^{i+1})\|_F^2$ where π denotes the permutation chosen for the k th epoch and $\mathbf{X}_\pi^j \in \mathbb{R}^{d \times B}$ and $\mathbf{Y}_\pi^j \in \mathbb{R}^{p \times B}$ consist of the $(jB - B + 1, \dots, jB)$ th columns of $\pi \circ \mathbf{X}$ and $\pi \circ \mathbf{Y}$, respectively. Since this notation is a bit lengthy, we simplify it to $\mathcal{L}(\mathbf{X}_\pi^j; \Theta) \triangleq \mathcal{L}(f(\mathbf{X}_\pi^j; \Theta), \mathbf{Y}_\pi^j)$ for any $j \in [m]$. Here, we can also view the mini-batch risk as a function of $\mathbf{M} = \mathbf{W}\mathbf{\Gamma}$, so we will sometimes abuse notation and write

$$\begin{aligned} \mathcal{L}(\mathbf{X}_\pi^j; \mathbf{M}) &\triangleq \|\mathbf{Y}_\pi^j - \mathbf{M} \text{BN}(\mathbf{X}_\pi^j)\|_F^2, \\ \nabla_{\mathbf{M}} \mathcal{L}(\mathbf{X}_\pi^j; \mathbf{M}) &\triangleq -(\mathbf{Y}_\pi^j - \mathbf{M} \text{BN}(\mathbf{X}_\pi^j)) \text{BN}(\mathbf{X}_\pi^j)^\top. \end{aligned}$$

For SS, we work with a fixed permutation $\pi \in \mathbb{S}_n$ and input dataset (\mathbf{X}, \mathbf{Y}) . Recall that we defined $\overline{\mathbf{X}}_\pi \triangleq \text{BN}_\pi(\mathbf{X})$ from Section 3, i.e., the column-wise concatenation of all batches after batch normalization: $\overline{\mathbf{X}}_\pi = [\text{BN}(\mathbf{X}_\pi^1) \cdots \text{BN}(\mathbf{X}_\pi^m)]$. When the context of parameters $\Theta = (\mathbf{W}, \mathbf{\Gamma})$ and permutation $\pi \in \mathbb{S}_n$ chosen by SS are clear, we denote the collection of outputs over the dataset by $\hat{\mathbf{Y}}_\pi \triangleq \mathbf{W}\mathbf{\Gamma}\overline{\mathbf{X}}_\pi$. Also recall that the distorted SS risk $\mathcal{L}_\pi(\Theta)$ we set out to optimize is defined to be $\mathcal{L}_\pi(\Theta) = \mathcal{L}_\pi(\mathbf{W}, \mathbf{\Gamma}) = \|\mathbf{Y}_\pi - \mathbf{W}\mathbf{\Gamma}\overline{\mathbf{X}}_\pi\|_F^2$. With $\mathbf{M} \triangleq \mathbf{W}\mathbf{\Gamma}$, we also abuse notation and write

$$\begin{aligned} \mathcal{L}_\pi(\mathbf{M}) &\triangleq \|\mathbf{Y}_\pi - \mathbf{M}\overline{\mathbf{X}}_\pi\|_F^2, \\ \nabla_{\mathbf{M}} \mathcal{L}_\pi(\mathbf{M}) &\triangleq -(\mathbf{Y}_\pi - \mathbf{M}\overline{\mathbf{X}}_\pi)\overline{\mathbf{X}}_\pi^\top. \end{aligned}$$

We will use big- O notation throughout to simplify the presentation of the proofs. When we write $O(\eta_k^t)$ for some exponent $t \geq 1$, we hide constants that depend on m , $\|\overline{\mathbf{X}}_\pi\|_F$, and various absolute constants defined explicitly below. These constants have at most polynomial dependence on these parameters and absolute constants.

A.1. Proof of convergence for SS

Let us first prove Theorem 3.2.2. First, we draw the reader's attention to some standard properties in optimization theory that allow us to prove global convergence. We then sketch out the proof in Appendix A.1.2 and flesh out the details in subsequent sections.

A.1.1. OPTIMIZATION PROPERTIES

It is profitable to keep in mind the general idea behind proving global convergence of SGD for a function $\mathcal{L}(\Theta)$, which has been exploited in Ahn, Yun, and Sra (2020); Zhou and Liang (2017); Nguyen, Tran-Dinh, Phan, Nguyen, and van Dijk (2021). The following two properties of the optimization problem are critical in such approaches:

Property 1 (Smoothness). G -smoothness of \mathcal{L} , i.e., the gradients of \mathcal{L} are G -Lipschitz. In particular, it implies the following two standard properties:

- (i) $\mathcal{L}(\Theta) \leq \mathcal{L}(\Theta') + \langle \nabla_{\Theta} \mathcal{L}(\Theta'), \Theta - \Theta' \rangle + \frac{G}{2} \|\Theta' - \Theta\|^2$ for all Θ, Θ' in the domain of \mathcal{L} .
- (ii) The Hessian $\mathbf{H} = \nabla_{\Theta}^2 \mathcal{L}(\Theta)$ satisfies $\|\mathbf{H}\|_2 \leq G$ for all Θ in the domain of \mathcal{L} .

Property 2 (PŁ condition). *The loss function \mathcal{L} satisfies the α -Polyak-Łojasiewicz condition, i.e., $\|\nabla\mathcal{L}(\Theta)\|^2 \geq 2\alpha(\mathcal{L}(\Theta) - \mathcal{L}^*)$ for all Θ in the domain of \mathcal{L} .*

In our case, we can use global smoothness and strong convexity (which implies the PŁ condition) of \mathcal{L}_π with respect to $M = W\Gamma$, but these global properties do not hold with respect to our optimization variables $\Theta = (W, \Gamma)$. Importantly, unlike the analyses of Ahn, Yun, and Sra (2020); Nguyen, Tran-Dinh, Phan, Nguyen, and van Dijk (2021), we cannot directly leverage the global smoothness and strong convexity as is, because we do not directly perform gradient updates on M . Instead, we effectively use a “dynamic” PŁ condition which depends on Γ . The subtlety in the analysis is to show that such behavior can be controlled to ensure convergence in the end.

Finally, a third property — which is often exploited to prove convergence results for linear neural networks — is the notion of an (approximate) invariance property satisfied by the layers of the neural network. Indeed, in the continuous time case, i.e., when we minimize $\mathcal{L}_\pi(\Theta(t))$ with gradient flow $\dot{\Theta}(t) = -\nabla_{\Theta}\mathcal{L}_\pi(\Theta(t))$, such an invariance can be directly shown by the differential equations, see Wu, Wang, and Ma (2019) for instance. To that end, define the following quantity

$$D \triangleq I + \text{diag}(W^\top W - \Gamma^2), \quad (1)$$

which we refer to as the invariance matrix. For each iterate Θ_i^k of SS, the corresponding D_i^k can also be naturally defined. In gradient flow, $D(t)$ actually remains invariant with time $t \in [0, \infty)$. We quickly prove this property here, and later prove that an approximate version holds in the discrete and stochastic case, although the bounds are messier.

Fact A.1.1. *In the gradient flow formulation, we have $\frac{d}{dt}D(t) = 0$. Moreover, in both the gradient flow and discrete time formulation, we have*

$$\text{diag}(W^\top \nabla_W \mathcal{L}_\pi) = (\nabla_\Gamma \mathcal{L}_\pi)\Gamma. \quad (2)$$

Proof. For the proof, we write out the (full) gradients of \mathcal{L}_π with respect to W and Γ for reference:

$$\nabla_W \mathcal{L}_\pi = -(\mathbf{Y}_\pi - \hat{\mathbf{Y}}_\pi)\bar{\mathbf{X}}_\pi^\top \Gamma, \quad (3)$$

$$\nabla_\Gamma \mathcal{L}_\pi = -\text{diag}(W^\top (\mathbf{Y}_\pi - \hat{\mathbf{Y}}_\pi)\bar{\mathbf{X}}_\pi^\top). \quad (4)$$

A direct calculation shows that $\text{diag}(W^\top \nabla_W \mathcal{L}_\pi) = (\nabla_\Gamma \mathcal{L}_\pi)\Gamma$. Due to the gradient flow formulation $\dot{\Theta}(t) = -\nabla_{\Theta}\mathcal{L}_\pi(\Theta(t))$ we have $\frac{d}{dt}W(t) = -\nabla_W \mathcal{L}_\pi$ and $\frac{d}{dt}\Gamma(t) = -\nabla_\Gamma \mathcal{L}_\pi$, so it follows from Equation (2) that $\frac{d}{dt}D(t) = 0$. \square

We now formally state the smoothness and PŁ guarantees for our setup.

Lemma A.1.2 (Smoothness with respect to M). *The SS risk \mathcal{L}_π is G_π -smooth with respect to $M = W\Gamma$, where*

$$G_\pi \triangleq \|\bar{\mathbf{X}}_\pi\|_2^2.$$

Proof. We directly check the Lipschitz gradient condition. Indeed, we have

$$\begin{aligned} & \|\nabla_M \mathcal{L}_\pi(M) - \nabla_M \mathcal{L}_\pi(M')\|_2 \\ &= \left\| (\mathbf{Y}_\pi - M\bar{\mathbf{X}}_\pi)\bar{\mathbf{X}}_\pi^\top - (\mathbf{Y}_\pi - M'\bar{\mathbf{X}}_\pi)\bar{\mathbf{X}}_\pi^\top \right\|_2 \\ &= \left\| (M - M')\bar{\mathbf{X}}_\pi\bar{\mathbf{X}}_\pi^\top \right\|_2 \leq \|\bar{\mathbf{X}}_\pi\|_2^2 \|M - M'\|_2, \end{aligned}$$

Note that the same inequality holds (with the same value of G_π) if we instead used the Frobenius norm, due to the fact that $\|AB\|_F \leq \|B\|_2 \|A\|_F$ in the last line. \square

Lemma A.1.3 (Strong convexity with respect to M). *Under Assumption 1(a), SS risk \mathcal{L}_π is α_π -strongly convex with respect to $M = W\Gamma$, where*

$$\alpha_\pi \triangleq \sigma_{\min}(\bar{\mathbf{X}}_\pi\bar{\mathbf{X}}_\pi^\top).$$

Hence, \mathcal{L}_π is also α_π -PŁ with respect to M .

Proof. Take the Hessian of $\mathcal{L}_\pi(M)$ with respect to the vectorized version $\text{vec}(M)$ of M to obtain $\nabla_{\text{vec}(M)}^2 \mathcal{L}_\pi(M) = \bar{\mathbf{X}}_\pi\bar{\mathbf{X}}_\pi^\top \otimes I_p$, where \otimes denotes the Kronecker product. Then evidently $\nabla_{\text{vec}(M)}^2 \mathcal{L}_\pi(M) \succeq \sigma_{\min}(\bar{\mathbf{X}}_\pi\bar{\mathbf{X}}_\pi^\top)I_p$. Owing to Assumption 1(a), this proves the claim. \square

A.1.2. PROOF SKETCH OF CONVERGENCE

Proof sketch of Theorem 3.2.2. The high level idea is this: we want to prove that $\mathcal{L}_\pi(\mathbf{M}_0^k) \rightarrow \mathcal{L}_\pi^*$ as $k \rightarrow \infty$. However, we will instead show the much stronger statement that $\mathcal{L}_\pi(\mathbf{M}_i^k) \rightarrow \mathcal{L}_\pi^*$ for all $i \in [m]$. Our high level approach is heavily inspired by the proof strategies in Wu et al. (2019); Ahn et al. (2020). Indeed, many of the technical lemmas in Appendix A.1.4 are analogous to ones proved in Wu et al. (2019), and the motivation for unrolling shuffling mini-batch updates to an epoch update with additional noise comes from Ahn et al. (2020).

As a necessary ingredient of the proof, we will demonstrate that for sufficiently small chosen η_k , we have an update equation that roughly looks like (modulo constants and noise terms)

$$\mathcal{L}_\pi(\mathbf{M}_i^{k+1}) - \mathcal{L}_\pi^* \lesssim (1 - \eta_k)(\mathcal{L}_\pi(\mathbf{M}_i^k) - \mathcal{L}_\pi^*) + O(\eta_k^2) \quad \text{for all } 0 \leq i \leq m - 1. \quad (5)$$

Remark A.1.4. Note that it is not necessarily the case that

$$\mathcal{L}_\pi(\mathbf{M}_{i+1}^k) - \mathcal{L}_\pi^* \lesssim (1 - \eta_k)(\mathcal{L}_\pi(\mathbf{M}_i^k) - \mathcal{L}_\pi^*) + O(\eta_k^2)$$

That is, the SS excess risk \mathcal{L}_π does not necessarily “decrease” from one iterate to the next; however, we can instead guarantee that the per-epoch progress bound (Equation (5)) holds for any fixed iteration index $i \in [m]$ after every epoch.

We impose an ordering relation on pairs (a, b) in the natural way: we say $(a, b) \leq (i, k)$ if $k = b$ and $a \leq i$, or if $b < k$. This is just tracking whether the iteration index (a, b) (the a th iterate of the b th epoch) is seen before the iterate (i, k) . To complete the induction on an iterate $(i, k + 1)$ we need three inductive hypotheses $L[a, b]$, $D[a, b]$, and $R[a, b]$ to hold for all $(a, b) < (i, k + 1)$. We define them formally below.

Hypothesis 1 (Loss stays bounded by an absolute constant). For all a, b satisfying $0 \leq a \leq m - 1$ and $b \geq 1$, the inductive property $L[a, b]$ states $\mathcal{L}_\pi(\Theta_a^b) \leq C_L$, for some appropriately chosen absolute constant C_L .

In particular, we can set $C_L \triangleq \max \{ \mathcal{L}_\pi(\Theta_t^1) : 0 \leq t \leq m - 1 \}$. Since we only look at the loss values for the first epoch, C_L is indeed an absolute constant depending on π .

Hypothesis 2 (Loss satisfies one-epoch inequality). For all a, b satisfying $0 \leq a \leq m - 1$ and $b > 1$, the inductive property $R[a, b]$ states that

$$\mathcal{L}_\pi(\mathbf{M}_a^b) - \mathcal{L}_\pi^* \leq \left(1 - \frac{\alpha_\pi \eta_k}{2}\right) (\mathcal{L}_\pi(\mathbf{M}_a^{b-1}) - \mathcal{L}_\pi^*) + O(\eta_k^2),$$

where the constant hidden in the $O(\eta_k^2)$ does not depend on k .

Hypothesis 3 (Approximate invariances hold). For all a, b satisfying $0 \leq a \leq m - 1$ and $b \geq 1$, the inductive property $D[a, b]$ states that

$$\|\mathbf{D}_a^b\|_2 \leq \begin{cases} C_D \sum_{t=1}^{b-1} \eta_t^2 \leq \frac{1}{2} & \text{if } a = 0, \\ C_D \sum_{t=1}^b \eta_t^2 \leq \frac{1}{2} & \text{otherwise,} \end{cases}$$

where C_D is an appropriately chosen absolute constant which does not depend on a or b .

Since the first iterate of the k th epoch Θ_0^k is the same as the last iterate of the $(k - 1)$ th epoch Θ_m^k , the same convention applies to inductive hypotheses; for example, by $L[m, k - 1]$ we mean $L[0, k]$.

In particular, the inductive hypotheses imply the following claims.

- (i) By Corollary A.1.8, $L[a, b]$ implies that $\|\mathbf{M}_a^b\|_2 \leq \frac{C_L^{1/2} + \|\mathbf{Y}_\pi\|_F}{\sigma_{\min}(\mathbf{X}_\pi)} \triangleq \xi$.
- (ii) Also by Corollary A.1.8, $D[a, b]$ and $L[a, b]$ together imply that we have $\|\mathbf{W}_a^b\|_2^2 \leq d^2(\frac{1}{2} + \xi)$ and $\|\mathbf{\Gamma}_a^b\|_2^2 \leq \frac{3}{2} + d^2(\frac{1}{2} + \xi)$. For the sake of notational convenience we will write $C_w \triangleq \sqrt{\frac{3}{2} + d^2(\frac{1}{2} + \xi)}$, so that $\max \{ \|\mathbf{W}_a^b\|_2, \|\mathbf{\Gamma}_a^b\|_2 \} \leq C_w$.
- (iii) By Corollary A.1.13, $D[a, b]$ implies that $\sigma_{\min}(\mathbf{\Gamma}_a^b)^2 \geq 1/2$.

(iv) By Proposition A.1.23, if $R[a, b]$ holds for all (a, b) , then for appropriately chosen η_k , the risk $\mathcal{L}_\pi(\mathbf{M}_a^b)$ converges to \mathcal{L}_π^* at a sublinear rate.

We will explain at a high level how these statements together allow us to conclude that $L[i, k+1]$, $D[i, k+1]$, and $R[i, k+1]$ hold. The idea, as in Ahn, Yun, and Sra (2020), is to accumulate the gradient updates in each epoch and isolate the signal and noise components of each gradient update. For clarity of exposition, we assume for now that $i = 0$. Here are a couple subtleties which we spell out explicitly, including how to generalize to $i > 0$.

- We are not directly performing gradient updates on \mathbf{M} ; we instead perform gradient updates on \mathbf{W} and $\mathbf{\Gamma}$. Nevertheless, the *effective* gradient signal for \mathbf{M} can still be extracted, and we term the remaining noise the *mismatched gradient noise*. For every iterate (j, k) , this will formally be denoted by \mathbf{q}_j^k .
- We are not taking a full batch gradient step from \mathbf{M}_0^k to \mathbf{M}_0^{k+1} . Rather, we are taking mini-batch updates which induce path dependency. Nevertheless, as previous works have shown, even at iterate (j, k) , we can still extract the full-batch gradient signal evaluated at \mathbf{M}_0^k , and we term the remaining noise the *path dependent noise*. For every iterate (j, k) , this will formally be denoted by \mathbf{e}_j^k .
- If $i > 0$, then the stepsize changes from η_k to η_{k+1} in the middle of our pass through the entire dataset. Nevertheless, it's not hard to see that this noise should be relatively small, of order $\eta_{k+1} - \eta_k$ — which is $O(\eta_k^2)$, as $\eta_k = \Omega(1/k)$. We will call this the *stepsize noise*, the accumulation of which for an epoch update starting from iterate (i, k) to $(i, k+1)$ will be denoted by $\mathbf{s}_{(i, k+1)}^{(i, k)}$.

We can accumulate these noise terms across the update across epoch k to form a composite noise term \mathbf{r}^k . The *full-batch update signal* for \mathbf{M} starting from \mathbf{M}_0^k will be denoted by $\tilde{\mathbf{g}}^k$. We emphasize that $\tilde{\mathbf{g}}^k \neq \nabla_{\mathbf{M}} \mathcal{L}_\pi(\mathbf{M}_0^k)$ because we only perform direct gradient updates on the component layers \mathbf{W} and $\mathbf{\Gamma}$. Then as we will show in Appendix A.1.3, we can write

$$\mathbf{M}_0^{k+1} = \mathbf{M}_0^k - \eta_k \tilde{\mathbf{g}}^k + \eta_k^2 \mathbf{r}^k. \quad (6)$$

Next, as seen in Lemma A.1.2, \mathcal{L}_π is globally G_π -smooth with respect to \mathbf{M} for some absolute constant G_π which depends on π . Thus, using the smoothness inequality as in Property 1, we obtain

$$\mathcal{L}_\pi(\mathbf{M}_0^{k+1}) - \mathcal{L}_\pi(\mathbf{M}_0^k) \leq \langle \nabla_{\mathbf{M}} \mathcal{L}_\pi(\mathbf{M}_0^k), \mathbf{M}_0^{k+1} - \mathbf{M}_0^k \rangle_F + \frac{G_\pi}{2} \|\mathbf{M}_0^{k+1} - \mathbf{M}_0^k\|_F^2.$$

The main idea is that we have the following inequality (proved in Lemma A.1.14) that shows that even though $\tilde{\mathbf{g}}^k \neq \nabla_{\mathbf{M}} \mathcal{L}_\pi(\mathbf{M}_0^k)$, it is nonetheless correlated to the “correct” gradient update $\nabla_{\mathbf{M}} \mathcal{L}_\pi(\mathbf{M}_0^k)$:

$$\langle \nabla_{\mathbf{M}} \mathcal{L}_\pi(\mathbf{M}_0^k), \tilde{\mathbf{g}}^k \rangle_F \geq \sigma_{\min}(\mathbf{\Gamma}_0^k)^2 \|\nabla_{\mathbf{M}} \mathcal{L}_\pi(\mathbf{M}_0^k)\|_F^2 \geq \frac{1}{2} \|\nabla_{\mathbf{M}} \mathcal{L}_\pi(\mathbf{M}_0^k)\|_F^2,$$

due to the inductive hypothesis $D[0, k]$.

For the stated stepsizes η_k , one can then plug in the gradient update Equation (6) and massage the inequalities a bit to obtain that

$$\mathcal{L}_\pi(\mathbf{M}_0^{k+1}) - \mathcal{L}_\pi(\mathbf{M}_0^k) \leq -\frac{\eta_k}{4} \|\nabla_{\mathbf{M}} \mathcal{L}_\pi(\mathbf{M}_0^k)\|_F^2 + O(\eta_k^2), \quad (7)$$

where the constant hidden by the big- O notation is $\text{poly}(m, C_w, C_L, \|\bar{\mathbf{X}}_\pi\|_F)$.

We now use α_π -strong convexity of \mathcal{L}_π with respect to \mathbf{M} (and hence α_π -PEL) shown in Lemma A.1.3 to obtain

$$\mathcal{L}_\pi(\mathbf{M}_0^{k+1}) - \mathcal{L}_\pi^* \leq \left(1 - \frac{\alpha_\pi \eta_k}{2}\right) (\mathcal{L}_\pi(\mathbf{M}_0^k) - \mathcal{L}_\pi^*) + O(\eta_k^2). \quad (8)$$

Note that this is precisely the statement of $R[0, k+1]$.

Provided that we can appropriately bound the noise terms \mathbf{r}^k to get the asserted $O(\eta_k^2)$ term above, this will imply $R[0, k+1]$. For sufficiently small stepsizes η_k , we can also use Equation (8) to prove $L[0, k+1]$.

On the other hand, to prove $D[0, k + 1]$, we can directly bound the update $\|D_0^{k+1} - D_{m-1}^k\|_2 \leq O(\eta_k^2)$ and combine this with the inductive hypothesis $D[m - 1, k]$ using the triangle inequality. If the stepsize $\eta_k = O(1/k^\beta)$ for $1/2 < \beta < 1$, then $\sum_{k \geq 1} \eta_k^2 < \infty$, so the absolute constant C_D can be picked such that $\|D_0^{k+1}\|_2 \leq \frac{1}{2}$.

Hence, $R[0, k]$, as stated in Equation (8), holds for all k by induction. We can thus unroll the inequality and conclude that $\mathcal{L}_\pi(M_0^k)$ converges to \mathcal{L}_π^* under the stated stepsize assumptions, as desired. \square

We now outline the structure of the proceeding sections, which fill in the details of the above proof sketch. In Appendix A.1.3, we explicitly write out the accumulation of gradient updates across an entire epoch, decomposing into signal and noise components. In Appendix A.1.4, we prove some technical lemmas controlling the singular values and norms of various weight matrices and gradients via the approximate invariance matrix D and the inductive hypotheses. In Appendix A.1.5 we leverage the norm bounds developed in Appendix A.1.4 to demonstrate that the accumulated noise terms defined in Appendix A.1.3 are negligible. Using these results, we are able to establish the $R[i, k + 1]$ and $L[i, k + 1]$ in Appendix A.1.6. We then turn to bounding the approximate invariances to establish $D[i, k + 1]$ in Appendix A.1.7. The stray details of the induction are spelled out in Appendix A.1.8.

A.1.3. REWRITING SS EPOCH GRADIENT UPDATES

To show that $L[0, k + 1]$ holds, we need to accumulate gradients from M_0^k to M_0^{k+1} .

First, we look at a single iterate update. For every $j < m$ we have

$$M_{j+1}^k = (W_j^k - \eta_k \nabla_{\mathbf{W}} \mathcal{L}(X_\pi^{j+1}; \Theta_j^k)) (\Gamma_j^k - \eta_k \nabla_{\Gamma} \mathcal{L}(X_\pi^{j+1}; \Theta_j^k)) \quad (9)$$

$$= M_j^k - \eta_k g_j^k + \eta_k^2 q_j^k, \quad (10)$$

where we have defined

$$g_j^k \triangleq \nabla_{\mathbf{W}} \mathcal{L}(X_\pi^{j+1}; \Theta_j^k) \Gamma_j^k + W_j^k \nabla_{\Gamma} \mathcal{L}(X_\pi^{j+1}; \Theta_j^k), \quad (11)$$

which is the gradient of the $(j + 1)$ th batch of \bar{X}_π evaluated on the j th iterate on epoch k , and

$$q_j^k \triangleq \nabla_{\mathbf{W}} \mathcal{L}(X_\pi^{j+1}; \Theta_j^k) \nabla_{\Gamma} \mathcal{L}(X_\pi^{j+1}; \Theta_j^k), \quad (12)$$

which is the mismatched gradient noise term associated with the fact that we performed gradient updates on W and Γ rather than M directly.

The key observation here is that

$$g_j^k = \nabla_M \mathcal{L}(X_\pi^{j+1}; M_j^k) (\Gamma_j^k)^2 + W_j^k \text{diag}((W_j^k)^\top \nabla_M \mathcal{L}(X_\pi^{j+1}; M_j^k)).$$

In other words, g_j^k is correlated to the “true” mini-batch gradient $\nabla_M \mathcal{L}(X_\pi^{j+1}; M_j^k)$ with respect to M through the “interaction terms” Γ_j^k and W_j^k .

We show in Lemma A.1.16 that we can control the size of the noise terms q_j^k which arise from the fact that we are not truly taking gradient updates with respect to M . More specifically, Lemma A.1.16 implies that $\|q_j^k\|_F = O(1)$.

Next, we actually accumulate gradients. The main obstacle we have to deal with is that the mini-batch updates prevent the gradient accumulation from being exactly equal to the full-batch update starting at M_0^k . Inspired by the approach in Ahn et al. (2020, Theorem 1), we separate out the gradient update g_j^k into a signal term \tilde{g}_j^k and noise term e_j^k . Specifically, we write

$$M_{j+1}^k = M_j^k - \eta_k \tilde{g}_j^k + \eta_k^2 e_j^k + \eta_k^2 q_j^k, \quad (13)$$

where

$$\tilde{g}_j^k \triangleq \nabla_{\mathbf{W}} \mathcal{L}(X_\pi^{j+1}; \Theta_0^k) \Gamma_0^k + W_0^k \nabla_{\Gamma} \mathcal{L}(X_\pi^{j+1}; \Theta_0^k), \quad (14)$$

is the signal of the gradient update of the $(j + 1)$ th batch evaluated with parameter values Θ_0^k (instead of Θ_j^k) and

$$e_j^k \triangleq \frac{\tilde{g}_j^k - g_j^k}{\eta_k}. \quad (15)$$

In particular, in Lemma A.1.18 below we show that $\|e_j^k\|_F = O(1)$, so that indeed the noise term is negligible with respect to the true gradient signal.

Taking this as given for now, when we accumulate the gradient updates across epoch k , we see that we can define

$$\tilde{\mathbf{g}}^k \triangleq \sum_{j=0}^{m-1} \tilde{\mathbf{g}}_j^k = \nabla_{\mathbf{W}} \mathcal{L}_\pi(\Theta_0^k) \Gamma_0^k + \mathbf{W}_0^k \nabla_{\Gamma} \mathcal{L}_\pi(\Theta_0^k), \quad (16)$$

so that the accumulation reads

$$\mathbf{M}_0^{k+1} = \mathbf{M}_0^k - \eta_k \tilde{\mathbf{g}}^k + \eta_k^2 \sum_{j=0}^{m-1} (\mathbf{e}_j^k + \mathbf{q}_j^k) \quad (17)$$

$$= \mathbf{M}_0^k - \eta_k \tilde{\mathbf{g}}^k + \eta_k^2 \mathbf{r}^k, \quad (18)$$

where we have additionally defined the composite noise term:

$$\mathbf{r}^k \triangleq \sum_{j=0}^{m-1} (\mathbf{e}_j^k + \mathbf{q}_j^k), \quad (19)$$

Note that if we instead start from $i > 0$, then the composite noise term \mathbf{r}^k will have an additional noise term $\mathbf{s}_{(i,k+1)}^{(i,k)}$, which we will address in Appendix A.1.5. In particular, we show there that the norm of $\mathbf{s}_{(i,k+1)}^{(i,k)}$ is $O(1)$. Combining this with Lemmas A.1.16 and A.1.18, we can conclude that $\|\mathbf{r}^k\|_F = O(1)$.

A.1.4. NORM AND SINGULAR VALUE BOUNDS BASED ON APPROXIMATE INVARIANCES

In this section, we prove several helper lemmas which help us bound noise terms in Appendix A.1.5 and the approximate invariances in Appendix A.1.7.

Upper bounds on the norms of \mathbf{W} and Γ . Much of Wu, Wang, and Ma (2019) is dedicated towards showing that the approximate invariances control the weight norms. The trouble with directly extending their strategy lies in the fact that in our setting the invariance \mathbf{D} is diagonal, which complicates the process of bounding various matrix norms. We first state the following technical lemma which involves the operator norm of Hadamard products.

Lemma A.1.5 (3.1f in Johnson (1990)). *Let $\mathbf{A}, \mathbf{B} \in \mathbb{R}^{d \times d}$ be matrices such that \mathbf{A} is positive definite. Then $\|\mathbf{A} \odot \mathbf{B}\|_2 \leq \|\mathbf{A}\|_2 \|\mathbf{B}\|_2$*

We leverage Lemma A.1.5 to prove the following useful helper lemma that relates bounds on $\|\mathbf{I} \odot \mathbf{W}^\top \mathbf{W}\|_2$ to $\|\mathbf{W}\|_2$.

