

# A Walk with SGD

Chen Xing<sup>\*†‡</sup> Devansh Arpit<sup>\*§</sup> Christos Tsirigotis<sup>¶‡</sup> Yoshua Bengio<sup>‡||</sup>

## Abstract

We present novel empirical observations regarding how stochastic gradient descent (SGD) navigates the loss landscape of over-parametrized deep neural networks (DNNs). These observations expose the qualitatively different roles of learning rate and batch-size in DNN optimization and generalization. Specifically we study the DNN loss surface along the trajectory of SGD by interpolating the loss surface between parameters from consecutive *iterations* and tracking various metrics during training. We find that the loss interpolation between parameters before and after each training iteration’s update is roughly convex with a minimum (*valley floor*) in between for most of the training. Based on this and other metrics, we deduce that for most of the training update steps, SGD moves in valley like regions of the loss surface by jumping from one valley wall to another at a height above the valley floor. This ‘bouncing between walls at a height’ mechanism helps SGD traverse larger distance for small batch sizes and large learning rates which we find play qualitatively different roles in the dynamics. While a large learning rate maintains a large height from the valley floor, a small batch size injects noise facilitating exploration. We find this mechanism is crucial for generalization because the valley floor has barriers and this exploration above the valley floor allows SGD to quickly travel far away from the initialization point (without being affected by barriers) and find flatter regions, corresponding to better generalization.

## 1 Introduction

Deep neural networks (DNNs) trained with algorithms based on stochastic gradient descent (SGD) are able to tune the parameters of massively over-parametrized models to reach small training loss with good generalization despite the existence of numerous bad minima. This is especially surprising given DNNs are capable of overfitting random data with almost zero training loss Zhang et al. (2016). This behavior has been studied by Arpit et al. (2017); Advani and Saxe (2017) where they suggest that deep networks generalize well because they tend to fit simple functions over training data before overfitting noise. It has been further discussed that model parameters that are in a region of flatter minima generalize better Hochreiter and Schmidhuber (1997); Keskar et al. (2016); Wu et al. (2017), and that SGD finds such minima when used with small batch size and large learning rate Keskar et al. (2016); Jastrzębski et al. (2017); Smith et al. (2017); Chaudhari and Soatto (2017). These recent papers frame SGD as a stochastic differential equation (SDE) under the assumption of using small learning rates. A main result of these papers is that the SDE dynamics remains the same as long as the ratio of learning rate to batch size remains unchanged. However, this view is limited due to its assumption and ignores the importance of the structure of SGD noise (i.e., the gradient covariance) and the qualitative roles of learning rate and batch size, which remain relatively obscure.

<sup>\*</sup>Equal Contribution

<sup>†</sup>College of Computer and Control Engineering, Nankai University, xingchen1113@gmail.com

<sup>‡</sup>MILA, Université de Montréal

<sup>§</sup>MILA, Université de Montréal, devansharpit@gmail.com

<sup>¶</sup>Aristotle University of Thessaloniki, Greece

<sup>||</sup>CIFAR Senior Fellow

On the other hand, various variants of SGD have been proposed for optimizing deep networks with the goal of addressing some of the common problems found in the high dimensional non-convex loss landscapes (Eg. saddle points, faster loss descent etc). Some of the popular algorithms used for training deep networks apart from vanilla SGD are SGD with momentum Polyak (1964); Sutskever et al. (2013), AdaDelta Zeiler (2012), RMSProp Tieleman and Hinton (2012), Adam Kingma and Ba (2014) etc. However, for any of these methods, currently there is little theory proving they help in improving generalization in DNNs (which by itself is currently not very well understood), although there have been some notable efforts (Eg. Hardt et al. (2015); Kawaguchi et al. (2017)). *This raises the question of whether optimization algorithms that are designed with the goal of solving the aforementioned high dimensional problems also help in finding minima that generalize well, or put differently, what attributes allow optimization algorithms to find such good minima in the non-convex setting.*

We take a step towards answering the above two questions in this paper for SGD (without momentum) through qualitative experiments. The main tool we use for studying the DNN loss surface along SGD's path is to interpolate the loss surface between parameters before and after each training update and track various metrics. Our findings about SGD's trajectory can be summarized as follows:

1. We observe that the loss interpolation between parameters before and after each *iteration's* update is roughly convex with a minimum (*valley floor*) in between. Thus, we deduce that SGD bounces off walls of a *valley-like-structure* at a height above the floor.
2. Learning rate controls the height at which SGD bounces above the valley floor while batch size controls gradient stochasticity which facilitates exploration (visible from larger parameter distance from initialization for small batch-size). In this way, learning rate and batch size exhibit different qualitative roles in SGD dynamics.<sup>7</sup>
3. The valley floor along SGD's path has many ups and downs (barriers) which may hinder exploration. Thus using a large learning rate helps avoid encountering barriers along SGD's path by maintaining a large height above valley floor (thus moving over the barriers instead of crossing<sup>8</sup> them).

Experiments are conducted on multiple data sets, architectures and hyper-parameter settings. The findings mentioned above hold true on all of them. We further find that stochasticity in SGD induced by mini-batches is needed both for better optimization and generalization. Conversely, artificially added isotropic noise in the absence of mini-batch induced stochasticity is bad for DNN optimization. We also discuss some striking similarities between our empirical findings about SGD's trajectory in DNNs and classical optimization theory in the quadratic setting.

## 2 Background and Related Work

Various algorithms have been proposed for optimizing deep neural networks that are designed from the view point of tackling various high dimensional optimization problems like oscillation during training (SGD with momentum Polyak (1964)), oscillations around minima (Nesterov momentum Nesterov (1983); Sutskever et al. (2013)), saddle points Dauphin et al. (2014), automatic decay of the learning rate (AdaDelta Zeiler (2012), RMSProp Tieleman and Hinton (2012) and ADAM Kingma and Ba (2014)), etc. However, currently there is insufficient theory to understand what kind of minima generalize better although empirically it has been observed that wider minima (that can be quantified by low Hessian norm) seem to have better generalization Keskar et al. (2016); Wu et al. (2017); Jastrzębski et al. (2017) due to their low complexity and are more likely to be reached under random initialization given their larger volumes Wu et al. (2017). *This argument raises the question of whether the intuitions behind the designs of the various optimization algorithms are really the reasons behind their success in deep learning or there are other underlying mechanisms that make them successful.*

To understand this aspect better, a number of (mostly) recent papers study SGD as a stochastic differential process Kushner and Yin (2003); Mandt et al. (2017); Chaudhari and Soatto (2017); Smith and Le (2017); Jastrzębski et al. (2017); Li et al. (2015) under the assumption (among

<sup>7</sup>This implies that except when using a reasonably small learning rate (which would make the SDE approximation hold), the effect of small batch size with a certain learning rate cannot be achieved by using a large batch size with a proportionally large learning rate (observed by Goyal et al. (2017)).

<sup>8</sup>A **barrier is crossed** when we see a point in the parameter space interpolated between the points just before and after an update step, such that the loss at the barrier point is higher than the loss at both the other points.

others) that the learning rate is reasonably small. Broadly, these papers show that the stochastic fluctuation in the stochastic differential equation **simulated by SGD is governed by the ratio of learning rate to batch size**. Hence according to this theoretical framework, the training dynamics of SGD should remain roughly identical when changing learning rate and batch size by the same factor. However, given DNNs (especially Resnet He et al. (2016) like architectures) are often trained with quite large learning rates, the small learning rate assumption may be a pitfall of this theoretical framework<sup>9</sup>. But this theory is nonetheless useful since learning rates do attain small values during training due to annealing or adaptive scheduling, so this framework may indeed apply during parts of training. *In this paper we attempt to go beyond these analyses to study the different qualitative roles of the noise induced by large learning rate versus the noise induced by a small batch size.*

There have also been work that consider SGD as a diffusion process where SGD is running a Brownian motion in the parameter space. Li et al. (2017) hypothesize this behavior of SGD and theoretically show that this diffusion process would allow SGD to cross barriers and thus escape sharp local minima. The authors use this theoretical result to support the findings of Keskar et al. (2016) who find that SGD with small mini-batch find wider minima. Hoffer et al. (2017) on the other hand make a similar hypothesis based on the evidence that the distance moved by SGD from initialization resembles a diffusion process, and make a similar claim about SGD crossing barriers during training. *Contrary to these claims, we find that interpolating the loss surface traversed by SGD on a per iteration basis suggests SGD almost never crosses any significant barriers for most of the training.*

There is also a long list of work towards understanding the loss surface geometry of DNNs from a theoretical standpoint. Dotsenko (1995); Amit et al. (1985); Choromanska et al. (2015) show that under certain assumptions, the DNN loss landscape can be modeled by a spherical spin glass model which is well studied in terms of its critical points. Safran and Shamir (2016) show that under **certain mild assumptions, the initialization is likely to be such that there exists a continuous monotonically decreasing path from the initial point to the global minimum**. Freeman and Bruna (2016) theoretically show that for DNNs with rectified linear units (ReLU), the level sets of the loss surface become more connected as network over-parametrization increases. This has also been justified by Sagun et al. (2017) who show that the hessian of deep ReLU networks is degenerate when the network is over-parametrized and hence the loss surface is flat along such degenerate directions. Goodfellow et al. (2014) empirically show that the convex interpolation of the loss surface from the initialization to the final parameters found by optimization algorithms do not cross any significant barriers, and that the landscape of loss

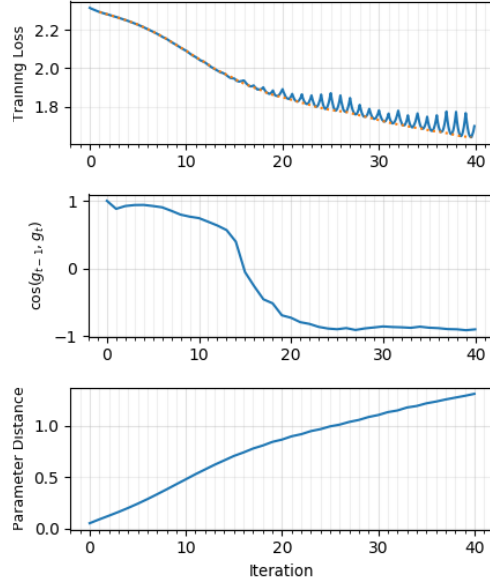


Figure 1: Plots for VGG-11 Epoch 1 trained using full batch **Gradient Descent (GD)** on CIFAR-10. **Top:** Training loss for the 1st 40 iterations of training. Between the training loss at every consecutive iteration (vertical gray lines), we uniformly sample 10 points between the parameters before and after a training update and calculate the loss at these points. Thus we take a slice of the loss surface between two iterations. These loss values are plotted between every consecutive training loss value from training updates. The dashed orange line connects the minimum of the loss interpolation between consecutive iterations (this minimum denotes the valley floor along the interpolation). **Middle:** Cosine of the angle between gradients from two consecutive iterations. **Bottom:** Parameter distance from initialization. **Gist:** The loss interpolation between consecutive iterations have a minimum for iterations where cosine is highly negative (close to  $-1$  after around 20 iterations meaning the consecutive gradients are almost along opposite directions), suggesting the optimization is oscillating along the walls of a valley like structure. The valley floor reduces monotonously.

