

---

# Asymmetric Valleys: Beyond Sharp and Flat Local Minima

---

Anonymous Author(s)

Affiliation

Address

email

## Abstract

Despite the non-convex nature of their loss functions, deep neural networks are known to generalize well when optimized with stochastic gradient descent (SGD). Recent work conjectures that SGD with proper configuration is able to find wide and flat local minima, which are correlated with good generalization performance. In this paper, we observe that local minima of modern deep networks are more than being flat or sharp. Instead, at a local minimum there exist many asymmetric directions such that the loss increases abruptly along one side, and slowly along the opposite side – we formally define such minima as *asymmetric valleys*. Under mild assumptions, we first prove that for asymmetric valleys, a solution biased towards the flat side generalizes better than the exact empirical minimizer. Then, we show that performing weight averaging along the SGD trajectory implicitly induces such biased solutions. This provides theoretical explanations for a series of intriguing phenomena observed in recent work [25, 5, 51]. Finally, extensive empirical experiments with modern deep networks are conducted to validate our assumptions and analyze the intriguing properties of asymmetric valleys.

## 1 Introduction

The loss landscape of neural networks has attracted great research interests in the deep learning community [9, 10, 32, 12, 15, 43, 36]. A deeper understanding of the loss landscape is important for designing better optimization algorithms, and helps to answer the question of when and how a deep network can achieve good generalization performance. One hypothesis that draws attention recently is that the local minima of neural networks can be characterized by their flatness, and it is conjectured that sharp minima tend to generalize worse than the flat ones [32]. A plausible explanation is that a flat minimizer of the training loss can achieve lower generalization error if the test loss is shifted from the training loss due to random perturbations. Figure 1(a) gives an illustration for this argument.

Although being supported by plenty of empirical observations [32, 25, 34], the definition of flatness was recently challenged in [11], which shows that one can construct arbitrarily sharp minima through weight re-parameterization without affecting the generalization performance. Moreover, recent evidences suggest that the minima of modern deep networks are connected with simple paths with low generalization error [12, 13]. It is empirically found that the minima found by large batch training and small batch training are shown to be connected by a path without any “bumps” [43]. In other words, a “sharp minimum” and a “flat minimum” may in fact belong to a same minimum in high dimensional space. Therefore, the notion of flat and sharp minima seems to be an oversimplification of the empirical loss surface.

In this paper, we expand the notion of flat and sharp minima by introducing the concept of *asymmetric valleys*. We observe that the loss surfaces of many neural networks are locally asymmetric. In specific, there exist many directions such that the loss increases abruptly along one side, and grows rather slowly along the opposite side (see Figure 1(b) as an illustration). We formally define this kind of

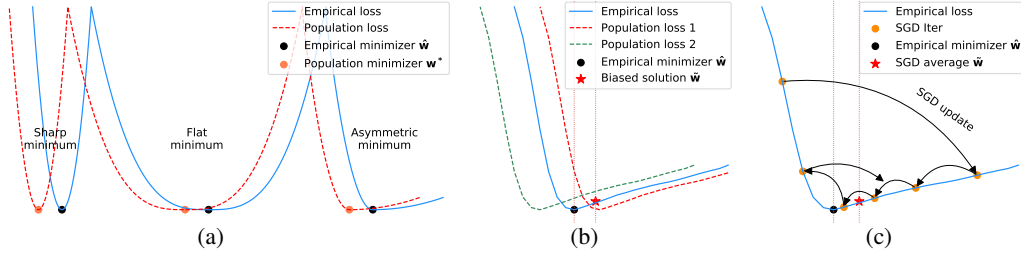


Figure 1: **(a)** An illustration of sharp, flat and asymmetric minima. If there exists a shift from empirical loss to population loss, flat minimum is more robust than sharp minimum. **(b)** For asymmetric valleys, if there exists a random shift, the solution  $\tilde{\mathbf{w}}$  biased towards the flat side is more robust than the minimizer  $\hat{\mathbf{w}}^*$ . **(c)** SGD tends to stay longer on the flat side of asymmetric valleys, therefore SGD averaging automatically produces a bias towards the flat side.

38 local minima as asymmetric valleys. As we will show in Section 6, asymmetric valleys generate  
 39 interesting illusions in high dimensional space. For example, located in the same valley shown in  
 40 Figure 1(b),  $\tilde{\mathbf{w}}$  may appear to be a wider and flatter minimum than  $\hat{\mathbf{w}}$  as the former is farther away  
 41 from the sharp side.

42 Asymmetric valleys also introduce novel insights to generalization. Folklore says when the exact  
 43 minimizer is flat, it tends to generalize better as it is more stable with respect to loss surface  
 44 perturbations [32]. Instead of following this argument, we show that in asymmetric valleys, the  
 45 solution biased towards the flat side of the valley generalizes better than the exact minimizer, under  
 46 mild assumptions. This result has at least two interesting implications: (1) converging to *which* local  
 47 minimum (if there are many) may not be critical for modern deep networks. However, it matters a  
 48 lot *where* the solution locates; and (2) the solution with lowest *a priori* generalization error is not  
 49 necessarily the minimizer of the training loss.

50 Given that a biased solution is preferred for asymmetric valleys, an immediate question is how we can  
 51 find such solutions in practice. It turns out that simply averaging the weights along the SGD trajectory,  
 52 naturally leads to the desired solutions. We give a theoretical analysis to support this argument, see  
 53 Figure 1(c) for an illustration. Our result nicely complements a series of recent empirical observations,  
 54 which demonstrated that averaged SGD has better performance over plain SGD, for various scenarios  
 55 including supervised/unsupervised/low-precision training [25, 5, 51].

56 In addition, we provide empirical analysis to verify our theoretical results and support our claims.  
 57 For example, we show that asymmetric valleys are indeed prevalent in modern deep networks, and  
 58 solutions with lower generalization error has bias towards the flat side of the valley.

## 59 2 Related Work

60 **Neural network landscape.** Neural network landscape analysis is an active and exciting area  
 61 [16, 34, 15, 40, 49, 10, 43]. For example, [12, 13] observed that essentially all local minima are  
 62 connected together with simple paths. In [22], cyclic learning rate was used to explore multiple local  
 63 optima along the training trajectory for model ensembling. There are also appealing visualizations  
 64 for the neural network landscape [34].

65 **Sharp and flat minima.** The discussion of sharp and flat local minima dates back to [20], and  
 66 recently regains its popularity. For example, Keskar et al. [32] proposed that large batch SGD finds  
 67 sharp minima, which leads to poor generalization. In [8], an entropy regularized SGD was introduced  
 68 to explicitly searching for flat minima. It was later pointed out that large batch SGD can yield  
 69 comparable performance when the learning rate or the number of training iterations are properly set  
 70 [21, 17, 47, 35, 46, 26]. Moreover, [11] showed that from a given flat minimum, one could construct  
 71 another minimum with arbitrarily sharp directions but equally good performance. In this paper, we  
 72 argue that the description of sharp or flat minima is an oversimplification. There may simultaneously  
 73 exist steep directions, flat directions, and asymmetric directions for the same minimum.

74 **SGD optimization and generalization.** As the de facto optimization tool for deep networks, SGD  
 75 and its variants are extensively studied in the literature. For example, it is shown that they could  
 76 escape saddle points or sharp local minima under reasonable assumptions [14, 28–30, 50, 1–3, 33].  
 77 For convex functions [41] or strongly convex but non-smooth functions [42], SGD averaging is shown  
 78 to give better convergence rate. In addition, it can also achieve higher generalization performance for

Lipschitz functions in theory [44, 7], or for deep networks in practice [22, 25, 5, 51]. Discussions on the generalization bound of neural networks can be found in [6, 39, 37, 31, 38, 4, 52]. We show that SGD averaging has implicit bias on the flat sides of the minima. Previously, it was shown that SGD has other kinds of implicit bias as well [48, 27, 18].

### 3 Asymmetric Valleys

In this section, we give a formal definition of asymmetric valley, and empirically show that it is prevalent in the loss landscape of modern deep neural networks.

**Preliminaries.** In supervised learning, we seek to optimize  $\mathbf{w}^* \triangleq \arg \min_{\mathbf{w} \in \mathbb{R}^d} \mathcal{L}(\mathbf{w})$ , where  $\mathcal{L}(\mathbf{w}) \triangleq \mathbb{E}_{\mathbf{x} \sim \mathcal{D}}[f(\mathbf{x}; \mathbf{w})] \in \mathbb{R}^d \rightarrow \mathbb{R}$  is the population loss,  $\mathbf{x} \in \mathbb{R}^m$  is the input sampled from distribution  $\mathcal{D}$ ,  $\mathbf{w} \in \mathbb{R}^d$  denotes the model parameter, and  $f \in \mathbb{R}^m \times \mathbb{R}^d \rightarrow \mathbb{R}$  is the loss function. Since the data distribution  $\mathcal{D}$  is usually unknown, instead of optimizing  $\mathcal{L}$  directly, we often use SGD to find the empirical risk minimizer  $\hat{\mathbf{w}}^*$  for a set of random samples  $\{\mathbf{x}_i\}_{i=1}^n$  from  $\mathcal{D}$  (a.k.a. training set):  $\hat{\mathbf{w}}^* \triangleq \arg \min_{\mathbf{w} \in \mathbb{R}^d} \hat{\mathcal{L}}(\mathbf{w})$ , where  $\hat{\mathcal{L}}(\mathbf{w}) \triangleq \frac{1}{n} \sum_{i=1}^n f(\mathbf{x}_i; \mathbf{w})$ .

In practice, it is numerically infeasible to find or test the exact local minimizer  $\hat{\mathbf{w}}^*$ . Fortunately, our theoretical results only depend on a good enough solution rather than an exact local minimum, as we will formally define in Section 4. For simplicity, we still refer to such solutions as “local minima”, although our analysis generalizes to “solutions found by SGD”.