Lemma A.1.6. *Suppose $\|\mathbf{I} \odot \mathbf{W}^\top \mathbf{W}\|_2 \leq \beta$, where $\mathbf{W} \in \mathbb{R}^{p \times d}$. Then $\|\mathbf{W}\|_2 \leq \sqrt{d\beta}$. Conversely, if $\|\mathbf{W}\|_2 \leq \beta$, then $\|\mathbf{I} \odot \mathbf{W}^\top \mathbf{W}\|_2 \leq \beta^2$.*

Proof. Note that $\mathbf{I} \odot \mathbf{W}^\top \mathbf{W}$ is a diagonal matrix with diagonal entries $\mathbf{W}_{:,1}^\top \mathbf{W}_{:,1}, \mathbf{W}_{:,2}^\top \mathbf{W}_{:,2}, \dots, \mathbf{W}_{:,d}^\top \mathbf{W}_{:,d}$, where $\mathbf{W}_{:,i}$ denotes the i th column of \mathbf{W} . Hence $\text{Tr}(\mathbf{I} \odot \mathbf{W}^\top \mathbf{W}) = \|\mathbf{W}\|_F^2$. Hence $\|\mathbf{W}\|_F^2 \leq d\beta$ (or tighter by replacing d with the rank of \mathbf{W}), from which it follows that $\|\mathbf{W}\|_2 \leq \sqrt{d\beta}$. For the other direction, we set $\mathbf{A} = \mathbf{I}$ and $\mathbf{B} = \mathbf{W}^\top \mathbf{W}$ in Lemma A.1.5, so $\|\mathbf{I} \odot \mathbf{W}^\top \mathbf{W}\|_2 \leq \|\mathbf{W}\|_2^2 \leq \beta^2$, as desired. \square

With Lemma A.1.6 in hand, we prove the following technical lemma which gives a uniform bound on the norms of Γ and \mathbf{W} based on $\xi = \|\mathbf{W}\Gamma\|_2$.

Lemma A.1.7. *If $\|\mathbf{D}\|_2 \leq \epsilon < 1$ and $\|\mathbf{W}\Gamma\|_2 \leq \xi$, we have*

$$\|\mathbf{W}\|_2 \leq d\sqrt{1 - \epsilon + \xi},$$

and

$$\|\Gamma^2\|_2 \leq 1 + \epsilon + d^2(1 - \epsilon + \xi).$$

Proof. We have from $\|\mathbf{W}\boldsymbol{\Gamma}\|_2 \leq \xi$ that

$$\|\mathbf{W}\boldsymbol{\Gamma}^2\mathbf{W}^\top\|_2 \leq \xi^2.$$

Next, our hypothesis that $\|\mathbf{D}\|_2 = \|\mathbf{I} + \text{diag}(\mathbf{W}^\top\mathbf{W}) - \boldsymbol{\Gamma}^2\|_2 \leq \epsilon$ implies that

$$\mathbf{W}\boldsymbol{\Gamma}^2\mathbf{W}^\top \succeq \mathbf{W}((1 - \epsilon)\mathbf{I} + \text{diag}(\mathbf{W}^\top\mathbf{W}))\mathbf{W}^\top.$$

Taking norms of both sides and applying the reverse triangle inequality, we obtain that

$$\xi^2 \geq \|\mathbf{W} \text{diag}(\mathbf{W}^\top\mathbf{W})\mathbf{W}^\top\|_2 - (1 - \epsilon)\|\mathbf{W}\|_2^2.$$

We now lower bound $\|\mathbf{W} \text{diag}(\mathbf{W}^\top\mathbf{W})\mathbf{W}^\top\|_2$. In particular, we expand out the matrix product. Note here that $\text{diag}(\mathbf{W}^\top\mathbf{W})_{i,i} = \|\mathbf{W}_{:,i}\|_2^2$. Thus we can write $\mathbf{W} \text{diag}(\mathbf{W}^\top\mathbf{W})\mathbf{W}^\top$ as

$$\begin{bmatrix} \mathbf{W}_{:,1} & \mathbf{W}_{:,2} & \cdots & \mathbf{W}_{:,d} \end{bmatrix} \begin{bmatrix} \|\mathbf{W}_{:,1}\|_2^2 & & & \\ & \|\mathbf{W}_{:,2}\|_2^2 & & \\ & & \ddots & \\ & & & \|\mathbf{W}_{:,d}\|_2^2 \end{bmatrix} \begin{bmatrix} \mathbf{W}_{:,1}^\top \\ \mathbf{W}_{:,2}^\top \\ \vdots \\ \mathbf{W}_{:,d}^\top \end{bmatrix},$$

from which we observe that the i th diagonal entry of $\mathbf{W} \text{diag}(\mathbf{W}^\top\mathbf{W})\mathbf{W}^\top$ is

$$(\mathbf{W} \text{diag}(\mathbf{W}^\top\mathbf{W})\mathbf{W}^\top)_{i,i} = \sum_{j=1}^d \|\mathbf{W}_{:,j}\|_2^2 W_{i,j}^2.$$

It follows that $\text{Tr}(\mathbf{W} \text{diag}(\mathbf{W}^\top\mathbf{W})\mathbf{W}^\top) = \sum_{j=1}^d \|\mathbf{W}_{:,j}\|_2^4$. Note that $\|\mathbf{A}\|_2 \geq \max_{i,j} |A_{i,j}|$ (the RHS is also known as the *max norm*). For our case we set $\mathbf{A} = \mathbf{W} \text{diag}(\mathbf{W}^\top\mathbf{W})\mathbf{W}^\top$ and note that the diagonal is nonnegative. So in fact in our case we obtain

$$\|\mathbf{W} \text{diag}(\mathbf{W}^\top\mathbf{W})\mathbf{W}^\top\|_2 \geq \frac{1}{d} \sum_{j=1}^d \|\mathbf{W}_{:,j}\|_2^4.$$

Now notice that $\sum_j \|\mathbf{W}_{:,j}\|_2^4 = \sum_j (\sum_i W_{i,j}^2)^2$. Applying Cauchy-Schwarz to the outer sum we find that

$$\sum_j \|\mathbf{W}_{:,j}\|_2^4 \geq \frac{(\sum_j \sum_i W_{i,j}^2)^2}{d},$$

but the RHS is equal to $\|\mathbf{W}\|_F^4$. Since $\|\mathbf{W}\|_F \geq \|\mathbf{W}\|_2$, we conclude that

$$\|\mathbf{W} \text{diag}(\mathbf{W}^\top\mathbf{W})\mathbf{W}^\top\|_2 \geq \frac{\|\mathbf{W}\|_2^4}{d^2}.$$

In summary, we have

$$\frac{\|\mathbf{W}\|_2^4}{d^2} - (1 - \epsilon)\|\mathbf{W}\|_2^2 - \xi^2 \leq 0.$$

Applying the quadratic formula, we find that

$$\|\mathbf{W}\|_2 \leq d\sqrt{1 - \epsilon + \xi}.$$

For the bound on $\|\boldsymbol{\Gamma}\|_2$, we start from the definition of \mathbf{D} and apply the reverse triangle inequality to obtain

$$|1 + \|\text{diag}(\mathbf{W}^\top\mathbf{W})\|_2 - \|\boldsymbol{\Gamma}^2\|_2| \leq \epsilon,$$

so we obtain

$$\|\boldsymbol{\Gamma}^2\|_2 \leq 1 + \epsilon + \|\mathbf{W}\|_2^2,$$

where we used $\|\mathbf{I} \odot \mathbf{W}^\top\mathbf{W}\|_2 \leq \|\mathbf{W}\|_2^2$ from Lemma A.1.6. From this, the conclusion directly follows. \square

Under the inductive hypotheses, Lemma A.1.7 implies that we can uniformly bound $\max \left\{ \|\mathbf{W}_j^k\|_2, \|\mathbf{\Gamma}_j^k\|_2 \right\}$. This is spelled out in the following corollary.

Corollary A.1.8 (Norms stay bounded). *Suppose that $L[j, k]$ and $D[j, k]$ hold. Define*

$$C_w \triangleq \sqrt{\frac{3}{2} + d^2 \left(\frac{1}{2} + \xi \right)},$$

with

$$\xi \triangleq \frac{C_L^{1/2} + \|\mathbf{Y}_\pi\|_F}{\sigma_{\min}(\overline{\mathbf{X}}_\pi^\top)}.$$

Here C_L was defined in Hypothesis 1. Then

$$\|\mathbf{M}_j^k\| \leq \xi,$$

and

$$\max \left\{ \|\mathbf{W}_j^k\|_2, \|\mathbf{\Gamma}_j^k\|_2 \right\} \leq C_w.$$

Proof. We have by triangle inequality that

$$\|\mathbf{M}_j^k \overline{\mathbf{X}}_\pi\|_2 \leq \|\mathbf{M}_j^k \overline{\mathbf{X}}_\pi\|_F \leq \|\mathbf{Y}_\pi - \mathbf{M}_j^k \overline{\mathbf{X}}_\pi\|_F + \|\mathbf{Y}_\pi\|_F \leq \mathcal{L}_\pi (\mathbf{M}_j^k)^{1/2} + \|\mathbf{Y}_\pi\|_F.$$

Since $L[j, k]$ holds, we have $\|\mathbf{Y}_\pi - \mathbf{M}_j^k \overline{\mathbf{X}}_\pi\|_F^2 \leq C_L$. Furthermore, as $n \geq d$, we know that $\|\mathbf{M}_j^k \overline{\mathbf{X}}_\pi\|_2 \geq \sigma_{\min}(\overline{\mathbf{X}}_\pi^\top) \|\mathbf{M}_j^k\|_2$ and by Item Assumption 1(a) we have $\sigma_{\min}(\overline{\mathbf{X}}_\pi^\top) > 0$. Hence we obtain

$$\|\mathbf{M}_j^k\|_2 \leq \frac{C_L^{1/2} + \|\mathbf{Y}_\pi\|_F}{\sigma_{\min}(\overline{\mathbf{X}}_\pi^\top)} = \xi.$$

It follows that ξ works as a bound on $\|\mathbf{M}_j^k\|_2$ for the application of Lemma A.1.7. Since $D[j, k]$ holds by assumption, this means that the hypothesis on \mathbf{D}_j^k is satisfied with $\epsilon = 1/2$. In summary, all the hypotheses of Lemma A.1.7 are satisfied. We can thus conclude that

$$\max \left\{ \|\mathbf{W}_j^k\|_2, \|\mathbf{\Gamma}_j^k\|_2 \right\} \leq C_w,$$

as desired. \square

The importance of these upper bounds on weight norms is that they allow us to upper bound the norms of gradients of \mathcal{L} with respect to various parameters.

Upper bounding the norms of gradients. The following lemma gives an upper bound on the norms of various gradients.

Lemma A.1.9. *For any $a \in [m]$ and $\Theta = (\mathbf{W}, \mathbf{\Gamma})$ we have*

$$\begin{aligned} \|\nabla_{\mathbf{W}} \mathcal{L}(\mathbf{X}_\pi^a; \Theta)\|_F^2 &\leq \|\mathbf{\Gamma}\|_2^2 \|\text{BN}(\mathbf{X}_\pi^a)\|_2^2 \mathcal{L}(\mathbf{X}_\pi^a; \Theta) \\ \|\nabla_{\mathbf{\Gamma}} \mathcal{L}(\mathbf{X}_\pi^a; \Theta)\|_F^2 &\leq \|\mathbf{W}\|_2^2 \|\text{BN}(\mathbf{X}_\pi^a)\|_2^2 \mathcal{L}(\mathbf{X}_\pi^a; \Theta) \\ \|\nabla_{\mathbf{M}} \mathcal{L}(\mathbf{X}_\pi^a; \mathbf{M})\|_F^2 &\leq \|\text{BN}(\mathbf{X}_\pi^a)\|_2^2 \mathcal{L}(\mathbf{X}_\pi^a; \mathbf{M}) \end{aligned}$$

Proof. First, we have by definition

$$\mathcal{L}(\mathbf{X}_\pi^a; \Theta) = \|\mathbf{W}\mathbf{\Gamma}\text{BN}(\mathbf{X}_\pi^a) - \mathbf{Y}_\pi^a\|_F^2.$$

Hence, the mini-batch gradients can be computed explicitly as

$$\nabla_{\mathbf{M}} \mathcal{L}(\mathbf{X}_\pi^a; \mathbf{M}) = -(\mathbf{Y}_\pi^a - \mathbf{M}\text{BN}(\mathbf{X}_\pi^a))\text{BN}(\mathbf{X}_\pi^a)^\top, \quad (20)$$

$$\nabla_{\mathbf{W}} \mathcal{L}(\mathbf{X}_\pi^a; \Theta) = \nabla_{\mathbf{M}} \mathcal{L}(\mathbf{X}_\pi^a; \mathbf{M})\mathbf{\Gamma}, \quad (21)$$

$$\nabla_{\mathbf{\Gamma}} \mathcal{L}(\mathbf{X}_\pi^a; \Theta) = \text{diag}(\mathbf{W}^\top \nabla_{\mathbf{M}} \mathcal{L}(\mathbf{X}_\pi^a; \mathbf{M})). \quad (22)$$

Since $\mathcal{L}(\mathbf{X}_\pi^a; \mathbf{M}) = \|\mathbf{Y}_\pi^a - \text{MBN}(\mathbf{X}_\pi^a)\|_F^2$ and $\|\mathbf{AB}\|_F \leq \|\mathbf{A}\|_2 \|\mathbf{B}\|_F$, Equation (20) gives

$$\|\nabla_{\mathbf{M}} \mathcal{L}(\mathbf{X}_\pi^a; \mathbf{M})\|_F^2 \leq \|\text{BN}(\mathbf{X}_\pi^a)\|_2^2 \mathcal{L}(\mathbf{X}_\pi^a; \mathbf{M}).$$

It thus follows from Equation (21) that

$$\|\nabla_{\mathbf{W}} \mathcal{L}(\mathbf{X}_\pi^a; \Theta)\|_F^2 \leq \|\Gamma\|_2^2 \|\nabla_{\mathbf{M}} \mathcal{L}(\mathbf{X}_\pi^a; \Theta)\|_F^2 \leq \|\Gamma\|_2^2 \|\text{BN}(\mathbf{X}_\pi^a)\|_2^2 \mathcal{L}(\mathbf{X}_\pi^a; \Theta).$$

Similarly, inspecting Equation (22), since $\|\text{diag}(\mathbf{A})\|_F^2 \leq \|\mathbf{A}\|_F^2$, we have

$$\|\nabla_{\Gamma} \mathcal{L}(\mathbf{X}_\pi^a; \Theta)\|_F^2 \leq \|\mathbf{W}^\top \nabla_{\mathbf{M}} \mathcal{L}(\mathbf{X}_\pi^a; \Theta)\|_F^2 \leq \|\mathbf{W}\|_2^2 \|\text{BN}(\mathbf{X}_\pi^a)\|_2^2 \mathcal{L}(\mathbf{X}_\pi^a; \Theta).$$

□

As a consequence of Corollary A.1.8, under the inductive hypotheses we can also bound the gradient norms by absolute constants.

Corollary A.1.10. *Assume $D[j, k]$ and $L[j, k]$ hold. Then, for any $a \in [m]$, we have*

$$\begin{aligned} \|\nabla_{\mathbf{M}} \mathcal{L}(\mathbf{X}_\pi^a; \mathbf{M}_j^k)\|_F^2 &\leq C_L \|\text{BN}(\mathbf{X}_\pi^a)\|_2^2, \\ \|\nabla_{\mathbf{W}} \mathcal{L}(\mathbf{X}_\pi^a; \Theta_j^k)\|_F^2 &\leq C_w^2 C_L \|\text{BN}(\mathbf{X}_\pi^a)\|_2^2, \\ \|\nabla_{\Gamma} \mathcal{L}(\mathbf{X}_\pi^a; \Theta_j^k)\|_F^2 &\leq C_w^2 C_L \|\text{BN}(\mathbf{X}_\pi^a)\|_2^2, \end{aligned}$$

where C_w was previously defined in Corollary A.1.8.

We now turn from upper bounds to lower bounds. The crux here is to start with bounding the minimum singular value of Γ away from zero. This in turns allows us to lower bound the correlation between $\tilde{\mathbf{g}}^k$ and $\nabla_{\mathbf{M}} \mathcal{L}_\pi(\mathbf{M}_0^k)$ away from zero. As we will see, we can also show similar correlation lower bounds for the cases $i > 0$.

Bounding the minimum singular value of Γ^2 . In order to bound $\sigma_{\min}(\Gamma_i^k)$ away from zero, we need to show that the approximate invariances prevent Γ from vanishing on any coordinate. To do so, we appeal to an alternate formulation of the Courant-Fisher theorem for singular values, which we restate below for completeness.

Theorem A.1.11 (Courant-Fisher). *Let $\mathbf{A}, \mathbf{B} \in \mathbb{R}^{m \times n}$. Then $|\sigma_k(\mathbf{A}) - \sigma_k(\mathbf{B})| \leq \|\mathbf{A} - \mathbf{B}\|_2$ for $k \in [\min\{m, n\}]$.*

With this in mind, we formally prove that the minimum singular value of Γ^2 is bounded away from zero.

Lemma A.1.12. *Suppose that $\|\mathbf{D}\|_2 = \|\mathbf{I} + \text{diag}(\mathbf{W}^\top \mathbf{W} - \Gamma^2)\|_2 \leq \epsilon$. Then we have*

$$\sigma_{\min}(\Gamma^2) \geq 1 - \epsilon.$$

Proof. Setting $\mathbf{A} \triangleq \mathbf{I} + \text{diag}(\mathbf{W}^\top \mathbf{W})$ and $\mathbf{B} \triangleq \Gamma^2$ in Courant-Fisher yields

$$|\sigma_d(\mathbf{I} + \text{diag}(\mathbf{W}^\top \mathbf{W})) - \sigma_d(\Gamma^2)| \leq \|\mathbf{I} + \text{diag}(\mathbf{W}^\top \mathbf{W}) - \Gamma^2\|_2.$$

Since the RHS is just \mathbf{D} , we obtain that

$$\sigma_{\min}(\Gamma^2) \geq 1 + \sigma_{\min}(\text{diag}(\mathbf{W}^\top \mathbf{W})) - \|\mathbf{D}\|_2.$$

The conclusion easily follows. □

Under the inductive hypothesis $D[i, k]$, i.e. $\|\mathbf{D}_i^k\|_2 \leq \frac{1}{2}$, this immediately implies the following corollary. We will see in the following section (in Corollary A.1.15) that this minimum singular value bound for Γ_i^k can be interpreted in the following manner. Although the effective PŁ condition evolves dynamically, the associated PŁ constant always stays bounded away from zero.

Corollary A.1.13 (PŁ bounded away from zero). *Assume $D[i, k]$ holds. Then we have*

$$\sigma_{\min}(\Gamma_i^k)^2 \geq \frac{1}{2}.$$

The accumulated gradient signal is correlated with the full-batch gradient signal.

Lemma A.1.14 (Correlation of $\tilde{\mathbf{g}}^k$ and $\nabla_{\mathbf{M}}\mathcal{L}_\pi(\mathbf{M}_0^k)$). *For all k , we have*

$$\langle \nabla_{\mathbf{M}}\mathcal{L}_\pi(\mathbf{M}_0^k), \tilde{\mathbf{g}}^k \rangle_F \geq \sigma_{\min}(\mathbf{\Gamma}_0^k)^2 \|\nabla_{\mathbf{M}}\mathcal{L}_\pi(\mathbf{M}_0^k)\|_F^2.$$

Proof. Recall that we previously defined

$$\tilde{\mathbf{g}}^k \triangleq \nabla_{\mathbf{W}}\mathcal{L}_\pi(\Theta_0^k)\mathbf{\Gamma}_0^k + \mathbf{W}_0^k \nabla_{\mathbf{\Gamma}}\mathcal{L}_\pi(\Theta_0^k).$$

Note that if we have $\mathbf{A}, \mathbf{\Lambda} \in \mathbb{R}^{n \times n}$, with $\mathbf{\Lambda} = \text{diag}(\lambda_1, \dots, \lambda_n)$ a diagonal matrix with nonnegative entries, then

$$\langle \mathbf{A}, \mathbf{A}\mathbf{\Lambda} \rangle_F = \langle \mathbf{A}\mathbf{\Lambda}^{1/2}, \mathbf{A}\mathbf{\Lambda}^{1/2} \rangle_F = \|\mathbf{A}\mathbf{\Lambda}^{1/2}\|_F^2 \geq \min_i \lambda_i \|\mathbf{A}\|_F^2.$$

Also, we have

$$\langle \mathbf{A}, \text{diag}(\mathbf{A}) \rangle_F = \langle \text{diag}(\mathbf{A}), \text{diag}(\mathbf{A}) \rangle_F = \|\text{diag}(\mathbf{A})\|_F^2 \geq 0.$$

Hence combining Equations (20) and (22) and the above inequalities, we have

$$\begin{aligned} \langle \nabla_{\mathbf{M}}\mathcal{L}_\pi(\mathbf{M}_0^k), \tilde{\mathbf{g}}^k \rangle_F &= \langle \nabla_{\mathbf{M}}\mathcal{L}_\pi(\mathbf{M}_0^k), \nabla_{\mathbf{W}}\mathcal{L}_\pi(\Theta_0^k)\mathbf{\Gamma}_0^k \rangle_F \\ &\quad + \langle \nabla_{\mathbf{M}}\mathcal{L}_\pi(\mathbf{M}_0^k), \mathbf{W}_0^k \nabla_{\mathbf{\Gamma}}\mathcal{L}_\pi(\Theta_0^k) \rangle_F \\ &= \langle \nabla_{\mathbf{M}}\mathcal{L}_\pi(\mathbf{M}_0^k), \nabla_{\mathbf{M}}\mathcal{L}_\pi(\mathbf{M}_0^k)(\mathbf{\Gamma}_0^k)^2 \rangle_F \\ &\quad + \langle (\mathbf{W}_0^k)^\top \nabla_{\mathbf{M}}\mathcal{L}_\pi(\mathbf{M}_0^k), \text{diag}((\mathbf{W}_0^k)^\top \nabla_{\mathbf{M}}\mathcal{L}_\pi(\mathbf{M}_0^k)) \rangle_F \\ &\geq \sigma_{\min}(\mathbf{\Gamma}_0^k)^2 \|\nabla_{\mathbf{M}}\mathcal{L}_\pi(\mathbf{M}_0^k)\|_F^2. \end{aligned}$$

□

We obtain the following corollary of the above lemma and Corollary A.1.13.

Corollary A.1.15. *Assume $D[0, k]$ holds. We have*

$$\langle \nabla_{\mathbf{M}}\mathcal{L}_\pi(\mathbf{M}_0^k), \tilde{\mathbf{g}}^k \rangle_F \geq \frac{1}{2} \|\nabla_{\mathbf{M}}\mathcal{L}_\pi(\mathbf{M}_0^k)\|_F^2. \quad (23)$$

More generally, assume $D[i, k]$ holds. We have

$$\langle \nabla_{\mathbf{M}}\mathcal{L}_\pi(\mathbf{M}_i^k), \tilde{\mathbf{g}}^{(i,k)} \rangle_F \geq \frac{1}{2} \|\nabla_{\mathbf{M}}\mathcal{L}_\pi(\mathbf{M}_i^k)\|_F^2, \quad (24)$$

where $\tilde{\mathbf{g}}^{(i,k)}$ is the analogous quantity to $\tilde{\mathbf{g}}^k$ for accumulating gradients starting at iterate (i, k) rather than $(0, k)$. It is defined more formally in Equation (26).

A.1.5. BOUNDING NOISE TERMS

We now turn to bounding the composite noise term \mathbf{r}^k . This is crucial to ensure the global convergence via Equation (8) and also to control the approximate invariances.

Mismatched gradient noise is negligible. As promised, we show that the mismatched gradient noise terms \mathbf{q}_j^k are negligible when we accumulate gradients from \mathbf{M}_0^k to \mathbf{M}_0^{k+1} .

Lemma A.1.16. *Assume that $L[j, k]$ and $D[j, k]$ hold for $j < m$. Then we have*

$$\|\mathbf{q}_j^k\|_F \leq C_w^2 C_L \|\text{BN}(\mathbf{X}_\pi^{j+1})\|_2^2.$$

Furthermore, for any $t < m$ we have

$$\sum_{j=0}^t \|\mathbf{q}_j^k\|_F \leq C_w^2 C_L \|\bar{\mathbf{X}}_\pi\|_F^2.$$

Proof. Recall the definition of \mathbf{q}_j^k , reproduced here for reference:

$$\mathbf{q}_j^k \triangleq \nabla_{\mathbf{W}} \mathcal{L}(\mathbf{X}_\pi^{j+1}; \Theta_j^k) \nabla_{\Gamma} \mathcal{L}(\mathbf{X}_\pi^{j+1}; \Theta_j^k).$$

Since $L[j, k]$ and $D[j, k]$ hold for (j, k) , we can apply Corollary A.1.10 to conclude that

$$\|\mathbf{q}_j^k\|_F \leq C_w^2 C_L \|\text{BN}(\mathbf{X}_\pi^{j+1})\|_2^2.$$

Since the inductive hypotheses hold for every $j < m$, when we accumulate the noise terms from $(0, k)$ to (t, k) , we can apply the above bound to conclude that

$$\begin{aligned} \sum_{j=0}^t \|\mathbf{q}_j^k\|_F &\leq \sum_{j=0}^t \|\text{BN}(\mathbf{X}_\pi^{j+1})\|_2^2 C_w^2 C_L \\ &\leq C_w^2 C_L \sum_{j=0}^{m-1} \|\text{BN}(\mathbf{X}_\pi^{j+1})\|_F^2 \\ &= C_w^2 C_L \|\overline{\mathbf{X}}_\pi\|_F^2, \end{aligned}$$

where the last equality used the definition of $\overline{\mathbf{X}}_\pi$. \square

Path dependent noise arising from mini-batch updates is negligible. In order to bound the noise term coming from mini-batch updates, we first prove the following auxiliary lemma that shows that the iterates don't move far within an epoch.

Lemma A.1.17. *Fix $t \leq m$ and assume $D[j, k]$ and $L[j, k]$ hold for all $j < t$. Then we have*

$$\|\mathbf{W}_t^k - \mathbf{W}_0^k\|_2 \leq \sqrt{t} \eta_k C_w C_L^{1/2} \|\overline{\mathbf{X}}_\pi\|_F.$$

The same inequality holds true if we replace \mathbf{W} with Γ .

We also have

$$\|\mathbf{M}_t^k - \mathbf{M}_0^k\|_2 \leq 2\sqrt{t} \eta_k C_w^2 C_L^{1/2} \|\overline{\mathbf{X}}_\pi\|_F + \eta_k^2 C_w^2 C_L \|\overline{\mathbf{X}}_\pi\|_F^2.$$

Proof. We have by definition that

$$\mathbf{W}_t^k = \mathbf{W}_0^k - \eta_k \sum_{j=0}^{t-1} \nabla_{\mathbf{W}} \ell(\mathbf{X}_\pi^{j+1}; \Theta_j^k).$$

Now, we have

$$\begin{aligned} \|\mathbf{W}_t^k - \mathbf{W}_0^k\|_2 &\leq \eta_k \sum_{j=0}^{t-1} \|\nabla_{\mathbf{W}} \ell(\mathbf{X}_\pi^{j+1}; \Theta_j^k)\|_2 \\ &\leq \eta_k C_w C_L^{1/2} \sum_{j=0}^t \|\text{BN}(\mathbf{X}_\pi^{j+1})\|_2 \\ &\leq \sqrt{t} \eta_k C_w C_L^{1/2} \|\overline{\mathbf{X}}_\pi\|_F \end{aligned}$$

where in the first line we have applied the triangle inequality, in the second line we have applied Corollary A.1.10, and in the last line we have applied Cauchy-Schwarz.

The same proof holds for Γ .

For \mathbf{M} , Equation (10) gives

$$\mathbf{M}_t^k = \mathbf{M}_0^k - \eta_k \sum_{j=0}^{t-1} \mathbf{g}_j^k + \eta_k^2 \sum_{j=0}^{t-1} \mathbf{q}_j^k.$$

Combining Equation (11) and Corollaries A.1.8 and A.1.10 yields

$$\begin{aligned} \|g_j^k\|_2 &\leq \|\Gamma_j^k\|_2 \|\nabla_{\mathbf{W}} \mathcal{L}(\mathbf{X}_\pi^{j+1}; \Theta_j^k)\|_2 + \|\mathbf{W}_j^k\|_2 \|\nabla_{\Gamma} \mathcal{L}(\mathbf{X}_\pi^{j+1}; \Theta_j^k)\|_2 \\ &\leq 2C_w^2 C_L^{1/2} \|\text{BN}(\mathbf{X}_\pi^{j+1})\|_2. \end{aligned}$$

Hence, summing up over j , using Cauchy-Schwarz, and applying the noise bound Lemma A.1.16, it follows that

$$\|\mathbf{M}_t^k - \mathbf{M}_0^k\|_2 \leq 2\sqrt{t}\eta_k C_w^2 C_L^{1/2} \|\overline{\mathbf{X}}_\pi\|_F + \eta_k^2 C_w^2 C_L \|\overline{\mathbf{X}}_\pi\|_F^2.$$

□

Now we show that the noise term $\|e_j^k\|_2$ is $O(1)$.