<sup>9</sup>For instance Goyal et al. (2017) investigate that increasing learning rate linearly with batch size helps to a certain extent but breaks down for very large learning rates.

surface near SGD’s trajectory has a valley-like 2D projection. *Broadly these studies analyze DNN loss surfaces (either theoretically or empirically) in isolation from the optimization dynamics.*

In our work we do not study the loss surface in isolation, but rather analyze it through the lens of SGD. In other words, we study the DNN loss surface along the trajectory of SGD and track various metrics while doing so, from which we deduce both how the landscape relevant to SGD looks like, and how the hyperparameters of SGD (learning rate and batch size) help SGD maneuver through it.

### 3 A Walk with SGD

We now begin our analysis of studying the loss surface of DNNs along the trajectory of optimization updates. Specifically, consider that the parameters  $\theta$  of a DNN are initialized to a value  $\theta_0$ . When using an optimization method to update these parameters, the  $t^{th}$  update step takes the parameter from  $\theta_t$  to  $\theta_{t+1}$  using estimated gradient  $\mathbf{g}_t$  as,

$$\theta_{t+1} = \theta_t - \eta \mathbf{g}_t \quad (1)$$

where  $\eta$  is the learning rate. Notice the  $t^{th}$  update step implies the  $t^{th}$  epoch *only* in the case when using the full batch gradient descent (GD– gradient computed using the whole dataset). In the case of stochastic gradient descent, one iteration is an update from gradient computed from a mini-batch. We then interpolate the DNN loss between the convex combination of  $\theta_t$  and  $\theta_{t+1}$  by considering parameter vectors  $\theta_t^\alpha = (1 - \alpha)\theta_t + \alpha\theta_{t+1}$ , where  $\alpha \in [0, 1]$  is chosen such that we obtain 10 samples uniformly placed between these two parameter points. Simultaneously, we also keep track of two metrics– the cosine of the angle between two consecutive gradients  $\cos(\mathbf{g}_{t-1}, \mathbf{g}_t) := \mathbf{g}_{t-1}^T \mathbf{g}_t / (\|\mathbf{g}_{t-1}\|_2 \|\mathbf{g}_t\|_2)$ , and the distance of the current parameter  $\theta_t$  from the initialization  $\theta_0$  given by  $\|\theta_t - \theta_0\|_2$ . As it will become apparent in a bit, these two metrics along with the interpolation curve help us make deductions about how the optimization *interacts* with the loss surface during its trajectory.

We perform experiments on MNIST Lecun and Cortes, CIFAR-10 Krizhevsky (2009) and a subset of the tiny Imagenet dataset Russakovsky et al. (2015) using multi-layer perceptrons (MLP), VGG-11 Simonyan and Zisserman (2014) and Resnet-56 He et al. (2016) architectures with various batch sizes and learning rates. In the main text, we mostly show results for CIFAR-10 using VGG-11 architecture with a batch size of 100 and a fixed learning rate of 0.1 due to space limitations. Experiments on all the other datasets, architectures and hyper-parameter settings can be found in the appendix. All the claims are consistent across them.

#### 3.1 Optimization Trajectory

We first experiment with full batch gradient descent (GD) to study its behavior before jumping to the analysis of SGD to isolate the confounding factor of mini-batch induced stochasticity. The plot of training loss interpolation between consecutive iterations (referred in the figure as *training loss*),  $\cos(\mathbf{g}_{t-1}, \mathbf{g}_t)$ , and *parameter distance*  $\|\theta_t - \theta_0\|_2$  for CIFAR-10 on VGG-11 architecture optimized using full batch gradient descent is shown in Figure 1 for the first 40 iterations of training. To be clear, the x-axis is calibrated by the number of iterations, and there are 10 interpolated loss values between each consecutive iterations (vertical gray lines) in the *training loss* plot which is as described

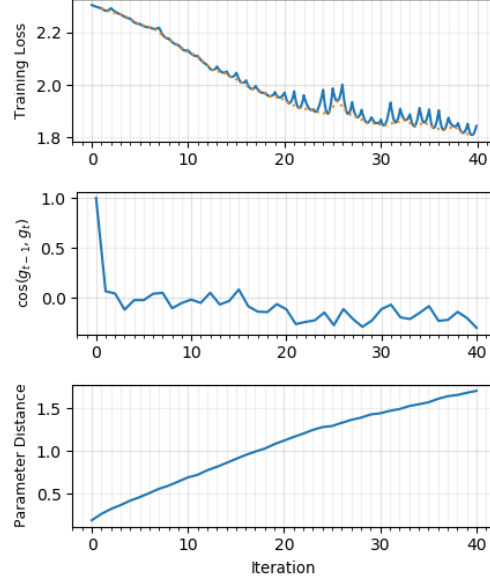


Figure 2: Plots for VGG-11 Epoch 1 trained using **SGD** on CIFAR-10. The descriptions of the plots are same as in Figure 1. **Gist:** The loss interpolation between consecutive iterations have a minimum for iterations and cosine is less negative compared with GD, suggesting the optimization is oscillating along the walls of a valley like structure but doing more exploration compared with GD. This is verified by the larger distance traveled by SGD compared with GD in Figure 1. The valley floor (dashed orange line) has many ups and downs showing barriers along SGD’s path which do not affect its dynamics because SGD travels at a height above the floor.

above (the cosine and parameter distance plots do not have any interpolations). This figure shows that the interpolated loss between every consecutive parameters from GD optimization update after iteration 15 appears to be a quadratic-like structure with a minimum in between.

Additionally, the cosine of the angle between consecutive gradients after iteration 15 is going negative and finally very close to  $-1$ , which means the consecutive gradients are almost along opposite directions. These two observations together suggest that the GD iterate is bouncing between walls of a valley-like landscape. For the iterations where there is a minimum in the interpolation between two iterations, we refer to this minimum as the *floor* of the valley (these valley floors are connected by dashed orange line in figure 1 for clarity). We will see GD behavior shows lack of exploration for better minima by comparing its parameter distance from initialization with that of SGD. Take note that the parameter distance of GD during these 40 iterations reaches  $\sim 1.4$ .

Now we perform the same analysis for SGD. Notice that even though the updates are performed using mini-batches for SGD, the training loss values used in the plot are computed using the full dataset to visualize the actual loss landscape. We show these plots for epoch 1 (Figure 2) in the main text and epoch 25 (Figure 16) and epoch 100 (Figure 17) in appendix. We find that while the loss interpolation also shows a quadratic-like structure with a minimum in between (similar to GD), there are some qualitative differences compared with GD. We see that the cosine of the angle between gradients from consecutive iterations are significantly less negative, suggesting that instead of oscillating in the same region, SGD is quickly moving away from its previous position. This can be verified by the parameter distance from initialization. We see that the distance after 40 iterations is  $\sim 1.7$ , which is larger than the distance moved by GD<sup>10</sup>. Finally and most interestingly, we see that the height of valley floor has many ups and downs for consecutive iterations in contrast with that of GD (emphasized by the dashed orange line in figure 2), which means that there is a rough terrain or barriers along the path of SGD that could hinder exploration if the optimization was traveling too close to the valley floor.

Arch\Epochs	1	10	25	100
VGG-11	0	0	5	13
Resnet-56	0	0	2	23
MLP	0	3	5	-

Table 1: Number of barriers crossed during training of one epoch (450 iterations) for VGG-11 and Resnet-56 on CIFAR-10 and MLP on MNIST. We say a barrier is crossed during an update step if there exists a point interpolated between the parameters before and after an update which has a loss value higher than the loss at either points. For most parts of the training, we find that SGD does not cross any significant number of barriers.

1 for VGG-11 and Resnet-56 trained on CIFAR-10 (trained for 100 epochs) and an MLP trained on MNIST (trained for 40 epochs). As we see, no barriers are crossed for most parts of the training. We further compute the number of barriers crossed for the first 40 epochs for VGG-11 on CIFAR-10 shown in Figure 9 in the appendix: no barriers are crossed for most of the epochs and even for the barriers that are crossed towards the end, we find that their heights are substantially smaller compared with the loss value at the corresponding point during training, meaning they are not significant.

Finally, we track the spectral norm of the Hessian along with the validation accuracy while the model is being trained<sup>11</sup>. This plot is shown in Figure 3 for VGG-11 (and Figure 6 for Resnet-56 in the appendix). We find that the spectral norm reduces as training progresses (hence SGD finds flatter regions) but starts increasing towards the end. This is mildly correlated with a drop in validation accuracy towards the end. Regarding this correlation, while Dinh et al. (2017) discuss that sharper

A similar analysis for Resnet-56 on CIFAR-10, MLP on MNIST, VGG-11 on tiny ImageNet trained using GD for the first epoch are shown in Figures 10, 18, 21 respectively in appendix. The same experiments analysis for SGD on different datasets, architectures under different hyper-parameters are also shown in section 1 in appendix. The observations and rules we discovered and described here are all consistent for all these experiments.

In order to show that the claim about optimization not crossing barriers extends to the whole training instead of only a few iterations we’ve shown, we quantitatively measure for the entire epoch in different phase of training if barriers are crossed. This result is shown in table

<sup>10</sup>In general, after the same number of updates, GD traverses a smaller distance compared with SGD, see Hoffer et al. (2017)

<sup>11</sup>Note that we track the spectral norm in the train mode of batchnorm; we observe that in validation mode the values are significantly larger. Tracking the value in train mode is fair because this is what SGD *experiences* during training. Additionally, we track spectral norm because it captures the largest eigenvalue of the Hessian in contrast with Frobenius norm which can be misleading because the Hessian may have negative eigenvalues and the Frobenius norm sums the square of all eigenvalues.



minima can perform as well as wider ones, it is empirically known that flatter minima generalize better than sharper ones with SGD (Keskar et al., 2016; Jastrzebski et al., 2017). This may be explained by Neyshabur et al. (2017); Achille and Soatto (2017) that discuss that minima that are both wide and have small norm may explain generalization in over-parametrized deep networks.

### 3.2 Qualitative Roles of Learning Rate and Batch Size

We now focus in more detail on how the learning rate and batch size play qualitatively different roles during SGD optimization. As an extreme case, we already saw in the last section that when using GD vs SGD, the cosine of the angle between gradients from two consecutive iterations  $\cos(\mathbf{g}_{t-1}, \mathbf{g}_t)$  is significantly closer to -1 (180 degrees) in the case of GD in contrast with SGD. Now we show on a more granular scale that changing the batch size gradually (keeping the learning rate fixed) changes  $\cos(\mathbf{g}_{t-1}, \mathbf{g}_t)$ , while changing the learning rate gradually (keeping batch size fixed) does not. It is shown in Figure 4 for VGG-11 trained on CIFAR-10 and in Figures 31, 32 and 33 for Resnet-56 with CIFAR-10, MLP with MNIST and VGG-11 with tiny ImageNet separately in the appendix. Notice the cosine is significantly more negative for larger batch sizes, implying that for larger batch sizes, the optimization is bouncing more within the same region instead of traveling farther along the valley as the case of small batch sizes. This behavior is verified by the smaller distance of parameters from initialization during training for larger batch size which is also discussed by Hoffer et al. (2017). This suggests that the noise from a small mini-batch size facilitates exploration that may lead to better minima and that this is hard to achieve by changing the learning rate.