#### 3.1 Definition of asymmetric valley

Before formally introducing asymmetric valleys, we first define asymmetric directions.

**Definition 1** (Asymmetric direction). *Given constants  $p > 0, \bar{r} > \underline{r} > 0, c > 1$ , a direction  $\mathbf{u}$  is  $(\bar{r}, \underline{r}, p, c)$ -asymmetric with respect to point  $\mathbf{w} \in \mathbb{R}^d$  and loss function  $\hat{\mathcal{L}}$ , if  $\nabla_l \hat{\mathcal{L}}(\mathbf{w} + l\mathbf{u}) < p$ , and  $\nabla_l \hat{\mathcal{L}}(\mathbf{w} - l\mathbf{u}) > cp$  for  $l \in (\underline{r}, \bar{r})$ .*

In the above definition,  $\mathbf{u} \in \mathbb{R}^d$  is a unit vector representing a direction such that the points on this direction passing  $\mathbf{w} \in \mathbb{R}^d$  can be written as  $\mathbf{w} + l\mathbf{u}$  for  $l \in (-\infty, \infty)$ . Intuitively, the loss landscape in the interval  $(-\bar{r}, -\underline{r})$  is “sharp”, while it is “flat” in the region  $(\underline{r}, \bar{r})$ . Note that we purposely leave out the region  $(-\underline{r}, \underline{r})$  without making further assumptions on it to comply with the fact that the second order derivatives of the loss function is usually continuous. It is impractical to assume the slope of the loss function change abruptly at the point  $l = 0$ .

As a concrete example, Figure 2 shows an asymmetric direction for a local minimum in ResNet-110 trained on the CIFAR-10 dataset. We verified that it is a  $(2.0, 0.6, 0.03, 15)$ -asymmetric direction, which means in the region  $(-2.0, -0.6) \cup (0.6, 2.0)$  the gradients are asymmetric with a relative ratio of  $c = 15$ .

With this Definition 1, we now formally define the *asymmetric valley*<sup>1</sup>.

**Definition 2** (Asymmetric valley). *Given constants  $p, \bar{r} > \underline{r} > 0, c > 1$ , a solution  $\hat{\mathbf{w}}^*$  of  $\hat{\mathcal{L}} \in \mathbb{R}^d \rightarrow \mathbb{R}$  is a  $(\bar{r}, \underline{r}, p, c)$ -asymmetric valley, if there exists at least one direction  $\mathbf{u}$  such that  $\mathbf{u}$  is  $(\bar{r}, \underline{r}, p, c)$ -asymmetric with respect to  $\hat{\mathbf{w}}^*$  and  $\hat{\mathcal{L}}$ .*

#### 3.2 Asymmetric valleys in deep networks

Empirically, by taking random directions with value  $(0, 1)$  in each dimension, we can find an asymmetric direction for a given solution  $\mathbf{w}^*$  with decent probability. We perform experiments with widely used deep networks, i.e., ResNet-56, ResNet-110, ResNet-164 [19], VGG-16 [45] and DenseNet-100 [23], on the CIFAR-10, CIFAR-100, SVHN and STL-10 image classification datasets. For each model on each dataset, we conduct 5 independent runs. The results show that we can *always* find asymmetric directions with certain specification  $(\bar{r}, \underline{r}, p, c)$  with  $c > 2$ , which means all the

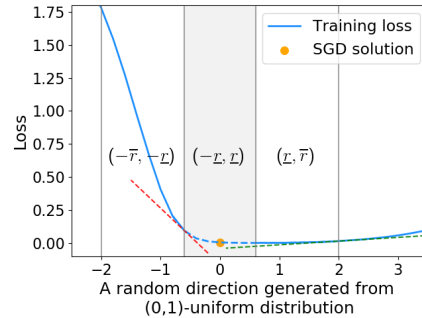


Figure 2: An asymmetric direction of a solution on the loss landscape of ResNet-110 trained on CIFAR-10.

<sup>1</sup>Here we abuse the name “valley”, since  $\hat{\mathbf{w}}^*$  is essentially a point at the center of a valley.

solutions that SGD found are located in asymmetric valleys. Asymmetric valleys widely exist in both simple and complex models, see Appendix A, Appendix E and Appendix F.

## 4 Bias and Generalization

As we show in the previous section, in the context of deep learning most local minima in practice are *asymmetric*, i.e., they might be sharp on one direction, but flat on the opposite direction. Therefore, it is interesting to investigate the generalization ability of a solution  $\mathbf{w}$  in this scenario, which may lead to different results as those obtained under the common symmetric assumption. In this section, we prove that a *biased* solution on the flat side of an asymmetric valley yields lower generalization error than the exact empirical minimizer  $\hat{\mathbf{w}}^*$  in that valley.

### 4.1 Theoretical analysis

Before presenting our theorem, we first introduce two mild assumptions. We will show that they empirically hold on modern deep networks in Section 4.2.

The first assumption (Assumption 1) states that there exists a shift between the empirical loss and true population loss. This is a common assumption in the previous works, e.g., [32], but was usually presented in an informal way. Here we define the “shift” in formally. Without loss of generality, we will compare the empirical loss  $\hat{\mathbf{L}}$  with  $\mathbf{L}' \triangleq \mathbf{L} - \min_{\mathbf{w}} \mathbf{L}(\mathbf{w}) + \min_{\mathbf{w}} \hat{\mathbf{L}}(\mathbf{w})$  to remove the “vertical difference” between  $\hat{\mathbf{L}}$  and  $\mathbf{L}$ . Notice that  $\min_{\mathbf{w}} \mathbf{L}(\mathbf{w})$  and  $\min_{\mathbf{w}} \hat{\mathbf{L}}(\mathbf{w})$  are constants and do not affect our generalization guarantee.

**Definition 3** ( $(\delta, R)$ -shift gap). *For  $\xi \geq 0$ ,  $\delta \in \mathbb{R}^d$ , and fixed functions  $\mathbf{L}$  and  $\hat{\mathbf{L}}$ , we define the  $(\delta, R)$ -shift gap between  $\mathbf{L}$  and  $\hat{\mathbf{L}}$  with respect to a point  $\mathbf{w}$  as*

$$\xi_{\delta}(\mathbf{w}) = \max_{\mathbf{v} \in \mathbb{B}(R)} |\mathbf{L}'(\mathbf{w} + \mathbf{v} + \delta) - \hat{\mathbf{L}}(\mathbf{w} + \mathbf{v})|$$

where  $\mathbf{L}'(\mathbf{w}) \triangleq \mathbf{L}(\mathbf{w}) - \min_{\mathbf{w}} \mathbf{L}(\mathbf{w}) + \min_{\mathbf{w}} \hat{\mathbf{L}}(\mathbf{w})$ , and  $\mathbb{B}(R)$  is the  $d$ -dimensional ball with radius  $R$  centered at  $\mathbf{0}$ .

From the above definition, we know that the two functions match well after the shift  $\delta$  if  $\xi_{\delta}(\mathbf{w})$  is very small. For example,  $\xi_{\delta}(\mathbf{w}) = 0$  means  $\mathbf{L}$  is locally identical to  $\hat{\mathbf{L}}$  after the shift  $\delta$ . Since  $\hat{\mathbf{L}}$  is computed on a set of random samples from  $\mathcal{D}$ , the actual shift  $\delta$  between  $\hat{\mathbf{L}}$  and  $\mathbf{L}$  is a random variable, ideally with zero expectation<sup>2</sup>.

**Assumption 1** (Random shift assumption). *For a given population loss  $\mathbf{L}$  and a random empirical loss  $\hat{\mathbf{L}}$ , constants  $R > 0, \bar{r} \geq \underline{r} > 0, \xi \geq 0$ , a vector  $\bar{\delta} \in \mathbb{R}^d$  with  $\bar{r} \geq \bar{\delta}_i \geq \underline{r}$  for all  $i \in [d]$ , a minimizer  $\hat{\mathbf{w}}^*$ , we assume that there exists a random variable  $\delta \in \mathbb{R}^d$  correlated with  $\hat{\mathbf{L}}$  such that  $\Pr(\delta_i = \bar{\delta}_i) = \Pr(\delta_i = -\bar{\delta}_i) = \frac{1}{2}$  for all  $i \in [d]$ , and the  $(\delta, R)$ -shift gap between  $\mathbf{L}$  and  $\hat{\mathbf{L}}$  with respect to  $\hat{\mathbf{w}}^*$  is bounded by  $\xi$ .*

Clearly,  $\delta$  has  $2^d$  possible values for a given shift vector  $\bar{\delta}$ , each with probability  $2^{-d}$ . Notice that Assumption 1 does not say that the difference between  $\mathbf{L}$  and  $\hat{\mathbf{L}}$  can only be one of the  $2^d$  possible  $\delta$ . Instead, it says after applying the shift  $\delta$ , the two functions have bounded  $L_{\infty}$  distance, which is a much milder assumption. It is also worth noting that our Definition 1 can mask out the central interval  $(-\underline{r}, \underline{r})$  because we have  $\bar{r} \geq \bar{\delta}_i \geq \underline{r}$  in Assumption 1. Therefore,  $\underline{r}$  cannot be arbitrarily large, otherwise Assumption 1 does not hold. Our second assumption stated below can be seen as an extension of Definition 2.