Lemma A.1.18. *If $L[j, k]$ and $D[j, k]$ both hold for all $j < m$ then we have for each j that*

$$\|e_j^k\|_2 \leq 4\sqrt{j} C_w^2 C_L^{1/2} \|\overline{\mathbf{X}}_\pi\|_F^2 (C_L^{1/2} + C_w^2 \|\overline{\mathbf{X}}_\pi\|_F) + O(\eta_k).$$

Hence, we also have

$$\sum_{i=0}^{m-1} \|e_j^k\|_2 \leq 4m^{3/2} C_w^2 C_L^{1/2} \|\overline{\mathbf{X}}_\pi\|_F^2 (C_L^{1/2} + C_w^2 \|\overline{\mathbf{X}}_\pi\|_F) + O(\eta_k).$$

Proof. Inspecting the definition of e_j^k (Equation (15)), let us bound the quantity

$$\begin{aligned} \eta_k e_j^k &= \underbrace{\nabla_{\mathbf{W}} \mathcal{L}(\mathbf{X}_\pi^{j+1}; \Theta_0^k) \Gamma_0^k - \nabla_{\mathbf{W}} \mathcal{L}(\mathbf{X}_\pi^{j+1}; \Theta_j^k) \Gamma_j^k}_{\text{(I)}} \\ &\quad + \underbrace{\mathbf{W}_0^k \nabla_{\Gamma} \mathcal{L}(\mathbf{X}_\pi^{j+1}; \Theta_0^k) - \mathbf{W}_j^k \nabla_{\Gamma} \mathcal{L}(\mathbf{X}_\pi^{j+1}; \Theta_j^k)}_{\text{(II)}}. \end{aligned}$$

First, we have by triangle inequality and the identity Equation (21) that the norm of (I) is at most

$$\begin{aligned} &\|\nabla_{\mathbf{W}} \mathcal{L}(\mathbf{X}_\pi^{j+1}; \Theta_0^k) \Gamma_0^k - \nabla_{\mathbf{W}} \mathcal{L}(\mathbf{X}_\pi^{j+1}; \Theta_j^k) \Gamma_j^k\|_2 \\ &\quad + \|\nabla_{\mathbf{W}} \mathcal{L}(\mathbf{X}_\pi^{j+1}; \Theta_0^k) \Gamma_0^k - \nabla_{\mathbf{W}} \mathcal{L}(\mathbf{X}_\pi^{j+1}; \Theta_j^k) \Gamma_j^k\|_2 \\ &\leq \|\Gamma_0^k (\Gamma_j^k - \Gamma_0^k)\|_2 \|\nabla_{\mathbf{M}} \mathcal{L}(\mathbf{X}_\pi^{j+1}; \mathbf{M}_0^k)\|_2 \\ &\quad + \|\Gamma_j^k\|_2 \|\nabla_{\mathbf{M}} \mathcal{L}(\mathbf{X}_\pi^{j+1}; \mathbf{M}_0^k) \Gamma_0^k - \nabla_{\mathbf{M}} \mathcal{L}(\mathbf{X}_\pi^{j+1}; \mathbf{M}_j^k) \Gamma_j^k\|_2 \\ &\leq \|\Gamma_0^k (\Gamma_j^k - \Gamma_0^k)\|_2 \|\nabla_{\mathbf{M}} \mathcal{L}(\mathbf{X}_\pi^{j+1}; \mathbf{M}_0^k)\|_2 \\ &\quad + \|\Gamma_j^k\|_2 \left(\|\nabla_{\mathbf{M}} \mathcal{L}(\mathbf{X}_\pi^{j+1}; \mathbf{M}_0^k) \Gamma_0^k - \nabla_{\mathbf{M}} \mathcal{L}(\mathbf{X}_\pi^{j+1}; \mathbf{M}_0^k) \Gamma_j^k\|_2 \right. \\ &\quad \left. + \|\nabla_{\mathbf{M}} \mathcal{L}(\mathbf{X}_\pi^{j+1}; \mathbf{M}_0^k) \Gamma_j^k - \nabla_{\mathbf{M}} \mathcal{L}(\mathbf{X}_\pi^{j+1}; \mathbf{M}_j^k) \Gamma_j^k\|_2 \right) \\ &\leq (\|\Gamma_0^k\|_2 + \|\Gamma_j^k\|_2) \|\Gamma_j^k - \Gamma_0^k\|_2 \|\nabla_{\mathbf{M}} \mathcal{L}(\mathbf{X}_\pi^{j+1}; \mathbf{M}_0^k)\|_2 \\ &\quad + \|\Gamma_j^k\|_2^2 \|\nabla_{\mathbf{M}} \mathcal{L}(\mathbf{X}_\pi^{j+1}; \mathbf{M}_0^k) - \nabla_{\mathbf{M}} \mathcal{L}(\mathbf{X}_\pi^{j+1}; \mathbf{M}_j^k)\|_2. \end{aligned}$$

Applying the weight bounds in Corollaries A.1.8 and A.1.10 and Lemma A.1.17 to the first term yields an upper bound of

$$2C_w \cdot (\sqrt{j}\eta_k C_w C_L^{1/2} \|\overline{\mathbf{X}}_\pi\|_F) \cdot (\|\text{BN}(\mathbf{X}_\pi^{j+1})\|_2 C_L^{1/2}) \leq 2\sqrt{j}\eta_k C_w^2 C_L \|\overline{\mathbf{X}}_\pi\|_F^2.$$

Turning to the second term, we can apply the smoothness bound in Lemma A.1.2 and the inductive bounds in Corollaries A.1.8 and A.1.10 and Lemma A.1.17 to obtain an upper bound of

$$C_w^2 \|\text{BN}(\mathbf{X}_\pi^{j+1})\|_2^2 \|\mathbf{M}_j^k - \mathbf{M}_0^k\|_2 \leq 2\sqrt{j}\eta_k C_w^4 C_L^{1/2} \|\overline{\mathbf{X}}_\pi\|_F^3 + \eta_k^2 C_w^4 C_L \|\overline{\mathbf{X}}_\pi\|_F^4.$$

Putting it together, we have

$$(I) \leq 2\sqrt{j}\eta_k C_w^2 C_L^{1/2} \|\bar{\mathbf{X}}_\pi\|_F^2 (C_L^{1/2} + C_w^2 \|\bar{\mathbf{X}}_\pi\|_F) + \eta_k^2 C_w^4 C_L \|\bar{\mathbf{X}}_\pi\|_F^4.$$

Similarly, for (II) we have the exact same bound since we can apply Lemma A.1.5 to remove the diagonal operator and uniformly bound $\|\mathbf{W}_j^k\|_2$ and $\|\Gamma_j^k\|_2$ by C_w .

Finally, combining (I) and (II) and dividing through by η_k , we can conclude that

$$\|e_j^k\|_2 \leq 4\sqrt{j}C_w^2 C_L^{1/2} \|\bar{\mathbf{X}}_\pi\|_F^2 (C_L^{1/2} + C_w^2 \|\bar{\mathbf{X}}_\pi\|_F) + O(\eta_k).$$

Summing up $\|e_j^k\|_2$ over all j and crudely bounding $\sum_{j=0}^{m-1} \sqrt{j} \leq m^{3/2}$, we see that

$$\sum_{i=0}^{m-1} \|e_j^k\|_2 \leq 4m^{3/2} C_w^2 C_L^{1/2} \|\bar{\mathbf{X}}_\pi\|_F^2 (C_L^{1/2} + C_w^2 \|\bar{\mathbf{X}}_\pi\|_F) + O(\eta_k).$$

□

Stepsize noise is negligible for $i > 0$. We now quickly show that the effect of generalizing the induction to $i > 0$ is negligible. In particular, we can carry out the same proof, except we will have to redefine the per-epoch update so it can account for gradient updates starting at an arbitrary iterate (i, k) rather than $(0, k)$. We explicitly redefine these terms below by quickly revisiting the signal-noise decomposition in Appendix A.1.3. Recall that a single-iterate update at iteration (a, b) can be written as (Equation (10))

$$\mathbf{M}_{a+1}^b = \mathbf{M}_a^b - \eta_b \mathbf{g}_a^b + \eta_b^2 \mathbf{q}_a^b,$$

where we defined

$$\begin{aligned} \mathbf{g}_a^b &\triangleq \nabla_{\mathbf{W}} \mathcal{L}(\mathbf{X}_\pi^{a+1}; \Theta_a^b) \Gamma_a^b + \mathbf{W}_a^b \nabla_{\Gamma} \mathcal{L}(\mathbf{X}_\pi^{a+1}; \Theta_a^b), \\ \mathbf{q}_a^b &\triangleq \nabla_{\mathbf{W}} \mathcal{L}(\mathbf{X}_\pi^{a+1}; \Theta_a^b) \nabla_{\Gamma} \mathcal{L}(\mathbf{X}_\pi^{a+1}; \Theta_a^b). \end{aligned}$$

Consider carrying out the same accumulation as in Appendix A.1.3, but this time choosing (i, k) instead $(0, k)$ as the ‘‘pivot.’’ For this purpose, we will change our notational convention a little bit and use superscripts to denote the pivot or the starting point (i, k) . As the redefinitions of the ‘‘signal’’ $\tilde{\mathbf{g}}_a^b$ (Equation (14)) and path dependent noise e_a^b (Equation (15)), we define

$$\begin{aligned} \tilde{\mathbf{g}}_a^{(i,k)} &\triangleq \nabla_{\mathbf{W}} \mathcal{L}(\mathbf{X}_\pi^{a+1}; \Theta_i^k) \Gamma_i^k + \mathbf{W}_i^k \nabla_{\Gamma} \mathcal{L}(\mathbf{X}_\pi^{a+1}; \Theta_i^k), \\ e_{(a,b)}^{(i,k)} &\triangleq \frac{\tilde{\mathbf{g}}_a^{(i,k)} - \mathbf{g}_a^b}{\eta_k}, \end{aligned}$$

for indices (a, b) satisfying $(i, k) \leq (a, b) \leq (i-1, k+1)$.

This way, the accumulation of updates on \mathbf{M} from iteration (i, k) to $(i-1, k+1)$ can be represented as

$$\begin{aligned} \mathbf{M}_i^{k+1} &= \mathbf{M}_i^k - \eta_k \sum_{j=i}^{m-1} \mathbf{g}_j^k + \eta_k^2 \sum_{j=i}^{m-1} \mathbf{q}_j^k - \eta_{k+1} \sum_{j=0}^{i-1} \mathbf{g}_j^{k+1} + \eta_{k+1}^2 \sum_{j=0}^{i-1} \mathbf{q}_j^{k+1} \\ &= \mathbf{M}_i^k - \eta_k \sum_{(a,b)=(i,k)}^{(i-1,k+1)} \mathbf{g}_a^b + \eta_k^2 \sum_{(a,b)=(i,k)}^{(i-1,k+1)} \mathbf{q}_a^b \\ &\quad - (\eta_{k+1} - \eta_k) \sum_{j=0}^{i-1} \mathbf{g}_j^{k+1} + (\eta_{k+1}^2 - \eta_k^2) \sum_{j=0}^{i-1} \mathbf{q}_j^{k+1} \\ &= \mathbf{M}_i^k - \eta_k \tilde{\mathbf{g}}^{(i,k)} + \eta_k^2 \sum_{(a,b)=(i,k)}^{(i-1,k+1)} \left(e_{(a,b)}^{(i,k)} + \mathbf{q}_a^b \right) + \eta_k^2 \mathbf{s}_{(i,k+1)}^{(i,k)} \end{aligned} \tag{25}$$

where in the last equality we used $\mathbf{g}_a^b = \tilde{\mathbf{g}}_a^{(i,k)} - \eta_k \mathbf{e}_{(a,b)}^{(i,k)}$ and also defined the accumulated signal $\tilde{\mathbf{g}}^{(i,k)}$ (a redefinition of $\tilde{\mathbf{g}}^k$ from Equation (16)) and the stepsize noise $\mathbf{s}_{(i,k+1)}^{(i,k)}$ as

$$\tilde{\mathbf{g}}^{(i,k)} \triangleq \sum_{(a,b)=(i,k)}^{(i-1,k+1)} \tilde{\mathbf{g}}_a^{(i,k)} = \nabla_{\mathbf{W}} \mathcal{L}_\pi(\Theta_i^k) \Gamma_i^k + \mathbf{W}_i^k \nabla_{\Gamma} \mathcal{L}_\pi(\Theta_i^k), \quad (26)$$

$$\mathbf{s}_{(i,k+1)}^{(i,k)} \triangleq -\frac{\eta_{k+1} - \eta_k}{\eta_k^2} \sum_{j=0}^{i-1} \mathbf{g}_j^{k+1} + \frac{\eta_{k+1}^2 - \eta_k^2}{\eta_k^2} \sum_{j=0}^{i-1} \mathbf{q}_j^{k+1}. \quad (27)$$

As a sanity check, we can quickly see that the stepsize noise $\mathbf{s}_{(i,k+1)}^{(i,k)}$ is zero if $i = 0$.

Now notice that Equation (25) can be thought of as a generalization of the per-epoch update (Equation (17)) originally obtained for $i = 0$. For now, suppose we ignore the last term of Equation (25) involving the stepsize noise $\mathbf{s}_{(i,k+1)}^{(i,k)}$. Then, if we carry out the above analysis for bounding the remaining terms in Equation (25), there is no difference in the argument up to reindexing; we can consider this as using η_k for the stepsize *even for the iterates* (a, b) where $b = k + 1$. In particular, the lemmas of the previous sections all hold up to reindexing notation.

Therefore, it now suffices to show that the stepsize noise $\mathbf{s}_{(i,k+1)}^{(i,k)}$ is of the same order as the other noise terms; in particular, $\left\| \mathbf{s}_{(i,k+1)}^{(i,k)} \right\|_F = O(1)$.

Lemma A.1.19. *Assume $L[j, k + 1]$ and $D[j, k + 1]$ hold for all $j \leq i - 1$. Suppose that*

$$\eta_k = O\left(\frac{1}{k^\beta}\right),$$

for some $1/2 < \beta < 1$. Then the stepsize noise

$$\left\| \mathbf{s}_{(i,k+1)}^{(i,k)} \right\|_F = O(1).$$

Proof. Since $D[j, k + 1]$ holds, Lemma A.1.16 demonstrates that $\left\| \mathbf{q}_j^{k+1} \right\|_2 = O(1)$. On the other hand, Corollaries A.1.8 and A.1.10 show that $\left\| \mathbf{g}_j^{k+1} \right\|_2 = O(1)$. Since $\eta_k = O(1/k^\beta)$ and $\beta \leq 1$, we have $\eta_{k+1} - \eta_k = O(1/k^{\beta+1}) = O(\eta_k^2)$. Similarly $\eta_{k+1}^2 - \eta_k^2 = O(1/k^{2\beta+1}) = O(\eta_k^3)$. Plugging these into Equation (27), we conclude that $\left\| \mathbf{s}_{(i,k+1)}^{(i,k)} \right\|_F = O(1)$, as desired. \square

Composite noise term is negligible. Now that we have formally defined the stepsize noise that arise for $i > 0$, we also redefine the composite noise term \mathbf{r}^k (Equation (19)) originally defined for $i = 0$. The updated definition is simply

$$\mathbf{r}^{(i,k)} \triangleq \sum_{(a,b)=(i,k)}^{(i-1,k+1)} \left(\mathbf{e}_{(a,b)}^{(i,k)} + \mathbf{q}_a^b \right) + \mathbf{s}_{(i,k+1)}^{(i,k)}, \quad (28)$$

which allows us to rewrite the epoch update spelled out in Equation (25) as

$$\mathbf{M}_i^{k+1} = \mathbf{M}_i^k - \eta_k \tilde{\mathbf{g}}^{(i,k)} + \eta_k^2 \mathbf{r}^{(i,k)}, \quad (29)$$

which is a generalization of Equation (18).

It is left to show formally that the composite noise term $\mathbf{r}^{(i,k)}$ defined in Equation (28), obtained from combining the mismatched gradient noise terms \mathbf{q}_a^b , the path dependent noise $\mathbf{e}_{(a,b)}^{(i,k)}$ for $(i, k) \leq (a, b) \leq (i - 1, k + 1)$, and the stepsize noise $\mathbf{s}_{(i,k+1)}^{(i,k)}$, is indeed $O(1)$.

Proposition A.1.20 (Composite noise term). *Suppose $L[a, b]$ and $D[a, b]$ hold for $(i, k) \leq (a, b) \leq (i - 1, k + 1)$, and $\eta_k = O(1/k^\beta)$ for some $1/2 < \beta < 1$. Then the composite noise term $\mathbf{r}^{(i,k)}$ satisfies*

$$\left\| \mathbf{r}^{(i,k)} \right\|_F \leq \text{poly}(m, C_w, C_L, \|\bar{\mathbf{X}}_\pi\|_F) + O(\eta_k).$$

Proof. Combining Lemmas A.1.16, A.1.18 and A.1.19, taking care to analyze the constants (which are all $\text{poly}(m, C_w, C_L, \|\bar{\mathbf{X}}_\pi\|_F)$) hidden by the big- O notation, yields the desired result. \square

A.1.6. ACCUMULATED LOSS UPDATE

In this section, we formally account for the noise terms and prove that an accumulated loss inequality holds. More precisely, we can use the gradient update spelled out in Equation (29) and the noise bounds in Appendix A.1.5 to obtain a single epoch loss update. In other words, this section prove the inductive step that the hypotheses $R[i, k + 1]$ and $L[i, k + 1]$ holds. For the sake of simplicity, in this section we focus on the case $i = 0$. However, as discussed above, the case $i > 0$ only adds a negligible stepsize error and the same arguments go through.

We start with the following proposition, which proves the hypothesis $R[0, k + 1]$.

Proposition A.1.21. *Assume that $L[j, k]$ and $D[j, k]$ hold for all $j < m$. Consider optimizing the linear+BN network with stepsize satisfying*

$$\eta_k \leq \frac{1}{2k^\beta},$$

for $1/2 < \beta < 1$.

Then

$$\mathcal{L}_\pi(\mathbf{M}_0^{k+1}) - \mathcal{L}_\pi^* \leq \left(1 - \frac{\alpha_\pi \eta_k}{2}\right) (\mathcal{L}_\pi(\mathbf{M}_0^k) - \mathcal{L}_\pi^*) + \text{poly}(m, C_w, C_L, \|\bar{\mathbf{X}}_\pi\|_F) \eta_k^2, \quad (30)$$

where the $\text{poly}(m, C_w, C_L, \|\bar{\mathbf{X}}_\pi\|_F)$ term is independent of k and has constant degree.

Proof. First, we use the G_π -smoothness of \mathcal{L}_π with respect to \mathbf{M} guaranteed by Lemma A.1.2 to obtain

$$\mathcal{L}_\pi(\mathbf{M}_0^{k+1}) - \mathcal{L}_\pi(\mathbf{M}_0^k) \leq \langle \nabla_{\mathbf{M}} \mathcal{L}_\pi(\mathbf{M}_0^k), \mathbf{M}_0^{k+1} - \mathbf{M}_0^k \rangle_F + \frac{G_\pi}{2} \|\mathbf{M}_0^{k+1} - \mathbf{M}_0^k\|_F^2.$$

Using the gradient update Equation (18) and Cauchy-Schwarz, we can upper bound the RHS by

$$\begin{aligned} & \langle \nabla_{\mathbf{M}} \mathcal{L}_\pi(\mathbf{M}_0^k), -\eta_k \tilde{\mathbf{g}}^k + \eta_k^2 \mathbf{r}^k \rangle_F + \frac{G_\pi}{2} \|\mathbf{M}_0^{k+1} - \mathbf{M}_0^k\|_F^2 \\ & \leq -\eta_k \langle \nabla_{\mathbf{M}} \mathcal{L}_\pi(\mathbf{M}_0^k), \tilde{\mathbf{g}}^k \rangle_F + \eta_k^2 \|\mathbf{r}^k\|_F \|\nabla_{\mathbf{M}} \mathcal{L}_\pi(\mathbf{M}_0^k)\|_F + \frac{G_\pi}{2} \|\mathbf{M}_0^{k+1} - \mathbf{M}_0^k\|_F^2. \end{aligned}$$

Next, we can use Lemma A.1.17, with $t = m$, together with the inequality $(a + b)^2 \leq 2a^2 + 2b^2$, to obtain an upper bound of

$$\begin{aligned} & \leq -\eta_k \langle \nabla_{\mathbf{M}} \mathcal{L}_\pi(\mathbf{M}_0^k), \tilde{\mathbf{g}}^k \rangle_F + \eta_k^2 \|\mathbf{r}^k\|_F \|\nabla_{\mathbf{M}} \mathcal{L}_\pi(\mathbf{M}_0^k)\|_F \\ & \quad + \frac{G_\pi}{2} (4m\eta_k^2 C_w^4 C_L \|\bar{\mathbf{X}}_\pi\|_F^2 + \eta_k^4 C_w^4 C_L^2 \|\bar{\mathbf{X}}_\pi\|_F^4). \end{aligned}$$

Then, because the inductive hypotheses apply we can apply Corollary A.1.10 to bound gradients, Corollary A.1.15 to bound the inner product $\langle \nabla_{\mathbf{M}} \mathcal{L}_\pi(\mathbf{M}_0^k), \tilde{\mathbf{g}}^k \rangle_F$. Moreover, since $\eta_k = O(1/k^\beta)$, we can use Proposition A.1.20 to bound $\|\mathbf{r}^k\|_F$. This yields an upper bound of

$$\begin{aligned} & -\frac{\eta_k}{2} \|\nabla_{\mathbf{M}} \mathcal{L}_\pi(\mathbf{M}_0^k)\|_F^2 + \eta_k^2 \|\nabla_{\mathbf{M}} \mathcal{L}_\pi(\mathbf{M}_0^k)\|_F \|\mathbf{r}^k\|_F + \text{poly}(m, C_w, C_L, \|\bar{\mathbf{X}}_\pi\|_F) \eta_k^2 \\ & \leq \left(-\frac{\eta_k}{2} + \frac{\eta_k^2}{2}\right) \|\nabla_{\mathbf{M}} \mathcal{L}_\pi(\mathbf{M}_0^k)\|_F^2 + \text{poly}(m, C_w, C_L, \|\bar{\mathbf{X}}_\pi\|_F) \eta_k^2 \\ & \leq -\frac{\eta_k}{4} \|\nabla_{\mathbf{M}} \mathcal{L}_\pi(k)\|_F^2 + \text{poly}(m, C_w, C_L, \|\bar{\mathbf{X}}_\pi\|_F) \eta_k^2 \end{aligned}$$

In the second line, we have used $ab \leq \frac{1}{2}(a^2 + b^2)$, and throughout, we have used the assumption $\eta_k \leq \frac{1}{2}$ to reduce higher order terms of η_k .

Putting it together, we find that

$$\mathcal{L}_\pi(\mathbf{M}_0^{k+1}) - \mathcal{L}_\pi(\mathbf{M}_0^k) \leq -\frac{\eta_k}{4} \|\nabla_{\mathbf{M}} \mathcal{L}_\pi(\mathbf{M}_0^k)\|_F^2 + \text{poly}(m, C_w, C_L, \|\bar{\mathbf{X}}_\pi\|_F) \eta_k^2$$

We now use α_π -strong convexity with respect to M (and hence α_π -PL) guaranteed by Lemma A.1.3 to obtain

$$\mathcal{L}_\pi(\mathbf{M}_0^{k+1}) - \mathcal{L}_\pi^* \leq \left(1 - \frac{\alpha_\pi \eta_k}{2}\right) (\mathcal{L}_\pi(\mathbf{M}_0^k) - \mathcal{L}_\pi^*) + \text{poly}(m, C_w, C_L, \|\bar{\mathbf{X}}_\pi\|_F) \eta_k^2.$$

□

One consequence of Proposition A.1.21 is that if the stepsize η_k is small enough, we can guarantee that the loss decreases from $\mathcal{L}_\pi(\mathbf{M}_i^k)$ to $\mathcal{L}_\pi(\mathbf{M}_i^{k+1})$.

Next, from Proposition A.1.21, we can prove the other inductive hypothesis, namely $L[0, k+1]$.

Corollary A.1.22. *Suppose $L[j, k]$ and $D[j, k]$ hold for all $j < m$ and the stepsize satisfies*

$$\eta_k \leq \frac{1}{k^\beta} \min \left\{ \frac{1}{2}, \frac{\alpha_\pi C_L}{\text{poly}(m, C_w, C_L, \|\bar{\mathbf{X}}_\pi\|_F)} \right\},$$

for some $1/2 < \beta < 1$.

Then $L[0, k+1]$ holds, i.e.

$$\mathcal{L}_\pi(\mathbf{M}_0^{k+1}) \leq C_L.$$

Proof. Since $\eta_k \leq \frac{1}{2k^\beta}$, we can apply Proposition A.1.21 to conclude that Equation (30) holds. Then, for the bound $\mathcal{L}_\pi(\mathbf{M}_0^{k+1}) \leq C_L$ to hold, it suffices to show that

$$\left(1 - \frac{\alpha_\pi \eta_k}{2}\right) C_L + \text{poly}(m, C_w, C_L, \|\bar{\mathbf{X}}_\pi\|_F) \eta_k^2 \leq C_L.$$

Equivalently,

$$\text{poly}(m, C_w, C_L, \|\bar{\mathbf{X}}_\pi\|_F) \eta_k \leq \frac{\alpha_\pi}{2} C_L.$$

Clearly this holds for the stated assumption on η_k . □

Finally, we show that by inductively unrolling the inequality in Proposition A.1.21, we can show that $\mathcal{L}_\pi(\mathbf{M}_i^k)$ converges to \mathcal{L}_π^* at a sublinear rate.

Proposition A.1.23. *Assume we are in the same setup as Proposition A.1.21. Suppose that the stepsize satisfies*

$$\eta_k = \frac{c}{k^\beta},$$

for some absolute constant c such that $c \leq \min \left\{ \frac{1}{2}, \frac{2}{\alpha_\pi} \right\}$ and $1/2 < \beta < 1$. Further suppose that $R[0, b]$ holds for every $b \in [k+1]$. Then if $\beta < 1$ we have

$$\mathcal{L}_\pi(\Theta_0^{k+1}) - \mathcal{L}_\pi^* \leq (\mathcal{L}_\pi(\Theta_0^1) - \mathcal{L}_\pi^*) \exp \left(\frac{c \alpha_\pi}{2(1-\beta)} (2 - k^{1-\beta}) \right) + \frac{c^2 \text{poly}(m, C_w, C_L, \|\bar{\mathbf{X}}_\pi\|_F) \log k}{k^\beta},$$

Proof. Note that by inspecting the proof of Proposition A.1.21 and Proposition A.1.20, the term $\text{poly}(m, C_w, C_L, \|\bar{\mathbf{X}}_\pi\|_F)$ has no dependence on k . So for simplicity we will assume that this term is bounded by some absolute constant A . Since $\eta_k \leq \frac{1}{2k^\beta}$, Proposition A.1.21 implies that

$$\mathcal{L}_\pi(\mathbf{M}_0^{k+1}) - \mathcal{L}_\pi^* \leq \left(1 - \frac{\alpha_\pi \eta_k}{2}\right) (\mathcal{L}_\pi(\mathbf{M}_0^k) - \mathcal{L}_\pi^*) + A \eta_k^2 \quad (31)$$

We can unroll the recurrence to obtain

$$\mathcal{L}_\pi(\mathbf{M}_0^{k+1}) - \mathcal{L}_\pi^* \leq (\mathcal{L}_\pi(\mathbf{M}_0^1) - \mathcal{L}_\pi^*) \prod_{t=1}^k \left(1 - \frac{\alpha_\pi \eta_t}{2}\right) + A \sum_{t=1}^k \eta_t^2 \left(\prod_{j=t+1}^k \left(1 - \frac{\alpha_\pi \eta_j}{2}\right) \right). \quad (32)$$

We have for any $c_t \leq 1$ that

$$\begin{aligned} \prod_{t=a}^b (1 - c_t) &\leq \exp\left(\sum_{t=a}^b \log(1 - c_t)\right) \\ &\leq \exp\left(-\sum_{t=a}^b c_t\right), \end{aligned}$$

where we have used $\log(1 - x) \leq -x$ for $x \leq 1$. For $1/2 < \beta < 1$ we have

$$\sum_{t=a}^b \frac{1}{t^\beta} \geq \int_a^b \frac{1}{t^\beta} dt = \frac{b^{1-\beta} - a^{1-\beta}}{1-\beta}.$$

Hence, since we assumed $\eta_k = \frac{c}{k^\beta}$, and $\frac{\alpha_\pi c}{2} \leq 1$, we have

$$\prod_{t=a}^b \left(1 - \frac{\alpha_\pi \eta_t}{2}\right) \leq \exp\left(-\frac{c\alpha_\pi}{2(1-\beta)}(b^{1-\beta} - a^{1-\beta})\right)$$

Now we can bound

$$A \sum_{t=1}^k \eta_t^2 \left(\prod_{j=t+1}^k \left(1 - \frac{\alpha_\pi \eta_j}{2}\right) \right) \leq \frac{c^2 A}{k^{2\beta}} + \sum_{t=1}^{k-1} \frac{c^2 A}{t^{2\beta}} \exp\left(-\frac{c\alpha_\pi}{2(1-\beta)}(k^{1-\beta} - (t+1)^{1-\beta})\right)$$

Define $T \triangleq k - Ck^\beta \log k$, where $C > 0$ is an absolute constant to be picked later. We can split up the sum into $t < T$ and $t \geq T$. For the terms $t < T$ we can use concavity to deduce that

$$k^{1-\beta} - (t+1)^{1-\beta} \geq (1-\beta)(k-t-1)k^{-\beta} \geq \Theta(\log k),$$

where the constant hidden in $\Theta(\log k)$ increases with C . Hence we pick C so that for $t < T$ we have

$$\frac{c\alpha_\pi}{2(1-\beta)}(k^{1-\beta} - (t+1)^{1-\beta}) \geq \beta \log k.$$

Then,

$$\begin{aligned} \sum_{t < T} \frac{c^2 A}{t^{2\beta}} \exp\left(-\frac{c\alpha_\pi}{2(1-\beta)}(k^{1-\beta} - (t+1)^{1-\beta})\right) &\leq \exp(-\beta \log k) \sum_{t < T} \frac{c^2 A}{t^{2\beta}} \\ &\leq O\left(\frac{c^2 A}{k^\beta}\right). \end{aligned}$$

On the other hand, for the terms $t \geq T$ we can naively bound the exponential term by 1 and obtain

$$\begin{aligned} \sum_{t \geq T} \frac{c^2 A}{t^{2\beta}} \exp\left(-\frac{c\alpha_\pi}{2(1-\beta)}(k^{1-\beta} - (t+1)^{1-\beta})\right) &\leq \sum_{t \geq T} \frac{c^2 A}{t^{2\beta}} \\ &\leq \Theta(k^\beta \log k) \frac{c^2 A}{T^{2\beta}} \\ &\leq \Theta\left(\frac{c^2 A \log k}{k^\beta}\right). \end{aligned}$$

Hence we have that

$$\mathcal{L}_\pi(\mathbf{M}_0^{k+1}) - \mathcal{L}_\pi^* \leq (\mathcal{L}_\pi(\mathbf{M}_0^1) - \mathcal{L}_\pi^*) \exp\left(\frac{c\alpha_\pi}{2(1-\beta)}(2 - k^{1-\beta})\right) + \Theta\left(\frac{c^2 A \log k}{k^\beta}\right), \quad (33)$$

and the inequality in the proposition statement holds by recalling that $A = \text{poly}(m, C_w, C_L, \|\bar{\mathbf{X}}_\pi\|_F)$. \square

A.1.7. BOUNDING APPROXIMATE INVARIANCES

Armed with Corollaries A.1.8 and A.1.10, we slog through the arduous task of inductively bounding the approximate invariance. As a reminder, these corollaries tell us that assuming the inductive hypotheses $L[j, k]$ and $D[j, k]$ hold, all weight norms and losses for iterate (j, k) can be bounded by uniform constants.