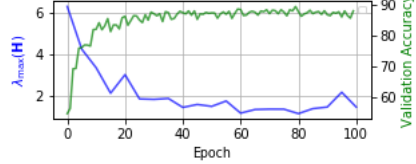


Figure 3: Spectral norm of the Hessian and validation accuracy for VGG-11 trained on CIFAR-10 using a fixed learning rate of 0.1 and batch size 100. Spectral norm and validation accuracy are roughly inversely correlated.

On the other hand, we find that the learning rate controls the height from the valley floor at which the optimization oscillates along the valley walls which is important for avoiding barriers along SGD’s path. Specifically, to quantify the height at which the optimization is bouncing above the valley floor, we make the following computations. Suppose at iterations  $t$  and  $t + 1$  of training, the parameters are given by  $\theta_t$  and  $\theta_{t+1}$  respectively, and from the 10 points sampled uniformly between  $\theta_t$  and  $\theta_{t+1}$  given by  $\theta_t^\alpha = (1 - \alpha)\theta_t + \alpha\theta_{t+1}$  for different values of  $\alpha \in [0, 1]$ , we define  $\theta_t^{\min} := \theta_t^{\arg \min_\alpha \mathcal{L}(\theta_t^\alpha)}$ , where  $\mathcal{L}(\theta)$  denotes the DNN loss at parameter  $\theta$  using the whole training set. Then we define the height of the iterate from the valley floor at iteration  $t$  as  $\frac{\mathcal{L}(\theta_t) + \mathcal{L}(\theta_{t+1}) - 2\mathcal{L}(\theta_t^{\min})}{2}$ . We then separately compute the average height for all iterations of epochs 1, 10, 25 and 100.

VGG-11	Epoch 1	Epoch 10	Epoch 25
<b>LR 0.1</b>	0.0625 $\pm$ 1.2e-3	0.0199 $\pm$ 4.8e-4	0.0104 $\pm$ 2.2e-5
<b>LR 0.05</b>	0.0102 $\pm$ 3.8e-5	0.0050 $\pm$ 2.8e-5	0.0035 $\pm$ 1.7e-5
Resnet-56	Epoch 1	Epoch 10	Epoch 25
<b>LR 0.3</b>	0.0380 $\pm$ 5.9e-4	0.0131 $\pm$ 4.5e-4	0.0094 $\pm$ 1.3e-5
<b>LR 0.15</b>	0.0084 $\pm$ 5.2e-5	0.0034 $\pm$ 3.2e-5	0.0020 $\pm$ 7.2e-6

Table 2: Average height of SGD above valley floor across the iterations in one epoch for different epochs during the training of VGG-11 and Resnet-56 using SGD on CIFAR-10. Here the height at iteration  $t$  is defined by  $\frac{\mathcal{L}(\theta_t) + \mathcal{L}(\theta_{t+1}) - 2\mathcal{L}(\theta_t^{\min})}{2}$ . Results of Epoch 100 and of VGG-11 on Tiny-ImageNet can be found at Table 3 in the appendix.

poor generalization Smith et al. (2017); Jastrzebski et al. (2017). Based on our observations on the role of learning rate and batch-size, we empirically study learning rate schedules in appendix C.

## 4 Importance of SGD Noise Structure

The gradient  $\mathbf{g}_{SGD}(\theta)$  from mini-batch SGD at a parameter value  $\theta$  is expressed as,  $\mathbf{g}_{SGD}(\theta) = \bar{\mathbf{g}}(\theta) + \frac{1}{\sqrt{B}}\mathbf{n}(\theta)$ , where  $\mathbf{n}(\theta) \sim \mathcal{N}(0, \mathbf{C}(\theta))$ ,  $\bar{\mathbf{g}}(\theta)$  denotes the expected gradient using all training samples,  $B$  is the mini-batch size and  $\mathbf{C}(\theta)$  is the gradient covariance matrix at  $\theta$ . In the previous

These values are shown in table 2 for VGG-11 and Resnet-56 trained on CIFAR-10 and in table 3 in appendix for VGG-11 on Tiny-ImageNet. They show that for almost all epochs, a smaller learning rate leads to a smaller height from the valley floor. Since the floor has barriers, it would increase the risk of hindering exploration for flatter minima. This has been corroborated by the recent empirical observations that smaller learning rates lead to sharper minima and

section we discussed how mini-batch induced stochasticity plays a crucial role in SGD based optimization. This stochasticity due to SGD has historically been attributed to helping the optimization escape local minima in DNNs. However, the importance of the structure of the gradient covariance matrix  $\mathbf{C}(\theta)$  is often neglected in these claims. To better understand its importance, we study the training dynamics of full batch gradient with artificially added isotropic noise. Specifically, we treat isotropic noise as our null hypothesis to confirm that the structure of noise induced by mini-batches in SGD is important.

In this experiment, we first train our models with gradient descent (GD), meaning there is no noise sampled from  $\mathcal{N}(0, \mathbf{C}(\theta))$  of SGD. For gradient descent with isotropic noise, we add isotropic noise at every iteration on  $\bar{\mathbf{g}}(\theta)$ . The noise is sampled from a normal distribution with variance calculated by multiplying the maximum gradient variance of the model at the initialization with a factor of 0.1 and 0.05. We train all models until their training losses saturate and monitor training loss, validation accuracy, cosine of the angle between gradients from two consecutive iterations and the parameter distance from initialization.

Figure 5 shows the results for VGG-11 and Figures 34 and 35 in the appendix shows the results for Resnet-56 on CIFAR-10 and MLP on MNIST. From the training loss and validation accuracy curves, we can see that adding even a small isotropic noise makes both the convergence<sup>12</sup> and generalization worse compared with the model trained with GD. The cosine of the angle between gradients of two consecutive iterations is close to  $-1$  for 1500 iterations for GD, which means two consecutive gradients are almost along opposite directions. It is an extra evidence that GD makes the optimization bounce off valley walls, which is what we discussed in section 3. The parameter distance from initialization shows that models trained with isotropic noise travel farther away compared with the model trained using noiseless GD. These distances are much larger even compared with models trained with SGD (not shown here) for the same number of updates.

To gain more intuitions into this behavior, we also calculate the norm of the final parameters found by GD, SGD and the isotropic noise cases. The parameter norms are 87 and 82 for GD and SGD respectively, and 369 and 443 for the 0.014 and 0.028 isotropic variance case. These numbers corroborate the generally discussed notion that SGD finds solutions with small  $\ell^2$  norm Zhang et al. (2016) compared with GD, and the fact that isotropic noise solutions have much larger norms and get stuck suggests that isotropic noise both hinders optimization and is bad for generalization.

Neelakantan et al. (2015) suggest adding isotropic noise to gradients and report performance improvement on a number of tasks. However, notice the crucial difference between our claim in this section and their setup is that they add isotropic noise on top of the noise due to the mini-batch induced stochasticity, while we add isotropic noise to the full dataset gradient (hence no noise is sampled from the gradient covariance matrix).

To gain insights into why the noise sampled from the gradient covariance matrix  $\mathbf{C}(\theta)$  helps SGD, we note that there is a relationship between the covariance  $\mathbf{C}(\theta)$  and the Hessian  $\mathbf{H}(\theta)$  of the loss surface at parameter  $\theta$  which is revealed by the generalized Gauss Newton decomposition (see Sagun et al. (2017)) when using the cross-entropy (or negative log likelihood) loss. Let  $p_i(\theta)$  denote the predicted probability output (of the correct class in the classification setting for instance) of a DNN parameterized by  $\theta$  for the  $i^{th}$  data sample (in total  $N$  samples). Then the negative log likelihood loss for the  $i^{th}$  sample is given by,  $\mathcal{L}_i(\theta) = -\log(p_i(\theta))$ . The relation between the Hessian  $\mathbf{H}(\theta)$  and

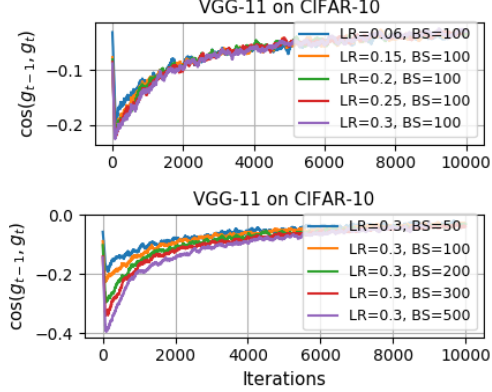


Figure 4: Changing batch size changes the cosine of angle between consecutive gradients while changing learning rate does not have any significant effect on the cosine. This shows batch size has a qualitatively different role compared with learning rate. Note that the curves are smoothed for visual clarity.

<sup>12</sup>We additionally find that the model trained with isotropic noise gets stuck because we find that neither reducing learning rate, nor switching to GD at this point leads to reduction in training loss. However, switching to SGD makes the loss go down.

the gradient covariance  $\mathbf{C}(\theta)$  for negative log likelihood loss is,

$$\mathbf{H}(\theta) = \mathbf{C}(\theta) + \bar{\mathbf{g}}(\theta)\bar{\mathbf{g}}(\theta)^T + \frac{1}{N} \sum_{i=1}^N \frac{\partial \mathcal{L}_i(\theta)}{\partial p_i(\theta)} \cdot \frac{\partial^2 p_i(\theta)}{\partial \theta^2}$$

The derivation can be found in section D of the appendix. Thus we find that the Hessian and covariance at any point  $\theta$  are related, and are almost equal near minima where the second term tends to zero. This relationship would imply that the mini-batch induced noise is roughly aligned with sharper directions of the loss landscape (empirically confirmed concurrently by Zhu et al. (2018)). This would prevent the optimization from converging along such directions unless a wider region is found, which could explain why SGD finds wider minima without relying on the stochastic differential equation framework (previous work) which assumes a reasonably small learning rate.