**Assumption 2** (Locally asymmetric). *For a given population loss  $\hat{\mathbf{L}}$ , and a minimizer  $\hat{\mathbf{w}}^*$ , there exist orthogonal directions  $\mathbf{u}^1, \dots, \mathbf{u}^k \in \mathbb{R}^d$  s.t.  $\mathbf{u}^i$  is  $(\bar{r}, \underline{r}, p_i, c_i)$ -asymmetric with respect to  $\hat{\mathbf{w}}^* + \mathbf{v} - \langle \mathbf{v}, \mathbf{u}^i \rangle \mathbf{u}^i$  for all  $\mathbf{v} \in \mathbb{B}(R')$  and  $i \in [k]$ .*

Assumption 2 states that if  $\mathbf{u}^i$  is an asymmetric direction at  $\hat{\mathbf{w}}^*$ , then the point  $\hat{\mathbf{w}}^* + \mathbf{v} - \langle \mathbf{v}, \mathbf{u}^i \rangle \mathbf{u}^i$  that deviates from  $\hat{\mathbf{w}}^*$  along the perpendicular direction of  $\mathbf{u}^i$ , is also asymmetric along the direction of  $\mathbf{u}^i$ . In other words, the *neighborhood* around  $\hat{\mathbf{w}}^*$  is an asymmetric valley.

<sup>2</sup>It may not be zero, as we are talking about the shift between two loss functions, rather than the difference between empirical/population loss values.

Under the above assumptions, we are ready to state our theorem, which says the empirical minimizer is not necessarily the optimal solution, and a biased solution leads to better generalization. We defer the proof to Appendix B.

**Theorem 1** (Bias leads to better generalization). *For any  $\mathbf{l} \in \mathbb{R}^k$ , if Assumption 1 holds for  $R = \|\mathbf{l}\|_2$ , Assumption 2 holds for  $R' = \|\bar{\delta}\|_2 + \|\mathbf{l}\|_2$ , and  $\frac{4\xi}{(c_i-1)p_i} < \mathbf{l}_i \leq \max\{\bar{r} - \bar{\delta}_i, \bar{\delta}_i - \underline{r}\}$ , then we have*

$$\mathbb{E}_{\delta} \mathcal{L}(\hat{\mathbf{w}}^*) - \mathbb{E}_{\delta} \mathcal{L}\left(\hat{\mathbf{w}}^* + \sum_{i=1}^k \mathbf{l}_i \mathbf{u}^i\right) \geq \sum_{i=1}^k (c_i - 1) \mathbf{l}_i p_i / 2 - 2k\xi > 0$$

**Remark on Theorem 1.** It is widely known that the empirical minimizer is usually different from the true optimum. However, in practice it is difficult to know how the training loss shifts from the population loss. Therefore, the best we could do to minimize the empirical loss function (with some regularizers). However, Theorem 1 states that in the asymmetric case, we should pick a biased solution even if the shift is unknown. This insight can be distilled into practical algorithms to achieve better generalization, as we will discuss in Section 5.

## 4.2 Validating assumptions

We conducted a series of experiments with modern deep networks to show that the two assumptions introduced above are generally valid.

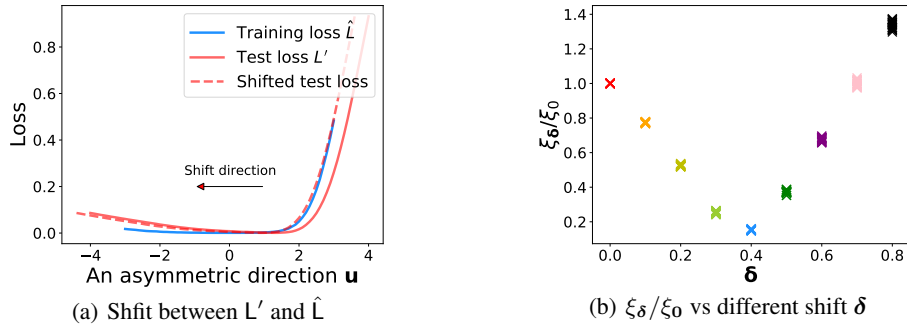


Figure 3: Shift exists between empirical loss and population loss for ResNet-110 on CIFAR-10.

**Verification of Assumption 1.** We show that a shift between  $L$  and  $\hat{L}$  is quite common in practice, by taking a ResNet-110 trained on CIFAR-10 as an example. Notice that we use test loss to represent  $L$  in practice. Since we could not visualize a shift in a high dimensional space, we randomly sample an asymmetric direction  $\mathbf{u}$  (more results are shown Appendix C) at the SGD solution  $\hat{\mathbf{w}}^*$ . The blue and red curves shown in Figure 3(a) are obtained by calculating  $\hat{L}(\hat{\mathbf{w}}^* + l\mathbf{u})$  and  $L'(\hat{\mathbf{w}}^* + l\mathbf{u})$  for  $l \in [-3, 3]$ , which correspond to the training and test loss, respectively.

We then try different shift values of  $\delta$  to “match” the two curves. As shown in Figure 3(a), after applying a horizontal shift  $\delta = 0.4$  to the test loss, the two curves overlap almost perfectly. Quantitatively, we can use the *shift gap* defined in Definition 3 to evaluate how well the two curves match each other after shifting. It turns out that  $\xi_{\delta=0.4} = 0.03$ , which is much lower than  $\xi_{\delta=0} = 0.22$  before shifting ( $\delta$  has only one dimension here). In Figure 3(b), we plot  $\xi_{\delta}/\xi_0$  as a function of  $\delta$ . Clearly, there exists a  $\delta$  that minimizes this ratio, indicating a good match.

We conducted the same experiments for different directions, models and datasets, and similar observations were made. Please refer to Appendix C for more results.

**Verification of Assumption 2.** This is a mild assumption that can be verified empirically. For example, we take a SGD solution of ResNet-110 on CIFAR-10 as  $\hat{\mathbf{w}}^*$ , and specify an asymmetric direction  $\mathbf{u}$  for  $\hat{\mathbf{w}}^*$ . We then randomly sample 100 different local adjustments for  $\mathbf{v} \in \mathbb{B}(25)$ . Based on these adjustments, we present the mean loss curves and standard variance zone on the asymmetric direction  $\mathbf{u}$  for all the points  $\hat{\mathbf{w}}^* + \mathbf{v} - \langle \mathbf{v}, \mathbf{u} \rangle \mathbf{u}$  in Figure 4. As we can see, the variance of these curves are very small, which means all of them are similar to each other. Moreover, we verified that  $\mathbf{u}$  is (4, 2, 0.1, 5.22)-asymmetric with respect to all neighboring points.

## 5 Averaging Generates Good Bias

In the previous section, we show that when the loss landscape of a local minimum is asymmetric, a solution with bias towards the flat side of the valley has better generalization performance. One immediate question is that how can we obtain such a solution via practical algorithms? Below we show that it can be achieved by simply taking the average of SGD iterates during the course of training. We first analyze the one dimensional case in Section 5.1, and then extend the analysis to the high dimensional case in Section 5.2.

Note that weight averaging is a classical algorithm in optimization [41], and recently regained its popularity in the context of deep learning [25, 5, 51]. Our following analysis can be viewed as a theoretical justification of recent algorithms that based on SGD iterates averaging.

### 5.1 One dimensional case

For asymmetric functions, as long as the learning rate is not too small, SGD will oscillate between the flat side and the sharp side. Below we focus on one round of oscillation, and show that the average of the iterates in each round has a bias on the flat side. Consequently, by aggregating all rounds of oscillation, averaging SGD iterates leads to a bias as well.

For each individual round  $i$ , we assume that it starts from the iteration when SGD goes from sharp side to flat side (denoted as  $w_0^i$ ), and ends at the iteration exactly before the iteration that SGD goes from sharp side to flat side again (denoted as  $w_{T_i}^i$ ). Here  $T_i$  denotes the number of iterations in the  $i$ -th rounds. The average iterate in the  $i$ -th round can be written as  $\bar{w} \triangleq \frac{1}{T_i} \sum_{j=0}^{T_i} w_j^i$ . For notational simplicity, we will omit the super script  $i$  on  $w_j^i$ .

The following theorem shows that the expectation of the average has bias on the flat side. To get a formal lower bound on  $\bar{w}$ , we consider the asymmetric case where  $\underline{r} = 0$ , and also assume lower bounds for the gradients on the function. We defer the proof to Appendix D.

**Theorem 2** (SGD averaging generates a bias). *Assume that a local minimizer  $w^* = 0$  is a  $(r, 0, a_+, c)$ -asymmetric valley, where  $b_- \leq \nabla L(w) \leq a_- < 0$  for  $w < 0$ , and  $0 < b_+ \leq \nabla L(w) \leq a_+$  for  $w \geq 0$ . Assume  $-a_- = ca_+$  for a large constant  $c$ , and  $\frac{-(b_- - \nu)}{b_+} = c' < \frac{e^{c/3}}{6}$ . The SGD updating rule is  $w_{t+1} = w_t - \eta(\nabla L(w) + \omega_t)$  where  $\omega_t$  is the noise and  $|\omega_t| < \nu$ , and assume  $\nu \leq a_+$ . Then we have*

$$\mathbb{E}[\bar{w}] > c_0 > 0,$$

where  $c_0$  is a constant that only depends on  $\eta, a_+, a_-, b_+, b_-$  and  $\nu$ .

Theorem 2 can be intuitively explained by Figure 5. If we run SGD on this one dimensional function, it will stay at the flat side for more iterations as the magnitude of the gradient on this side is much smaller. Therefore, the average of the locations is biased towards the flat side.

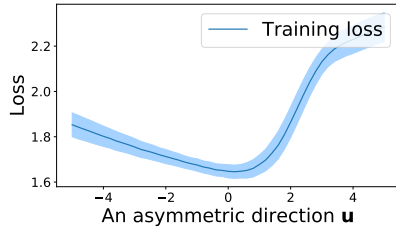


Figure 4: Training loss mean and variance for the neighborhood of  $\hat{w}^*$  at the direction of  $\mathbf{u}$ .

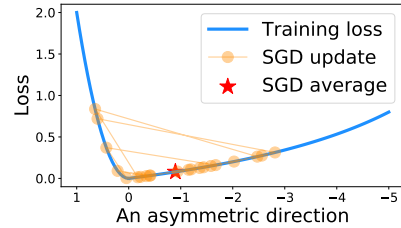


Figure 5: SGD iterates and their average on an asymmetric function.