Lemma A.1.24. *Suppose $L[j, k]$ and $D[j, k]$ hold for some $j < m$. We have*

$$\|D_{j+1}^k - D_j^k\|_2 \leq 2C_w^2 C_L \|\text{BN}(\mathbf{X}_\pi^{j+1})\|_2^2 \eta_k^2.$$

Hence, if $D[t, k]$ holds for all $t \leq j$, we have

$$\|D_{j+1}^k\|_2 \leq 2C_w^2 C_L \|\bar{\mathbf{X}}_\pi\|_F^2 \sum_{t=1}^k \eta_t^2.$$

Proof. We have

$$\begin{aligned} (\mathbf{W}_{j+1}^k)^\top \mathbf{W}_{j+1}^k &= (\mathbf{W}_j^k - \eta_k \nabla_{\mathbf{W}} \mathcal{L}(\mathbf{X}_\pi^{j+1}; \Theta_j^k))^\top (\mathbf{W}_j^k - \eta_k \nabla_{\mathbf{W}} \mathcal{L}(\mathbf{X}_\pi^{j+1}; \Theta_j^k)) \\ &= (\mathbf{W}_j^k)^\top \mathbf{W}_j^k - \eta_k [\nabla_{\mathbf{W}} \mathcal{L}(\mathbf{X}_\pi^{j+1}; \Theta_j^k)^\top \mathbf{W}_j^k + (\mathbf{W}_j^k)^\top \nabla_{\mathbf{W}} \mathcal{L}(\mathbf{X}_\pi^{j+1}; \Theta_j^k)] \\ &\quad + \eta_k^2 [\nabla_{\mathbf{W}} \mathcal{L}(\mathbf{X}_\pi^{j+1}; \Theta_j^k)^\top \nabla_{\mathbf{W}} \mathcal{L}(\mathbf{X}_\pi^{j+1}; \Theta_j^k)] \end{aligned}$$

Similarly, we have

$$\begin{aligned} (\Gamma_{j+1}^k)^2 &= (\Gamma_j^k - \eta_k \nabla_{\Gamma} \mathcal{L}(\mathbf{X}_\pi^{j+1}; \Theta_j^k)) (\Gamma_j^k - \eta_k \nabla_{\Gamma} \mathcal{L}(\mathbf{X}_\pi^{j+1}; \Theta_j^k)) \\ &= (\Gamma_j^k)^2 - \eta_k [\nabla_{\Gamma} \mathcal{L}(\mathbf{X}_\pi^{j+1}; \Theta_j^k) \Gamma_j^k + \Gamma_j^k \nabla_{\Gamma} \mathcal{L}(\mathbf{X}_\pi^{j+1}; \Theta_j^k)] \\ &\quad + \eta_k^2 [\nabla_{\Gamma} \mathcal{L}(\mathbf{X}_\pi^{j+1}; \Theta_j^k) \nabla_{\Gamma} \mathcal{L}(\mathbf{X}_\pi^{j+1}; \Theta_j^k)] \end{aligned}$$

The gradient invariance in Fact A.1.1 cancels out the η_k term in $D_{j+1}^k = \mathbf{I} + \text{diag}((\mathbf{W}_{j+1}^k)^\top \mathbf{W}_{j+1}^k - (\Gamma_{j+1}^k)^2)$. Hence, if we take the operator norm of $D_{j+1}^k - D_j^k$ and use Lemma A.1.5, we can ignore the diagonal operator. Then, since the inductive hypotheses hold, we can apply the inductive gradient bound (Corollary A.1.10) to obtain

$$\begin{aligned} \|D_{j+1}^k - D_j^k\|_2 &\leq \eta_k^2 \left[\|\nabla_{\mathbf{W}} \mathcal{L}(\mathbf{X}_\pi^{j+1}; \Theta_j^k)\|_2^2 + \|\nabla_{\Gamma} \mathcal{L}(\mathbf{X}_\pi^{j+1}; \Theta_j^k)\|_2^2 \right] \\ &\leq 2C_w^2 C_L \|\text{BN}(\mathbf{X}_\pi^{j+1})\|_2^2 \eta_k^2. \end{aligned}$$

To conclude, we apply triangle inequality and the inductive hypothesis on $D[t, k]$, yielding

$$\begin{aligned} \|D_{j+1}^k\|_2 &\leq \|D_0^k\|_2 + 2C_w^2 C_L \eta_k^2 \sum_{t=0}^j \|\text{BN}(\mathbf{X}_\pi^{t+1})\|_2^2 \\ &\leq \|D_0^k\|_2 + 2C_w^2 C_L \eta_k^2 \|\bar{\mathbf{X}}_\pi\|_F^2 \\ &\leq 2C_w^2 C_L \|\bar{\mathbf{X}}_\pi\|_F^2 \sum_{t=1}^k \eta_t^2. \end{aligned}$$

□

As a corollary, we see that for stepsizes of the form $\eta_k = c/k^\beta$ for $1/2 < \beta < 1$, we can select c to guarantee that $\|D_j^k\|_2 \leq 1/2$ for all (j, k) .

Corollary A.1.25. *Assume that $D[t, k]$ holds for all $t \leq j < m$ and $\eta_k = \frac{c}{k^\beta}$ for $1/2 < \beta < 1$. If*

$$c^2 \leq \frac{1}{4(1 + \frac{1}{2\beta-1})C_w^2 C_L \|\bar{\mathbf{X}}_\pi\|_F^2},$$

then

$$\|\mathbf{D}_{j+1}^k\|_2 \leq 2C_w^2 C_L \|\overline{\mathbf{X}}_\pi\|_F^2 \sum_{t=1}^k \eta_t^2 \leq \frac{1}{2}.$$

In other words, $D[j+1, k]$ holds.

Proof. Lemma A.1.24 implies that

$$\|\mathbf{D}_{j+1}^k\|_2 \leq 2C_w^2 C_L \|\overline{\mathbf{X}}_\pi\|_F^2 \sum_{t=1}^k \eta_t^2 \leq 2C_w^2 C_L \|\overline{\mathbf{X}}_\pi\|_F^2 \sum_{t=1}^{\infty} \eta_t^2.$$

For $1/2 < \beta < 1$, we have

$$\sum_{k=1}^{\infty} \frac{1}{k^{2\beta}} \leq 1 + \int_1^{\infty} \frac{1}{t^{2\beta}} dt = 1 + \frac{1}{1-2\beta} t^{1-2\beta} \Big|_1^{\infty} = 1 + \frac{1}{2\beta-1}.$$

Hence, if

$$c^2 \leq \frac{1}{4(1 + \frac{1}{2\beta-1})C_w^2 C_L \|\overline{\mathbf{X}}_\pi\|_F^2},$$

then evidently $\|\mathbf{D}_{j+1}^k\|_2 \leq \frac{1}{2}$, as desired. \square

A.1.8. COMPLETING THE INDUCTION

With all the pieces in place, we formally state the theorem for SS convergence.

Theorem A.1.26 (Formal statement of convergence for SS). *Let $\pi \in \mathbb{S}_n$ be such that Assumption 1(a) holds. Let $f(\cdot; \Theta) = \mathbf{W}\Gamma\mathbf{B}\mathbf{N}(\cdot)$ be a 2-layer linear+BN network initialized at $\Theta_0^1 = (\mathbf{W}_0^1, \Gamma_0^1) = (\mathbf{0}, \mathbf{I})$. Consider optimizing f using SS with permutation π and decreasing stepsize*

$$\eta_k = \frac{1}{k^\beta} \cdot \min \left\{ \frac{1}{2}, \frac{2}{\alpha_\pi}, \frac{\alpha_\pi C_L}{\text{poly}(m, C_w, C_L, \|\overline{\mathbf{X}}_\pi\|_F)}, \sqrt{\frac{1}{4(1 + \frac{1}{2\beta-1})C_w^2 C_L \|\overline{\mathbf{X}}_\pi\|_F^2}} \right\}$$

for any $1/2 < \beta < 1$. Then the SS risk satisfies

$$\mathcal{L}_\pi(\Theta_0^{k+1}) - \mathcal{L}_\pi^* \leq \frac{\text{poly}(m, d, C_L, \|\overline{\mathbf{X}}_\pi\|_F, \frac{1}{\sigma_{\min}(\overline{\mathbf{X}}_\pi^\top)}) \log k}{k^\beta}.$$

In particular, the SS risk converges to the global optimal risk \mathcal{L}_π^* .

Proof of Theorem A.1.26. We proceed by induction on the epoch. We restate the key inductive statements to prove: (with appropriate selection of η_k):

- $L[j, k]$: $\mathcal{L}_\pi(\Theta_j^k)$ stays bounded above by some uniform constant $C_L \geq \|\mathbf{Y}_\pi\|_F^2$ — this is the content of Corollary A.1.22.
- $R[j, k]$: $\mathcal{L}_\pi(\Theta_j^k)$ satisfies the per-epoch loss inequality — this is the content of Proposition A.1.21.
- $D[j, k]$: The approximate invariances stay bounded in norm away from 1. More precisely, $\|\mathbf{D}_j^k\|_2 \leq 2C_w^2 C_L \|\overline{\mathbf{X}}_\pi\|_F^2 \sum_{t=1}^{\infty} \eta_t^2 \leq \frac{1}{2}$ — this is the content of Corollary A.1.25.

Notice that the assumptions on η_k exactly satisfy the hypotheses of Corollaries A.1.22 and A.1.25 and Proposition A.1.23.

The base cases follows from the initialization. Recall that we set

$$C_L = \max \{ \mathcal{L}_\pi(\Theta_t^1) : 0 \leq t \leq m-1 \}.$$

Since we only look at the loss values for the first epoch, C_L is indeed an absolute constant depending on π . So $L[j, 1]$, as defined in Hypothesis 1, holds for all $j < m$. Next, we have from the initialization that $\mathbf{D}_0^1 = \mathbf{0}$. Then, applying Lemma A.1.24, we can conclude $D[j, 1]$ also holds for all $j < m$, so Hypothesis 3 holds. Finally, there is no need to check Hypothesis 2 because it is only defined for $k > 1$.

Now our inductive hypotheses are that $L[j, k]$, $R[j, k]$, $D[j, k]$ hold for all $j < m$. For the inductive step, we want to prove $L[0, k+1]$, $R[0, k+1]$, $D[0, k+1]$. By construction of η_k , the hypotheses of Corollary A.1.22 is satisfied, so $L[0, k+1]$ holds. Moreover, the hypotheses of Proposition A.1.21 are satisfied, so $R[0, k+1]$ holds. Finally, the hypotheses of Lemma A.1.24 are satisfied, so that $D[0, k+1]$ holds.

As asserted earlier, the above argument is robust up to reindexing if we want to prove the statement for $i > 0$. In particular, all of the results of the previous section go through, as we showed in Proposition A.1.20 that the stepsize noise is negligible. Hence the induction is completed for all (i, k) and so the unrolled update equation in Proposition A.1.23 holds for all k . This gives the formal rate of convergence for the stated stepsize. In particular, we see that the SS risk \mathcal{L}_π converges to its global minimum. \square

Proof of Theorem 3.2.2. All we need to do is convert the stepsize requirements. Examining the stepsize requirements in Theorem A.1.26, they depend on C_L , C_w , and $\|\overline{\mathbf{X}}_\pi\|_F$.

Now recall the definition of C_w in Corollary A.1.8:

$$C_w^2 \leq \frac{3}{2} + d^2 \left(\frac{1}{2} + \frac{C_L^{1/2} + \|\mathbf{Y}_\pi\|_F}{\sigma_{\min}(\overline{\mathbf{X}}_\pi^\top)} \right).$$

Hence $C_w = \text{poly}(d, C_L, \|\mathbf{Y}_\pi\|_F, 1/\sigma_{\min}(\overline{\mathbf{X}}_\pi^\top))$. Finally, since $\|\mathbf{X}_\pi\|_F^2 = dn$, $\sigma_{\min}(\overline{\mathbf{X}}_\pi \overline{\mathbf{X}}_\pi^\top) = \sigma_{\min}(\overline{\mathbf{X}}_\pi^\top)^2$, and C_L is an absolute constant, we can convert the stepsize requirements into

$$\eta_k = \frac{1}{k^\beta} \cdot \min \left\{ \frac{1}{2}, \frac{2}{\sigma_{\min}(\overline{\mathbf{X}}_\pi \overline{\mathbf{X}}_\pi^\top)}, \frac{\sqrt{2\beta-1} \text{poly}(\sigma_{\min}(\overline{\mathbf{X}}_\pi^\top))}{\text{poly}(n, d, \|\mathbf{Y}_\pi\|_F)} \right\},$$

which directly implies the stepsize requirements in Theorem 3.2.2. \square

A.2. Proof of convergence for RR

In this section, we prove Theorem 3.2.3. With RR, we randomly resample permutation $\pi_k \in \mathbb{S}_n$ on epoch k . Hence, it is natural to seek a convergence bound in expectation. We briefly comment on the complications that arise in this setting.

Since we want to prove convergence an expectation, an inductive approach that controls the approximate invariances and loss evolution is complicated by the necessity for bounds on these quantities that are stronger than merely being in expectation. This is precisely why we need Assumption 3.

Additional notation. As introduced in Section 2, we can view RR as optimizing the risk

$$\mathcal{L}_{\text{RR}}(\Theta) \triangleq \frac{1}{n!} \sum_{\pi \in \mathbb{S}_n} \mathcal{L}_\pi(\Theta)$$

via with-replacement SGD on an *epoch level* (i.e., \mathcal{L}_π is sampled uniformly with replacement at every epoch), albeit with noise terms due to the shuffling algorithm. Motivated by the setup in Section 2, we can also write \mathcal{L}_{RR} in an equivalent form to \mathcal{L}_{GD} as follows:

$$\begin{aligned} \mathcal{L}_{\text{RR}}(\Theta) &\triangleq \frac{1}{n!} \sum_{\pi \in \mathbb{S}_n} \mathcal{L}(f(\mathbf{X}_\pi; \Theta), \mathbf{Y}_\pi) = \frac{1}{n!} \mathcal{L}(\mathbf{W}\Gamma\text{BN}_{\text{RR}}(\mathbf{X}), \mathbf{Y}_{\text{RR}}), \\ \text{where } \text{BN}_{\text{RR}}(\mathbf{X}) &\triangleq [\overline{\mathbf{X}}_{\pi_1} \quad \dots \quad \dots \quad \overline{\mathbf{X}}_{\pi_{n!}}] \in \mathbb{R}^{d \times (n \cdot n!)}, \\ \mathbf{Y}_{\text{RR}} &\triangleq [\mathbf{Y}_{\pi_1} \quad \dots \quad \mathbf{Y}_{\pi_{n!}}] \in \mathbb{R}^{p \times (n \cdot n!)}. \end{aligned}$$

For notational convenience, we also write $\overline{\mathbf{X}}_{\text{RR}} \triangleq \text{BN}_{\text{RR}}(\mathbf{X})$, and reiterate that overlines indicate the presence of batch normalization.

Just as in the SS case, we will abuse notation and refer to the RR risk function as a function of $M = W\Gamma$ by writing $\mathcal{L}_{\text{RR}}(M) \triangleq \frac{1}{n!} \sum_{\pi \in \mathbb{S}_n} \mathcal{L}_\pi(M)$. Similarly, we will often refer to the gradient of the RR risk with respect to M as $\nabla_M \mathcal{L}_{\text{RR}}(M) \triangleq \frac{1}{n!} \sum_{\pi \in \mathbb{S}_n} \nabla_M \mathcal{L}_\pi(M)$. We will find it helpful to use the notation $\bar{X}_{\max,2} \triangleq \arg \max_{\pi \in \mathbb{S}_n} \|\bar{X}_\pi\|_2$. Similarly we will denote the maximum Frobenius norm batch normalized dataset by $\bar{X}_{\max,F} \triangleq \arg \max_{\pi \in \mathbb{S}_n} \|\bar{X}_\pi\|_F$. Furthermore, it follows from the unit variance constraint in the definition of BN that $\|\bar{X}_{\max,2}\|_2 \leq \|\bar{X}_{\max,F}\|_F \leq \sqrt{dn}$.

At a high level, the RR proof of convergence closely follows the SS proof of convergence. Indeed, most of the technical legwork has already been fleshed out in the SS case — most of the results port over immediately, taking care to replace π with π_k . However, we will be careful in accounting for where we need to deviate from the SS logic.

A.2.1. CHECKING OPTIMIZATION PROPERTIES

We first check the smoothness and strong convexity property with respect to the merged matrix M that we heavily relied on for the proof of convergence for SS.

Fact A.2.1 (Smoothness of RR). *Define*

$$G_{\text{RR}} \triangleq \frac{1}{n!} \sum_{\pi \in \mathbb{S}_n} G_\pi.$$

Then $\mathcal{L}_{\text{RR}}(M)$ is G_{RR} -smooth with respect to M .

Proof. The statement follows from combining Lemma A.1.2 with the fact that if f_i are G_i -smooth, then $\sum_{i=1}^n f_i$ is $\sum_{i=1}^n G_i$ -smooth. \square

As before, we cannot directly use a PL inequality on W or Γ ; we must instead bootstrap this from the strong convexity of the risk with respect to $M = W\Gamma$.

Fact A.2.2 (Strong convexity). *Suppose that Assumption 1(b) holds for some $\pi \in \mathbb{S}_n$. Then the loss function $\mathcal{L}_{\text{RR}}(M) = \frac{1}{n!} \sum_{\pi \in \mathbb{S}_n} \mathcal{L}_\pi(M)$ is α_{RR} -strongly convex with respect to M with $\alpha_{\text{RR}} \triangleq \frac{1}{n!} \sum_{\pi} \sigma_{\min}(\bar{X}_\pi \bar{X}_\pi^\top) = \frac{1}{n!} \sum_{\pi} \alpha_\pi$.*

Proof. Take the Hessian of \mathcal{L}_{RR} with respect to $\text{vec}(M)$ to obtain $\nabla_{\text{vec}(M)}^2 \mathcal{L}_{\text{RR}}(M) = \frac{1}{n!} \sum_{\pi \in \mathbb{S}_n} \bar{X}_\pi \bar{X}_\pi^\top \otimes I$. Hence we have $\nabla_{\text{vec}(M)}^2 \mathcal{L}_{\text{RR}}(M) \succeq \frac{1}{n!} \sum_{\pi} \sigma_{\min}(\bar{X}_\pi \bar{X}_\pi^\top)$. Since we assumed Assumption 1(b) holds, the sum is strictly positive, so \mathcal{L}_{RR} is indeed strongly convex. \square

A.2.2. PROOF SKETCH OF RR CONVERGENCE

We first state the modified inductive hypothesis for the one-epoch risk update, which replaces $R[j, k]$.

Hypothesis 4. *For $k > 1$, the inductive hypothesis $S[k]$ states that*

$$\begin{aligned} \mathbb{E}_{\pi_{k-1}}[\mathcal{L}_{\text{RR}}(M_0^k) | \mathcal{F}_{k-1}] - \mathcal{L}_{\text{RR}}^* &\leq \left(1 - \frac{\alpha_{\text{RR}} \eta_{k-1}}{2}\right) (\mathcal{L}_{\text{RR}}(M_0^{k-1}) - \mathcal{L}_{\text{RR}}^*) \\ &\quad + \eta_{k-1}^2 \text{poly}(m, A_w, A_L, \|\bar{X}_{\max,F}\|_F), \end{aligned}$$

where $\mathbb{E}_{\pi_{k-1}}$ denotes the expectation over random draws of the permutation π_{k-1} .

Proof sketch of Theorem 3.2.3. As before, we start by writing out the smoothness inequality with respect to M :

$$\mathcal{L}_{\text{RR}}(M_0^{k+1}) \leq \mathcal{L}_{\text{RR}}(M_0^k) + \langle \nabla_M \mathcal{L}_{\text{RR}}(M_0^k), M_0^{k+1} - M_0^k \rangle + \frac{G_{\text{RR}}}{2} \|M_0^{k+1} - M_0^k\|^2. \quad (34)$$

Next, we have the same gradient update

$$M_0^{k+1} = M_0^k - \eta_k \tilde{g}^k + \eta_k^2 \mathbf{r}^k,$$

but now all quantities involving π turn into π_k . For example, we redefine

$$\tilde{g}^k \triangleq \sum_{t=0}^{m-1} \tilde{g}_t^k = \nabla_W \mathcal{L}_{\pi_k}(M_0^k) \Gamma_0^k + W_0^k \nabla_\Gamma \mathcal{L}_{\pi_k}(M_0^k).$$

Since we want to prove convergence in expectation, it is standard to consider the natural filtration \mathcal{F}_k of which permutations we have picked up to (but not including) epoch k . Formally, $\mathcal{F}_k = \sigma(\pi_1, \dots, \pi_{k-1})$, where $\sigma(Z)$ denotes the σ -algebra generated by the random variable Z .

Noting the identity

$$\mathbb{E}_\pi[\nabla_M \mathcal{L}_\pi(\mathbf{M}_0^k)] = \nabla_M \mathcal{L}_{\text{RR}}(\mathbf{M}_0^k),$$

it follows that if we can apply Corollary A.1.15, then

$$\mathbb{E}_{\pi_k} [\langle \nabla_M \mathcal{L}_{\text{RR}}(\mathbf{M}_0^k), \tilde{\mathbf{g}}^k \rangle_F] = \langle \nabla_M \mathcal{L}_{\text{RR}}(\mathbf{M}_0^k), \mathbb{E}_{\pi_k}[\tilde{\mathbf{g}}^k] \rangle_F \geq \frac{1}{2} \|\nabla_M \mathcal{L}_{\text{RR}}(\mathbf{M}_0^k)\|_F^2.$$

Hence, we can follow the same argument for upper bounding the smoothness inequality for SS and take a conditional expectation over π_k conditioned on \mathcal{F}_k . Assuming for now that the weight norms are bounded by some absolute constant A_w and the relevant losses are bounded by an absolute constant A_L , this yields

$$\mathbb{E}_{\pi_k} [\mathcal{L}_{\text{RR}}(\mathbf{M}_0^{k+1}) | \mathcal{F}_k] \leq \mathcal{L}_{\text{RR}}(\mathbf{M}_0^k) - \frac{\eta k}{4} \|\nabla_M \mathcal{L}_{\text{RR}}(\mathbf{M}_0^k)\|_F^2 + \eta_k^2 \mathbb{E}_{\pi_k} [\text{poly}(m, A_w, A_L, \|\mathbf{X}_{\pi_k}\|_F)].$$

Noting that we can upper bound $\|\mathbf{X}_{\pi_k}\|_F$ uniformly by $\|\overline{\mathbf{X}}_{\max, F}\|_F$, which does not depend on k , this shows that we have

$$\mathbb{E}_{\pi_k} [\mathcal{L}_{\text{RR}}(\mathbf{M}_0^{k+1}) | \mathcal{F}_k] \leq \mathcal{L}_{\text{RR}}(\mathbf{M}_0^k) - \frac{\eta k}{4} \|\nabla_M \mathcal{L}_{\text{RR}}(\mathbf{M}_0^k)\|_F^2 + \eta_k^2 \text{poly}(m, A_w, A_L, \|\overline{\mathbf{X}}_{\max, F}\|_F).$$

Next, α_{RR} -strong convexity yields

$$\mathbb{E}_{\pi_k} [\mathcal{L}_{\text{RR}}(\mathbf{M}_0^{k+1}) | \mathcal{F}_k] - \mathcal{L}_{\text{RR}}^* \leq \left(1 - \frac{\alpha_{\text{RR}} \eta k}{2}\right) (\mathcal{L}_{\text{RR}}(\mathbf{M}_0^k) - \mathcal{L}_{\text{RR}}^*) + \eta_k^2 \text{poly}(m, A_w, A_L, \|\overline{\mathbf{X}}_{\max, F}\|_F).$$

This is exactly the statement of $S[k+1]$. To proceed, we must fill in the following details. First, we must show that the relevant losses are bounded by A_L — see Corollary A.2.6. Then, we show that the weight norms are bounded by A_w — this is shown in Corollary A.2.7. Finally, we must show that $D[j, k]$ holds, i.e., bound the approximate invariances — this is the content of Corollary A.2.9. Once we address these technicalities, an inductive argument similar to the SS version proves the theorem. Note that the SS inductive hypotheses $L[j, k]$ and $R[j, k]$ are not active in this proof. \square

A.2.3. WEIGHT BOUNDS

In this section we elucidate the connection between weight norms and the loss evolution. In particular, we show that bounds on the weight norms confer a bound on the loss function value. First, we state the following inequality, which follows from a quick application of Cauchy-Schwarz.

Fact A.2.3. *We have $\frac{1}{n!} \sum_{\pi \in \mathbb{S}_n} (\mathcal{L}_\pi(\mathbf{M}))^{1/2} \leq \mathcal{L}_{\text{RR}}(\mathbf{M})^{1/2}$.*

With this identity in hand, we derive the following corollaries about weight and gradient bounds.

Proposition A.2.4. *We have*

$$\|\nabla_M \mathcal{L}_{\text{RR}}(\mathbf{M})\| \leq \|\overline{\mathbf{X}}_{\max, 2}\|_2 \mathcal{L}_{\text{RR}}(\mathbf{M})^{1/2}.$$

Proof. We have $\|\nabla_M \mathcal{L}_{\text{RR}}(\mathbf{M})\|_2 \leq \frac{1}{n!} \sum_{\pi} \|\nabla_M \mathcal{L}_\pi(\mathbf{M})\|_2$ by the triangle inequality. Applying the individual gradient bounds in Lemma A.1.9 and uniformly bounding $\|\overline{\mathbf{X}}_\pi\|$ by $\|\overline{\mathbf{X}}_{\max, 2}\|$, we see that

$$\|\nabla_M \mathcal{L}_{\text{RR}}(\mathbf{M})\|_2 \leq \|\overline{\mathbf{X}}_{\max, 2}\|_2 \cdot \frac{1}{n!} \sum_{\pi} \mathcal{L}_\pi^{1/2}(\mathbf{M}).$$

Applying Fact A.2.3 yields the desired result. \square

As promised, we quantify the relationship between \mathcal{L}_{RR} , \mathcal{L}_π , and $\|\mathbf{M}\|_2$ with the following proposition.

Proposition A.2.5. Let $\sigma_0 \triangleq \frac{1}{n!} \sum_{\pi} \sigma_{\min}(\overline{\mathbf{X}}_{\pi}^{\top})$. We have

$$\frac{\mathcal{L}_{\text{RR}}(\mathbf{M}) - 2\|\mathbf{Y}\|_F^2}{2\|\overline{\mathbf{X}}_{\max, F}\|_F^2} \leq \|\mathbf{M}\|_2^2 \leq \left(\frac{\|\mathbf{Y}\|_F + \mathcal{L}_{\text{RR}}(\mathbf{M})^{1/2}}{\sigma_0} \right)^2.$$

Similarly, for any $\pi \in \mathbb{S}_n$, we have

$$\frac{\mathcal{L}_{\pi}(\mathbf{M}) - 2\|\mathbf{Y}\|_F^2}{2\|\overline{\mathbf{X}}_{\pi}\|_F^2} \leq \|\mathbf{M}\|_2^2 \leq \left(\frac{\|\mathbf{Y}\|_F + \mathcal{L}_{\pi}(\mathbf{M})^{1/2}}{\sigma_{\min}(\overline{\mathbf{X}}_{\pi}^{\top})} \right)^2.$$

Proof. First, we have

$$\begin{aligned} \frac{1}{n!} \sum_{\pi \in \mathbb{S}_n} \|\mathbf{M}\overline{\mathbf{X}}_{\pi}\|_2 &\leq \|\mathbf{Y}\|_F + \frac{1}{n!} \sum_{\pi} \|\mathbf{Y}_{\pi} - \mathbf{M}\overline{\mathbf{X}}_{\pi}\|_F \\ &= \|\mathbf{Y}\|_F + \frac{1}{n!} \sum_{\pi} \mathcal{L}_{\pi}(\mathbf{M})^{1/2}, \end{aligned}$$

where we have used in the first line the fact that $\|\mathbf{Y}_{\pi}\|_F = \|\mathbf{Y}\|_F$ for all π . Therefore by Fact A.2.3 and using $\|\mathbf{M}\overline{\mathbf{X}}_{\pi}\|_2 \geq \sigma_{\min}(\overline{\mathbf{X}}_{\pi}^{\top})\|\mathbf{M}\|_2$, we find that

$$\frac{1}{n!} \sum_{\pi} \sigma_{\min}(\overline{\mathbf{X}}_{\pi}^{\top})\|\mathbf{M}\|_2 \leq \|\mathbf{Y}\|_F + \mathcal{L}_{\text{RR}}(\mathbf{M})^{1/2}.$$

It follows that

$$\|\mathbf{M}\|_2 \leq \frac{\|\mathbf{Y}\|_F + \mathcal{L}_{\text{RR}}(\mathbf{M})^{1/2}}{\sigma_0}.$$

For the other direction, note that

$$\mathcal{L}_{\pi}(\mathbf{M}) = \|\mathbf{Y}_{\pi} - \mathbf{M}\overline{\mathbf{X}}_{\pi}\|_F^2 \leq 2\|\mathbf{Y}\|_F^2 + 2\|\mathbf{M}\|_2^2\|\overline{\mathbf{X}}_{\pi}\|_F^2.$$

Averaging over all $\pi \in \mathbb{S}_n$ gives us

$$\mathcal{L}_{\text{RR}}(\mathbf{M}) \leq 2\|\mathbf{Y}\|_F^2 + \frac{2\|\mathbf{M}\|_2^2}{n!} \sum_{\pi} \|\overline{\mathbf{X}}_{\pi}\|_F^2.$$

Uniformly bounding $\|\overline{\mathbf{X}}_{\pi}\|_F^2$ by $\|\overline{\mathbf{X}}_{\max, F}\|_F^2$ and rearranging yields the desired result.