## 5 Discussion

We presented qualitative results to understand how GD and SGD interact with the DNN loss surface and avoided assumptions in order to rely instead on empirical evidence. We now draw similarities between the optimization trajectory in DNNs that we have empirically found, with those in quadratic loss optimization (see section 5 of LeCun et al. (1998)). Based on our empirical evidence, we deduce that both GD and SGD move in a valley like landscape by bouncing off valley walls. This is reminiscent of optimization in a quadratic loss setting with a non-isotropic positive semi-definite Hessian, where the optimal learning rate  $\eta$  causes under-damping without divergence along eigenvectors of the Hessian which have eigenvalues  $\lambda_i$  such that  $\lambda_i^{-1} < \eta < 2\lambda_i^{-1}$ . On the other hand, in the case of DNNs trained with GD, we find that even though the training loss oscillates between valley walls during consecutive iterations, the valley floor decreases smoothly (see Figure 1). This is similar to the quadratic loss optimization with over-damped convergence along the eigenvectors corresponding to eigenvalues  $\lambda_i$  such that  $\eta < \lambda_i^{-1}$ .

On a different note, it is commonly conjectured that when training DNNs, SGD crosses barriers to escape local minima. Contrary to this commonly held intuition, we find that SGD almost never crosses any significant barriers along its path. More interestingly, when training with a certain learning rate (see figure 2), we find barriers at the floor of the valley but SGD avoids them by traveling at a height above the floor (due to large learning rate). Hence if we use a small learning rate, SGD should encounter such barriers and likely cross them. But we found in our experiments that this was not the case. This suggests that while in theory SGD is capable of crossing barriers (due to stochasticity), it does not do so because probably there exist other directions in such regions along which SGD can continue to optimize without crossing barriers. But since small learning rates empirically correlate with bad generalization, this suggests that moving over such barriers instead of crossing them by using a large learning rate is a good mechanism for exploration for good regions.

Finally, much of what we have discussed is based on the loss landscape of specific datasets and architectures along with network parameterization choices like rectified linear activation units (ReLUs) and batch normalization Ioffe and Szegedy (2015). These conclusions may differ depending on these choices. In these cases analysis similar to ours can be performed to see if similar dynamics hold or not. Studying these dynamics may provide more practical guidelines for setting optimization hyperparameters.

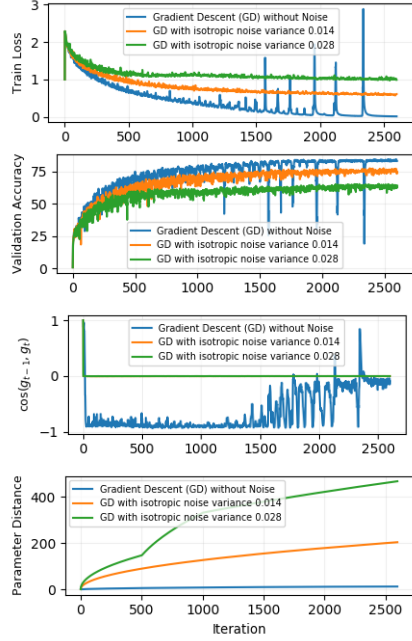


Figure 5: Plots for VGG-11 trained by GD (without noise) and GD with artificial isotropic noise sampled from Gaussian distribution with different variances. Models trained using GD with added isotropic noise get stuck in terms of training loss and have worse validation performance compared with the model trained with GD.



## References

- Alessandro Achille and Stefano Soatto. On the emergence of invariance and disentangling in deep representations. *arXiv preprint arXiv:1706.01350*, 2017.
- Madhu S Advani and Andrew M Saxe. High-dimensional dynamics of generalization error in neural networks. *arXiv preprint arXiv:1710.03667*, 2017.
- Daniel J Amit, Hanoch Gutfreund, and Haim Sompolinsky. Spin-glass models of neural networks. *Physical Review A*, 32(2):1007, 1985.
- Devansh Arpit, Stanisław Jastrzębski, Nicolas Ballas, David Krueger, Emmanuel Bengio, Maxinder S Kanwal, Tegan Maharaj, Asja Fischer, Aaron Courville, Yoshua Bengio, et al. A closer look at memorization in deep networks. In *International Conference on Machine Learning*, pages 233–242, 2017.
- Pratik Chaudhari and Stefano Soatto. Stochastic gradient descent performs variational inference, converges to limit cycles for deep networks. *arXiv preprint arXiv:1710.11029*, 2017.
- Anna Choromanska, Mikael Henaff, Michael Mathieu, Gérard Ben Arous, and Yann LeCun. The loss surfaces of multilayer networks. In *Artificial Intelligence and Statistics*, pages 192–204, 2015.
- Yann N Dauphin, Razvan Pascanu, Caglar Gulcehre, Kyunghyun Cho, Surya Ganguli, and Yoshua Bengio. Identifying and attacking the saddle point problem in high-dimensional non-convex optimization. In *Advances in neural information processing systems*, pages 2933–2941, 2014.
- Laurent Dinh, Razvan Pascanu, Samy Bengio, and Yoshua Bengio. Sharp minima can generalize for deep nets. *arXiv e-prints*, 1703.04933, March 2017.
- Viktor Dotsenko. *An introduction to the theory of spin glasses and neural networks*, volume 54. World Scientific, 1995.
- C Daniel Freeman and Joan Bruna. Topology and geometry of half-rectified network optimization. *arXiv preprint arXiv:1611.01540*, 2016.
- Ian J Goodfellow, Oriol Vinyals, and Andrew M Saxe. Qualitatively characterizing neural network optimization problems. *arXiv preprint arXiv:1412.6544*, 2014.
- Priya Goyal, Piotr Dollár, Ross Girshick, Pieter Noordhuis, Lukasz Wesolowski, Aapo Kyrola, Andrew Tulloch, Yangqing Jia, and Kaiming He. Accurate, large minibatch sgd: training imagenet in 1 hour. *arXiv preprint arXiv:1706.02677*, 2017.
- Moritz Hardt, Benjamin Recht, and Yoram Singer. Train faster, generalize better: Stability of stochastic gradient descent. *arXiv preprint arXiv:1509.01240*, 2015.
- Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 770–778, 2016.
- Sepp Hochreiter and Jürgen Schmidhuber. Flat minima. *Neural Computation*, 9(1):1–42, 1997.
- Elad Hoffer, Itay Hubara, and Daniel Soudry. Train longer, generalize better: closing the generalization gap in large batch training of neural networks. In *Advances in Neural Information Processing Systems*, pages 1729–1739, 2017.
- Sergey Ioffe and Christian Szegedy. Batch normalization: Accelerating deep network training by reducing internal covariate shift. In *International conference on machine learning*, pages 448–456, 2015.
- Stanisław Jastrzębski, Zachary Kenton, Devansh Arpit, Nicolas Ballas, Asja Fischer, Yoshua Bengio, and Amos Storkey. Three factors influencing minima in sgd. *arXiv preprint arXiv:1711.04623*, 2017.
- Kenji Kawaguchi, Leslie Pack Kaelbling, and Yoshua Bengio. Generalization in deep learning. *arXiv preprint arXiv:1710.05468*, 2017.

- Nitish Shirish Keskar, Dheevatsa Mudigere, Jorge Nocedal, Mikhail Smelyanskiy, and Ping Tak Peter Tang. On large-batch training for deep learning: Generalization gap and sharp minima. *arXiv preprint arXiv:1609.04836*, 2016.
- Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014.
- Alex Krizhevsky. Learning Multiple Layers of Features from Tiny Images. Technical report, 2009.
- Harold Kushner and G George Yin. *Stochastic approximation and recursive algorithms and applications*, volume 35. Springer Science & Business Media, 2003.
- Yann Lecun and Corinna Cortes. The MNIST database of handwritten digits. URL <http://yann.lecun.com/exdb/mnist/>.
- Yann LeCun, Léon Bottou, Genevieve B Orr, and Klaus-Robert Müller. Efficient backprop. In *Neural networks: Tricks of the trade*, pages 9–50. Springer, 1998.
- Chris Junchi Li, Lei Li, Junyang Qian, and Jian-Guo Liu. Batch size matters: A diffusion approximation framework on nonconvex stochastic gradient descent. *arXiv preprint arXiv:1705.07562*, 2017.
- Qianxiao Li, Cheng Tai, et al. Stochastic modified equations and adaptive stochastic gradient algorithms. *arXiv preprint arXiv:1511.06251*, 2015.
- Stephan Mandt, Matthew D Hoffman, and David M Blei. Stochastic gradient descent as approximate bayesian inference. *arXiv preprint arXiv:1704.04289*, 2017.
- Arvind Neelakantan, Luke Vilnis, Quoc V Le, Ilya Sutskever, Lukasz Kaiser, Karol Kurach, and James Martens. Adding gradient noise improves learning for very deep networks. *arXiv preprint arXiv:1511.06807*, 2015.
- Yurii Nesterov. A method of solving a convex programming problem with convergence rate  $o(1/k^2)$ . In *Soviet Mathematics Doklady*, volume 27, pages 372–376, 1983.
- Behnam Neyshabur, Srinadh Bhojanapalli, David McAllester, and Nati Srebro. Exploring generalization in deep learning. In *Advances in Neural Information Processing Systems*, pages 5949–5958, 2017.
- Boris T Polyak. Some methods of speeding up the convergence of iteration methods. *USSR Computational Mathematics and Mathematical Physics*, 4(5):1–17, 1964.
- Olga Russakovsky, Jia Deng, Hao Su, Jonathan Krause, Sanjeev Satheesh, Sean Ma, Zhiheng Huang, Andrej Karpathy, Aditya Khosla, Michael Bernstein, Alexander C. Berg, and Li Fei-Fei. ImageNet Large Scale Visual Recognition Challenge. *International Journal of Computer Vision (IJCV)*, 115(3):211–252, 2015. doi: 10.1007/s11263-015-0816-y.
- Itay Safran and Ohad Shamir. On the quality of the initial basin in overspecified neural networks. In *International Conference on Machine Learning*, pages 774–782, 2016.
- Levent Sagun, Utku Evci, V Ugur Guney, Yann Dauphin, and Leon Bottou. Empirical analysis of the hessian of over-parametrized neural networks. *arXiv preprint arXiv:1706.04454*, 2017.
- Karen Simonyan and Andrew Zisserman. Very deep convolutional networks for large-scale image recognition. *arXiv preprint arXiv:1409.1556*, 2014.
- Leslie N Smith. Cyclical learning rates for training neural networks. In *Applications of Computer Vision (WACV), 2017 IEEE Winter Conference on*, pages 464–472. IEEE, 2017.
- Leslie N Smith and Nicholay Topin. Super-convergence: Very fast training of residual networks using large learning rates. *arXiv preprint arXiv:1708.07120*, 2017.
- Samuel L Smith and Quoc V Le. Understanding generalization and stochastic gradient descent. *arXiv preprint arXiv:1710.06451*, 2017.