### 5.2 High dimensional case

For high dimensional functions, the analysis on averaging SGD iterates would be more complicated compared to that given in the previous subsection. However, if we only care about the bias on a specific direction  $\mathbf{u}$ , we could directly apply Theorem 2 with one additional assumption. Specifically, if the projections of the loss function onto  $\mathbf{u}$  along the SGD trajectory satisfy the assumptions in Theorem 2, i.e., being asymmetric and the gradient on both sides have upper and lower bounds, then the claim of Theorem 2 directly applies. This is because only the gradient along the direction  $\mathbf{u}$  will affect the SGD trajectory projected onto  $\mathbf{u}$ , and we could safely omit all other directions.

We find that this assumption holds empirically. For a given SGD solution, we fix a random asymmetric direction  $\mathbf{u} \in \mathbb{R}^d$ , and sample the loss surface on direction  $\mathbf{u}$  that passes the  $t$ -th epoch of SGD trajectory (denoted as  $\mathbf{w}_t$ ), i.e., evaluate  $\hat{L}(\mathbf{w}_t + l\mathbf{u})$ , for  $0 \leq t \leq 200$  and  $l \in [-15, 15]$ . As shown in the Figure 6, after the first 40 epochs, the projected loss surfaces becomes relatively stable. Therefore, we can directly apply Theorem 2 to the direction  $\mathbf{u}$ .

As we will see in Section 6.1, compared with SGD solutions, SGD averaging indeed creates bias along different asymmetric directions, as predicted by our theory.

## 6 Experimental Observations

In this section, we empirically show that asymmetric valleys create interesting illusions when visualizing high dimensional loss landscape in low dimensional space. In addition, as a refinement of judging the generalization performance by the sharpness/flatness of a local minimum, we show that *where* the solution locates at a local minimum basin is important. We also find that batch normalization [24] seems to be a major cause for asymmetric valleys in deep networks, but the results are deferred to Appendix H due to space limit.

### 6.1 Experiments with weight averaging

Recently, Izmailov et al. [25] proposed the stochastic weight averaging (SWA) algorithm, which explicitly takes the average of SGD iterates to achieve better generalization. Inspired by their observation that “SWA leads to solutions corresponding to wider optima than SGD”, we provide a more refined explanation in this subsection. That is, averaging weights leads to “biased” solutions in an asymmetric valley, which correspond to better generalization.

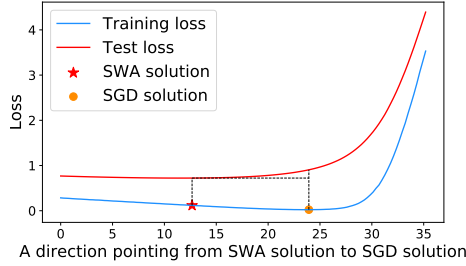


Figure 7: SWA solution and SGD solution interpolation (ResNet-164 on CIFAR-100)

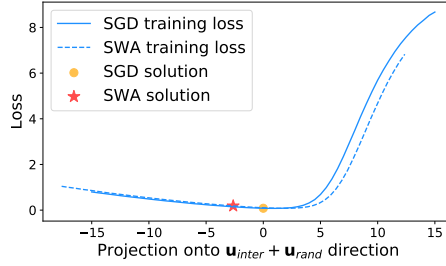


Figure 8: The average of SGD has a bias on flat side (ResNet-110 on CIFAR-100)

Specifically, we run the SWA algorithm (with decreasing learning rate) with popular deep networks, including ResNet-56, ResNet-110, ResNet-164, VGG-16, and DenseNet-100, on various datasets including CIFAR-10, CIFAR-100, SVHN and STL-10, following the configurations in [25]. Then we run SGD with small learning rate *from the SWA solutions* to find a solution located in the same basin (denoted as SGD).

In Figure 7, We draw an interpolation between the solutions obtained by SWA and SGD<sup>3</sup>. One can observe that there is no “bump” between these two solutions, meaning they are located in the same basin. Clearly, the SWA solution is biased towards the flat side, which verifies our theoretical analysis in Section 5. Further, we notice that although the biased SWA solution has higher training loss than the solution found by SGD, it indeed yields lower test loss. This verifies our analysis in Section 4. Similar observations are made on other networks and other datasets, which we present in Appendix E.

Table 1: Training and test accuracy on CIFAR-100.

Network	CIFAR-100	
	train	test
ResNet-110-SWA	94.98%	78.94%
ResNet-110-SGD	97.52%	78.29%
ResNet-164-SWA	97.48%	80.69%
ResNet-164-SGD	99.12%	76.56%

<sup>3</sup>Izmailov et al. [25] have done a similar experiment.



To further support our claim, we list our result in Table 1, from which we can observe that SGD solutions always have higher training accuracy, but worse test accuracy, compared to SWA solutions. This supports our claim in Theorem 1, which states that a bias towards the flat sides of asymmetric valleys could help improve generalization, although it yields higher training error.

**Verifying Theorem 2.** We further verify that averaging SGD solutions creates a bias towards the flat side in expectation for many other asymmetric directions, not just for the specific direction we discussed above.

We take a ResNet-110 trained on CIFAR-100 as an example. Denote  $\mathbf{u}_{inter}$  as the unit vector pointing from the SGD solution to the SWA solution,  $\mathbf{u}_{rand}$  as another unit random direction, and the direction  $\mathbf{u}_{inter} + \mathbf{u}_{rand}$  is used to explore the asymmetric landscape.

The results are shown in Figure 8, from which we can observe that SWA has a bias on the flat side compared with the SGD solution. We create 10 different random vectors for each network and each dataset, and similar observations can be made (see more examples in Appendix F).

**Batch size effect** In addition to SWA algorithm, we also observe similar trend when training with different batch sizes. The results are deferred to Appendix G.

## 6.2 Illusions created by asymmetric valleys

We further point out that visualizing the “width” of a given solution  $\mathbf{w}$  in a low-dimensional space may lead to illusive results. For example, one visualization technique used in [25] is to show how the loss changes along many random directions  $\mathbf{v}_i$ ’s drawn from the  $d$ -dimensional Gaussian distribution.

We take the large batch and small batch solutions from the previous subsection as an example. Figure 9 visualizes the “width” of the two solutions using the method described above. From the figure, one may draw the conclusion that small batch training leads to a wider minimum compared to large batch training. However, these two solutions are in fact from the *same* basin (see the discussion in Appendix G). In other words, the loss curvature near the two solutions looks different because they are located at *different locations* in a same asymmetric valley, instead of being located at *different local minima*. Similar observation holds for SWA and SGD solutions, see Figure 10<sup>4</sup>.

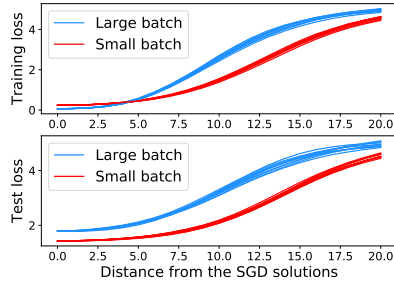


Figure 9: Random ray of large batch and small batch solution.

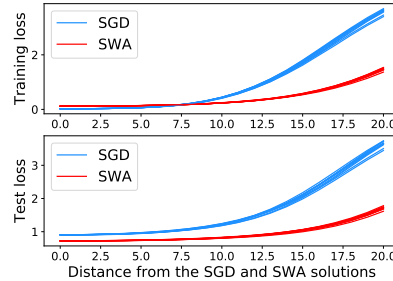


Figure 10: Random ray of SGD and SWA solution

## 7 Conclusion

In this paper, we introduced the notion of asymmetric valley to characterize the loss landscape of deep networks, expanding the current research that simply categorizes local minima by sharpness/flatness. This notion allowed us to analyze and understand the geometry of loss landscape from a new perspective. For example, based on a formal definition of asymmetric valley, we showed that a biased solution lying on the flat side of the valley generalizes better than the exact empirical minimizer. Further, it is proved that by averaging the weights obtained along the SGD trajectory naturally leads to such biased solution. We also conducted extensive experiments with state-of-the-art deep models to analyze the properties of asymmetric valleys. It is showed that due to the existence of asymmetric valleys, intriguing illusions can be created when visualizing high dimensional loss surface in the 1D space. We hope this work will deepen our understanding on the loss landscape of deep neural networks, and inspire new theories and algorithms that further improve generalization.

<sup>4</sup>Similar observations were made by Izmailov et al. [25] as well.