The set of inequalities for π also follow by a similar argument. \square

As a corollary of Proposition A.2.5, it follows from Assumption 3 that each of the losses $\mathcal{L}_{\pi}(\mathbf{M}_i^k)$ stay bounded by an absolute constant throughout training.

Corollary A.2.6 (Uniform bound on SS losses). *Under Assumption 3, for every $\pi \in \mathbb{S}_n$ we have*

$$\mathcal{L}_{\pi}(\mathbf{M}_i^k) \leq A_L,$$

where

$$A_L \triangleq 2\|\mathbf{Y}\|_F^2 + 2A_{\text{RR}}^2\|\overline{\mathbf{X}}_{\max, F}\|_F^2.$$

Here, A_{RR} was previously defined in Assumption 3.

Finally, Assumption 3 implies the following inductive statement about the weight norms.

Corollary A.2.7 (Uniform bound on weight norms). *Assume Assumption 3 and $D[j, k]$ holds. Then for RR, we have*

$$\max \{ \|\mathbf{W}_i^k\|_2, \|\mathbf{\Gamma}_i^k\|_2 \} \leq A_w,$$

where

$$A_w^2 \triangleq \frac{3}{2} + d^2 \left(\frac{1}{2} + A_{\text{RR}} \right).$$

Here, A_{RR} was previously defined in Assumption 3.

A.2.4. BOUNDING APPROXIMATE INVARIANCES

In this section we formally bound the approximate invariances throughout RR training under Assumption 3. In particular, as a consequence of Corollaries A.2.6 and A.2.7, the following two claims follow almost directly upon inspection of the proofs of Lemma A.1.24 and Corollary A.1.25.

Lemma A.2.8. *Suppose $D[j, k]$ holds for some $j < m$. We have*

$$\|\mathbf{D}_{j+1}^k - \mathbf{D}_j^k\|_2 \leq 2A_w^2 A_L \|\text{BN}(\mathbf{X}_{\pi_k}^{j+1})\|_2^2 \eta_k^2.$$

Hence, if $D[t, k]$ holds for all $t \leq j$, we also have

$$\|\mathbf{D}_{j+1}^k\|_2 \leq 2A_w^2 A_L \|\bar{\mathbf{X}}_{\max, F}\|_F^2 \sum_{t=1}^k \eta_t^2.$$

Corollary A.2.9. *Assume that $D[t, k]$ holds for all $t \leq j < m$ and $\eta_k = \frac{c}{k^\beta}$ for $1/2 < \beta < 1$. If*

$$c^2 \leq \frac{1}{4(1 + \frac{1}{2\beta-1})A_w^2 A_L \|\bar{\mathbf{X}}_{\max, F}\|_F^2},$$

then

$$\|\mathbf{D}_{j+1}^k\|_2 \leq 2A_w^2 A_L \|\bar{\mathbf{X}}_{\max, F}\|_F^2 \sum_{t=1}^k \eta_t^2 \leq \frac{1}{2}.$$

In other words, $D[j+1, k]$ holds.

A.2.5. COMPLETING THE PROOF

With the connection between the RR loss function and weight bounds in hand, we can complete the proof of convergence for RR. Finally, we formally state the RR convergence result.

Theorem A.2.10 (Formal statement of convergence for RR). *Suppose Assumption 1(b) and Assumption 3 hold. Let $f(\cdot; \Theta) = \text{WTBN}(\cdot)$ be a 2-layer linear+BN network initialized at $\Theta_0^1 = (\mathbf{W}_0^1, \mathbf{\Gamma}_0^1) = (\mathbf{0}, \mathbf{I})$. Consider optimizing f using RR with decreasing stepsize*

$$\eta_k = \frac{1}{k^\beta} \cdot \min \left\{ \frac{1}{2}, \frac{2}{\alpha_{\text{RR}}}, \sqrt{\frac{1}{4(1 + \frac{1}{2\beta-1})A_w^2 A_L \|\bar{\mathbf{X}}_{\max, F}\|_F^2}} \right\}$$

for any $1/2 < \beta < 1$. Then the RR risk satisfies

$$\mathbb{E}[\mathcal{L}_{\text{RR}}(\Theta_0^{k+1})] - \mathcal{L}_{\text{RR}}^* \leq \frac{\text{poly}(n, d, A_{\text{RR}}, \|\mathbf{Y}\|_F) \log k}{k^\beta}.$$

In other words, the RR risk converges to the global optimal risk $\mathcal{L}_{\text{RR}}^*$.

Proof of Theorem A.2.10. We inductively prove that $D[j, k]$ and $S[k]$ hold. There is no need to check $S[1]$ because Hypothesis 4 is defined for $k > 1$. The proof that the base cases $D[j, 0]$ all hold follows the same proof as that in the SS case.

Now suppose for the sake of induction that $D[j, k]$ and $S[k]$ hold. We will show that $D[j+1, k]$ holds. Once we show that $D[j, k]$ holds for all $j < m$, we can then show that $S[k+1]$ holds.

In particular, by the assumption on η_k , Corollary A.2.9 implies that $D[j+1, k]$ holds. Hence by induction $D[j, k]$ holds for all $j < m$. We see that assuming Assumption 3 simplified the proof strategy significantly, as we did not have to go through the trouble of proving $L[j, k]$.

Next, let's understand what happens to the per-epoch loss bound in $S[k+1]$. Explicitly, we can follow the same steps as in the proof sketch — which only required $\eta_k \leq \frac{c}{k^\beta}$ where $c \leq \min \left\{ \frac{1}{2}, \frac{2}{\alpha_{\text{RR}}} \right\}$ — to see that

$$\mathbb{E}_{\pi_k}[\mathcal{L}_{\pi_k}(\mathbf{M}_0^{k+1}) | \mathcal{F}_k] - \mathcal{L}_\pi^* \leq \left(1 - \frac{\alpha_\pi \eta_k}{2}\right) (\mathcal{L}_{\text{RR}}(\mathbf{M}_0^k) - \mathcal{L}_{\text{RR}}^*) + \text{poly}(m, A_w, A_L, \|\bar{\mathbf{X}}_{\max, F}\|_F) \eta_k^2.$$

Indeed, under Assumption 3 and $D[j, k]$ for $j < m$, we can apply Corollaries A.2.6 and A.2.7 to rigorously bound all of the analogous noise terms a.s.. We can then follow the same argument as in Proposition A.1.23 to unroll the recurrence, using iterated expectation to obtain a total expectation in the end. We can thus conclude that with $\eta_k = \frac{c}{k^\beta}$ for $1/2 < \beta < 1$ and the constant c chosen as the theorem statement that

$$\begin{aligned} \mathbb{E}[\mathcal{L}_{\text{RR}}(\Theta_0^{k+1})] - \mathcal{L}_{\text{RR}}^* &\leq (\mathcal{L}_{\text{RR}}(\Theta_0^1) - \mathcal{L}_\pi^*) \exp\left(\frac{c\alpha_{\text{RR}}}{2(1-\beta)}(2 - k^{1-\beta})\right) \\ &\quad + \frac{\text{poly}(m, A_w, A_L, \|\overline{\mathbf{X}}_{\max, F}\|_F) \log k}{k^\beta} \\ &\leq (\mathcal{L}_{\text{RR}}(\Theta_0^1) - \mathcal{L}_\pi^*) \exp\left(\frac{c\alpha_{\text{RR}}}{2(1-\beta)}(2 - k^{1-\beta})\right) \\ &\quad + \frac{\text{poly}(n, d, A_{\text{RR}}, \|\mathbf{Y}\|_F) \log k}{k^\beta}. \end{aligned}$$

In the last line we used the fact that $A_L = \text{poly}(A_{\text{RR}}, \|\mathbf{Y}\|_F, \|\overline{\mathbf{X}}_{\max, F}\|_F)$ from Corollary A.2.6, $A_w = \text{poly}(d, A_{\text{RR}})$ from Corollary A.2.7, and $\|\overline{\mathbf{X}}_{\max, F}\|_F^2 = dn$. The desired claim immediately follows.

We can also immediately see how the stepsize requirements match that of Theorem 3.2.3. \square

B. Proofs for classification results

In this section we lay out the groundwork for formally proving our main results Theorems B.2.1 and B.3.1 about the separability decomposition of SS+BN and RR+BN (cf. Theorems 4.1.3 and 4.1.4). At a high level, we show that the separability decomposition is closely linked to the presence of monochromatic batches (Lemma B.1.1) and the dimensionality of the batch normalized dataset (Lemma B.1.2). In Appendix B.4, we formally characterize the optimal directions of linear+BN classifiers. We defer the proofs of the more technical lemmas to Appendix C.

ADDITIONAL NOTATION AND SETUP

We lay out some additional notation that will aid in our discussion of classification. Division of two vectors should be interpreted in a coordinatewise fashion, so $\frac{\mu}{\sigma} \in \mathbb{R}^d$ with k th coordinate μ_k/σ_k . For a matrix $\mathbf{A} \in \mathbb{R}^{d \times n}$, we define $\|\mathbf{A}\|_{2, \infty} = \max_{i \in [n]} \|\mathbf{A}_{:, i}\|_2$, i.e., the maximum Euclidean norm of the columns.

We remind the reader of some notation introduced previously, with an important redefinition for $\overline{\mathbf{X}}_{\text{RR}}$. For a dataset $\mathbf{Z} = (\mathbf{X}, \mathbf{Y})$, we write $\mathbf{Z}^+ = (\mathbf{X}^+, \mathbf{Y}^+)$ and $\mathbf{Z}^- = (\mathbf{X}^-, \mathbf{Y}^-)$ to denote the positive and negative examples, respectively. Recall that we write the dataset batch normalized under a permutation π as $\overline{\mathbf{Z}}_\pi \triangleq (\overline{\mathbf{X}}_\pi, \mathbf{Y}_\pi)$, where $\overline{\mathbf{X}}_\pi \triangleq \text{BN}_\pi(\mathbf{X}) \in \mathbb{R}^{d \times n}$ and $\mathbf{Y}_\pi = \pi \circ \mathbf{Y}$. We also use $\overline{\mathbf{X}}_\pi^+$ and $\overline{\mathbf{X}}_\pi^-$ to denote the submatrices of $\overline{\mathbf{X}}_\pi$ containing its columns corresponding to the positive and negative examples, respectively.

For the sake of analyzing the rank of $\overline{\mathbf{X}}_{\text{RR}}$, we will redefine it as follows by throwing out redundant batches. Let $\binom{[n]}{B}$ denote the set of all $\binom{[n]}{B}$ unique (up to permutation) batches of size B that can be created from choosing the columns of $\mathbf{X} \in \mathbb{R}^{d \times n}$. Fix an arbitrary labelling of these $\binom{[n]}{B}$ batches, and let $\mathbf{B}^j \in \mathbb{R}^{d \times B}$ refer to the j th such batch. Then

$$\overline{\mathbf{X}}_{\text{RR}} \triangleq \text{BN}_{\text{RR}}(\mathbf{X}) \triangleq \left[\text{BN}(\mathbf{B}^1) \quad \dots \quad \text{BN}(\mathbf{B}^{\binom{[n]}{B}}) \right] \in \mathbb{R}^{d \times B \binom{[n]}{B}}$$

Note that the rank of $\overline{\mathbf{X}}_{\text{RR}}$ is the same as the rank of the original definition, since all we did was throw out redundant batches for the purposes of analyzing the rank.

We now turn to laying down some of the background necessary to introduce our technical results. As a motivating step, recall Proposition 4.1.2. It states that SS with permutation π can cause divergence of the GD risk if $\overline{\mathbf{Z}}_\pi$ is PLS or LS, but not if it is SC. Hence, determining sufficient conditions for when $\overline{\mathbf{Z}}_\pi$ is SC is of primary interest. Intuitively, $\overline{\mathbf{Z}}_\pi$ being SC should be related to some notion of genericity — the convex hulls of positive and negative features should be full dimensional. To formalize this intuition, we take a quick detour and recall several standard notions in convex analysis, defined for example in Boyd et al. (2004).

For $S \subseteq \mathbb{R}^d$, its interior $\text{int}(S)$ denotes the largest open set contained in S . We say that S is *affine* if for any $\mathbf{x}_1, \mathbf{x}_2 \in S$, the

line $\lambda \mathbf{x}_1 + (1-\lambda)\mathbf{x}_2 \subseteq S$. An *affine combination* of k points $\mathbf{x}_1, \dots, \mathbf{x}_k \in \mathbb{R}^d$ is given by $\sum_{i=1}^k \lambda_i \mathbf{x}_i$ where $\sum_{i=1}^k \lambda_i = 1$.⁶ The affine hull of S is the set of all affine combinations of S , and is denoted by $\text{aff}(S)$, and clearly $\text{aff}(S)$ is an affine set. Similarly, the relative interior $\text{relint}(S)$ denotes the largest open subset of $\text{aff}(S)$ contained in S . For a matrix \mathbf{A} , we slightly abuse notation and write $\text{aff}(\mathbf{A})$ to denote the affine hull of its columns. Similarly, we use $\text{conv}(\mathbf{A})$ to denote the convex hull of its columns.

Let $\mathbf{x}_0 \in \text{aff}(S)$ be any element of the affine hull of S . It is not hard to see that $\text{aff}(S) = \mathbf{x}_0 + V$, where V is a linear subspace of \mathbb{R}^d . Furthermore, this V is uniquely determined by S . One can think of \mathbf{x}_0 as an offset and V as the space of valid directions to move in to stay in $\text{aff}(S)$. We define $\dim(S) \triangleq \dim(V)$.

Definition 3. A set $S \subseteq \mathbb{R}^d$ is called *full dimensional* if $\dim(S) = d$. This definition is equivalent to saying that $\text{int}(\text{conv}(S))$ is nonempty. Similarly, for any matrix $\mathbf{A} \in \mathbb{R}^{d \times n}$, we say \mathbf{A} is *full dimensional* if the set of its columns is full dimensional.

This formal definition of full dimensional allows us to identify sufficient conditions for a dataset $\mathbf{Z} = (\mathbf{X}, \mathbf{Y})$ to be SC.

B.1. Preliminary results on separability decomposition

In this section, we introduce the technical results that help us analyze the separation between shuffling SGD and GD. A unifying theme is to understand the effect of monochromatic batches and rank on the separability decomposition — and thus, divergence. In particular, we show that having monochromatic batches and being full-rank prevents divergence in the underparameterized regime. In a later section (Appendix B.4), we also show that these properties also significantly influence the optimal directions under the logistic loss.

The following lemma formalizes how monochromatic batches affect the separability decomposition.

Lemma B.1.1. *Given a permutation π , suppose there are two batches \mathbf{Z}_π^1 and \mathbf{Z}_π^2 such that \mathbf{Z}_π^1 consists entirely of positive examples and \mathbf{Z}_π^2 consists entirely of negative examples. Then, if we consider the resulting $\bar{\mathbf{Z}}_\pi = (\bar{\mathbf{X}}_\pi, \mathbf{Y}_\pi)$, the submatrices $\bar{\mathbf{X}}_\pi^+$ and $\bar{\mathbf{X}}_\pi^-$ of $\bar{\mathbf{X}}_\pi$ satisfy*

$$\text{relint}(\text{conv}(\bar{\mathbf{X}}_\pi^+)) \cap \text{relint}(\text{conv}(\bar{\mathbf{X}}_\pi^-)) \neq \emptyset.$$

Consequently, $\bar{\mathbf{Z}}_\pi$ is not LS.

Proof. Batch normalization ensures that the batch normalized features $\bar{\mathbf{X}}_\pi^1$ and $\bar{\mathbf{X}}_\pi^2$ are mean-zero. But this implies that $\mathbf{0}$ is in the convex hulls of each batch, which implies that $\text{conv}(\bar{\mathbf{X}}_\pi^1)$ intersects $\text{conv}(\bar{\mathbf{X}}_\pi^2)$. In fact, $\mathbf{0}$ is in the intersection of their relative interiors as well. To see this, we prove that the mean $\boldsymbol{\mu}$ of a batch $\mathbf{B} = \{\mathbf{x}_1, \dots, \mathbf{x}_B\}$ lies in the relative interior of $\text{conv}(\mathbf{B})$.

This can be shown via an inductive argument on the batch size. If $B = 2$, then $\boldsymbol{\mu}$ is the midpoint between \mathbf{x}_1 and \mathbf{x}_2 , which is in the relative interior of $\text{conv}(\mathbf{B})$. Now assume it's true for all possible batches of size B . When we add \mathbf{x}_{B+1} , we get a new batch \mathbf{B}' , with mean $\boldsymbol{\mu}' = \frac{B}{B+1}\boldsymbol{\mu} + \frac{1}{B+1}\mathbf{x}_{B+1}$. Hence if $\mathbf{x}_{B+1} \in \text{conv}(\mathbf{B})$, clearly $\boldsymbol{\mu}' \in \text{relint}(\text{conv}(\mathbf{B}))$ by convexity. If $\mathbf{x}_{B+1} \notin \text{conv}(\mathbf{B})$, then \mathbf{x}_{B+1} is one of the vertices of $\text{conv}(\mathbf{B}')$. Since $\boldsymbol{\mu} \in \text{relint}(\text{conv}(\mathbf{B}))$ and convexity, the segment between $\boldsymbol{\mu}$ and \mathbf{x}_{B+1} must stay in the relative interior of $\text{conv}(\mathbf{B}')$ except at the endpoints. Since $\boldsymbol{\mu}'$ is in the relative interior of this segment, the conclusion follows.

Hence $\text{relint}(\text{conv}(\bar{\mathbf{X}}_\pi^1))$ intersects $\text{relint}(\text{conv}(\bar{\mathbf{X}}_\pi^2))$. Since we have $\bar{\mathbf{X}}_\pi^+ \supseteq \bar{\mathbf{X}}_\pi^1$ and $\bar{\mathbf{X}}_\pi^- \supseteq \bar{\mathbf{X}}_\pi^2$, the relative interiors of the larger convex hulls intersect as well.

Finally, suppose $\bar{\mathbf{Z}}_\pi$ was LS. By definition there must exist a strict separating hyperplane for the two hulls. But the hulls intersect, so this is a contradiction. We conclude that $\bar{\mathbf{Z}}_\pi$ is not LS. \square

As Lemma B.1.1 establishes, monochromatic batches lead to $\bar{\mathbf{Z}}_\pi$ being PLS or SC. The following lemma synthesizes nicely with the above result; it demonstrates that if \mathbf{X} is full dimensional, and the relative interiors of the convex hulls of \mathbf{X}^+ and \mathbf{X}^- intersect, then \mathbf{Z} is SC.

Lemma B.1.2. *Let $\mathbf{Z} = (\mathbf{X}, \mathbf{Y})$ be such that*

⁶Here, unlike the definition of a convex combination, the λ_i 's are allowed to be negative.

- (1) \mathbf{X} is full dimensional.
- (2) $\text{relint}(\text{conv}(\mathbf{X}^+))$ intersects $\text{relint}(\text{conv}(\mathbf{X}^-))$.

Then \mathbf{Z} is SC.

Proof. Consider any halfspace in \mathbb{R}^d . First, since the relative interiors of the hull of positive points and negative points intersect, then so too do the hulls themselves. So there is no hyperplane that separates \mathbf{X}^+ from \mathbf{X}^- , i.e. \mathbf{Z} is either PLS or SC.

Suppose that \mathbf{Z} is PLS, i.e., there exists a hyperplane \mathbf{v} such that $y_i \mathbf{v}^\top \mathbf{x}_i \geq 0$ for every $(\mathbf{x}_i, y_i) \in \mathbf{Z}$. Since \mathbf{X} is full dimensional, the hyperplane orthogonal to \mathbf{v} cannot pass through every point of \mathbf{X} — hyperplanes are affine subspaces of dimension at most $d - 1$. So $y_i \mathbf{v}^\top \mathbf{x}_i > 0$ for some i and also $y_j \mathbf{v}^\top \mathbf{x}_j = 0$ for some $i \neq j$; otherwise, \mathbf{Z} is LS, which is a contradiction. Hence $\text{conv}(\mathbf{X}^+)$ and $\text{conv}(\mathbf{X}^-)$ touch only at the hyperplane defined by \mathbf{v} , which contradicts the assumption that the relative interiors intersect. Hence \mathbf{Z} is SC. \square

Lemmas B.1.1 and B.1.2 taken together show that to identify sufficient conditions for $\overline{\mathbf{Z}}_\pi$ to be SC, one should look for conditions under which $\overline{\mathbf{X}}_\pi$ is full dimensional and monochromatic batches are present. We already answered the former question in the main text with Proposition 3.2.1, which we restate for reference. Its proof is deferred to Appendix C.3.

Proposition B.1.3. *Assume that the original features $\mathbf{X} \in \mathbb{R}^{d \times n}$ satisfies Assumption 2 and $B > 2$. Then if we batch normalize and remove one datapoint from each normalized batch, to form a $d \times (B - 1) \frac{n}{B}$ matrix in the SS case and a $d \times (B - 1) \binom{n}{B}$ matrix in the RR case, the dataset is full-rank almost surely, regardless of which datapoints we remove. In particular, we have $\text{rank}(\overline{\mathbf{X}}_\pi) = \min \{d, (B - 1) \frac{n}{B}\}$ and $\text{rank}(\overline{\mathbf{X}}_{\text{RR}}) = \min \{d, (B - 1) \binom{n}{B}\}$ almost surely.*

Let us now consider the other question about the presence of monochromatic batches. Intuitively, under Assumption 4(a), there should be many monochromatic batches w.h.p. as long as B is small. The following lemma formalizes this intuition; its proof is contained in Appendix C.1.

Lemma B.1.4. *Assume Assumption 4(a). If $B = o(\log n)$ then there are $\Omega(n)$ monochromatic batches w.h.p.. If $B = \Omega(\log n)$, then there are no monochromatic batches w.h.p..*

The upshot of Lemmas B.1.1, B.1.2 and B.1.4 and Proposition B.1.3 is that small batch sizes naturally prevent divergence. However, there is a natural tradeoff here: small batch sizes also entail significant variance in the batch statistics, so they can lead to large (but non-diverging) values of the GD risk anyway when we train the network with SS.

Remark B.1.5 (Multiclass classification). *In the multiclass case with $K > 2$ different classes, one can directly generalize the above analysis to look at all $\binom{K}{2}$ pairwise combinations of classes. Lemma B.1.1 generalizes by requiring the existence of a monochromatic batch for each class. Hence as K increases the batch size must shrink to ensure that $\overline{\mathbf{Z}}_\pi$ is SC w.h.p., opening up a wider range of batch sizes for SS divergence.*

Finally, we formally define a robust notion of the separability decomposition that will prove helpful for quantifying the effects of increasing the batch size. To do so, we rely on the notion of the margin of a linearly separable dataset and the so-called penetration depth of overlapping convex hulls.

Definition 4 (Margin and penetration depth). *Let $\mathbf{Z} = (\mathbf{X}, \mathbf{Y})$ be a dataset, and let $\mathbf{X}^+, \mathbf{X}^-$ denote the positive and negative features, respectively. If \mathbf{Z} is linearly separable, then the margin of \mathbf{Z} is defined to be the ℓ_2 margin corresponding to the maximum margin classifier for \mathbf{Z} .*

If \mathbf{Z} is not linearly separable, then by definition $\text{conv}(\mathbf{X}^+)$ intersects $\text{conv}(\mathbf{X}^-)$. The penetration depth (Agarwal et al., 2000) of \mathbf{Z} is defined as the smallest Euclidean distance of translation of $\text{conv}(\mathbf{X}^+)$ such that the resulting convex body still intersects $\text{conv}(\mathbf{X}^-)$. In words, this quantifies the smallest perturbation we need to make \mathbf{Z} linearly separable.

Definition 5 (γ -robust separability decomposition). *Let $\mathbf{Z} = (\mathbf{X}, \mathbf{Y})$ be a dataset and $\mathbf{S}_{\text{GD}}^{\text{LS}} \sqcup \mathbf{S}_{\text{GD}}^{\text{SC}}$ be the separability decomposition of $\overline{\mathbf{Z}}_{\text{GD}} \triangleq (\overline{\mathbf{X}}_{\text{GD}}, \mathbf{Y}) \triangleq (\text{BN}(\mathbf{X}), \mathbf{Y})$.*

For $\gamma > 0$, we say that \mathbf{Z} is γ -robust if the following conditions hold.

- (1) One of the following:

- (a) $\bar{\mathbf{Z}}_{\text{GD}}$ is LS and $\mathbf{S}_{\text{GD}}^{\text{LS}}$ has margin at least γ ; or,
 (b) $\bar{\mathbf{Z}}_{\text{GD}}$ is SC and $\mathbf{S}_{\text{GD}}^{\text{SC}}$ has penetration depth at least γ .
- (2) Let σ_k , a_k , and b_k denote the standard deviation, min, and max of the k th feature of \mathbf{X} . Then $\min_{k \in [d]} \frac{\sigma_k}{b_k - a_k} = \Omega(1)$.
- (3) $\|\bar{\mathbf{X}}_{\text{GD}}\|_{2,\infty} = O(\sqrt{d})$, i.e., the maximum Euclidean norm of datapoints in $\text{BN}(\mathbf{X})$ is $O(\sqrt{d})$.

Definition 5 formalizes the informal statement in Theorem 4.1.3 that “ $\bar{\mathbf{Z}}_{\text{GD}}$ ’s separability decomposition can change with small perturbations.” If a dataset is robust, its separability decomposition cannot change easily by small perturbations (e.g., due to batch normalization) on datapoints. Notice here that PLS datasets can never be robust; to see why, note that \mathbf{Z} is PLS whenever $\bar{\mathbf{Z}}_{\text{GD}}$ is PLS. It follows that if \mathbf{Z} is PLS then it is not γ -robust for any $\gamma > 0$.

This definition of robustness is also natural from the perspective of concentration, since it provides an immediate link between concentration of batch statistics and the separability decomposition of $\bar{\mathbf{Z}}_{\pi}$. In order to estimate $\bar{\mathbf{X}}_{\text{GD}}$ well, we need to estimate σ to within a multiplicative factor, which explains (2). Moreover, the degree to which the SS datapoints concentrate around the GD datapoints also depend multiplicatively on the size of GD datapoints, which explains (3).

The following proposition provides a sufficient condition for the datapoints concentrating within a distance γ of the GD datapoints (cf. the definition of γ -robustness); its proof is deferred to Appendix C.2.

Proposition B.1.6. *Suppose \mathbf{Z} is γ -robust and $B = \Omega(d \log(nd)/\gamma^2)$. Then with probability at least $1 - 1/\text{poly}(n, d)$ over the choice of π , we have $\|\bar{\mathbf{X}}_{\pi} - \pi \circ \bar{\mathbf{X}}_{\text{GD}}\|_{2,\infty} = O(\gamma)$.*

With this setup in hand, we are ready to present our main theorems characterizing the separability decompositions of SS+BN and RR+BN.

B.2. Separability decomposition for SS

Theorem B.2.1 (Separability decomposition for SS, formal). *Throughout this theorem, assume that $B > 2$ and Assumption 4(a) and Assumption 2 hold.*

Suppose that $d \leq (B - 1)\frac{n}{B}$. Then the following hold.

- (1) $\bar{\mathbf{Z}}_{\pi}$ is SC w.h.p. for $B = o(\log n)$. If we relax Assumption 2, then $\bar{\mathbf{Z}}_{\pi}$ can be PLS as well.
 (2) Suppose further that \mathbf{Z} is γ -robust. Then $\bar{\mathbf{Z}}_{\pi}$ has the same separability decomposition as $\bar{\mathbf{Z}}_{\text{GD}}$ w.h.p. for $B = \Omega(d \log(nd)/\gamma^2)$.

Now suppose $d > (B - 1)\frac{n}{B}$. Then the following hold.

- (3) $\bar{\mathbf{Z}}_{\pi}$ is PLS w.h.p. for $B = o(\log n)$.
 (4) $\bar{\mathbf{Z}}_{\pi}$ is LS w.h.p. for $B = \Omega(\log n)$.

Proof. We consider cases based on whether the batch size is large or small.

Small batch size. By Lemma B.1.4, when $B = o(\log n)$, w.h.p. we get a batch comprised solely of positive examples and a batch comprised solely of negative examples. By Lemma B.1.1, this implies that $\bar{\mathbf{Z}}_{\pi}$ is PLS or SC, and the relative interior of the positive features intersects the relative interior of the negative features. By Lemma B.1.2, if $\bar{\mathbf{X}}_{\pi}$ is full dimensional, then $\bar{\mathbf{Z}}_{\pi}$ is SC. Hence it suffices to analyze the rank of $\bar{\mathbf{X}}_{\pi}$.

The maximal rank of $\bar{\mathbf{X}}_{\pi}$ is $\min\{d, (B - 1)\frac{n}{B}\}$ because of the mean zero constraint enforced in every batch. Because we assumed that Assumption 2 holds, Proposition 3.2.1 implies that $\bar{\mathbf{X}}_{\pi}$ achieves this upper bound almost surely. Putting it all together, we conclude that when $d \leq (B - 1)\frac{n}{B}$, $\bar{\mathbf{X}}_{\pi}$ is full-dimensional. It follows that $\bar{\mathbf{Z}}_{\pi}$ is SC w.h.p. over the choice of π , which proves (1).

On the other hand, when $d > (B - 1)\frac{n}{B}$, there always exists a hyperplane that passes through all of the monochromatic batches of $\bar{\mathbf{Z}}_{\pi}$ and perfectly classifies non-monochromatic batches (see Appendix B.4). It follows that $\bar{\mathbf{Z}}_{\pi}$ is PLS, which proves (3).