- Samuel L Smith, Pieter-Jan Kindermans, and Quoc V Le. Don't decay the learning rate, increase the batch size. *arXiv preprint arXiv:1711.00489*, 2017.
- Ilya Sutskever, James Martens, George Dahl, and Geoffrey Hinton. On the importance of initialization and momentum in deep learning. In *International conference on machine learning*, pages 1139–1147, 2013.
- Tijmen Tieleman and Geoffrey Hinton. Lecture 6.5-rmsprop: Divide the gradient by a running average of its recent magnitude. *COURSERA: Neural networks for machine learning*, 4(2):26–31, 2012.
- Lei Wu, Zhanxing Zhu, et al. Towards understanding generalization of deep learning: Perspective of loss landscapes. *arXiv preprint arXiv:1706.10239*, 2017.
- Matthew D Zeiler. Adadelta: an adaptive learning rate method. *arXiv preprint arXiv:1212.5701*, 2012.
- Chiyuan Zhang, Samy Bengio, Moritz Hardt, Benjamin Recht, and Oriol Vinyals. Understanding deep learning requires rethinking generalization. *arXiv preprint arXiv:1611.03530*, 2016.
- Zhanxing Zhu, Jingfeng Wu, Lei Wu, Jinwen Ma, and Bing Yu. The regularization effects of anisotropic noise in stochastic gradient descent. *arXiv preprint arXiv:1803.00195*, 2018.

Table 3: Average height of SGD above valley floor across the iterations in one epoch for different epochs during the training of VGG-11 and Resnet-56 using SGD on CIFAR-10 and VGG-11 on Tiny-ImageNet. Here the height at iteration  $t$  is defined by  $\frac{\mathcal{L}(\theta_t) + \mathcal{L}(\theta_{t+1}) - 2\mathcal{L}(\theta_t^{\min})}{2}$ .

VGG-11 on CIFAR-10	<b>Epoch 1</b>	<b>Epoch 10</b>	<b>Epoch 25</b>	<b>Epoch 100</b>
<b>LR 0.1</b>	0.0625 $\pm$ 1.2e-3	0.0199 $\pm$ 4.8e-4	0.0104 $\pm$ 2.2e-5	0.0025 $\pm$ 3.1e-6
<b>LR 0.05</b>	0.0102 $\pm$ 3.8e-5	0.0050 $\pm$ 2.8e-5	0.0035 $\pm$ 1.7e-5	0.0011 $\pm$ 1.0e-6
Resnet-56 on CIFAR-10	<b>Epoch 1</b>	<b>Epoch 10</b>	<b>Epoch 25</b>	<b>Epoch 100</b>
<b>LR 0.3</b>	0.0380 $\pm$ 5.9e-4	0.0131 $\pm$ 4.5e-4	0.0094 $\pm$ 1.3e-5	0.0017 $\pm$ 1.0e-5
<b>LR 0.15</b>	0.0084 $\pm$ 5.2e-5	0.0034 $\pm$ 3.2e-5	0.0020 $\pm$ 7.2e-6	0.0013 $\pm$ 6.7e-6
VGG-11 on Tiny-ImageNet	<b>Epoch 1</b>	<b>Epoch 10</b>	<b>Epoch 25</b>	<b>Epoch 100</b>
<b>LR 0.5</b>	0.028 $\pm$ 1.0e-3	0.213 $\pm$ 1.5e-3	0.187 $\pm$ 1.9e-3	9.8e-5 $\pm$ 2.0e-9
<b>LR 0.1</b>	0.0039 $\pm$ 5.2e-5	0.163 $\pm$ 2.64e-5	0.116 $\pm$ 0.013	1.1e-5 $\pm$ 3.6e-11

## Appendix

### A Optimization Trajectory

This is a continuation of section 3.1 in the main text. Here we show further experiments on other datasets, architectures and hyper-parameter settings. The analysis of GD training for Resnet-56 on CIFAR-10, MLP on MNIST and VGG-11 on tiny ImageNet are shown in figures 10, 18 and 21 respectively. Similarly, the analysis of SGD training for Resnet-56 on CIFAR-10 dataset with batch size of 100 and learning rate 0.1 for epochs 1, 2, 25 and 100 are shown in figures 11, 12, 13 and 14 respectively. The analysis of SGD training for VGG-11 on CIFAR-10 with the batch size of 100 and learning rate 0.1 on epochs 2, 25, 100 are shown in figures 15, 16 and 17. The analysis of SGD training for MLP on MNIST for epochs 1 and 2 are shown in figures 19 and 20. The analysis of SGD training for VGG-11 on tiny ImageNet for epochs 1 is shown in figure 22. We also conducted the same experiment and analysis on various batch sizes and learning rates for every architecture. Results of VGG-11 can be found in figures 23, 24, 25 and 26. Results of Resnet-56 can be found in figures 27, 28, 29 and 30. The observations and rules we discovered and described in section 3 are all consistent for all these experiments. Specifically, for the interpolation of SGD for VGG-11 on tiny ImageNet, the valley-like trajectory is weird-looking but even so, according to our quantitative evaluation there is no barrier between any two consecutive iterations.

We track the spectral norm of the Hessian along with the validation accuracy while the model is being trained. This is shown in figure 6 for Resnet-56 trained on CIFAR-10.

### B Qualitative Roles of Learning Rate and Batch Size

This is a continuation of section 3.2 in the main text. In this section we show further experiments for the analysis of different roles of learning rate and batch size during training on various architectures and data sets. Figures 31, 32 and 33 shows the results for Resnet-56 on CIFAR-10, MLP on MNIST and VGG-11 on Tiny-Imagenet. In all of the experiments, training the model with smaller batch size will make the angle between gradients of two consecutive iterations larger, which means for smaller batch size, instead of oscillating within the same region, the optimization travels farther along the valley, as we described in section 3.2. For all architectures, changing learning rate doesn't change the angles.

### C Learning Rate Schedule

We observe from table 2 that the optimization oscillates at a lower height as training progresses (which is likely because SGD finds flatter regions as training progresses, see Figure 3). As we discussed based on Figure 2, the floor of the DNN valley is highly non-linear with many barriers. Based on these two observations, it seems that it should be advantageous for SGD to maintain a large height from the floor of the valley to facilitate further exploration without getting hindered by barriers as it may allow the optimization to find flatter regions. Hence, this line of thought suggests that we should increase the learning rate as training progresses (of course eventually it needs to be annealed for convergence to a minimum). Smith (2017); Smith and Topin (2017) propose a cyclical learning

rate (CLR) schedule which partially has this property. It involves linearly increasing the learning rate every iteration until a certain number of iterations, then similarly linearly reducing it, and repeat this process in a cycle. We now empirically show that multiple cycles of CLR are redundant, and simply increasing the learning rate until a certain point, and then annealing it leads to similar or better performance. Specifically, to rule out the need for cycles, as a null hypothesis, we increase the learning rate as in the first cycle of CLR, then keep it flat, then linearly anneal it (we call it the *trapezoid schedule*). For fairness, we also plot the widely used step-wise learning rate annealing schedule. *In our experiments, we find that methods which increase learning rate during training may be considered slightly better.* The learning curves are shown in figures 8 in main text and 7 in appendix (with other details). We leave an extensive study of learning rate schedule design based on the proposed guideline as future work.

We run the same experiment as described above on Resnet-56 with CIFAR-10 and it shows the rule for CLR, trapezoid schedule and SGD with stepwise annealing. Plots can be seen at figure 7. All schedules are tuned to their best performance with a hyperparameter grid search. For both Resnet-56 and VGG-11, we use batch size 100 for all models. The learning rate schedules are apparent from the figures themselves.

## D Importance of SGD Noise Structure

Here we derive in detail the relation between the Hessian and gradient covariance using the fact that for the negative log likelihood loss  $\mathcal{L}_i(\theta) = -\log(p_i(\theta))$ . Note we use the fact that for this particular loss function,  $\frac{\partial \mathcal{L}_i(\theta)}{\partial p_i(\theta)} = -\frac{1}{p_i(\theta)}$ , and  $\frac{\partial^2 \mathcal{L}_i(\theta)}{\partial p_i(\theta)^2} = \frac{1}{p_i^2(\theta)}$ , which yields  $\frac{\partial^2 \mathcal{L}_i(\theta)}{\partial p_i(\theta)^2} = \left(\frac{\partial \mathcal{L}_i(\theta)}{\partial p_i(\theta)}\right)^2$ .

$$\mathbf{H}(\theta) = \frac{1}{N} \sum_{i=1}^N \frac{\partial^2 \mathcal{L}_i(\theta)}{\partial \theta^2} \quad (2)$$

$$= \frac{1}{N} \sum_{i=1}^N \frac{\partial}{\partial \theta} \left( \frac{\partial \mathcal{L}_i(\theta)}{\partial p_i(\theta)} \cdot \frac{\partial p_i(\theta)}{\partial \theta} \right) \quad (3)$$

$$= \frac{1}{N} \sum_{i=1}^N \frac{\partial^2 \mathcal{L}_i(\theta)}{\partial p_i(\theta)^2} \cdot \frac{\partial p_i(\theta)}{\partial \theta} \frac{\partial p_i(\theta)}{\partial \theta}^T + \frac{\partial \mathcal{L}_i(\theta)}{\partial p_i(\theta)} \cdot \frac{\partial^2 p_i(\theta)}{\partial \theta^2} \quad (4)$$

$$= \frac{1}{N} \sum_{i=1}^N \left( \frac{\partial \mathcal{L}_i(\theta)}{\partial p_i(\theta)} \right)^2 \cdot \frac{\partial p_i(\theta)}{\partial \theta} \frac{\partial p_i(\theta)}{\partial \theta}^T + \frac{\partial \mathcal{L}_i(\theta)}{\partial p_i(\theta)} \cdot \frac{\partial^2 p_i(\theta)}{\partial \theta^2} \quad (5)$$

$$= \frac{1}{N} \sum_{i=1}^N \frac{\partial \mathcal{L}_i(\theta)}{\partial \theta} \frac{\partial \mathcal{L}_i(\theta)}{\partial \theta}^T + \frac{\partial \mathcal{L}_i(\theta)}{\partial p_i(\theta)} \cdot \frac{\partial^2 p_i(\theta)}{\partial \theta^2} \quad (6)$$

$$= \mathbf{C}(\theta) + \bar{\mathbf{g}}(\theta) \bar{\mathbf{g}}(\theta)^T + \frac{1}{N} \sum_{i=1}^N \frac{\partial \mathcal{L}_i(\theta)}{\partial p_i(\theta)} \cdot \frac{\partial^2 p_i(\theta)}{\partial \theta^2} \quad (7)$$

where  $\bar{\mathbf{g}}(\theta) = \frac{1}{N} \sum_{i=1}^N \frac{\partial \mathcal{L}_i(\theta)}{\partial \theta}$ .

## E Discussion

In the main text, we talk about converge in the quadratic setting depending on the value of learning rate relative to the largest eigenvalue of the Hessian. The convergence in this setting has been visualized in 36.



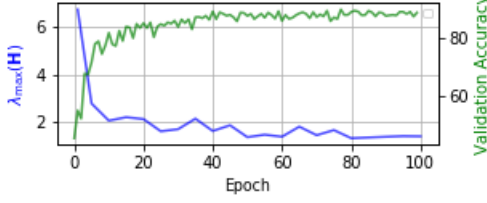


Figure 6: Max eigenvalue (spectral norm) of the Hessian and validation accuracy for Resnet-56 trained on CIFAR-10 using a fixed learning rate of 0.3 and batch size 100. The spectral norm roughly decreases with training but starts increasing slightly towards the end. Similarly validation accuracy roughly improves throughout training but drops towards the end.