## References

- [1] Allen-Zhu, Z. How to make the gradients small stochastically: Even faster convex and nonconvex SGD. In *NeurIPS*, pp. 1165–1175, 2018.
- [2] Allen-Zhu, Z. Natasha 2: Faster non-convex optimization than SGD. In *NeurIPS*, pp. 2680–2691, 2018.
- [3] Allen-Zhu, Z. and Li, Y. NEON2: finding local minima via first-order oracles. In *NeurIPS*, pp. 3720–3730, 2018.
- [4] Arora, S., Ge, R., Neyshabur, B., and Zhang, Y. Stronger generalization bounds for deep nets via a compression approach. In *ICML*, volume 80 of *Proceedings of Machine Learning Research*, pp. 254–263. PMLR, 2018.
- [5] Athiwaratkun, B., Finzi, M., Izmailov, P., and Wilson, A. G. There are many consistent explanations of unlabeled data: Why you should average. In *International Conference on Learning Representations*, 2019. URL <https://openreview.net/forum?id=rkgKBhA5Y7>.
- [6] Bartlett, P. L., Foster, D. J., and Telgarsky, M. J. Spectrally-normalized margin bounds for neural networks. In *NIPS*, pp. 6241–6250, 2017.
- [7] Cesa-Bianchi, N., Conconi, A., and Gentile, C. On the generalization ability of on-line learning algorithms. In *NIPS*, pp. 359–366. MIT Press, 2001.
- [8] Chaudhari, P., Choromanska, A., Soatto, S., LeCun, Y., Baldassi, C., Borgs, C., Chayes, J. T., Sagun, L., and Zecchina, R. Entropy-sgd: Biasing gradient descent into wide valleys. In *ICLR*. OpenReview.net, 2017.
- [9] Choromanska, A., LeCun, Y., and Arous, G. B. Open problem: The landscape of the loss surfaces of multilayer networks. In *COLT*, volume 40 of *JMLR Workshop and Conference Proceedings*, pp. 1756–1760. JMLR.org, 2015.
- [10] Cooper, Y. The loss landscape of overparameterized neural networks. *CoRR*, abs/1804.10200, 2018.
- [11] Dinh, L., Pascanu, R., Bengio, S., and Bengio, Y. Sharp minima can generalize for deep nets. In *ICML*, volume 70 of *Proceedings of Machine Learning Research*, pp. 1019–1028. PMLR, 2017.
- [12] Dräxler, F., Veschgini, K., Salmhofer, M., and Hamprecht, F. A. Essentially no barriers in neural network energy landscape. In *ICML*, volume 80 of *Proceedings of Machine Learning Research*, pp. 1308–1317. PMLR, 2018.
- [13] Garipov, T., Izmailov, P., Podoprikin, D., Vetrov, D. P., and Wilson, A. G. Loss surfaces, mode connectivity, and fast ensembling of dnns. In *NeurIPS*, pp. 8803–8812, 2018.
- [14] Ge, R., Huang, F., Jin, C., and Yuan, Y. Escaping from saddle points - online stochastic gradient for tensor decomposition. In *COLT*, volume 40 of *JMLR Workshop and Conference Proceedings*, pp. 797–842. JMLR.org, 2015.
- [15] Ge, R., Lee, J. D., and Ma, T. Learning one-hidden-layer neural networks with landscape design. In *ICLR*. OpenReview.net, 2018.
- [16] Goodfellow, I. J. and Vinyals, O. Qualitatively characterizing neural network optimization problems. In *ICLR*, 2015.
- [17] Goyal, P., Dollár, P., Girshick, R. B., Noordhuis, P., Wesolowski, L., Kyrola, A., Tulloch, A., Jia, Y., and He, K. Accurate, large minibatch SGD: training imagenet in 1 hour. *CoRR*, abs/1706.02677, 2017.
- [18] Gunasekar, S., Lee, J., Soudry, D., and Srebro, N. Characterizing implicit bias in terms of optimization geometry. In *ICML*, volume 80 of *Proceedings of Machine Learning Research*, pp. 1827–1836. PMLR, 2018.

- [19] He, K., Zhang, X., Ren, S., and Sun, J. Deep residual learning for image recognition. In *CVPR*, pp. 770–778. IEEE Computer Society, 2016.
- [20] Hochreiter, S. and Schmidhuber, J. Simplifying neural nets by discovering flat minima. In *NIPS*, pp. 529–536. MIT Press, 1994.
- [21] Hoffer, E., Hubara, I., and Soudry, D. Train longer, generalize better: closing the generalization gap in large batch training of neural networks. In *NIPS*, pp. 1729–1739, 2017.
- [22] Huang, G., Li, Y., Pleiss, G., Liu, Z., Hopcroft, J. E., and Weinberger, K. Q. Snapshot ensembles: Train 1, get M for free. In *ICLR*. OpenReview.net, 2017.
- [23] Huang, G., Liu, Z., van der Maaten, L., and Weinberger, K. Q. Densely connected convolutional networks. In *CVPR*, pp. 2261–2269. IEEE Computer Society, 2017.
- [24] Ioffe, S. and Szegedy, C. Batch normalization: Accelerating deep network training by reducing internal covariate shift. In *ICML*, volume 37 of *JMLR Workshop and Conference Proceedings*, pp. 448–456. JMLR.org, 2015.
- [25] Izmailov, P., Podoprikin, D., Garipov, T., Vetrov, D. P., and Wilson, A. G. Averaging weights leads to wider optima and better generalization. In *UAI*, pp. 876–885. AUAI Press, 2018.
- [26] Jastrzebski, S., Kenton, Z., Arpit, D., Ballas, N., Fischer, A., Bengio, Y., and Storkey, A. J. Three factors influencing minima in SGD. *CoRR*, abs/1711.04623, 2017.
- [27] Ji, Z. and Telgarsky, M. Risk and parameter convergence of logistic regression. *CoRR*, abs/1803.07300, 2018.
- [28] Jin, C., Ge, R., Netrapalli, P., Kakade, S. M., and Jordan, M. I. How to escape saddle points efficiently. In *ICML*, volume 70 of *Proceedings of Machine Learning Research*, pp. 1724–1732. PMLR, 2017.
- [29] Jin, C., Liu, L. T., Ge, R., and Jordan, M. I. On the local minima of the empirical risk. In *NeurIPS*, pp. 4901–4910, 2018.
- [30] Jin, C., Netrapalli, P., and Jordan, M. I. Accelerated gradient descent escapes saddle points faster than gradient descent. In *COLT*, volume 75 of *Proceedings of Machine Learning Research*, pp. 1042–1085. PMLR, 2018.
- [31] Kawaguchi, K., Kaelbling, L. P., and Bengio, Y. Generalization in deep learning. *CoRR*, abs/1710.05468, 2017.
- [32] Keskar, N. S., Mudigere, D., Nocedal, J., Smelyanskiy, M., and Tang, P. T. P. On large-batch training for deep learning: Generalization gap and sharp minima. In *ICLR*. OpenReview.net, 2017.
- [33] Kleinberg, R., Li, Y., and Yuan, Y. An alternative view: When does SGD escape local minima? In *ICML*, volume 80 of *Proceedings of Machine Learning Research*, pp. 2703–2712. PMLR, 2018.
- [34] Li, H., Xu, Z., Taylor, G., Studer, C., and Goldstein, T. Visualizing the loss landscape of neural nets. In *NeurIPS*, pp. 6391–6401, 2018.
- [35] Masters, D. and Luschi, C. Revisiting small batch training for deep neural networks. *CoRR*, abs/1804.07612, 2018.
- [36] Mehta, D., Chen, T., Tang, T., and Hauenstein, J. D. The loss surface of deep linear networks viewed through the algebraic geometry lens. *CoRR*, abs/1810.07716, 2018.
- [37] Neyshabur, B., Bhojanapalli, S., McAllester, D., and Srebro, N. Exploring generalization in deep learning. In *NIPS*, pp. 5949–5958, 2017.
- [38] Neyshabur, B., Bhojanapalli, S., and Srebro, N. A pac-bayesian approach to spectrally-normalized margin bounds for neural networks. In *ICLR*. OpenReview.net, 2018.

- 415 [39] Neyshabur, B., Li, Z., Bhojanapalli, S., LeCun, Y., and Srebro, N. Towards understanding the  
416 role of over-parametrization in generalization of neural networks. *CoRR*, abs/1805.12076, 2018.
- 417 [40] Pennington, J. and Bahri, Y. Geometry of neural network loss surfaces via random matrix theory.  
418 In *ICML*, volume 70 of *Proceedings of Machine Learning Research*, pp. 2798–2806. PMLR,  
419 2017.
- 420 [41] Polyak, B. T. and Juditsky, A. B. Acceleration of stochastic approximation by averaging. *SIAM*  
421 *J. Control Optim.*, 30(4):838–855, July 1992. ISSN 0363-0129.
- 422 [42] Rakhlin, A., Shamir, O., and Sridharan, K. Making gradient descent optimal for strongly convex  
423 stochastic optimization. In *ICML*. icml.cc / Omnipress, 2012.
- 424 [43] Sagun, L., Evci, U., Güney, V. U., Dauphin, Y., and Bottou, L. Empirical analysis of the hessian  
425 of over-parametrized neural networks. In *ICLR (Workshop)*. OpenReview.net, 2018.
- 426 [44] Shalev-Shwartz, S., Shamir, O., Srebro, N., and Sridharan, K. Stochastic convex optimization.  
427 In *COLT*, 2009.
- 428 [45] Simonyan, K. and Zisserman, A. Very deep convolutional networks for large-scale image  
429 recognition. In *ICLR*, 2015.
- 430 [46] Smith, S. L. and Le, Q. V. A bayesian perspective on generalization and stochastic gradient  
431 descent. In *ICLR*. OpenReview.net, 2018.
- 432 [47] Smith, S. L., Kindermans, P., Ying, C., and Le, Q. V. Don’t decay the learning rate, increase the  
433 batch size. In *ICLR*. OpenReview.net, 2018.
- 434 [48] Soudry, D., Hoffer, E., Nacson, M. S., Gunasekar, S., and Srebro, N. The implicit bias of  
435 gradient descent on separable data. *Journal of Machine Learning Research*, 19:70:1–70:57,  
436 2018.
- 437 [49] Wu, L., Zhu, Z., and E, W. Towards understanding generalization of deep learning: Perspective  
438 of loss landscapes. *CoRR*, abs/1706.10239, 2017.
- 439 [50] Xu, Y., Rong, J., and Yang, T. First-order stochastic algorithms for escaping from saddle points  
440 in almost linear time. In *NeurIPS*, pp. 5535–5545, 2018.
- 441 [51] Yang, G., Zhang, T., Kirichenko, P., Bai, J., Wilson, A. G., and Sa, C. D. Swalp: Stochastic  
442 weight averaging in low precision training. In *ICML*, 2019.
- 443 [52] Zhou, W., Veitch, V., Austern, M., Adams, R. P., and Orbanz, P. Non-vacuous generaliza-  
444 tion bounds at the imagenet scale: a PAC-bayesian compression approach. In *International*  
445 *Conference on Learning Representations*, 2019.

## 446 A Additional Figures for Section 3.2: Asymmetric Directions

447 To show that asymmetric valley can be commonly observed, we conduct experiments ranging from  
 448 the simplest network to modern deep neural networks.