Large batch size. Now consider when $B = \Omega(\log n)$. In this regime, Lemma B.1.4 implies there are no monochromatic batches with high probability. Moreover, Proposition 3.2.1 implies that the features are full rank. It thus follows that when $d > (B - 1)\frac{n}{B}$, $\bar{\mathbf{Z}}_\pi$ is LS, which proves (4). Let us now consider the case $d \leq (B - 1)\frac{n}{B}$.

When \mathbf{Z} is γ -robust, we can directly apply Proposition B.1.6 to prove (2). Indeed, for $B = \Omega(d \log(nd)/\gamma^2)$, each SS datapoint is within distance $O(\gamma)$ of the corresponding GD datapoint with probability at least $1 - 1/\text{poly}(n, d)$. By increasing the batch size by at most a constant factor, we can guarantee that each SS datapoint is in fact within distance $\gamma/3$ of the corresponding GD datapoint.

If $\bar{\mathbf{Z}}_{\text{GD}}$ is LS, γ -robustness implies that the max-margin hyperplane for $\bar{\mathbf{Z}}_{\text{GD}}$ has margin at least γ . Since each SS datapoint moves at most $\gamma/3$ from the corresponding GD datapoint, this hyperplane still has margin at least $2\gamma/3$, implying that $\bar{\mathbf{Z}}_\pi$ is LS. If $\bar{\mathbf{Z}}_{\text{GD}}$ is SC, then γ -robustness implies that we need to translate the convex hulls of positive and negative points by at least γ to separate them. But we can only translate them by a total of $\gamma/3 + \gamma/3 = 2\gamma/3 < \gamma$, so the hulls stay strictly overlapping. This implies $\bar{\mathbf{Z}}_\pi$ is also SC. This concludes the proof of (2). \square

Remark B.2.2. If $\bar{\mathbf{Z}}_{\text{GD}}$ is not γ -robust for any $\gamma > 0$ (for example if $\bar{\mathbf{Z}}_{\text{GD}}$ is PLS), then it is not hard to construct examples where the separability decomposition of $\bar{\mathbf{Z}}_\pi$ is LS or PLS, and SS diverges. For a good example of this scenario, see Figures 4b and 4d. Even if $\bar{\mathbf{Z}}_{\text{GD}}$ is γ -robust, if $B = o(d \log(nd)/\gamma^2)$, then concentration can fail to hold in the worst case, and we can find analogous constructions where SS diverges.

B.3. Separability decomposition for RR

Theorem B.3.1 (Separability decomposition for RR, formal). *Suppose that $B > 2$ and Assumption 4(b) and Assumption 2 hold.*

If $d \leq (B - 1)\binom{n}{B}$, then $\bar{\mathbf{Z}}_{\text{RR}}$ is SC and $\bar{\mathbf{X}}_{\text{RR}}$ is full-rank almost surely, regardless of the separability decomposition of $\bar{\mathbf{Z}}_{\text{GD}}$.

Otherwise, if $d > (B - 1)\binom{n}{B}$, then $\bar{\mathbf{Z}}_{\text{RR}}$ is deterministically PLS, regardless of the separability decomposition of $\bar{\mathbf{Z}}_{\text{GD}}$.

Proof. When there are at least B positive and B negative examples, there exists a batch of all positive and a batch of all negative examples. Hence by Lemma B.1.1, $\bar{\mathbf{Z}}_{\text{RR}}$ is PLS or SC. By Proposition B.1.3, under Assumption 2, $\bar{\mathbf{X}}_{\text{RR}}$ attains the maximal rank of $\min\{d, (B - 1)\binom{n}{B}\}$ almost surely. So by Lemma B.1.2, if $d \leq (B - 1)\binom{n}{B}$, $\bar{\mathbf{Z}}_{\text{RR}}$ will be SC and $\bar{\mathbf{X}}_{\text{RR}}$ is full-rank almost surely. On the other hand, if $d > (B - 1)\binom{n}{B}$, then $\bar{\mathbf{Z}}_{\text{RR}}$ is PLS deterministically because there always exists a hyperplane that passes through all of the monochromatic batches of $\bar{\mathbf{Z}}_{\text{RR}}$ and perfectly classifies non-monochromatic batches. \square

B.4. Characterizing the optimal direction of classifiers

Thus far, we have primarily considered the separability decomposition of $\bar{\mathbf{Z}}_\pi$ and $\bar{\mathbf{Z}}_{\text{RR}}$. In fact, we can say more about the direction of optimal classifiers under the logistic loss. First, we prove Lemma 4.1.1, which constrains optimal directions via the separability decomposition. Next, we leverage overparameterization and the rank properties shown in Proposition B.1.3 to characterize the optimal direction under the logistic loss for data drawn from a density. Using these insights, we can prove our main result of this section, Proposition B.4.3.

We first restate and prove Lemma 4.1.1.

Lemma B.4.1. *Let $\mathbf{Z} = \mathbf{S}^{\text{LS}} \sqcup \mathbf{S}^{\text{SC}}$ be the separability decomposition of \mathbf{Z} . If \mathbf{v} is an optimal direction for \mathcal{L} , then $\mathbf{v}^\top \mathbf{x} = 0$ for all $\mathbf{x} \in \text{Span}(\mathbf{X}^{\text{SC}})$ and $y_i \mathbf{v}^\top \mathbf{x}_i > 0$ for every $(\mathbf{x}_i, y_i) \in \mathbf{S}^{\text{LS}}$.*

Proof. By definition of \mathbf{S}^{SC} , there exists some $(\mathbf{x}_i, y_i) \in \mathbf{S}^{\text{SC}}$ such that $y_i \langle \mathbf{v}, \mathbf{x}_i \rangle \leq 0$. If $\langle \mathbf{v}, \mathbf{x}_i \rangle \neq 0$, then $y_i \langle t\mathbf{v}, \mathbf{x}_i \rangle \rightarrow -\infty$ as $t \rightarrow \infty$, which contradicts the assumption that \mathbf{v} is an optimal direction. Similarly, if $y_i \langle \mathbf{v}, \mathbf{x}_i \rangle \leq 0$ for some $(\mathbf{x}_i, y_i) \in \mathbf{S}^{\text{LS}}$, then clearly $\mathbf{u} + t\mathbf{v}$ cannot infimize \mathcal{L} , as there exists a hyperplane which strictly separates \mathbf{S}^{LS} and is orthogonal to $\text{Span}(\mathbf{S}^{\text{SC}})$. \square

Next, we restate and prove Proposition 4.1.2.

Proposition B.4.2. *Suppose Assumption 1(a) holds. Assume that the iterates $\mathbf{v}_\pi(t)$ infimize \mathcal{L}_π , and their projections onto $\text{Span}(\overline{\mathbf{X}}_\pi^{\text{SC}})^\perp$ converge in direction to some optimal direction \mathbf{v}_π^* for \mathcal{L}_π . Then the GD risk \mathcal{L}_{GD} diverges if and only if $\overline{\mathbf{Z}}_\pi$ is PLS or LS and there exists some $(\mathbf{x}_i, y_i) \in \overline{\mathbf{Z}}_{\text{GD}}$ such that $y_i \mathbf{v}_\pi^{*\top} \mathbf{x}_i < 0$. The analogous statement holds true for $\overline{\mathbf{Z}}_{\text{RR}}$ under Assumption 1(b). Furthermore, the if part holds true for SS and RR without Assumption 1.*

Proof. First suppose Assumption 1(a) holds. Then if $\overline{\mathbf{Z}}_\pi$ is SC, Lemma B.4.1 implies that any optimal direction \mathbf{v}_π^* must be orthogonal to all of \mathbb{R}^d , so $\mathbf{v}_\pi^* = \mathbf{0}$. Hence the iterates $\mathbf{v}_\pi(t)$ converge to a finite optimum, which implies the GD risk cannot diverge. Note that if Assumption 1(a) doesn't hold, then the only difference is that being orthogonal to $\overline{\mathbf{X}}_\pi^{\text{SC}}$ does not imply that $\mathbf{v}_\pi^* = \mathbf{0}$.

Now suppose $\overline{\mathbf{Z}}_\pi$ is PLS or LS. Regardless of whether Assumption 1(a) holds, if $\mathbf{v}_\pi(t)$ infimizes \mathcal{L}_π , we necessarily have $\|\mathbf{v}_\pi(t)\|_2 \rightarrow +\infty$. Hence any mistake on $(\mathbf{x}_i, y_i) \in \overline{\mathbf{Z}}_{\text{GD}}$ implies divergence. And clearly, if there is no mistake on any (\mathbf{x}_i, y_i) then the GD risk does not diverge.

The same proof also goes through for $\overline{\mathbf{Z}}_{\text{RR}}$, so this proves the theorem. \square

At first glance, one might expect to be able to perfectly classify all the datapoints in the overparameterized regime. However, under the logistic risk, the optimal direction instead puts monochromatic batches *on the decision boundary*; the following proposition formalizes this notion.

Proposition B.4.3. *Suppose Assumption 2, $B > 2$, and $d > (B - 1)\frac{n}{B}$. Almost surely, for any $\pi \in \mathbb{S}_n$, there exists $\mathbf{v} \in \mathbb{R}^d$ which satisfies (1) for any non-monochromatic batch $\overline{\mathbf{Z}}_\pi^j$ we have $\text{sgn}(\mathbf{v}^\top \overline{\mathbf{X}}_\pi^j) = \mathbf{Y}_\pi^j$ and (2) for any monochromatic batch $\overline{\mathbf{Z}}_\pi^k$ we have $\mathbf{v}^\top \overline{\mathbf{X}}_\pi^k = \mathbf{0}^\top$. Furthermore, any optimal direction \mathbf{v}_π^* for the logistic risk \mathcal{L}_π necessarily satisfies both (1) and (2). The same conclusion holds for RR as well, with the requirement $d > (B - 1)\binom{n}{B}$.*

More precisely, our analysis is motivated by the fact that if $d \geq n$ and $\mathbf{X} \in \mathbb{R}^{d \times n}$ is full-rank, then given any $\mathbf{c} \in \mathbb{R}^n$ we can find some halfspace $\mathbf{v} \in \mathbb{R}^d$ such that $\mathbf{v}^\top \mathbf{X} = \mathbf{c}^\top$. In our binary classification setting, linear separability is implied by $\text{sgn}(\mathbf{c}^\top) = \mathbf{Y}$. However, we cannot directly apply this fact because BN actually prevents $\overline{\mathbf{X}}_\pi$ from being full-rank due to the mean zero constraint. However, it turns out that the following slightly weaker statement is true. We can always find a halfspace $\mathbf{v} \in \mathbb{R}^d$ that perfectly separates all the non-monochromatic batches. The following lemma proves this half of the proposition.

Lemma B.4.4. *Suppose Assumption 2, $B > 2$, and $d > (B - 1)\frac{n}{B}$. Then almost surely there exists $\mathbf{v} \in \mathbb{R}^d$ such that for every non-monochromatic batch $\overline{\mathbf{Z}}_\pi^j = (\overline{\mathbf{X}}_\pi^j, \mathbf{Y}_\pi^j)$, we have $\text{sgn}(\mathbf{v}^\top \overline{\mathbf{X}}_\pi^j) = \mathbf{Y}_\pi^j$. The same conclusion holds for RR as well, with the requirement $d > (B - 1)\binom{n}{B}$.*

Proof. We denote $\overline{\mathbf{X}}_\pi^j = \text{BN}(\mathbf{X}_\pi^j) = [\mathbf{x}_1 \ \cdots \ \mathbf{x}_B]$ and $\mathbf{Y}_\pi^j = [y_1 \ \cdots \ y_B]$. Since Assumption 2 holds, it also follows that Proposition 3.2.1 holds. Hence, within each batch, any $B - 1$ of the datapoints are linearly independent almost surely. This implies that we can find $\mathbf{v} \in \mathbb{R}^d$ such that for any batch $\overline{\mathbf{X}}_\pi^j$ and $\mathbf{c} \in \mathbb{R}^{B-1}$, we have $\mathbf{v}^\top \mathbf{x}_i = c_i$ for all $i \in [B - 1]$. In particular, we can pick c_i such that $\text{sgn}(c_i) = y_i$. Next, we show that \mathbf{v} can be picked such that $\text{sgn}(\mathbf{v}^\top \mathbf{x}_B) = y_B$. The mean zero constraint enforces that $\mathbf{v}^\top \mathbf{x}_B = -\sum_{i=1}^{B-1} c_i$. But if the labels are not monochromatic, we can just increase the absolute value of one of the c_i 's so that $\text{sgn}(-\sum_{i=1}^{B-1} c_i) = y_B$, as desired. \square

Hence, in the overparameterized regime minimizing the logistic risk will lead to a classifier which separates all the non-monochromatic batches. What happens to the monochromatic batches? It turns out that minimizing the logistic risk will lead to a classifier whose decision boundary *contains* all of the monochromatic batches.

Lemma B.4.5. *Assuming $d > (B - 1)\frac{n}{B}$, any optimal direction $\mathbf{v}_* \in \mathbb{R}^d$ for the logistic risk simultaneously puts all of the monochromatic batches on the decision boundary. More precisely, for any monochromatic batch $\overline{\mathbf{Z}}_\pi^j = (\overline{\mathbf{X}}_\pi^j, \mathbf{Y}_\pi^j)$ we have $\mathbf{v}_*^\top \overline{\mathbf{X}}_\pi^j = \mathbf{0}^\top$. Similarly, if $d > (B - 1)\binom{n}{B}$, the same conclusion holds for the RR dataset.*

Proof. Again, we denote the monochromatic batch by $\overline{\mathbf{Z}}_\pi^j$ by $\overline{\mathbf{X}}_\pi^j = [\mathbf{x}_1 \ \cdots \ \mathbf{x}_B]$ and $\mathbf{Y}_\pi^j = [y_1 \ \cdots \ y_B]$. The logistic loss on a single input \mathbf{x}_i for classifier $\mathbf{v} \in \mathbb{R}^d$ is $\ell(\mathbf{v}^\top \mathbf{x}_i, y_i) = -\log \rho(y_i \mathbf{v}^\top \mathbf{x}_i)$, where $\rho(t) = 1/(1 + \exp(-t))$

is the sigmoid function. WLOG suppose that $y_i = 1$ for all $i \in [B]$. Hence the minibatch risk is

$$-\sum_{i=1}^B \log \rho(\mathbf{v}^\top \mathbf{x}_i) = -\sum_{i=1}^{B-1} \log \rho(\mathbf{v}^\top \mathbf{x}_i) - \log \rho\left(-\mathbf{v}^\top \sum_{i=1}^{B-1} \mathbf{x}_i\right).$$

For each i we can look at the first order optimality condition for $s_i = \mathbf{v}^\top \mathbf{x}_i$. This yields $\rho(s_i) = \rho(-\sum_{i=1}^{B-1} s_i)$.

Note that this is satisfied when $s_i = 0$ for all $i \in [B]$, and by strict convexity this is the unique minimizer. And because of overparameterization, we can find a \mathbf{v}_* that will satisfy $s_i = 0$ for each monochromatic batch, i.e., the batch entirely lies on the decision boundary of the classifier defined by \mathbf{v}_* .

The proof carries over immediately to the RR setting, except in that setting, being overparameterized means $d > (B-1)\binom{n}{B}$. Hence the lemma is proved. \square

The geometric interpretation of Lemma B.4.5 is that in the overparameterized regime, if the dataset contains any monochromatic batches, then any optimal direction \mathbf{v}_* must be orthogonal to the subspace spanned by the monochromatic batches. Note, however, that the definition of overparameterized here depends on whether we look at \mathcal{Z}_π or \mathcal{Z}_{RR} . In the former case, overparameterized means $d > (B-1)\frac{n}{B}$, whereas in the latter case, overparameterized means $d > (B-1)\binom{n}{B}$. This insight motivates the construction of the toy datasets in Section 4.2. We conclude our characterization of the optimal direction in the overparameterized regime by proving Proposition B.4.3.

Proof of Proposition B.4.3. Proposition 3.2.1 implies that for $d > (B-1)\frac{n}{B}$ and assuming assumption 2, almost surely we have $\text{rank}(\overline{\mathbf{X}}_\pi) = (B-1)\frac{n}{B}$. We can lower bound the infimum of the SS logistic risk by the sum of the infima of the mini-batch SS logistic risks. Combining Lemmas B.4.4 and B.4.5, the claim follows. The same argument holds for $\overline{\mathcal{Z}}_{\text{RR}}$ assuming $d > (B-1)\binom{n}{B}$. \square

C. Proofs of technical lemmas

In this section we spell out the formal details of our guarantees for how permutations interact with BN. In Appendix C.1, we show that given a batch size $B = o(\log n)$ and a constant number of classes K , then w.h.p. over the choice of π there exists (many) monochromatic batches. In other words, for small batch sizes there are many batches consisting solely of positive or negative examples (in case of $K = 2$). Conversely, we show that above this threshold such monochromatic batches do not appear w.h.p.. In Appendix C.2, we show that there is a commensurate threshold above which the batch statistics themselves concentrate in the without-replacement setting. To do so we will appeal to the recent results of [Bardenet & Maillard \(2015\)](#) on the concentration of without-replacement estimators. Finally, in Appendix C.3, we prove that assuming the original features were drawn from a density, the batch normalized features have maximal rank almost surely. In other words, batch normalization preserves genericity of the original features modulo the mean zero constraint inside each batch.

C.1. Presence of monochromatic batches

In this section we formally prove Lemma B.1.4 via standard concentration arguments. One of the potential pitfalls of any mini-batch based algorithm is that its batches may not be representative of the entire dataset. More precisely, let's suppose we have a dataset \mathcal{Z} with labels coming from K classes. Suppose furthermore that the dataset is balanced — each class contains $n = Bm$ examples (here K is a constant). For any batch size B which is not sufficiently large, i.e. $B = o(\log n)$, then w.h.p. over the permutation π we will have $\Theta(n)$ batches which are *monochromatic* — batches which only consist of examples in the same class. This is closely related to the classic coupon collector problem, but we restate the guarantees here for the sake of completeness.

To prove this claim, we appeal to standard martingale concentration inequalities. Consider the batches $\overline{\mathcal{Z}}_\pi^i$ for $i \in [Km]$ and define the indicator variables $T_i = \mathbf{1}[\overline{\mathcal{Z}}_\pi^i \text{ is monochromatic}]$, and set $T \triangleq \sum_i T_i$, the total number of monochromatic batches.

Next, note that the sequence

$$\mathbb{E}[T], \mathbb{E}[T|T_1], \mathbb{E}[T|T_1, T_2], \dots, \mathbb{E}[T|T_1, T_2, \dots, T_{Km}]$$

is a Doob martingale. Indeed the martingale property follows from the tower property:

$$\mathbb{E}[\mathbb{E}[T|T_{1:k}] | \mathbb{E}[T|T_{1:(k-1)}]] = \mathbb{E}[T|T_{1:(k-1)}].$$

Note that $\mathbb{E}[T|T_1, \dots, T_{K_m}] = T$ and $\mathbb{E}[T_1] = K \binom{n}{B} / \binom{Kn}{B}$. Hence by linearity $\mathbb{E}[T] = \frac{K^2 n}{B} \frac{\binom{n}{B}}{\binom{Kn}{B}}$.

Let us now show that the martingale increments $\mathbb{E}[T|T_{1:k}] - \mathbb{E}[T|T_{1:(k-1)}]$ are bounded a.s.. In the worst case, the k th batch can only decrease this conditional expectation by at most K , since for any fixed class we can only remove at most one monochromatic batch from it. Hence the total number of potential monochromatic batches left can decrease by at most K . This worst case guarantee still holds under conditional expectation, so $|\mathbb{E}[T|T_{1:k}] - \mathbb{E}[T|T_{1:(k-1)}]| \leq K$ a.s..

Azuma-Hoeffding then tells us that for any $\epsilon > 0$ we have

$$\mathbb{P}[|T - \mathbb{E}[T]| \geq \epsilon] \leq 2 \exp\left(-\frac{B\epsilon^2}{2nK^3}\right).$$

This gives us the following fact which we use in both the regression and classification setting.

Fact C.1.1. *For any $\delta \in (0, 1)$, and a constant number of classes K with n datapoints each, we have with probability at least $1 - \delta$ that the total number of monochromatic batches T satisfies*

$$\left| T - \frac{K^2 n}{B} \frac{\binom{n}{B}}{\binom{Kn}{B}} \right| \leq \sqrt{\frac{2nK^3 \log(2/\delta)}{B}}.$$

For $K = O(1)$, we note that the above inequality guarantees that we can get within $O(\sqrt{n \log n})$ of the true expectation with at most $1/\text{poly}(n)$ failure probability. For the toy regression dataset in Section 3.3, we have $K = 2$ and $B = 2$, so we can use Fact C.1.1 to deduce that there will be $\Theta(n)$ monochromatic batches with high probability. We can also use Fact C.1.1 to prove Lemma B.1.4.

Proof of Lemma B.1.4. In the classification setting, we have $K = 2$ classes (positive and negative). Using the folklore inequalities

$$\left(\frac{n}{k}\right)^k \leq \binom{n}{k} \leq \left(\frac{en}{k}\right)^k,$$

we can deduce the lower bound on the expectation of T :

$$\mathbb{E}[T] \geq \frac{4n}{B(2e)^B}.$$

The lower bound is $\Omega(n^{1-\epsilon})$ for any $\epsilon > 0$ whenever $B = o(\log n)$, so indeed when $B = o(\log n)$ we have a positive number of monochromatic batches w.h.p.

Let us now upper bound the probability of obtaining any monochromatic batches. We have

$$\mathbb{P}[T_1 = 1] = 2 \frac{\binom{n}{B}}{\binom{2n}{B}} \leq 2 \frac{\prod_{k=0}^{B-1} (n-k)}{\prod_{k=0}^{B-1} (2n-k)} \leq 2 \frac{\prod_{k=0}^{B-1} (n-k)}{\prod_{k=0}^{B-1} 2(n-k)} \leq 2^{-B+1}.$$

Hence by union bound the probability that $T > 0$ is upper bounded by $\frac{4n}{B \cdot 2^{-B}}$. This is $1/\text{poly}(n)$ for some $B = \Omega(\log n)$, so indeed when $B = \Omega(\log n)$ we have no monochromatic batches with probability at least $1 - 1/\text{poly}(n)$. This concludes the proof. \square

C.2. Concentration of batch statistics for without-replacement sampling

In the following, we are generating B samples of $X_i \in \mathbb{R}$ without replacement from a population of size n , contained in $[a, b]$ a.s.. We let μ be the population mean and σ^2 be the population variance.

Lemma C.2.1 (Corollary 2.5 in (Bardenet & Maillard, 2015)). *Let $\hat{\mu}_B$ denote the sample mean for a sample of size B drawn without replacement from the overall population. For all $B \leq n$ and $\delta \in (0, 1)$ we have with probability at least $1 - \delta$ that*

$$|\hat{\mu}_B - \mu| \leq (b - a) \sqrt{\frac{\log(2/\delta)}{B}}.$$

Similarly, they prove the following result about concentration of the empirical variance

Lemma C.2.2 (Lemma 4.1 in (Bardenet & Maillard, 2015)). *Let $\hat{\sigma}_B^2 \triangleq \frac{1}{B} \sum_{i=1}^B (X_i - \hat{\mu}_B)^2$ be the (biased) empirical variance estimator and $\hat{\sigma}_B \triangleq \sqrt{\hat{\sigma}_B^2}$. Then for all $\delta \in (0, 1)$ we have with probability at least $1 - \delta$ that*

$$\hat{\sigma}_B \geq \sigma - 3(b - a) \sqrt{\frac{\log(3/\delta)}{2B}}.$$

We now prove the other side of this concentration inequality with a quick application of (Maurer, 2006, Theorem 1), following the same notation as in Bardenet & Maillard (2015).

Lemma C.2.3. *For all $\delta \in (0, 1)$ we have with probability at least $1 - \delta$ that*

$$\hat{\sigma}_B \leq \sigma + (b - a) \sqrt{\frac{\log(1/\delta)}{2B}}.$$

Proof. We take the self bounded random variable $Z = \frac{B}{(b-a)^2} \tilde{V}_B$, where

$$\tilde{V}_B \triangleq \frac{1}{B} \sum_{i=1}^B (X_i - \mu)^2$$

is computed with the samples X_i sampled *with replacement*. On the other hand,

$$V_B \triangleq \frac{1}{B} \sum_{i=1}^B (X'_i - \mu)^2$$

is computed with the samples X'_i sampled *without replacement*. We can relate the concentration of V_B to that of \tilde{V}_B ; the latter quantity is possible to analyze with the entropy method. Indeed, a routine modification of the proof of Bardenet & Maillard (2015, Lemma 3.3) (which uses essentially the same definition of Z) implies that

$$\mathbb{P} \left[V_B - \sigma^2 \geq \frac{(b-a)^2}{B} \epsilon \right] \leq \exp \left(- \frac{(b-a)^2 \epsilon^2}{2B\sigma^2} \right).$$

Solving for ϵ in terms of δ yields $\epsilon = \sqrt{\frac{2B\sigma^2}{(b-a)^2} \log(1/\delta)}$, so we obtain

$$\mathbb{P} \left[V_B - \sigma^2 \geq (b-a)\sigma \sqrt{\frac{2 \log(1/\delta)}{B}} \right] \leq \delta.$$

Since $V_B = (\hat{\mu}_B - \mu)^2 + \hat{\sigma}_B^2 \geq \hat{\sigma}_B^2$, we can complete the square to obtain

$$\mathbb{P} \left[\hat{\sigma}_B^2 \geq \left(\sigma + (b-a) \sqrt{\frac{\log(1/\delta)}{2B}} \right)^2 \right] \leq \delta.$$

So with probability at least $1 - \delta$, we have

$$\hat{\sigma}_B^2 \leq \left(\sigma + (b-a) \sqrt{\frac{\log(1/\delta)}{2B}} \right)^2,$$

and taking square roots implies the desired result. \square

Now, let us return to the question of concentration for batch norm with a randomly selected permutation π . The following proposition shows that assuming the batch size is large enough, the features of corresponding SS and GD datapoints are close to each other.

Proposition C.2.4. *If*

$$B = \Omega\left(\frac{\log(nd)}{\min_{k \in [d]} \left(\frac{\sigma_k}{b_k - a_k}\right)^2 \epsilon^2}\right),$$

then with probability at least $1 - 1/\text{poly}(n, d)$ we have

$$\|\bar{\mathbf{X}}_\pi - \pi \circ \bar{\mathbf{X}}_{\text{GD}}\|_{2, \infty} \leq \frac{\epsilon}{1 - \epsilon} \|\bar{\mathbf{X}}_{\text{GD}}\|_{2, \infty} + \frac{\epsilon\sqrt{d}}{1 - \epsilon}.$$

Here, we remind the reader that for a matrix $\mathbf{A} \in \mathbb{R}^{d \times n}$, we define $\|\mathbf{A}\|_{2, \infty} = \max_{i \in [n]} \|\mathbf{A}_{:, i}\|_2$, i.e. the maximum Euclidean norm of the columns.

Proof. WLOG consider the unnormalized first batch $\mathbf{X}_\pi^1 = \{\mathbf{x}_1, \dots, \mathbf{x}_B\}$. We initially handle concentration along its first coordinate $\{x_1, \dots, x_B\}$ for its first datapoint x_1 . Write $\hat{\boldsymbol{\mu}}, \hat{\boldsymbol{\sigma}} \in \mathbb{R}^d$ to denote the mini-batch mean and standard deviation, respectively, and $\hat{\mu}_k, \hat{\sigma}_k \in \mathbb{R}$ for the k th coordinate of these vectors. Similarly let $\boldsymbol{\mu}, \boldsymbol{\sigma} \in \mathbb{R}^d$ denote the full-batch mean and standard deviation, and let $\mu_k, \sigma_k \in \mathbb{R}$ to denote the k th coordinate of these vectors, respectively. Finally, let $\mathbf{a}, \mathbf{b} \in \mathbb{R}^d$ denote the coordinate-wise min and max of \mathbf{X} , with $a_k, b_k \in \mathbb{R}$ denoting the min and max for the k th coordinate of \mathbf{X} .

Hence, the first feature of the first datapoint in the normalized first batch $\bar{\mathbf{X}}_\pi^1$ is $\frac{x_1 - \hat{\mu}_1}{\hat{\sigma}_1}$, whereas the corresponding quantity in $\bar{\mathbf{X}}_{\text{GD}}$ is $\frac{x_1 - \mu_1}{\sigma_1}$.

Pick

$$B = \Omega\left(\frac{(\log(1/\delta))}{\left(\frac{\sigma_1}{b_1 - a_1}\right)^2 \epsilon^2}\right).$$

Then from Lemmas C.2.1 to C.2.3, we see that $|\hat{\mu}_1 - \mu_1| \leq \epsilon\sigma_1$ and $|\hat{\sigma}_1 - \sigma_1| \leq \epsilon\sigma_1$ with probability at least $1 - \delta$ for the first datapoint x_1 .