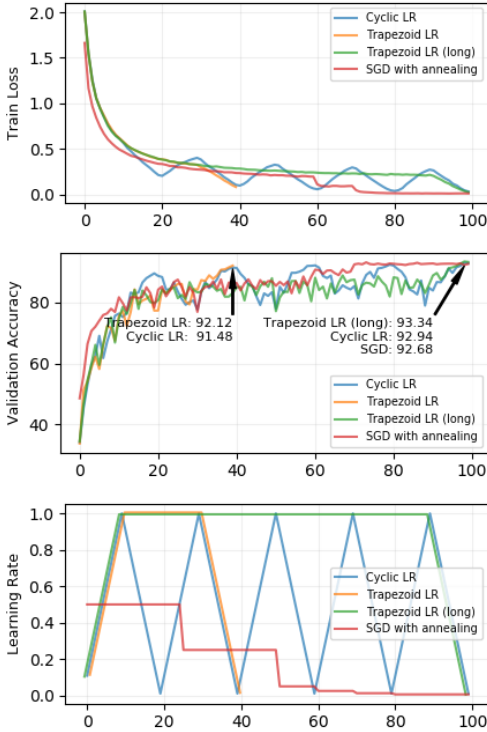


Figure 7: Plots for Resnet-56 on CIFAR-10 trained using Cyclic learning rate (CLR), SGD with stepsize annealing, and trapezoid schedule. Cycles in the CLR schedule are redundant, which is shown by the trapezoid schedule.

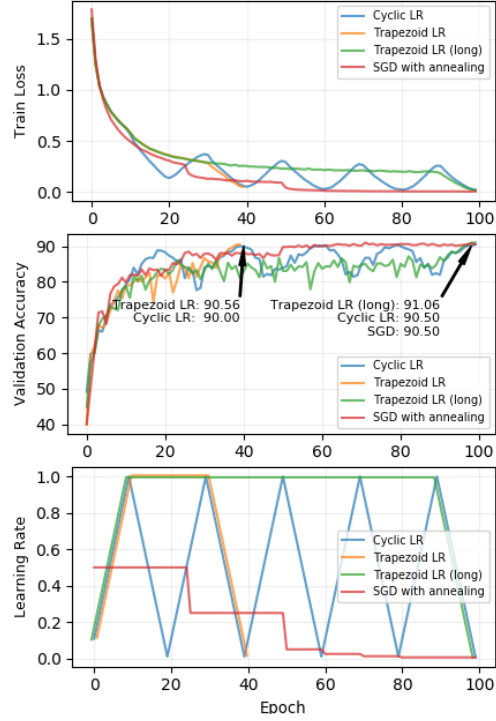


Figure 8: Plots for VGG-11 on CIFAR-10 trained using Cyclic learning rate (CLR), SGD with stepwise annealing, and trapezoid schedule. Cycles in the CLR schedule are redundant, which is shown by the trapezoid schedule.

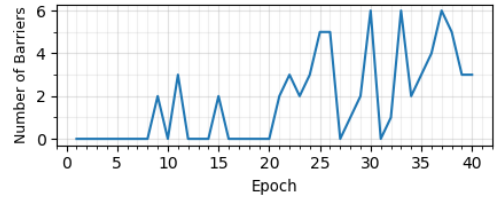


Figure 9: Numbers of barriers found during training loss interpolation for every epoch(450 iterations) for VGG-11 on CIFAR-10. We say a barrier exists during a training update step if there exists a point between the parameters before and after an update which has a loss value higher than the loss at either points. Note that even for these barriers, their heights (defined by  $\frac{\mathcal{L}(\theta_t) + \mathcal{L}(\theta_{t+1}) - 2\mathcal{L}(\theta_t^{\min})}{2}$ ) are substantially smaller compared with the value of loss at the corresponding iterations (not mentioned here), meaning they are not significant barriers.

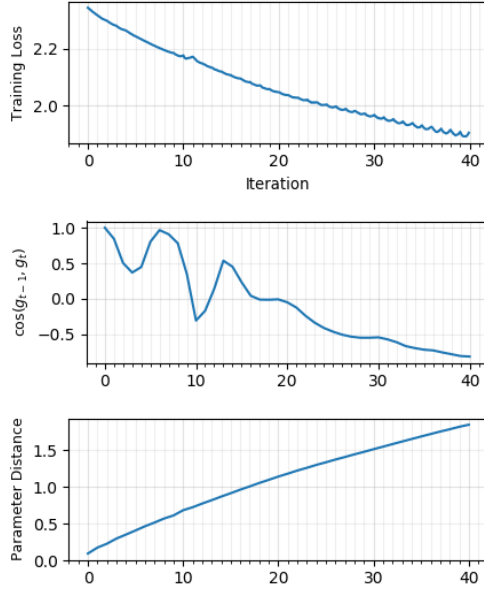


Figure 10: Plots for Resnet-56 Epoch 1 trained using full batch **Gradient Descent (GD)** on CIFAR-10. The descriptions of the plots are same as in figure 1.

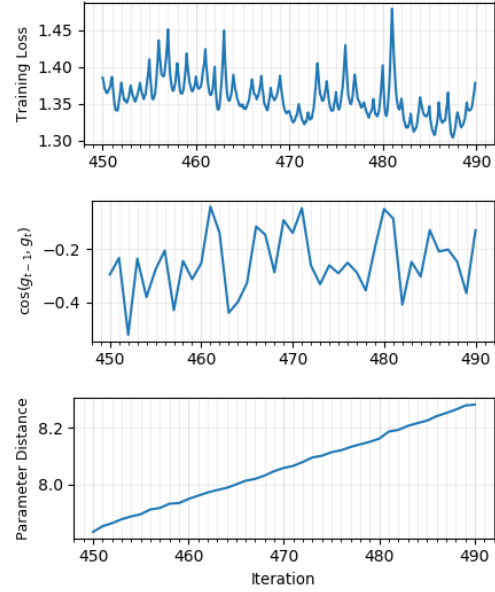


Figure 12: Plots for Resnet-56 Epoch 2 trained using **SGD** on CIFAR-10. The descriptions of the plots are same as in figure 1.

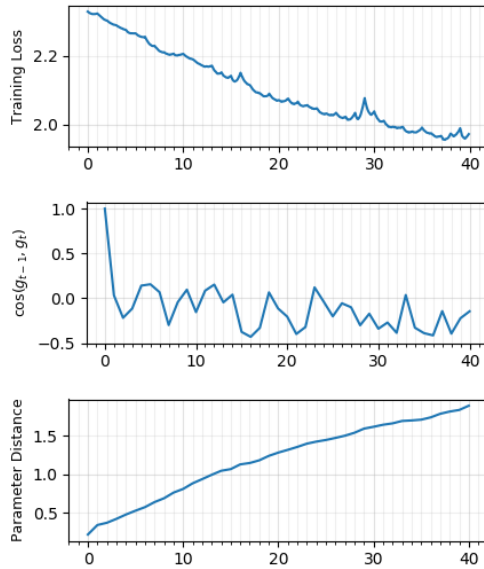


Figure 11: Plots for Resnet-56 Epoch 1 trained using **SGD** on CIFAR-10. The descriptions of the plots are same as in figure 1.

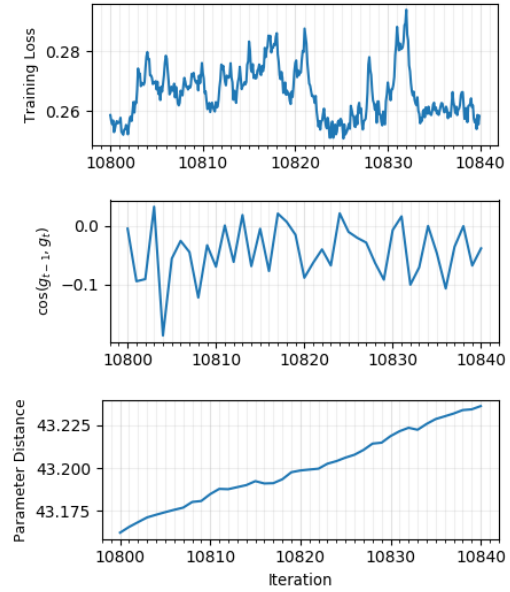


Figure 13: Plots for Resnet-56 Epoch 25 trained using **SGD** on CIFAR-10. The descriptions of the plots are same as in figure 1.

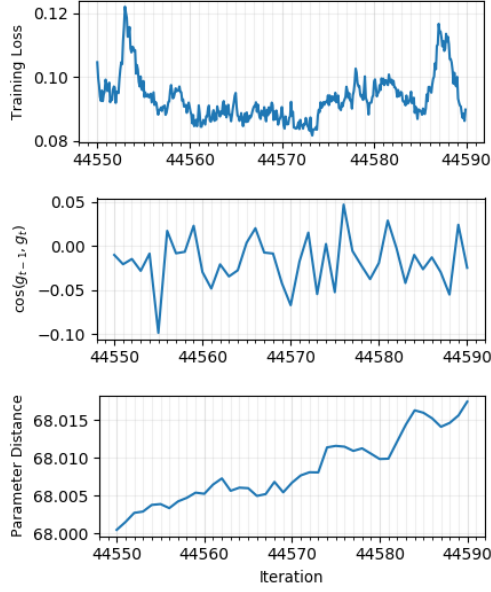


Figure 14: Plots for Resnet-56 Epoch 100 trained using **SGD** on CIFAR-10. The descriptions of the plots are same as in figure 1.

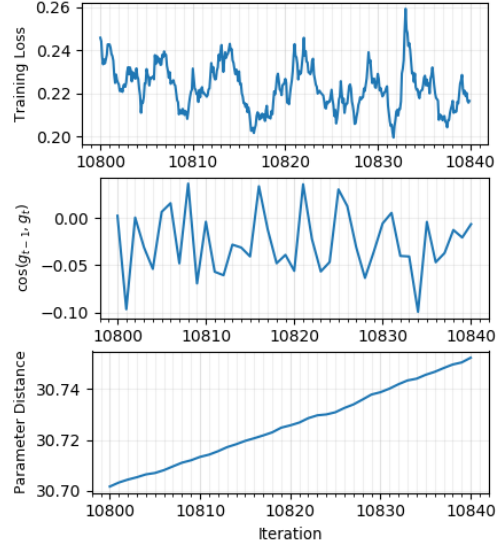


Figure 16: Plots for VGG-11 Epoch 25 trained using **SGD** on CIFAR-10. The descriptions of the plots are same as in figure 1.