449 **A simple case** First, we will show that asymmetric valley can be observed on a simple MLP (one  
 450 hidden layer with 10 hidden neurons) on a logistic regression task in Figure 11

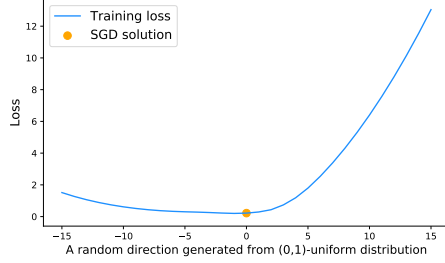


Figure 11: Asymmetric direction for a solution of MLP on logistic regression.  $(\bar{r}, r, p, c) = (10.0, 5.0, 0.11, 6.0)$ .

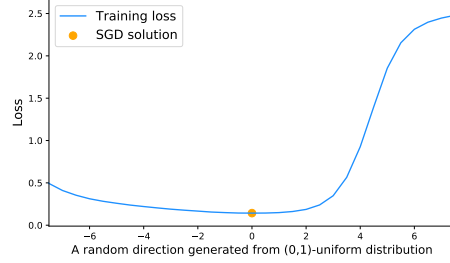


Figure 12: Asymmetric direction for a solution of ResNet-164 on CIFAR-10.  $(\bar{r}, r, p, c) = (4.0, 2.5, 0.033, 4.8)$ .

451 **Other datasets and networks** See Figure 12, Figure 13, Figure 14, Figure 15, and Figure 16.

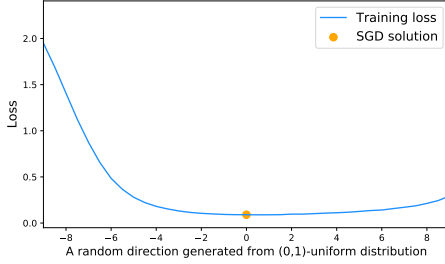


Figure 13: Asymmetric direction for a solution of DenseNet-100 on CIFAR-10.  $(\bar{r}, r, p, c) = (7.0, 5.0, 0.030, 4.8)$ .

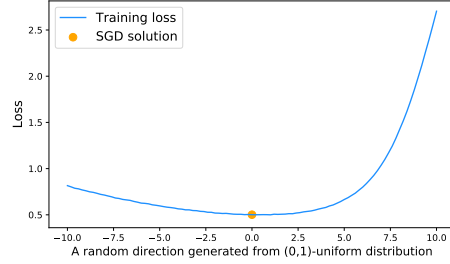


Figure 14: Asymmetric direction for a solution of ResNet-110 on CIFAR-100.  $(\bar{r}, r, p, c) = (7.0, 5.0, 0.039, 2.7)$ .

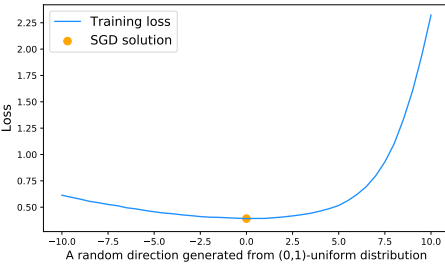


Figure 15: Asymmetric direction for a solution of ResNet-164 on CIFAR-100.  $(\bar{r}, r, p, c) = (7.0, 5.0, 0.031, 2.5)$ .

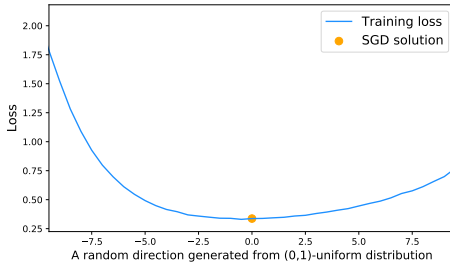


Figure 16: Asymmetric direction for a solution of DenseNet-100 on CIFAR-100.  $(\bar{r}, r, p, c) = (8.5, 6.5, 0.087, 2.1)$ .

## 452 B Proof for Theorem 1

453 *Proof.* Since  $\delta$  has  $2^d$  possible value for a given  $\bar{\delta}$ , we can use an integer  $j \in \{0, \dots, 2^d - 1\}$  to  
 454 represent each value. When writing  $j$  in binary, its  $i$ -th digit represents whether  $\delta_i = \bar{\delta}_i$  (equal to 1)  
 455 or  $\delta_i = -\bar{\delta}_i$  (equal to 0). We use  $j \wedge 2^i$  to represent the bitwise AND operator between  $j$  and  $2^i$ ,  
 456 which equals 0 if the  $i$ -th digit of  $j$  is 0.

457 To prove our theorem, it suffices to show that for any  $i \in [k]$ ,

$$\mathbb{E}_\delta \mathbb{L} \left( \hat{\mathbf{w}}^* + \sum_{i_0=1}^{i-1} \mathbf{l}_{i_0} \mathbf{u}_{i_0} \right) - \mathbb{E}_\delta \mathbb{L} \left( \hat{\mathbf{w}}^* + \sum_{i_0=1}^i \mathbf{l}_{i_0} \mathbf{u}_{i_0} \right) \geq (c_i - 1) \mathbf{l}_i p_i / 2 - 2\xi > 0 \quad (1)$$

458 If (1) is true, it suffices to take summation over  $i$  on both sides, and we will get our conclusion.  
 459 Therefore, below we will prove (1).

$$\begin{aligned} & \mathbb{E}_\delta \mathbb{L} \left( \hat{\mathbf{w}}^* + \sum_{i_0=1}^{i-1} \mathbf{l}_{i_0} \mathbf{u}_{i_0} \right) - \min_{\mathbf{w}} \mathbb{L}(\mathbf{w}) + \min_{\mathbf{w}} \hat{\mathbb{L}}(\mathbf{w}) \\ &= \mathbb{E}_\delta \mathbb{L}' \left( \hat{\mathbf{w}}^* + \sum_{i_0=1}^{i-1} \mathbf{l}_{i_0} \mathbf{u}_{i_0} \right) \stackrel{\textcircled{1}}{\geq} \frac{1}{2^d} \sum_{j=0}^{2^d-1} \hat{\mathbb{L}} \left( \hat{\mathbf{w}}^* + \sum_{i_0=1}^{i-1} \mathbf{l}_{i_0} \mathbf{u}_{i_0} + \delta^j \right) - \xi \\ &= \frac{1}{2^d} \sum_{\substack{j=0 \\ j \wedge 2^i = 0}}^{2^d-1} \left[ \hat{\mathbb{L}} \left( \hat{\mathbf{w}}^* + \sum_{i_0=1}^{i-1} \mathbf{l}_{i_0} \mathbf{u}_{i_0} + \delta^j \right) + \hat{\mathbb{L}} \left( \hat{\mathbf{w}}^* + \sum_{i_0=1}^{i-1} \mathbf{l}_{i_0} \mathbf{u}_{i_0} + \delta^{j+2^i} \right) \right] - \xi \end{aligned} \quad (2)$$

460 Where  $\textcircled{1}$  holds by Assumption 1, and the fact that  $\|\sum_{i_0=1}^{i-1} \mathbf{l}_{i_0} \mathbf{u}_{i_0}\|_2 \leq \|\mathbf{l}\|_2 = R$ . For every  $j$  s.t.  
 461  $j \wedge 2^i = 0$ ,

$$\begin{aligned} & \hat{\mathbf{w}}^* + \sum_{i_0=1}^i \mathbf{l}_{i_0} \mathbf{u}_{i_0} + \delta^j \\ &= \hat{\mathbf{w}}^* + \sum_{i_0=1}^i \mathbf{l}_{i_0} \mathbf{u}_{i_0} + \delta^j + \langle \delta^j, \mathbf{u}^i \rangle \mathbf{u}^i - \langle \delta^j, \mathbf{u}^i \rangle \mathbf{u}^i \\ &= \hat{\mathbf{w}}^* + \sum_{i_0=1}^{i-1} \mathbf{l}_{i_0} \mathbf{u}_{i_0} + \delta^j - \bar{\delta}_i \mathbf{u}^i - \langle \delta^j, \mathbf{u}^i \rangle \mathbf{u}^i + \mathbf{l}_i \mathbf{u}^i \\ &= \hat{\mathbf{w}}^* + \sum_{i_0=1}^{i-1} \mathbf{l}_{i_0} \mathbf{u}_{i_0} + \delta^j - \langle \delta^j, \mathbf{u}^i \rangle \mathbf{u}^i + (\mathbf{l}_i - \bar{\delta}_i) \mathbf{u}^i \end{aligned}$$

462 Since  $\|\sum_{i_0=1}^{i-1} \mathbf{l}_{i_0} \mathbf{u}_{i_0}\|_2 \leq \|\mathbf{l}\|_2$ ,  $\|\delta^j\|_2 = \|\bar{\delta}\|_2$ , we know that  $\forall j$ ,  $\sum_{i_0=1}^{i-1} \mathbf{l}_{i_0} \mathbf{u}_{i_0} + \delta^j \in \mathbb{B}(R')$ . By  
 463 Assumption 2, for every  $i \in [k]$ ,  $\mathbf{u}^i$  is asymmetric with respect to  $\hat{\mathbf{w}}^* + \sum_{i_0=1}^{i-1} \mathbf{l}_{i_0} \mathbf{u}_{i_0} + \delta^j - \langle \delta^j, \mathbf{u}^i \rangle \mathbf{u}^i$ .  
 464 Since  $\mathbf{l}_i \leq \bar{\delta}_i - r$ , we have  $\mathbf{l}_i - \bar{\delta}_i < -r$ . By the definition of asymmetric direction, we know

$$\hat{\mathbb{L}} \left( \hat{\mathbf{w}}^* + \sum_{i_0=1}^{i-1} \mathbf{l}_{i_0} \mathbf{u}_{i_0} + \delta^j \right) \geq \hat{\mathbb{L}} \left( \hat{\mathbf{w}}^* + \sum_{i_0=1}^i \mathbf{l}_{i_0} \mathbf{u}_{i_0} + \delta^j \right) + c_i \mathbf{l}_i p_i \quad (3)$$