Now, we have

$$\begin{aligned} \left| \frac{x_1 - \hat{\mu}_1}{\hat{\sigma}_1} - \frac{x_1 - \mu_1}{\sigma_1} \right| &= \left| \frac{x_1 - \mu_1}{\hat{\sigma}_1} - \frac{x_1 - \mu_1}{\sigma_1} + \frac{\mu_1 - \hat{\mu}_1}{\hat{\sigma}_1} \right| \\ &\leq \left| \frac{\sigma_1 - \hat{\sigma}_1}{\hat{\sigma}_1} \right| \left| \frac{x_1 - \mu_1}{\sigma_1} \right| + \frac{|\mu_1 - \hat{\mu}_1|}{\hat{\sigma}_1} \\ &\leq \frac{\epsilon\sigma_1}{(1 - \epsilon)\sigma_1} \left| \frac{x_1 - \mu_1}{\sigma_1} \right| + \frac{\epsilon\sigma_1}{(1 - \epsilon)\sigma_1} \\ &\leq \frac{\epsilon}{1 - \epsilon} \left| \frac{x_1 - \mu_1}{\sigma_1} \right| + \frac{\epsilon}{1 - \epsilon}. \end{aligned}$$

To aggregate this bound across features, we need to pick B such that $|\hat{\mu}_k - \mu_k| \leq \epsilon\sigma_k$ and $|\hat{\sigma}_k - \sigma_k| \leq \epsilon\sigma_k$ for every feature $k \in [d]$. Indeed, this is achieved via the union bound by picking $\delta = 1/\text{poly}(d)$ and

$$B = \Omega\left(\frac{\log(1/\delta)}{\min_{k \in [d]} \left(\frac{\sigma_k}{b_k - a_k}\right)^2 \epsilon^2}\right).$$

For this batch size, we see that

$$\left\| \frac{\mathbf{x}_1 - \hat{\boldsymbol{\mu}}}{\hat{\boldsymbol{\sigma}}} - \frac{\mathbf{x}_1 - \boldsymbol{\mu}}{\boldsymbol{\sigma}} \right\|_2 \leq \frac{\epsilon}{1 - \epsilon} \left\| \frac{\mathbf{x}_1 - \boldsymbol{\mu}}{\boldsymbol{\sigma}} \right\|_2 + \frac{\epsilon\sqrt{d}}{1 - \epsilon}.$$

Note that a similar inequality holds for all of $\mathbf{x}_1, \dots, \mathbf{x}_B$. Now, recalling the definition of $\|\cdot\|_{2, \infty}$, by applying union bound on all the batches, we have

$$\|\bar{\mathbf{X}}_\pi - \pi \circ \bar{\mathbf{X}}_{\text{GD}}\|_{2, \infty} \leq \frac{\epsilon}{1 - \epsilon} \|\bar{\mathbf{X}}_{\text{GD}}\|_{2, \infty} + \frac{\epsilon\sqrt{d}}{1 - \epsilon},$$

which occurs with probability at least $1 - 1/\text{poly}(n, d)$ when we appropriately choose $B = \Omega\left(\frac{\log(nd)}{\min_{k \in [d]} \left(\frac{\sigma_k}{b_k - a_k}\right)^2 \epsilon^2}\right)$. \square

We now use the above proposition to prove Proposition B.1.6.

Proof of Proposition B.1.6. Conditions (2) and (3) in the definition of γ -robustness ensures that $\min_{k \in [d]} \left(\frac{\sigma_k}{b_k - a_k}\right)^2 = \Omega(1)$ and $\|\bar{\mathbf{X}}_{\text{GD}}\|_{2, \infty} = O(\sqrt{d})$. Hence taking $B = \Omega\left(\frac{\log(nd)}{\epsilon^2}\right)$ as in Proposition C.2.4, we conclude that with probability at least $1 - 1/\text{poly}(n, d)$ we have

$$\|\bar{\mathbf{X}}_{\pi} - \pi \circ \bar{\mathbf{X}}_{\text{GD}}\|_{2, \infty} \leq O\left(\frac{\epsilon \sqrt{d}}{1 - \epsilon}\right).$$

Hence, by taking $\epsilon = O\left(\frac{\gamma}{\sqrt{d}}\right)$, we see that

$$\|\bar{\mathbf{X}}_{\pi} - \pi \circ \bar{\mathbf{X}}_{\text{GD}}\|_{2, \infty} \leq O(\gamma).$$

Plugging this choice of ϵ back into our definition of B , we conclude that when $B = \Omega\left(\frac{d \log(nd)}{\gamma^2}\right)$, $\bar{\mathbf{X}}_{\pi}$ concentrates around $\pi \circ \bar{\mathbf{X}}_{\text{GD}}$ within distance γ . \square

C.3. Rank of batch normalized features

In this section we prove Proposition 3.2.1, which states that for batch sizes greater than 2, Assumption 2 implies that the batch normalized dataset will be full-rank (and hence full-dimensional) almost surely.

One shift in perspective that is fruitful for proving linear independence is to view batch normalization as an operation that returns functions of the input dataset. Since BN operates independently on each of the d features, we first handle the case where the input is a batch of scalars. Let $\binom{[n]}{B}$ denote the set of all $\binom{n}{B}$ batches of size B that can be created from n datapoints contained in \mathbf{X} . Fix an arbitrary labelling of these $\binom{n}{B}$ batches, and let \mathbf{B}^j refer to the j th such batch. WLOG suppose that $\mathbf{B}^1 = \{x_1, \dots, x_B\} \in \mathbb{R}^B$.

Formally, let $\mathcal{F}^B \triangleq \{f : \mathbb{R}^B \setminus \{\mathbf{x} \in \mathbb{R}^B \mid x_1 = x_2 = \dots = x_B\} \rightarrow \mathbb{R}\}$ denote the space of real valued functions on batches of size B where BN is defined. On batch $\mathbf{B}^j = \{x_{j_1}, \dots, x_{j_B}\}$, BN is an operation that maps this batch to the set of B functions $\left\{g_i^j(\mathbf{B}^j)\right\}_{i=1}^B$, where

$$g_i^j(\mathbf{B}^j) \triangleq \frac{x_{j_i} - \mu^j}{\sigma^j} \in \mathcal{F}^B$$

where μ^j and σ^j are the empirical mean and standard deviation, respectively of \mathbf{B}^j . If j is clear from context, we may drop the superscript j without chance of confusion. We also sometimes abuse notation and write g_i^j as a function of \mathbf{X} , since \mathbf{X} contains all datapoints in \mathbf{B}^j . From this perspective, BN_{π} maps a dataset of n datapoints to n functions.

We first show that, within a batch, the functions have rank $B - 1$ over \mathbb{R} .

Lemma C.3.1. *Viewed as functions, any subset of $B - 1$ functions of the batch normalized outputs $\{g_i^1\} = \left\{\frac{x_1 - \mu}{\sigma}, \dots, \frac{x_B - \mu}{\sigma}\right\}$ have rank $B - 1$ over \mathbb{R} .*

Proof. First, we note that the functions x_i are linearly independent over \mathbb{R} for $i \in [n]$. WLOG take the subset of $B - 1$ functions to be $\{g_i^1\}_{i=1}^{B-1}$. Suppose that there is a dependence relationship

$$\sum_{i=1}^{B-1} c_i \frac{x_i - \mu}{\sigma} = 0.$$

Rearranging, we see that

$$\sum_{i=1}^{B-1} c_i x_i = \frac{1}{B} \left(\sum_{t=1}^B x_t \right) \left(\sum_{i=1}^{B-1} c_i \right)$$

and because of linear independence of the x_i 's as functions over \mathbb{R} and only the right-hand side contains x_B , the only way this can happen is if $\sum_{i=1}^{B-1} c_i = 0$. But then we have a dependence relationship between the x_i 's for $i \leq B - 1$. Linear independence of the x_i 's thus implies that $c_i = 0$ for each i . \square

With the above lemma in hand, we can show that as functions, any collection of batch normalized outputs are essentially full rank. The caveat is we need to throw out one function in each batch, as each batch is trivially dependent because of the zero mean constraint.

Proposition C.3.2. *Let $B > 2$. Consider the $B \binom{n}{B}$ functions that are the batch normalized outputs of all $\binom{n}{B}$ batches $B^j \in \binom{[n]}{B}$. If we take any subset of $(B-1) \binom{n}{B}$ of these functions obtained by removing one function from each of the $\binom{n}{B}$ batches, then the rank of these functions over \mathbb{R} is $(B-1) \binom{n}{B}$. In particular, the rank of the functions corresponding to any π is $(B-1) \frac{n}{B}$.*

Proof. Consider the batch $B^1 = \{x_1, \dots, x_B\}$. By Lemma C.3.1, we know that the functions $\{g_i^1\}_{i=1}^{B-1} = \left\{ \frac{x_i - \mu}{\sigma} \right\}_{i=1}^{B-1}$ are linearly independent. Consider a dependence relation amongst any subset of the $(B-1) \binom{n}{B}$ described in the theorem statement. WLOG we can suppose that this was formed by throwing out g_B^j from each batch B^j and consider a dependence relation between the $(B-1) \binom{n}{B}$ remaining functions. The dependence relation reads

$$\sum_{i=1}^{B-1} c_{i,1} \frac{x_i - \mu}{\sigma} = \sum_{i=1}^{B-1} \sum_{j>1} c_{i,j} g_i^j. \quad (35)$$

We show that some setting of the input x_i for $i \in [n]$ yields a contradiction unless $c_{i,1} = 0$ for $i \in [B]$. The main insight is that BN has jump discontinuities at the points where it is undefined, i.e., where the entire batch is equal to the same thing.

More formally, for $i > B$ we set the other datapoints x_i to be arbitrary pairwise distinct positive real numbers. Suppose for the sake of contradiction that $c_{i,1} \neq 0$ for some $i \in [B]$; WLOG we can assume that $i = 1$. Since the functions on the LHS are linearly independent by Lemma C.3.1, the LHS is not identically zero. We show that the LHS of Equation (35) exhibits discontinuous behavior in the punctured neighborhood around $(x_1, \dots, x_B) = \mathbf{0}$, whereas the RHS of Equation (35) is continuous on the same neighborhood; this yields a contradiction.

Indeed, set $(x_1, x_2, x_3, \dots, x_B) = (\epsilon, \epsilon^2, -\epsilon - \epsilon^2, 0, \dots, 0)$, where $\epsilon \neq 0$. We have

$$\frac{x_1 - \mu}{\sigma} = \frac{\sqrt{B}\epsilon}{\sqrt{\epsilon^2 + \epsilon^4 + (\epsilon + \epsilon^2)^2}}.$$

If $\epsilon \rightarrow 0^+$, then the first normalized coordinate approaches $+\sqrt{\frac{B}{2}}$. However, if $\epsilon \rightarrow 0^-$, then the first normalized coordinate approaches $-\sqrt{\frac{B}{2}}$. This is a contradiction, since the RHS of Equation (35) is continuous as a function of ϵ . We conclude that $c_{i,1} = 0$ for all $i \in [B-1]$.

The same argument holds if we replace the LHS with any batch B^j . Hence, $c_{i,j} = 0$ for all $i \in [B-1]$ and $j \in \left[\binom{n}{B}\right]$. We conclude the rank of these $(B-1) \binom{n}{B}$ functions is $(B-1) \binom{n}{B}$, as desired. Notice also that this argument also shows that the functions corresponding to the batches in π also have rank $(B-1) \frac{n}{B}$. \square

Note that the assumption that $B > 2$ is critical for the construction in the proof. If $B = 2$, then actually $\frac{x_i - \mu}{\sigma} \in \{\pm 1\}$, and the proof breaks down. The batch normalized dataset will be a Boolean matrix; hence, its rank cannot be analyzed by using density arguments.

The following lemma establishes that the zero set of any nontrivial linear combination of the g_i^j 's is a measure zero subset of \mathbb{R}^n .

Lemma C.3.3. *Suppose that $c_{i,j}$ are not identically zero. Then for $B > 2$, the zero set of $\sum_{i=1}^{B-1} \sum_{j=1}^{\binom{n}{B}} c_{i,j} g_i^j$ is a measure zero subset of \mathbb{R}^n .*

Proof. Since $B > 2$, Proposition C.3.2 implies that $f(\mathbf{X}) \triangleq \sum_{i=1}^{B-1} \sum_{j=1}^{\binom{n}{B}} c_{i,j} g_i^j$ is not identically zero. Furthermore, $f(\mathbf{X})$ is real analytic on finitely many connected components of \mathbb{R}^n , since each of the functions g_i^j are real analytic on finitely many connected components of \mathbb{R}^n . The claim follows by applying Proposition 0 in Mityagin (2015). \square

Having established these results, we can finally prove Proposition 3.2.1.

Proof of Proposition 3.2.1. Denote

$$\mathbf{X} = \begin{bmatrix} \mathbf{X}_1 \\ \vdots \\ \mathbf{X}_d \end{bmatrix} \in \mathbb{R}^{d \times n},$$

and the $(B-1)\binom{n}{B}$ functions

$$\mathbf{g}_i^j(\mathbf{X}) = \begin{bmatrix} g_{i,1}^j(\mathbf{X}_1) \\ \vdots \\ g_{i,d}^j(\mathbf{X}_d) \end{bmatrix} \in \mathbb{R}^d$$

We can assemble these vector valued functions into a matrix $\mathbf{G}(\mathbf{X})$ of size $d \times (B-1)\binom{n}{B}$ with the k th row consisting of the scalar valued functions $g_{i,k}^j(\mathbf{X})$. Now consider the determinant of any $\min\{d, (B-1)\frac{n}{B}\} \times \min\{d, (B-1)\frac{n}{B}\}$ submatrix of these scalar functions, which is itself a function of \mathbf{X} .

To prove the claim, it suffices to show that this determinant — which is a function of $\mathbf{X} \in \mathbb{R}^{d \times n}$ — is analytic almost everywhere and not identically zero; then Lemma C.3.3 implies that it vanishes on a measure zero set of $\mathbb{R}^{d \times n}$.

We prove the above claim by induction on d . For $d=1$, note that for any i, j , the scalar function $g_{i,1}^j(\mathbf{X})$ is not identically zero and is analytic on finitely connected components of $\mathbb{R}^{d \times n}$. Now suppose the claim is true for d , we prove the claim for $d+1$. When we use cofactor expansion along the first row of \mathbf{G} , which has functions $g_{i,1}^j$ which crucially only depend on the \mathbf{X}_1 , we obtain

$$\sum_{i=1}^{B-1} \sum_{j=1}^{\binom{n}{B}} (-1)^{i+j} g_{i,1}^j(\mathbf{X}_1) \det(\mathbf{M}_{i,j}(\mathbf{X})),$$

where $\mathbf{M}_{i,j}$ denotes the minor corresponding to the i th function of batch j . Note that these minors are not functions of \mathbf{X}_1 , so they can be treated as constants with respect to \mathbf{X}_1 . By induction these constants are nonzero almost surely. Then Proposition C.3.2 implies that the determinant, which is a linear combination of the functions $g_{i,1}^j$ is itself not zero identically. Also, the determinant is manifestly piecewise analytic on finitely connected components of $\mathbb{R}^{(d+1) \times n}$, being a polynomial of such functions. Hence the induction is completed.

We can now finish off the proof of the proposition. When Assumption 2 holds, the probability that the data falls in measure zero set is a probability zero event. In other words, almost surely any such $d \times (B-1)\binom{n}{B}$ matrix constructed by batch normalizing and throwing out one normalized point in each batch is full rank. The same argument holds for the $(B-1)\frac{n}{B}$ functions that correspond to the batch normalized outputs for a permutation π . Hence the proposition is proved. \square

D. Additional experiments on real data

In this section we provide detailed explanations of our experimental setup and present our additional experiments for regression and classification. All experiments were implemented in PyTorch 1.12.

D.1. Experiment details

We first define the architectures used in our real data experiments outlined in Section 4.3.

For the linear+BN networks, the 1-layer network is

$$\mathbf{X} \mapsto \mathbf{W}_1 \Gamma_1 \text{BN}(\mathbf{X}),$$

On the other hand, the 2-layer network is

$$\mathbf{X} \mapsto \mathbf{W}_2 \Gamma_2 \text{BN}(\mathbf{W}_1 \mathbf{X})$$

and the 3-layer network is

$$\mathbf{X} \mapsto \mathbf{W}_3 \Gamma_3 \text{BN}(\mathbf{W}_2 \Gamma_2 \text{BN}(\mathbf{W}_1 \mathbf{X})).$$

Hence the difference between the 1-layer network and deeper network is that the deeper networks have tunable parameters inside of BN.

For the MLP experiments, the 3-layer network is

$$\mathbf{X} \mapsto \mathbf{W}_3 \text{ReLU}(\Gamma_3 \text{BN}(\mathbf{W}_2 \text{ReLU}(\Gamma_2 \text{BN}(\mathbf{W}_1 \mathbf{X}))))).$$

For the ResNet18 experiments, we used the ResNet18 architecture available through PyTorch, using `ResNet18_Weights.DEFAULT` pretrained weights.

The linear and BN layers were all initialized using the default PyTorch initialization. For the linear+BN networks, the linear layers were instantiated with a width of 512 and `bias=False`. For the 3 layer MLP, the linear layers were instantiated with a width of 512 and `bias=True`. The BN layers were instantiated with `track_running_stats=False`. As alluded to in Section 1.2, to evaluate the training GD risk \mathcal{L}_{GD} in the eval loop, we passed in the entire dataset as a single batch, thus avoiding EMA altogether. Except for the ResNet18 experiments, the images in the dataset were flattened into vectors.

Except for in the respective batch size and momentum ablation study (Figures 11 and 12), we used batch size $B = 128$ and no momentum. Note that for all of the datasets, $\log_2 n \approx 16$, which suggests that we are in the asymptotic regime where divergence can happen as stated in Theorem 4.1.3.

We now explain the difference in divergence behavior between the 1-layer and deeper linear+BN networks for SS. As suggested by Theorem 4.1.3, divergence can happen if the separability decomposition of $\bar{\mathbf{Z}}_{\text{GD}}$ is not robust to perturbation. In the 1-layer case, the data remains far from being linearly separable. Meanwhile, in the deeper case, the network is incentivized to train the parameters inside BN such that the final features (e.g., $\text{BN}_\pi(\mathbf{W}_1 \mathbf{X})$) are closer to being LS. But the nonlinearity of BN is not enough to make $\text{BN}_\pi(\mathbf{W}_1 \mathbf{X})$ robustly LS. This also explains why introducing nonlinear activations prevents the divergence phenomenon.

We also note that in reality RR is run for T epochs. Thus, if $T < \frac{\binom{n}{B}}{B}$, the optimization routine only sees a proper subset of $\bar{\mathbf{Z}}_{\text{RR}}$. However, there are other forces that help ensure that the subset of $\bar{\mathbf{Z}}_{\text{RR}}$ actually seen during optimization is SC and satisfies Assumption 1(b). For example, it is likely that the algorithm sees non-monochromatic batches that also cause the hulls to overlap. One way this can happen is if $\mathbf{0}$ is in the relative interiors of the convex hulls of the monochromatic portion of each batch. Moreover, with extremely high probability we never see a repeat batch, so by Proposition B.1.3 the rank of the subset of $\bar{\mathbf{Z}}_{\text{RR}}$ is w.h.p. equal to $T \cdot \min\{d, (B-1)\frac{n}{B}\}$. For us, $d = 10 \cdot 512$, $B = 128$, and $n \geq 50000$, so the rank of the subset of $\bar{\mathbf{Z}}_{\text{RR}}$ we see easily outstrips the dimensionality of the final linear layers. This ensures that Assumption 1(b) holds.

In Figure 11, we see that divergence on 3 layer linear+BN networks generally occurs for large batch sizes, which corroborates Theorem 4.1.3. These batch sizes were picked because they are common choices for batch sizes in practice. For the largest batch size ($B = 128$), there does not appear to be divergence within 1000 epochs, which we address below.

In Figure 12, we see that the presence of momentum preserves divergence for SS, and in some cases accelerates it. Note that for each stepsize we used the same permutation. For the $\eta = 10^{-4}$ experiment, although the 0 momentum run did not start to diverge within 1000 epochs, the 0.9 and 0.99 momentum runs started to diverge. This further lends evidence to the claim that the apparent reason for no divergence for $\eta = 10^{-4}$ without momentum is that the small learning rate leads to slower convergence to an optimal direction for \mathcal{L}_π .

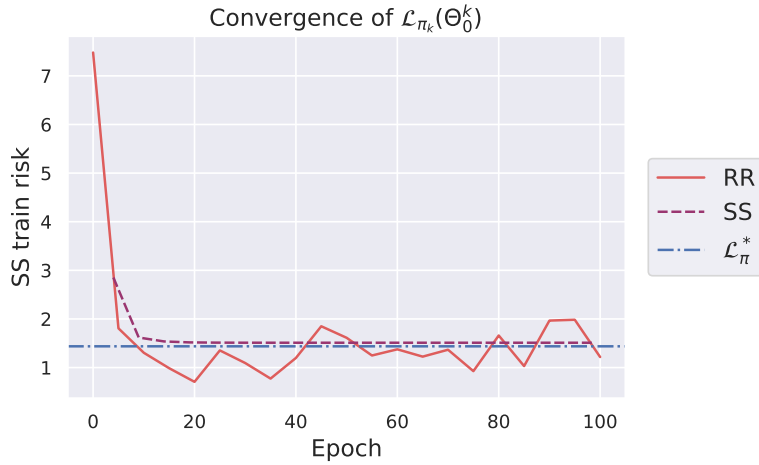


Figure 7. Loss evolution for $\ell_{\pi_k}(\mathbf{M}(k))$ for experiment described in Section 3.4. Note how SS converges to \mathcal{L}_{π}^* ; RR oscillates due to resampling π_k but appears to converge to a value close to \mathcal{L}_{π}^* . This corroborates the convergence results Theorems 3.2.2 and 3.2.3 to distorted optima.

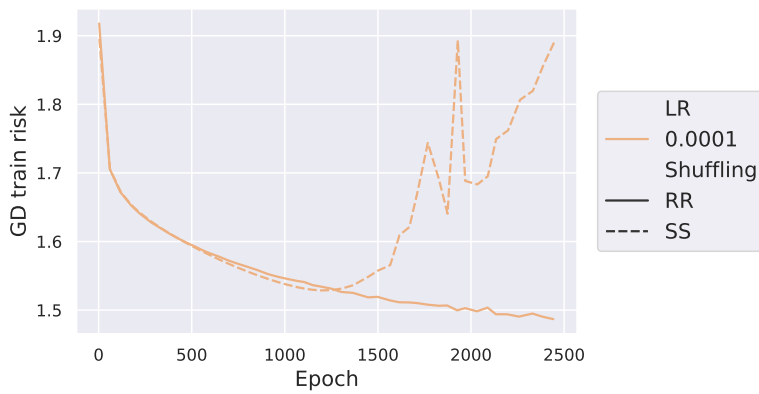


Figure 8. Eventual separation between SS and RR for 3 layer linear+BN network for $\eta = 10^{-4}$ after around epoch 1200 on CIFAR10. The color is consistent with the original Figure 3b.

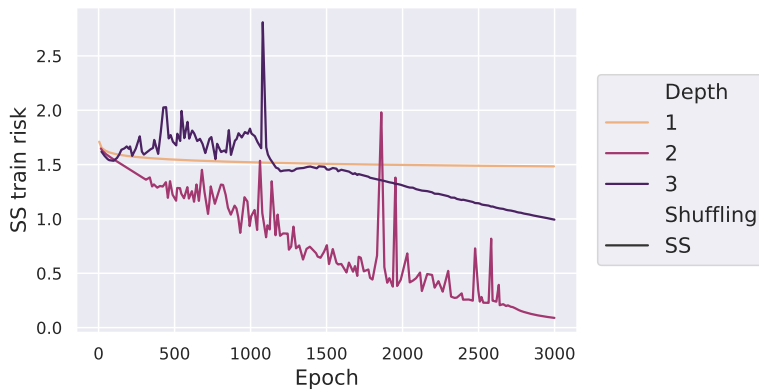


Figure 9. Evidence that the 2 layer and 3 layer linear+BN actually had diverging GD risks when trained with SS. On the y -axis we plot the value of $\mathcal{L}_{\pi}(\Theta_0^k)$; divergence occurs when the SS features $\bar{\mathbf{X}}_{\pi}$ are LS. As the SS risk for both 2 and 3 layer networks continues to decrease, this supports $\bar{\mathbf{X}}_{\pi}$ being LS. However, the SS risk for the 1 layer network seems to plateau, suggesting $\bar{\mathbf{X}}_{\pi}$ is SC rather than LS.

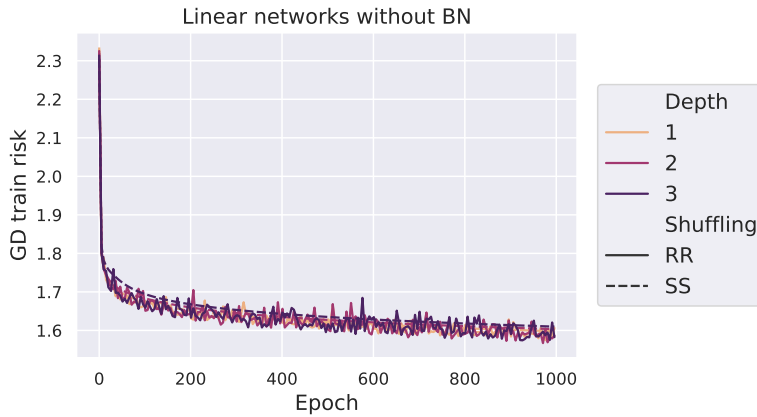


Figure 10. No divergence occurs for 1, 2, and 3 layer linear networks without BN trained with SS. This supports the theory that the strange training behavior occurs when using SS and BN in combination, rather than being an intrinsic failing of SS (cf. Figure 3b). Here we picked $\eta = 10^{-2}$, because we generally observed faster divergence with larger learning rates.

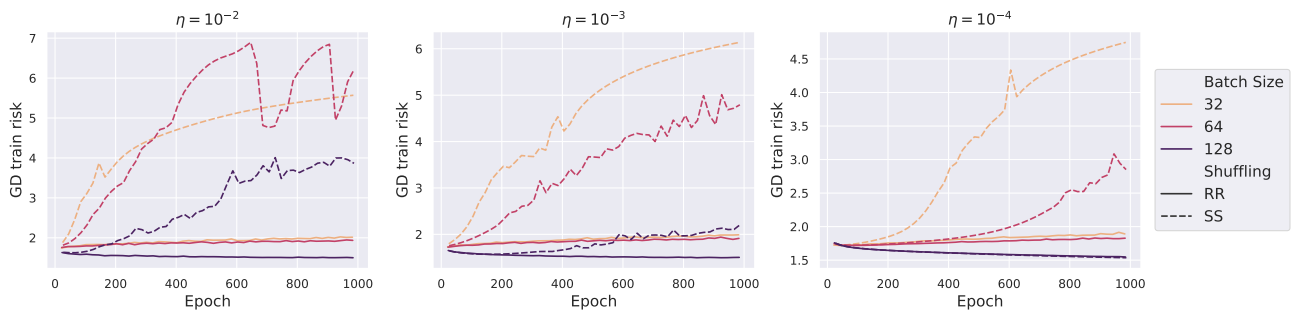


Figure 11. Batch size ablation. These experiments were done on 3 layer linear+BN networks; each subfigure shows the results of training with different stepsizes $\eta \in \{10^{-2}, 10^{-3}, 10^{-4}\}$. All experiments were performed on CIFAR10. Since $\log n = \log 50000 \approx 16$, all batch sizes are in the regime where divergence can happen according to Theorem 4.1.3.

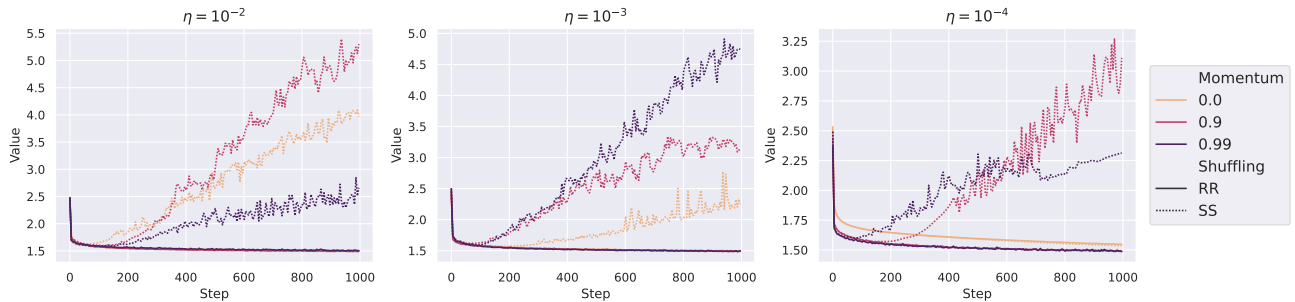


Figure 12. Momentum ablation. These experiments were done on 3 layer linear+BN networks; each subfigure shows the results of training with different stepsizes $\eta \in \{10^{-2}, 10^{-3}, 10^{-4}\}$. All experiments were performed on CIFAR10. For each subfigure, the experiment was run with the same random permutation. Generally, momentum doesn't prevent divergence, and in some cases even hastens it.

E. Calculations for toy datasets

In this section we do the detailed calculations and provide additional figures to understand the toy datasets we introduced in Sections 3.3 and 4.2. Since both constructions use $B = 2$, we remind the reader that the batch normalization of any two distinct real numbers is $(-1, 1)$. It follows that if we batch normalize with $B = 2$, we obtain $\overline{\mathbf{X}}_\pi \in \{-1, 1\}^{d \times n}$. Recall the distorted risk \mathcal{L}_π reflects the single-shuffle batch normalized dataset $\overline{\mathbf{Z}}_\pi$. This \mathcal{L}_π is a *random* quantity, and it is over this source of randomness (the construction of the batches) that we show that there is a gap between SS, RR, and GD.

E.1. Regression toy dataset

Proposition E.1.1. *There exists a regression dataset $\mathbf{Z} = (\mathbf{X}, \mathbf{Y}) \in [-1, 1]^{1 \times 16n} \times [-1, 1]^{1 \times 16n}$ such that the following statements hold with batch size $B = 2$:*

- (1) $M_{\text{GD}}^* = M_{\text{RR}}^* = 0$.
- (2) $M_\pi^* \neq 0$ with probability at least $1 - O(\frac{1}{\sqrt{n}})$.
- (3) $|M_\pi^*| = \Omega(\frac{1}{\sqrt{n}})$ with constant probability.

We first describe the construction and then prove that the dataset satisfies the properties outlined above.

Formal construction of regression dataset: Construct the dataset as follows. Take $\mathbf{A} \in \mathbb{R}^{1 \times 4n}$ to be $4n$ equally spaced points in the interval $(\frac{3}{4}, 1)$. Define

$$\begin{aligned} \mathbf{X}^1 &= \mathbf{A}; & \mathbf{X}^2 &= -\mathbf{A} \\ \mathbf{X}^3 &= -\mathbf{A} + \frac{1}{2}\mathbf{1}^\top; & \mathbf{X}^4 &= \mathbf{A} - \frac{1}{2}\mathbf{1}^\top \end{aligned}$$

and

$$\begin{aligned} \mathbf{Y}^1 &= \mathbf{1}^\top; & \mathbf{Y}^2 &= \mathbf{1}^\top \\ \mathbf{Y}^3 &= -\mathbf{1}^\top; & \mathbf{Y}^4 &= -\mathbf{1}^\top. \end{aligned}$$

Notice that the indices also match which quadrant the cluster of points are in. Visually, these four groups of points $\mathbf{Z}^i \triangleq (\mathbf{X}^i, \mathbf{Y}^i)$ are clusters with the i th cluster in the i th quadrant. For brevity, we also refer to these clusters by their index i only, so cluster 1 refers to the cluster $\mathbf{Z}^1 = (\mathbf{X}^1, \mathbf{Y}^1)$, and so on. These definitions are consistent with the clusters depicted in Figure 13a.