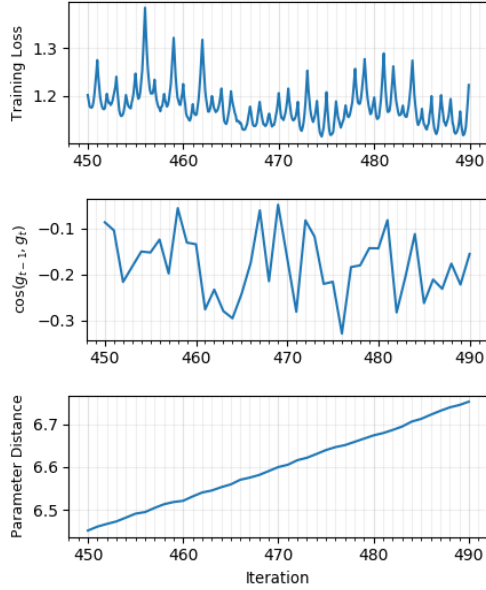


Figure 15: Plots for VGG-11 Epoch 2 trained using **SGD** on CIFAR-10. The descriptions of the plots are same as in figure 1.

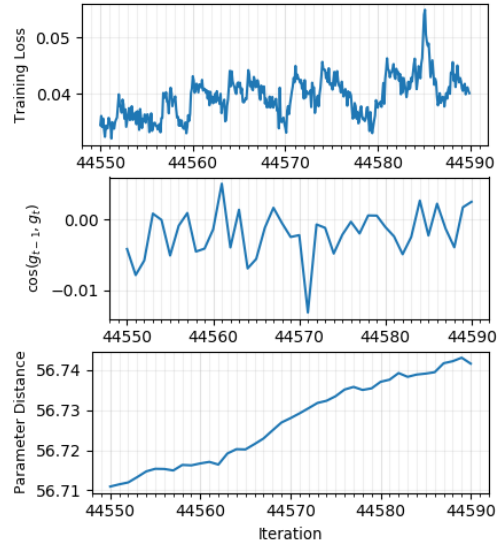


Figure 17: Plots for VGG-11 Epoch 100 trained using **SGD** on CIFAR-10. The descriptions of the plots are same as in figure 1.

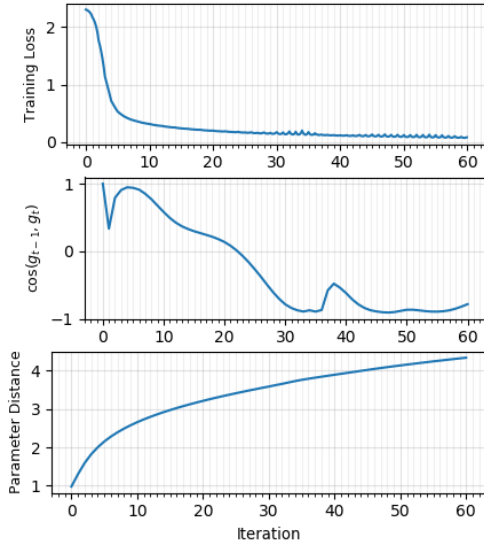


Figure 18: Plots for MLP Epoch 1 trained using full batch **Gradient Descent (GD)** on MNIST. The descriptions of the plots are same as in figure 1.

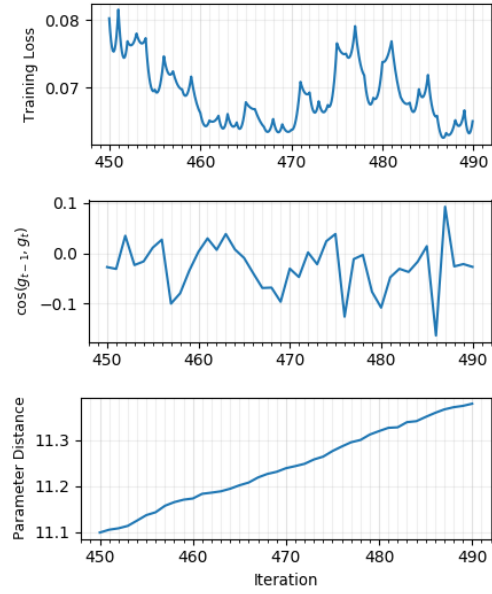


Figure 20: Plots for MLP Epoch 2 trained using **SGD** on MNIST. The descriptions of the plots are same as in figure 1.

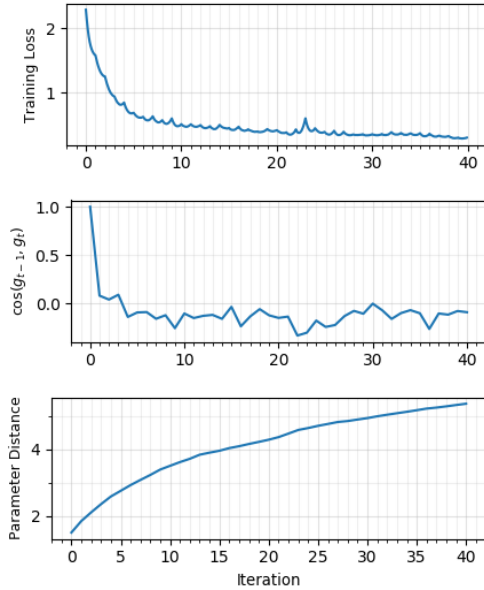


Figure 19: Plots for MLP Epoch 1 trained using **SGD** on MNIST. The descriptions of the plots are same as in figure 1.

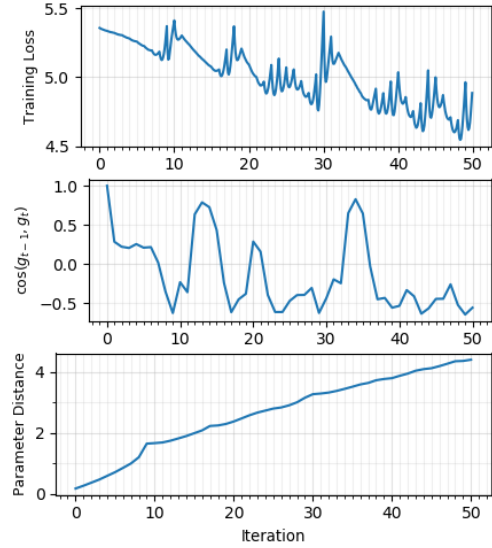


Figure 21: Plots for VGG-11 Epoch 1 trained using full batch **Gradient Descent (GD)** on Tiny-ImageNet. The descriptions of the plots are same as in figure 1.

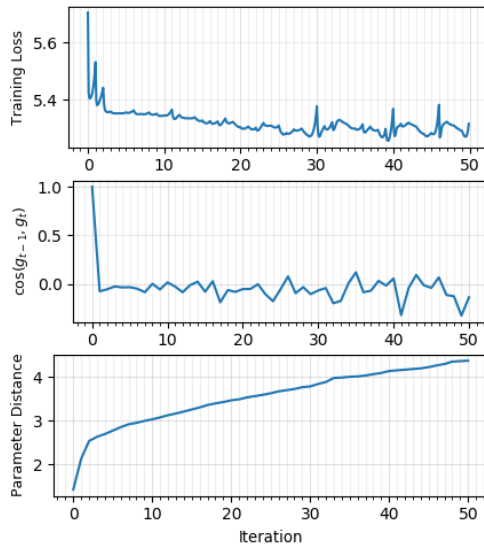


Figure 22: Plots for VGG-11 Epoch 1 trained using **SGD** on Tiny-ImageNet. The descriptions of the plots are same as in figure 1.

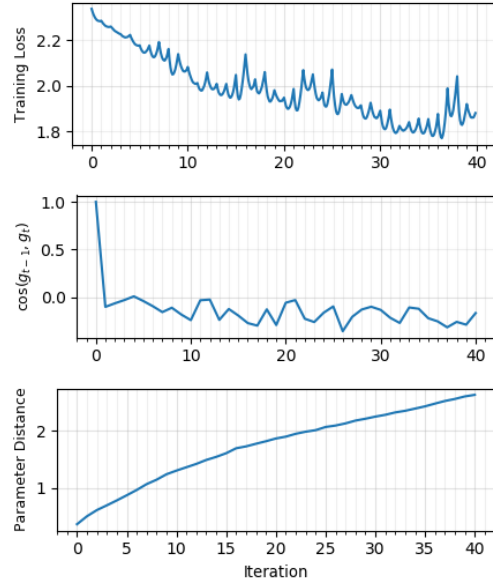


Figure 24: Plots for VGG-11 Epoch 1 trained using learning rate 0.2 batch size 100 on CIFAR-10.

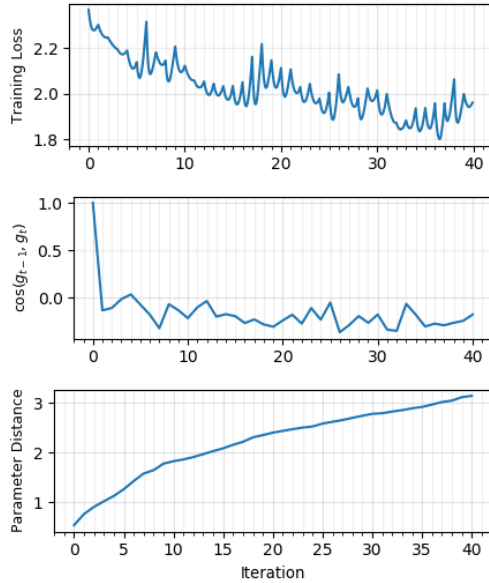


Figure 23: Plots for VGG-11 Epoch 1 trained using learning rate 0.3 batch size 100 on CIFAR-10.

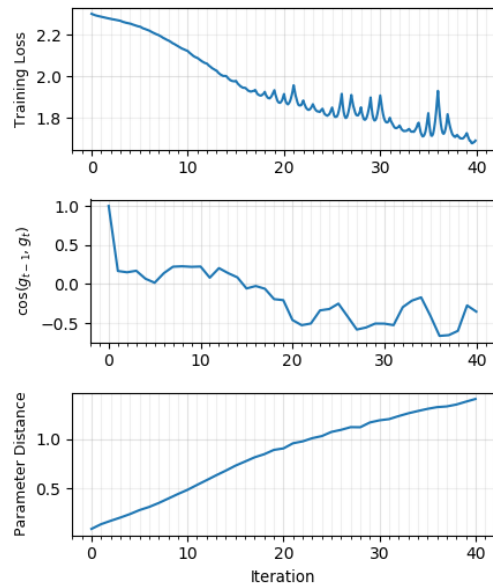


Figure 25: Plots for VGG-11 Epoch 1 trained using learning rate 0.1 batch size 500 on CIFAR-10.