465 Similarly,

$$\begin{aligned} & \hat{\mathbf{w}}^* + \sum_{i_0=1}^i \mathbf{l}_{i_0} \mathbf{u}_{i_0} + \delta^{j+2^i} \\ &= \hat{\mathbf{w}}^* + \sum_{i_0=1}^{i-1} \mathbf{l}_{i_0} \mathbf{u}_{i_0} + \delta^{j+2^i} + \langle \delta^{j+2^i}, \mathbf{u}^i \rangle \mathbf{u}^i - \langle \delta^{j+2^i}, \mathbf{u}^i \rangle \mathbf{u}^i + \mathbf{l}_i \mathbf{u}^i \\ &= \hat{\mathbf{w}}^* + \sum_{i_0=1}^{i-1} \mathbf{l}_{i_0} \mathbf{u}_{i_0} + \delta^{j+2^i} - \langle \delta^{j+2^i}, \mathbf{u}^i \rangle \mathbf{u}^i + (\bar{\delta}_i + \mathbf{l}_i) \mathbf{u}^i \end{aligned}$$

466 Since  $\mathbf{l}_i \leq r - \bar{\delta}_i$ , we have  $\bar{\delta}_i + \mathbf{l}_i \leq r$ . Therefore,

$$\hat{\mathbb{L}} \left( \hat{\mathbf{w}}^* + \sum_{i_0=1}^{i-1} \mathbf{l}_{i_0} \mathbf{u}_{i_0} + \delta^{j+2^i} \right) \geq \hat{\mathbb{L}} \left( \hat{\mathbf{w}}^* + \sum_{i_0=1}^i \mathbf{l}_{i_0} \mathbf{u}_{i_0} + \delta^{j+2^i} \right) - \mathbf{l}_i p_i \quad (4)$$

Combining (3) and (4), we have,

$$\begin{aligned}
(2) &\geq \frac{1}{2^d} \sum_{j=0}^{2^d-1} \left[ \hat{\mathcal{L}} \left( \hat{\mathbf{w}}^* + \sum_{i_0=1}^i \mathbf{l}_{i_0} \mathbf{u}_{i_0} + \delta^j \right) + c_i \mathbf{l}_i p_i + \hat{\mathcal{L}} \left( \hat{\mathbf{w}}^* + \sum_{i_0=1}^i \mathbf{l}_{i_0} \mathbf{u}_{i_0} + \delta^{j+2^i} \right) - \mathbf{l}_i p_i \right] - \xi \\
&= \frac{1}{2^d} \sum_{j=0}^{2^d-1} \left[ \hat{\mathcal{L}} \left( \hat{\mathbf{w}}^* + \sum_{i_0=1}^i \mathbf{l}_{i_0} \mathbf{u}_{i_0} + \delta^j \right) \right] + (c_i - 1) \mathbf{l}_i p_i / 2 - \xi \\
&\stackrel{\textcircled{2}}{\geq} \mathbb{E}_{\delta} \mathcal{L}' \left( \hat{\mathbf{w}}^* + \sum_{i_0=1}^i \mathbf{l}_{i_0} \mathbf{u}_{i_0} \right) + (c_i - 1) \mathbf{l}_i p_i / 2 - 2\xi \\
&= \mathbb{E}_{\delta} \mathcal{L} \left( \hat{\mathbf{w}}^* + \sum_{i_0=1}^i \mathbf{l}_{i_0} \mathbf{u}_{i_0} \right) - \min_{\mathbf{w}} \mathcal{L}(\mathbf{w}) + \min_{\mathbf{w}} \hat{\mathcal{L}}(\mathbf{w}) + (c_i - 1) \mathbf{l}_i p_i / 2 - 2\xi
\end{aligned}$$

Where ② holds by Assumption 1 and the fact that  $\|\sum_{i_0=1}^i \mathbf{l}_{i_0} \mathbf{u}_{i_0}\|_2 \leq \|\mathbf{l}\|_2 = R$ . That means,

$$\mathbb{E}_{\delta} \mathcal{L} \left( \hat{\mathbf{w}}^* + \sum_{i_0=1}^i \mathbf{l}_{i_0} \mathbf{u}_{i_0} \right) \geq \mathbb{E}_{\delta} \mathcal{L} \left( \hat{\mathbf{w}}^* + \sum_{i_0=1}^i \mathbf{l}_{i_0} \mathbf{u}_{i_0} \right) + (c_i - 1) \mathbf{l}_i p_i / 2 - 2\xi > 0$$

Where the last inequality holds as  $\mathbf{l}_i > \frac{4\xi}{(c_i-1)p_i}$ .

□

## C Additional Figures for Section 4.2: Shift Exists Empirically

See Figure 17, Figure 18, and Figure 19.

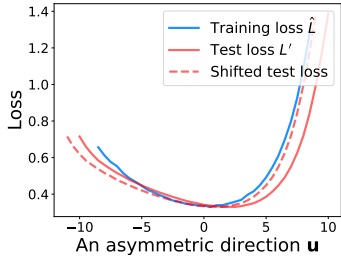


Figure 17: Shift on asymmetric direction (DenseNet-100 on CIFAR-100),  $\xi_{\delta=1} = 0.119$ ,  $\xi_{\delta=0} = 0.439$

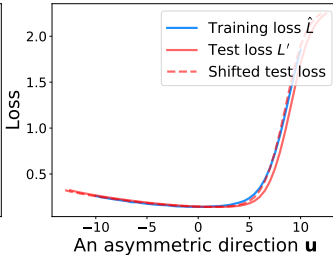


Figure 18: Shift on asymmetric direction (ResNet-164 on CIFAR-10),  $\xi_{\delta=0.5} = 0.0699$ ,  $\xi_{\delta=0} = 0.189$

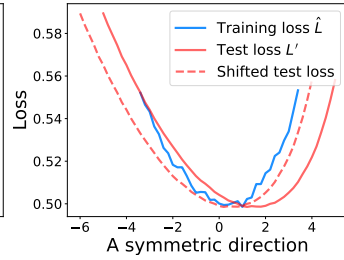


Figure 19: Shift on symmetric direction (ResNet-110 on CIFAR-100),  $\xi_{\delta=1} = 0.0197$ ,  $\xi_{\delta=0} = 0.0431$

## D Proof for Theorem 2

To prove Theorem 2, we will need the following concentration bound.

**Lemma 3** (Azuma's inequality). *Let  $X_1, X_2, X_3, \dots, X_n$  be independent random variables satisfying  $|X_i - \mathbb{E}[X_i]| \leq c_i$ , for  $1 \leq i \leq n$ . We have the following bound for  $X = \sum_{i=1}^n X_i$ :*

$$\Pr(|X - \mathbb{E}(X)| \geq \lambda) \leq 2e^{-\frac{\lambda^2}{2 \sum_{i=1}^n c_i^2}}$$

Let  $p_{\min} \triangleq -\eta(a_- + a_+ + 2\nu)$ ,  $p_{\max} \triangleq -\eta(b_- - \nu)$ . Since  $-a_- = ca_+$ , we know  $p_{\min} > (c-1)\eta a_+ - 2\eta\nu$ . First, we have the following bounds on the first step  $w_0$ .

**Lemma 4.** *For every  $i \in [h]$ ,  $w_0 \in [p_{\min}, p_{\max}]$ .*



478 *Proof.* Since  $w_0$  is the first step that SGD jumps from the flat side to the sharp side, denote the  
 479 previous location as  $w_{-1} < 0$ . Since  $w_{-1}$  is at the sharp side, we know that the gradient is  
 480  $\nabla L(w_{-1}) \leq a_-$ . Therefore, we have

$$w_0 = w_{-1} - \eta(\nabla L(w_{-1}) + \omega_{-1})$$

481 Where  $\omega_{-1}$  is the noise bounded by  $\nu$ .

482 At the time when SGD jump from the flat side to sharp side, denote the target position as  $w'_{-1}$ . We  
 483 know that  $w'_{-1} \in [-\eta(a_+ + \nu), 0]$ . Since the gradient on the sharp side is at most  $a_-$ , we know the  
 484 next step is lower bounded by  $-\eta(a_+ + 2\nu + a_-) = p_{\min} > 0$ . In other words, SGD stays at the  
 485 sharp side for only 1 iterations (this matches with our empirical observation, see e.g. Figure 5).

486 That means, the bound on  $w'_{-1}$  can be applied to  $w_{-1}$  as well, because they are the same iterate. By  
 487 applying the upper and lower bound on  $\nabla L(w_{-1})$ , we get:

$$w_0 \geq -\eta(a_+ + \nu) - \eta(a_- + \nu) = p_{\min}$$

488 and also

$$w_0 \leq 0 - \eta(b_- - \nu) = p_{\max}$$

489 □

490 Below we first define  $T_{\min} \triangleq \left( \frac{-\sqrt{2\nu} \log^{1/2}(2\tau) + \sqrt{2\nu^2 \log(2\tau) - 4a_+(a_- + a_+ + 2\nu)}}{2a_+} \right)^2$ , where  $\tau$  is a con-  
 491 stant with value to be set later.  $T_{\min}$  satisfies the following inequality.

492 **Lemma 5.**  $\forall t \leq T_{\min}, p_{\min} - t\eta a_+ - \sqrt{2t\eta\nu} \log^{1/2}(2\tau) \geq 0$ .