Then take $\mathbf{X} = [\mathbf{X}^1 \ \mathbf{X}^2 \ \mathbf{X}^3 \ \mathbf{X}^4] \in \mathbb{R}^{1 \times 16n}$ and $\mathbf{Y} = [\mathbf{Y}^1 \ \mathbf{Y}^2 \ \mathbf{Y}^3 \ \mathbf{Y}^4] \in \mathbb{R}^{1 \times 16n}$. After applying BN with permutation π and batch size 2, we obtain a dataset $\overline{\mathbf{X}}_\pi$ with every point being located in one of four SS clusters $\overline{\mathbf{Z}}_\pi^i \triangleq (\overline{\mathbf{X}}_\pi^i, \mathbf{Y}_\pi^i)$ for $i \in [4]$ located at $(\pm 1, \pm 1)$ with the same relative labelling: $\overline{\mathbf{Z}}_\pi^1$ is located at $(1, 1)$ and then labelling counterclockwise (see Figure 13b).

In Figure 13a, we visualize the construction with $n = 3$ (so the depicted dataset has $16n = 48$ datapoints). We plot the slopes of the M_{GD}^* (green solid line), M_{RR}^* (purple dash-dotted line), and typical values for M_π^* (yellow dotted line). In Figure 13b, we show what $\overline{\mathbf{Z}}_\pi$ looks like for a typical permutation π . The sizes of the points represent the number of points that end up in the corresponding cluster.

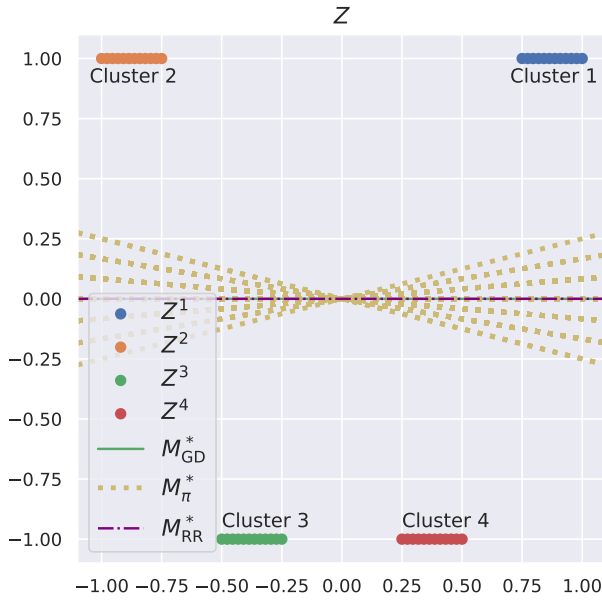
E.1.1. ANALYZING THE REGRESSION TOY DATASET

In order to prove this proposition, we will need the following standard technical estimate on the asymptotics of binomial coefficients (see e.g. Thomas & Joy (2006))

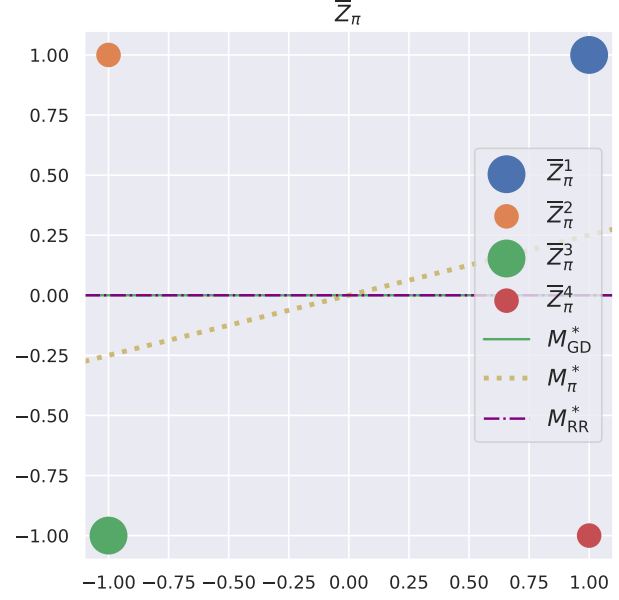
Lemma E.1.2. *For all n and k we have*

$$\sqrt{\frac{n}{8k(n-k)}} 2^{H(k/n)n} \leq \binom{n}{k} \leq \sqrt{\frac{n}{\pi k(n-k)}} 2^{H(k/n)n},$$

where $H(p) \triangleq -p \log_2 p - (1-p) \log_2 (1-p)$ is the binary entropy function.



(a) Unnormalized toy regression dataset \mathbf{Z} demonstrating distortion of SS with constant probability. Notice how the GD and RR lines are aligned with slope 0, but the SS lines are distorted away from 0.



(b) Toy regression dataset after BN with permutation π . The size of the point corresponding to $\bar{\mathbf{Z}}_\pi^i$ represents the number of normalized points in $\bar{\mathbf{Z}}_\pi^i$ (not to scale), and here we picked π such that $M_\pi^* = \frac{1}{4}$.

Now let us prove the proposition.

Proof of Proposition E.1.1. First, it is clear from the symmetry of \mathbf{X} that $\mathbf{Y}_{\text{GD}} \bar{\mathbf{X}}_{\text{GD}}^\top = 0$, so $M_{\text{GD}}^* = 0$. This proves the first half of (1).

Setup: For $i \in [4]$, let k_π^i denote the number of points that end up in cluster $\bar{\mathbf{Z}}_\pi^i$ after normalizing with permutation π . Since we have $M_\pi^* = \frac{1}{16n} \mathbf{Y}_\pi \bar{\mathbf{X}}_\pi^\top$, evidently the k_π^i completely determine M_π^* . In fact, more is true: we claim that $M_\pi^* = 0$ if and only if $k_\pi^1 = 4n$.

To see why this is true, first note that the label \mathbf{Y} is unaffected by BN. Hence, there are necessarily $8n$ points that end up with y coordinate 1. It follows that if there are k points in $\bar{\mathbf{Z}}_\pi^1$, then there are $8n - k$ points in $\bar{\mathbf{Z}}_\pi^2$. Similarly, if there are j in $\bar{\mathbf{Z}}_\pi^3$, then there are $8n - j$ in $\bar{\mathbf{Z}}_\pi^4$. On the other hand, recall that BN with $B = 2$ and $d = 1$ always sends one point in each batch to $x = +1$ and one point to $x = -1$. Hence, there are $8n$ points that end up with x coordinate 1, which means $k + 8n - j = 8n$, implying that $k = j$.

Referring back to the formula for M_π^* , we see that

$$\begin{aligned} M_\pi^* &= \frac{1}{16n} \sum_{i=1}^4 \mathbf{Y}_\pi^i (\bar{\mathbf{X}}_\pi^i)^\top \\ &= \frac{1}{16n} (k - (8n - k) + k - (8n - k)) \\ &= \frac{k - 4n}{4n}. \end{aligned}$$

Hence $M_\pi^* = 0$ if and only if $k_\pi^1 = 4n$. Referring back to Figure 13b, we see that the sizes of the clusters represent k_π^i , and the plotted M_π^* with slope $\frac{1}{4}$ corresponds to a π where $k_\pi^1 = 5n$.

To analyze k_π^1 , for $i, j \in [4]$ we introduce the random variables $T_\pi^{\{i,j\}}$ to denote the number of batches formed with one point from cluster i and cluster j with permutation π . Similarly let $U \triangleq \{1, 2\}$ denote the upper clusters and $L \triangleq \{3, 4\}$

denote the lower clusters. Define $T_\pi^{\{i,L\}} \triangleq T_\pi^{\{i,3\}} + T_\pi^{\{i,4\}}$, which represents the total number of batches with one point in cluster i and another point in L when using permutation π . Finally, define $T_\pi^{\{U,L\}} \triangleq T_\pi^{\{1,L\}} + T_\pi^{\{2,L\}}$, which represents the total number of batches with one point in cluster U and another in L with permutation π .

With this notation in hand, let us prove the claim. Evidently, $k_\pi^1 = T_\pi^{\{U,U\}} + T_\pi^{\{1,L\}}$. Similarly we have $k_\pi^2 = T_\pi^{\{U,U\}} + T_\pi^{\{2,L\}}$. Because we established earlier that $k_\pi^1 + k_\pi^2 = 8n$, then $k_\pi^1 = 4n$ if and only if $T_\pi^{\{1,L\}} = T_\pi^{\{2,L\}}$. In words, this means that $M_\pi^* = 0$ if and only if the number of $\{1, L\}$ batches is the same as the number of $\{2, L\}$ batches.

RR averages out distortion: For every π we can find π' such that $M_{\pi'}^* = -M_\pi^*$. This is because we can always find π' that swap $T_\pi^{\{1,L\}}$ and $T_\pi^{\{2,L\}}$, by turning all $\{1, L\}$ batches into $\{2, L\}$ batches and vice versa. Then Proposition 3.3.1 implies that $M_{\text{RR}}^* = 0$, which proves the second half of (1).

SS is distorted: Now, let us show that $\mathbb{P}[k_\pi^1 = 4n] = O(\frac{1}{\sqrt{n}})$. The main idea is that conditioned on $T_\pi^{\{U,L\}} = 2t$, we can compute the probability that $T_\pi^{\{1,L\}} = T_\pi^{\{2,L\}} = t$ exactly. Indeed, of the $4n$ points in cluster 1, we pick t of them to form batches with L , and similarly for cluster 2. This gives $\binom{4n}{t}^2$ ways for $T_\pi^{\{1,L\}} = T_\pi^{\{2,L\}} = t$. In total, there are $8n$ points in U and we picked $2t$ of them to match with L , which gives a denominator of $\binom{8n}{2t}$. Hence

$$\mathbb{P}[k_\pi^1 = 4n | T_\pi^{\{U,L\}} = 2t] = \frac{\binom{4n}{t}^2}{\binom{8n}{2t}}.$$

In order to obtain the $O(\frac{1}{\sqrt{n}})$ bound we desire, we need to use the fact that $T_\pi^{\{U,L\}}$ — the number of batches between U and L — concentrates tightly. In fact, if we color the 4 clusters corresponding to membership in U and L and slightly generalize the analysis leading to Fact C.1.1, we obtain that for some absolute constant C , we have $|T_\pi^{\{U,L\}} - 4n| \leq 2C\sqrt{n \log n}$ with probability at least $1 - 1/n$.

Applying Lemma E.1.2, we obtain for all t such that $|t - 2n| \leq 2C\sqrt{n \log n}$, we have

$$\mathbb{P}[k_\pi^1 = 4n | T_\pi^{\{U,L\}} = 2t] = O\left(\frac{\frac{n}{t(4n-t)} 2^{8H(\frac{t}{4n})n}}{\sqrt{\frac{n}{t(4n-t)} 2^{8H(\frac{t}{4n})n}}}\right) \quad (36)$$

$$= O\left(\sqrt{\frac{n}{t(4n-t)}}\right) \quad (37)$$

$$= O\left(\frac{1}{\sqrt{n}}\right). \quad (38)$$

Hence we have

$$\begin{aligned} \mathbb{P}[k_\pi^1 = 4n] &= \sum_{t=0}^{4n} \mathbb{P}[k_\pi^1 = 4n | T_\pi^{\{U,L\}} = 2t] \mathbb{P}[T_\pi^{\{U,L\}} = 2t] \\ &\leq \frac{1}{n} + \sum_{|t-2n| \leq C\sqrt{n \log n}} \mathbb{P}[k_\pi^1 = 4n | T_\pi^{\{U,L\}} = 2t] \mathbb{P}[T_\pi^{\{U,L\}} = 2t] \\ &\leq \frac{1}{n} + O\left(\frac{1}{\sqrt{n}}\right) \\ &\leq O\left(\frac{1}{\sqrt{n}}\right), \end{aligned}$$

where in the second line we have used the union bound along with the fact that $T_\pi^{\{U,L\}}$ concentrates, and in the third line we have used Equation (38). This proves (2).

Quantitative SS distortion bounds with constant probability: Finally, we show (3). Suppose that $k_\pi^1 = 4n + d$ for $|d| = O(\sqrt{n})$. The above analysis for the case of $d = 0$ immediately generalizes to show that, if $t - d > 0$,

$$\mathbb{P}[k_\pi^1 = 4n + d | T_\pi^{\{U,L\}} = 2t - d] = \frac{\binom{4n}{t} \binom{4n}{t-d}}{\binom{8n}{2t-d}}.$$

Notice that since $2t - d$ concentrates around $4n$, it suffices to only consider the high probability regime where $2t - d = 4n + O(\sqrt{n \log n})$. In particular, since $|d| = O(\sqrt{n})$, we have $t - d = O(t)$.

Thus, if we plug in Lemma E.1.2, we obtain in the regime where $t - d = O(t)$ that

$$\frac{\binom{4n}{t} \binom{4n}{t-d}}{\binom{8n}{2t-d}} = O\left(\frac{1}{\sqrt{t}}\right) 2^{4n[H(\frac{t}{4n}) + H(\frac{t-d}{4n}) - 2H(\frac{2t-d}{8n})]}.$$

Concavity of binary entropy implies that $H(\frac{t}{4n}) + H(\frac{t-d}{4n}) - 2H(\frac{2t-d}{8n}) \leq 0$. It follows that

$$\mathbb{P}[k_\pi^1 = 4n + d | T_\pi^{\{U,L\}} = 2t - d] = O\left(\frac{1}{\sqrt{t}}\right).$$

Following the same argument as in the $d = 0$ case, we have for $|d| = O(\sqrt{n})$ that

$$\begin{aligned} \mathbb{P}[k_\pi^1 = 4n + d] &= \sum_{2t-d=0}^{4n} \mathbb{P}[k_\pi^1 = 4n | T_\pi^{\{U,L\}} = 2t - d] \mathbb{P}[T_\pi^{\{U,L\}} = 2t - d] \\ &\leq \frac{1}{n} + \sum_{|t-\frac{d}{2}-2n| \leq C\sqrt{n \log n}} \mathbb{P}[k_\pi^1 = 4n + d | T_\pi^{\{U,L\}} = 2t - d] \mathbb{P}[T_\pi^{\{U,L\}} = 2t - d] \\ &\leq \frac{1}{n} + O\left(\frac{1}{\sqrt{n}}\right) \\ &\leq O\left(\frac{1}{\sqrt{n}}\right), \end{aligned}$$

From here, it follows that there exists some positive constant c such that

$$\mathbb{P}[|k_\pi^1 - 4n| > c\sqrt{n}] = \Omega(1),$$

which proves (3). □

E.2. Classification toy dataset

In this section, we motivate how we constructed our toy classification dataset parameterized by n . We then give a detailed construction and analysis of the dataset, parameterized by n . We plot the unnormalized \mathbf{Z} in Figure 14a and the normalized $\overline{\mathbf{Z}}_{\text{GD}}$ in Figure 14b for $n = 10$.

The main idea is that since $d = 2$ and the optimal direction is orthogonal to the SC portion of the separability decomposition (Lemma 4.1.1), we can fix the optimal directions of $\overline{\mathbf{Z}}_\pi$ and $\overline{\mathbf{Z}}_{\text{GD}}$ by carefully constraining $\text{Span}(\overline{\mathbf{X}}_\pi^{\text{SC}})$ and $\text{Span}(\overline{\mathbf{X}}_{\text{GD}}^{\text{SC}})$, respectively. For, $\overline{\mathbf{Z}}_{\text{GD}}$ we carefully select the boundary points \mathbf{X}_{bdr} which define the boundary of $\text{conv}(\overline{\mathbf{X}}_{\text{GD}}^+)$ and $\text{conv}(\overline{\mathbf{X}}_{\text{GD}}^-)$ so that $\dim(\text{Span}(\overline{\mathbf{X}}_{\text{GD}}^{\text{SC}})) = 1$. For $\overline{\mathbf{Z}}_\pi$, we need to ensure that $\text{Span}(\overline{\mathbf{X}}_\pi^{\text{SC}})$ is a one dimensional subspace of \mathbb{R}^2 which is close to orthogonal with $\text{Span}(\overline{\mathbf{X}}_{\text{GD}}^{\text{SC}})$. Although we can guarantee $\dim(\text{Span}(\overline{\mathbf{X}}_\pi^{\text{SC}})) \geq 1$ w.h.p., the main subtlety here is ensuring that equality holds with constant probability. Given the above, we are afforded the luxury of adding datapoints which are misclassified by \mathbf{v}_π^* , the optimal direction of \mathcal{L}_π .

Proposition E.2.1. *There exists a classification dataset $\mathbf{Z} = (\mathbf{X}, \mathbf{Y}) \in [-3, 3]^{2 \times (2n+6)} \times \{-1, 1\}^{2n+6}$ such that the following statements hold with batch size $B = 2$:*

- (1) $\overline{\mathbf{Z}}_{\text{GD}}$ is PLS and as $n \rightarrow \infty$, \mathbf{v}_{GD}^* converges in direction to $[1 \quad 2]^\top$.

(2) $\overline{\mathbf{Z}}_\pi$ is PLS with constant probability. If so, we have $\mathbf{v}_\pi^* = [1 \quad -1]^\top$.

(3) There exists points $(\mathbf{x}_i, y_i) \in \overline{\mathbf{Z}}_{\text{GD}}$ such that $y_i \langle \mathbf{v}_\pi^*, \mathbf{x}_i \rangle < 0$, i.e., GD points that \mathbf{v}_π^* misclassifies.

Hence, the GD risk \mathcal{L}_{GD} diverges with constant probability if we train with SS.

Proof. In our construction, we separate out $\mathbf{Z} = (\mathbf{X}, \mathbf{Y})$ into several groups of points: $\mathbf{X}_{\text{cor}}^+$, $\mathbf{X}_{\text{err}}^+$, $\mathbf{X}_{\text{bdr}}^+$, and their negative variants. Let the overlined version of these matrices denote the corresponding points after full-batch BN (i.e., after taking $\overline{\mathbf{X}}_{\text{GD}} \triangleq \text{BN}(\mathbf{X})$), and the overlined version with an extra π subscript denoting the corresponding points after BN with permutation π (i.e., $\overline{\mathbf{X}}_\pi \triangleq \text{BN}_\pi(\mathbf{X})$). For example, $\overline{\mathbf{X}}_{\text{cor}}^+$ denotes the features for the positive examples in $\overline{\mathbf{X}}_{\text{GD}}$ that are to be classified correctly by \mathbf{v}_π^* , whereas $\overline{\mathbf{X}}_{\text{cor},\pi}^+$ denotes those same datapoints in $\overline{\mathbf{X}}_\pi$ after batch normalization with π .

Setup: Let us first explain the semantic meanings of the different groups of points in our construction.

$\mathbf{X}_{\text{cor}}^+$: Unnormalized positive examples correctly classified by \mathbf{v}_π^* with positive margin

$\mathbf{X}_{\text{err}}^+$: Unnormalized positive example incorrectly classified by \mathbf{v}_π^*

$\mathbf{X}_{\text{bdr}}^+$: Unnormalized positive examples on the decision boundary of \mathbf{v}_π^*

$\overline{\mathbf{X}}_{\text{cor}}^+$: full-batch-normalized positive examples correctly classified by \mathbf{v}_π^* with positive margin

$\overline{\mathbf{X}}_{\text{err}}^+$: full-batch-normalized positive example incorrectly classified by \mathbf{v}_π^*

$\overline{\mathbf{X}}_{\text{bdr}}^+$: full-batch-normalized positive examples on the decision boundary of \mathbf{v}_{GD}^*

The semantic meanings of the negative versions of these points are completely analogous. We also define $\mathbf{X}_{\text{bdr}}^- = \mathbf{X}_{\text{bdr}}^+ \cup \mathbf{X}_{\text{bdr}}^-$, and the normalized quantity analogously.

We construct $\mathbf{X}_{\text{cor}}^+$, $\mathbf{X}_{\text{err}}^+$, and $\mathbf{X}_{\text{bdr}}^+$ as follows. Take $\mathbf{X}_{\text{cor}}^+$ to be n equally spaced points on the line segment of width $\frac{1}{n}$ centered at $[2 \quad 2]^\top$. For the sake of visual clarity, we increase the spacing of the points in the diagram Figure 14a. Define

$\mathbf{X}_{\text{err}}^+$ to be $[3 \quad 2.5]^\top$. Next, define $\mathbf{X}_{\text{bdr}}^+$ to be $\begin{bmatrix} -3 & 1 \\ 1.5 & -0.5 \end{bmatrix}$, lying on the line $y = -0.5x$. Finally, define

$$\begin{aligned} \mathbf{X}_{\text{cor}}^- &= -\mathbf{X}_{\text{cor}}^+ \\ \mathbf{X}_{\text{err}}^- &= -\mathbf{X}_{\text{err}}^+ \\ \mathbf{X}_{\text{bdr}}^- &= -\mathbf{X}_{\text{bdr}}^+ \end{aligned}$$

In Figure 14, we visualize this toy dataset. Note the visual similarity between \mathbf{Z} in Figure 14a and $\overline{\mathbf{Z}}_{\text{GD}}$ in Figure 14b; this is a feature of the construction. Indeed, as we'll see shortly, as $n \rightarrow \infty$, $\overline{\mathbf{X}}_{\text{GD}}$ approaches a uniform rescaling of \mathbf{X} in all coordinates. We also plotted the decision boundaries corresponding to \mathbf{v}_{GD}^* and \mathbf{v}_π^* . We highlight the fact that in Figure 14b, $\overline{\mathbf{X}}_{\text{err}}^+$ and $\overline{\mathbf{X}}_{\text{err}}^-$ are both on the wrong side of the decision boundary for \mathbf{v}_π^* .

GD is PLS: Evidently $\boldsymbol{\mu} = \mathbf{0}$ and one can compute that as $n \rightarrow \infty$ that $\boldsymbol{\sigma} \rightarrow [2 \quad 2]^\top$. Regardless, we see that $\overline{\mathbf{Z}}_{\text{GD}}$ is PLS with $\text{Span}(\overline{\mathbf{X}}_{\text{GD}}^{\text{SC}})$ corresponding to the one-dimensional subspace spanned by $\overline{\mathbf{X}}_{\text{bdr}}^-$ and \mathbf{v}_{GD}^* correctly classifies all of the other points. In the limit $n \rightarrow \infty$, we have that $\text{Span}(\overline{\mathbf{X}}_{\text{bdr}}^-) = \text{Span}([-2 \quad 1]^\top)$ and \mathbf{v}_{GD}^* is in the direction $[1 \quad 2]^\top$. This proves (1).

SS is PLS with constant probability: Now, let us compute what happens to $\overline{\mathbf{Z}}_\pi$. Because BN with $B = 2$ sends features to ± 1 , this implies that a normalized batch is either mapped to

(1) $\begin{bmatrix} -1 & +1 \\ -1 & +1 \end{bmatrix}$, i.e. the normalized batch lies in the direction $[+1 \quad +1]^\top$; or

(2) $\begin{bmatrix} -1 & +1 \\ +1 & -1 \end{bmatrix}$, i.e the normalized batch lies in the direction $[+1 \quad -1]^\top$.

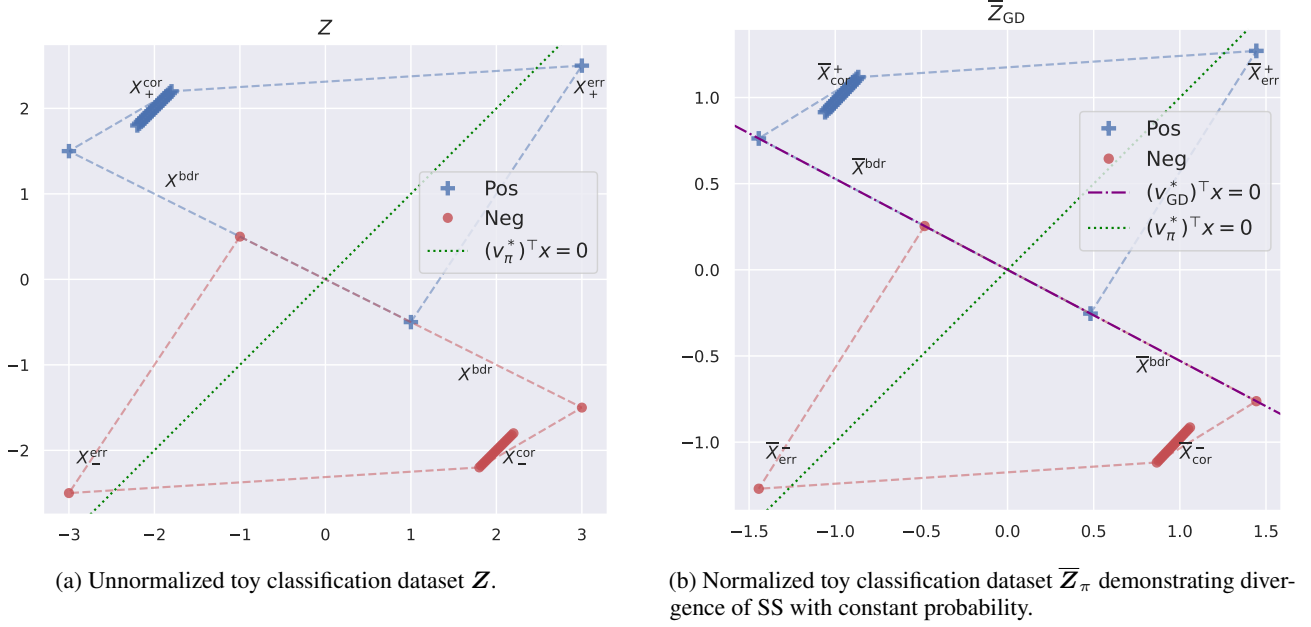


Figure 14. Toy classification dataset (a) before full-batch BN, i.e. Z (b) after full-batch BN, i.e. \bar{Z}_{GD} .

Note that due to Fact C.1.1, with high probability, there will be a batch drawn from X_{cor}^+ and a batch drawn from X_{cor}^- . These batches necessarily land in situation (1), so $[1 \ 1]^\top \in \text{Span}(\bar{X}_\pi^{\text{SC}})$ with high probability. On the other hand, with high probability there will be a batch with one point from X_{cor}^+ and one from X_{cor}^- , which lands in (2). Hence, as long as $\text{Span}(\bar{X}_\pi^{\text{SC}}) = \text{Span}([1 \ 1]^\top)$, this implies that \bar{Z}_π is PLS with optimal direction $[-1 \ +1]^\top$.

By inspection, to ensure that $\text{Span}(\bar{X}_\pi^{\text{SC}}) = \text{Span}([1 \ 1]^\top)$, there are two bad events we need to avoid:

- (a) We send a positive example to $[+1 \ -1]^\top$.
- (b) We send a negative example to $[-1 \ +1]^\top$.

We will use

$$\mathbb{P}[\text{avoid (a) and (b)}] = \mathbb{P}[\text{avoid (a)} \mid \text{avoid (b)}] \mathbb{P}[\text{avoid (b)}].$$

Note that by symmetry, $\mathbb{P}[\text{avoid (b)}] = \mathbb{P}[\text{avoid (a)}]$.

A little thought reveals that (a) can happen only if the positive boundary example (i.e. in X_{bdr}^+) located originally at $[1 \ -0.5]^\top$ is batched together with a point in $X_{\text{bdr}} \cup X_{\text{cor}}^+$. In turn, this event occurs with probability at most $\frac{2}{3}$. So $\mathbb{P}[\text{avoid (b)}] \geq \frac{1}{3}$. Also, notice avoiding (a) still happens with probability at least $\frac{1}{3}$ even after conditioning on avoiding (b). Hence, the probability that we avoid both (a) and (b) is at least $\frac{1}{3}^2 = \frac{1}{9}$. So with constant probability, $\text{Span}(\bar{X}_\pi^{\text{SC}}) = \text{Span}([1 \ 1]^\top)$. This proves (2).

Putting it all together, we see that with constant probability, $\text{Span}(\bar{X}_\pi^{\text{SC}}) = \text{Span}([1 \ 1]^\top)$, and \bar{Z}_π is PLS with optimal direction $\mathbf{v}_\pi^* = [+1 \ -1]^\top$. Recall that $\mathbf{v}_{GD}^* \rightarrow [1 \ 2]^\top$ as $n \rightarrow \infty$. So asymptotically we have

$$\frac{|\langle \mathbf{v}_\pi^*, \mathbf{v}_{GD}^* \rangle|}{\|\mathbf{v}_\pi^*\| \|\mathbf{v}_{GD}^*\|} = \frac{1}{\sqrt{10}}.$$

SS misclassifies GD points: On the constant probability event that \bar{Z}_π is PLS, we have \mathbf{v}_π^* is in the direction $[1 \ -1]^\top$. This misclassifies \bar{X}_{err}^+ and \bar{X}_{err}^- . So this proves (3). \square

Remark E.2.2. *Note that at the cost of visual clarity, we can modify the construction to obtain optimal classifiers \mathbf{v}_π^* and \mathbf{v}_{GD}^* which are asymptotically orthogonal. In this alternate construction, we take $\mathbf{X}_{\text{bdr}}^+ = \begin{bmatrix} -1 & 0.5 \\ 1 & -0.5 \end{bmatrix}$, lying on the line $y = -x$, and $\mathbf{X}_{\text{cor}}^+$ to be n equally spaced points on the line segment between $[-2 \ 2]^\top$ and $[-2 + \frac{1}{n} \ 2 + \frac{1}{n}]^\top$. Then $\mathbf{v}_{\text{GD}}^* \rightarrow [1 \ 1]^\top$, and $\mathbf{v}_\pi^* = [1 \ -1]^\top$ with constant probability, as desired.*