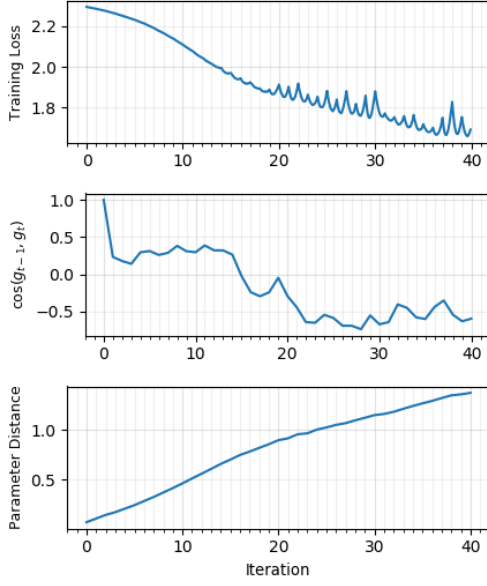


Figure 26: Plots for VGG-11 Epoch 1 trained using learning rate 0.1 batch size 1000 on CIFAR-10.

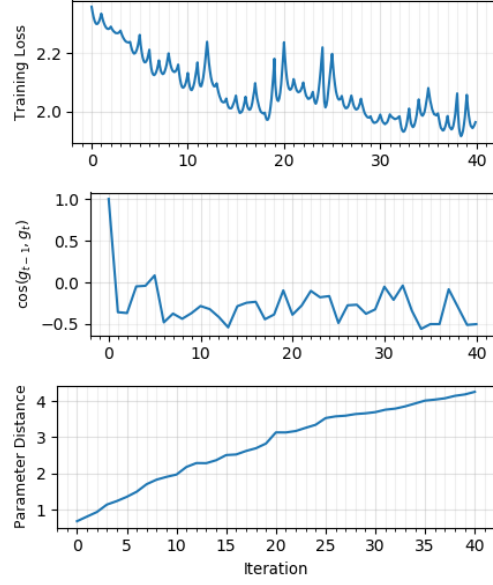


Figure 28: Plots for Resnet-56 Epoch 1 trained using learning rate 1 batch size 100 on CIFAR-10.

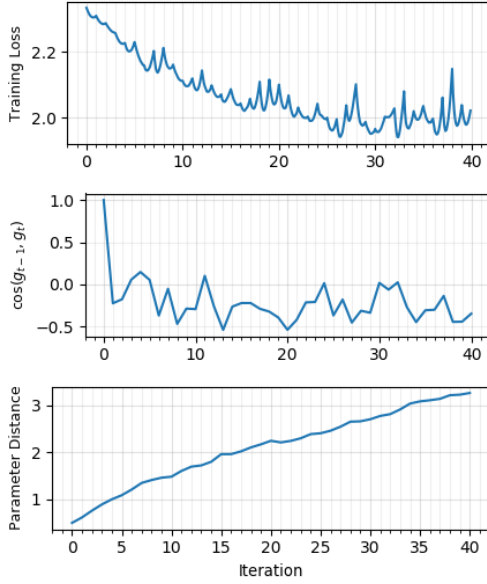


Figure 27: Plots for Resnet-56 Epoch 1 trained using learning rate 0.7 batch size 100 on CIFAR-10.

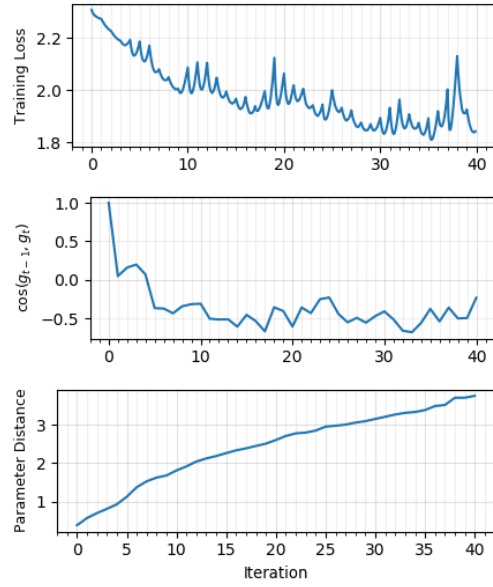


Figure 29: Plots for Resnet-56 Epoch 1 trained using learning rate 1 batch size 500 on CIFAR-10.

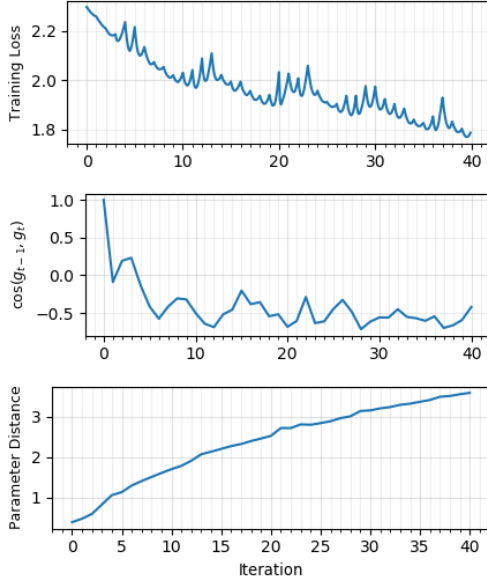


Figure 30: Plots for Resnet-56 Epoch 1 trained using learning rate 1 batch size 1000 on CIFAR-10.

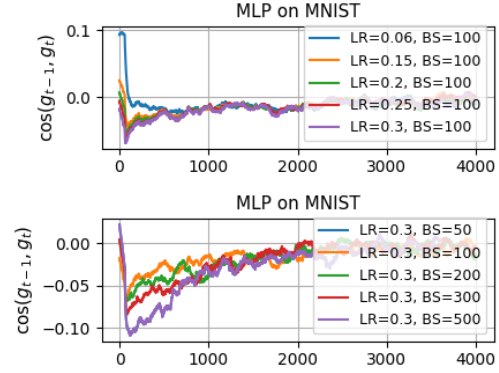


Figure 32: Changing batch size changes the cosine of angle between consecutive gradients while changing learning rate does not have any significant effect on the cosine. This shows batch size has a qualitatively different role compared with learning rate. Note that the curves are smoothened for visual clarity.

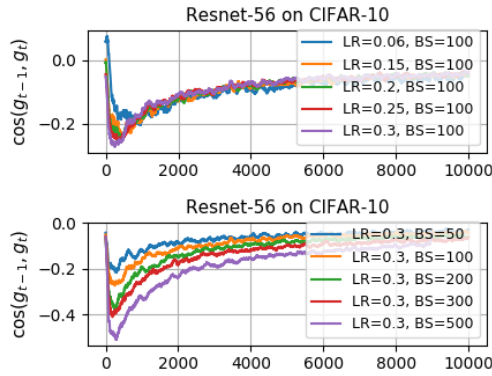


Figure 31: Changing batch size changes the cosine of angle between consecutive gradients while changing learning rate does not have any significant effect on the cosine. This shows batch size has a qualitatively different role compared with learning rate. Note that the curves are smoothened for visual clarity.

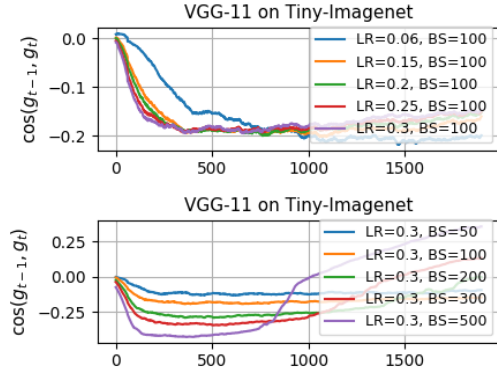


Figure 33: Changing batch size changes the cosine of angle between consecutive gradients while changing learning rate does not have any significant effect on the cosine. This shows batch size has a qualitatively different role compared with learning rate. Note that the curves are smoothened for visual clarity.

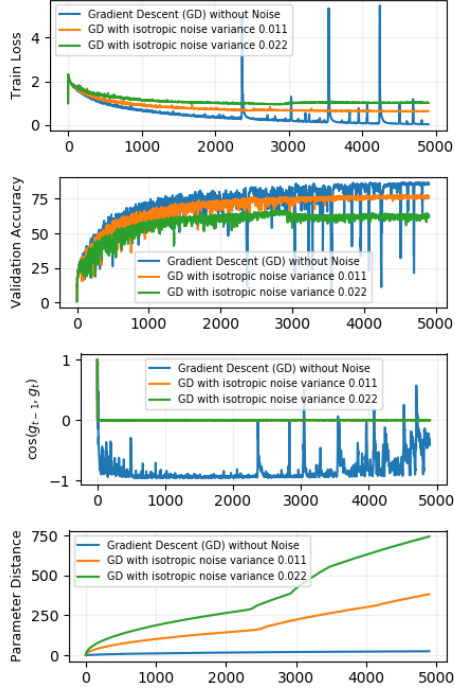


Figure 34: Plots for Resnet-56 trained by GD (without noise) and GD with artificial isotropic noise sampled from Gaussian distribution with different variances. Models trained using GD with added isotropic noise get stuck in terms of training loss and have worse validation performance compared with the model trained with GD.

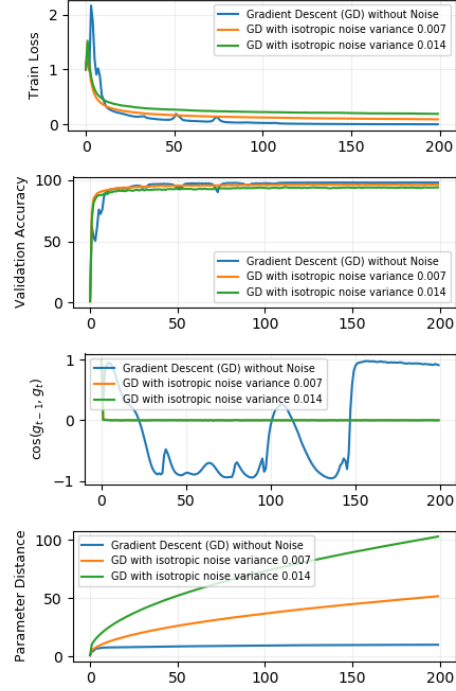


Figure 35: Plots for MLP trained by GD (without noise) and GD with artificial isotropic noise sampled from Gaussian distribution with different variances. Models trained using GD with added isotropic noise get stuck in terms of training loss and have worse validation performance compared with the model trained with GD.

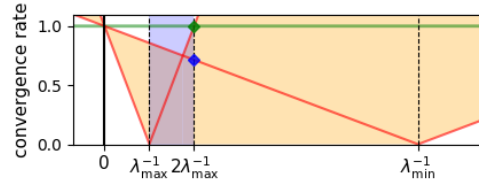


Figure 36: Graph of  $|\lambda_i \eta - 1|$ : convergence rates of gradient descent on  $\lambda_{\min}$ -strongly convex and  $\lambda_{\max}$ -smooth quadratic surfaces. The area which is shaded by orange contains possibly graphs of other eigenvalues of the hessian. In the range of learning rates shaded by blue, trajectory underdamps in the direction of the maximum eigenvalue. For a certain learning rate, while the trajectory oscillates in the direction of the maximum eigenvalue (green diamond), it overdamps in all the others (e.g. blue diamond - in a 'flat' direction).