493 *Proof.* By the definition of  $p_{\min}$ , we have

$$\begin{aligned} & -\eta(a_- + a_+ + 2\nu) - t\eta a_+ - \sqrt{2t\eta\nu} \log^{1/2}(2\tau) \geq 0 \\ \Leftrightarrow & (a_- + a_+ + 2\nu) + ta_+ + \sqrt{2t\eta\nu} \log^{1/2}(2\tau) \leq 0 \\ \Leftrightarrow & (a_- + a_+ + 2\nu) + \Delta^2 a_+ + \sqrt{2\Delta r} \log^{1/2}(2\tau) \leq 0 \quad (\Delta \triangleq \sqrt{t}) \\ \Leftrightarrow & \Delta \in \left[ 0, \frac{-\sqrt{2\nu} \log^{1/2}(2\tau) + \sqrt{2\nu^2 \log(2\tau) - 4a_+(a_- + a_+ + 2\nu)}}{2a_+} \right] \\ \Leftrightarrow & t \leq \left( \frac{-\sqrt{2\nu} \log^{1/2}(2\tau) + \sqrt{2\nu^2 \log(2\tau) - 4a_+(a_- + a_+ + 2\nu)}}{2a_+} \right)^2 \end{aligned} \quad \square$$

494 Now, we have the following theorem that says with decent probability, the minimum number of  
 495 iterates on the flat side in  $i$ -th round is at least  $T_{\min}$ .

496 **Theorem 6.** *If we start at  $w_0 \geq p_{\min}$ , for every fixed  $\tau > T_{\min}$ , with probability at least  $1 - \frac{T_{\min}}{\tau}$ ,  
 497 we have  $\forall t \leq T_{\min}, w_t > w_0 - t\eta a_+ - \sqrt{2t\eta\nu} \log^{1/2}(2\tau) \geq 0$ .*

498 *Proof.* Define filtration  $\mathcal{F}_t = \sigma\{\omega_0, \dots, \omega_{t-1}\}$ , where  $\sigma\{\cdot\}$  denotes the sigma field. Define the  
 499 event  $\mathfrak{E}_T = \{\forall t \leq T, w_t > w_0 - t\eta a_+ - \sqrt{2t\eta\nu} \log^{1/2}(2\tau)\}$  and define  $G_t = w_0 - w_t - t\eta a_+ + M$ ,  
 500 where  $M \triangleq (T_{\min} + 1)(w_0 + \nu + 2\eta a_+)$ . Since we only consider the case  $t \leq T_{\min}$ , we have

$$G_t = w_0 - w_t - t\eta a_+ + (T_{\min} + 1)(w_0 + \nu + 2\eta a_+) > w_0 - w_t - t\eta a_+ + w_t + t\eta a_+ > 0$$

501 Therefore,  $G_t$  is always positive. By SGD updating rule, we have

$$\begin{aligned} \mathbb{E}[G_{t+1} \mathbb{1}_{\mathfrak{E}_t} | \mathcal{F}_t] &= \mathbb{E}[(w_0 - w_{t+1} - (t+1)\eta a_+ + M) \mathbb{1}_{\mathfrak{E}_t} | \mathcal{F}_t] \\ &\leq \mathbb{E}[(w_0 - w_t + \eta \omega_t - t\eta a_+ + M) \mathbb{1}_{\mathfrak{E}_t} | \mathcal{F}_t] = w_0 - w_t - t\eta a_+ + M = G_t \mathbb{1}_{\mathfrak{E}_t} \end{aligned} \quad (5)$$

502 Since  $\mathbb{1}_{\mathfrak{E}_t} \leq \mathbb{1}_{\mathfrak{E}_{t-1}}$ , and  $G_t$  is always positive, we have

$$G_t \mathbb{1}_{\mathfrak{E}_t} \leq G_t \mathbb{1}_{\mathfrak{E}_{t-1}} \quad (6)$$

503 Combining (5) and (6) together, we know  $G_t \mathbb{1}_{\mathfrak{E}_{t-1}}$  is a supermartingale.

504 We can also bound the absolute value of the difference in every iteration:

$$\begin{aligned} & |G_{t+1} \mathbb{1}_{\mathfrak{E}_t} - \mathbb{E}[G_{t+1} \mathbb{1}_{\mathfrak{E}_t} | \mathcal{F}_t]| \\ &= |(w_0 - w_{t+1} - (t+1)\eta a_+ + M) - (w_0 - w_t - \nabla L(w_t) - (t+1)\eta a_+ + M) | \mathcal{F}_t] \\ &\leq \eta \nu \end{aligned}$$

505 By Azuma's inequality, we get:

$$\Pr(G_t \mathbb{1}_{\mathfrak{E}_{t-1}} - G_0 \geq \lambda) \leq 2e^{-\frac{\lambda^2}{2t\eta^2\nu^2}}$$

506 That gives,

$$\Pr(G_t \mathbb{1}_{\mathfrak{E}_{t-1}} - G_0 \geq \sqrt{2t\eta\nu} \log^{1/2}(2\tau)) \leq 1/\tau$$

507 That means, if  $\mathbb{1}_{\mathfrak{E}_{t-1}}$  holds, with probability at least  $1 - 1/\tau$ ,

$$w_0 - w_t - t\eta a_+ + M < \sqrt{2t\eta\nu} \log^{1/2}(2\tau) + G_0 = \sqrt{2t\eta\nu} \log^{1/2}(2\tau) + M$$

508 Which gives

$$w_t > w_0 - t\eta a_+ - \sqrt{2t\eta\nu} \log^{1/2}(2\tau)$$

509 In other words, that means if  $\mathbb{1}_{\mathfrak{E}_{t-1}}$  holds, then  $\mathbb{1}_{\mathfrak{E}_t}$  also holds with probability at least  $1 - 1/\tau$ .

510 Therefore, if we are running  $T_{\min}$  steps, we know that with probability at least  $1 - \frac{T_{\min}}{\tau}$ ,  $\mathbb{1}_{\mathfrak{E}_{T_{\min}}}$   
511 holds. Therefore, by Lemma 5,

$$\forall t \leq T_{\min}, w_t > w_0 - t\eta a_+ - \sqrt{2t\eta\nu} \log^{1/2}(2\tau) \geq p_{\min} - t\eta a_+ - \sqrt{2t\eta\nu} \log^{1/2}(2\tau) \geq 0 \quad \square$$

512 Similarly, we define  $T_{\max} \triangleq \left( \frac{-\sqrt{2\nu} \log^{1/2}(2\tau) + \sqrt{2\nu^2 \log(2\tau) - 4(b_- - \nu)b_+}}{2b_+} \right)^2$ , which satisfies the fol-  
513 lowing inequality.

514 **Lemma 7.**  $p_{\max} - T_{\max}\eta b_+ - \sqrt{2T_{\max}}\eta\nu \log^{1/2}(2\tau) < 0$ .

515 *Proof.* By the definition of  $p_{\max}$ , we want to show that

$$(b_- - \nu) + T_{\max}b_+ + \sqrt{2T_{\max}}\eta\nu \log^{1/2}(2\tau) \geq 0$$

516 Which holds by the definition of  $T_{\max}$ .  $\square$

517 The Theorem below shows with decent probability,  $T_{\max} - 1$  is an upper bound on the total number  
518 of iterates on the flat side in the  $i$ -th round.

519 **Theorem 8.** If  $w_0 \leq p_{\max}$ , with probability at least  $1 - \frac{T_{\max}}{\tau}$ ,  $w_{T_{\max}} < 0$ .

520 *Proof.* Define event  $\mathfrak{E}'_T = \{\forall t \leq T, w_t < w_0 - t\eta b_+ + \sqrt{2t\eta\nu} \log^{1/2}(2\tau)\}$ , and  $G'_t = w_t + t\eta b_+ >$   
521 0.

522 We have

$$\begin{aligned} & \mathbb{E}[G'_{t+1} \mathbb{1}_{\mathfrak{E}'_t} | \mathcal{F}_t] \\ &= \mathbb{E}[(w_{t+1} + (t+1)\eta b_+) \mathbb{1}_{\mathfrak{E}'_t} | \mathcal{F}_t] \\ &\leq \mathbb{E}[(w_t - \eta \omega_t + t\eta b_+) \mathbb{1}_{\mathfrak{E}'_t} | \mathcal{F}_t] \\ &= (w_t + t\eta b_+) \mathbb{1}_{\mathfrak{E}'_t} \\ &= G'_t \mathbb{1}_{\mathfrak{E}'_t} \end{aligned}$$

523 Moreover, we know  $\mathbb{1}_{\mathfrak{E}'_t} \leq \mathbb{1}_{\mathfrak{E}'_{t-1}}$ , which means  $G'_t \mathbb{1}_{\mathfrak{E}'_t} \leq G'_t \mathbb{1}_{\mathfrak{E}'_{t-1}}$ . So  $G'_t \mathbb{1}_{\mathfrak{E}'_{t-1}}$  is a supermartingale.

524 We can also bound the absolute value of the difference in every iteration:

$$\begin{aligned} & |G'_{t+1} \mathbb{1}_{\mathfrak{E}'_t} - \mathbb{E}[G'_{t+1} \mathbb{1}_{\mathfrak{E}'_t} | \mathcal{F}_t]| \\ &= |(w_{t+1} + (t+1)\eta b_+) - (w_t - \eta \nabla L(w_t) + (t+1)\eta b_+) | \mathcal{F}_t]| \\ &\leq \eta \nu \end{aligned}$$

525 Using Azuma inequality, we get

$$\Pr \left( G'_t \mathbb{1}_{\mathfrak{E}'_{t-1}} - G'_0 \geq \sqrt{2t\eta\nu} \log^{1/2}(2\tau) \right) \leq 2e^{-\frac{t\eta^2\nu^2 \log(2\tau)}{t\eta^2\nu^2}} = \frac{1}{\tau}$$

526 That means, if  $\mathbb{1}_{\mathfrak{E}'_{t-1}}$  holds, with probability at least  $1 - 1/\tau$ ,

$$w_t < w_0 - t\eta b_+ + \sqrt{2t\eta\nu} \log^{1/2}(2\tau)$$

527 In other words,  $\mathbb{1}_{\mathfrak{E}'_t}$  also holds. Therefore, if we are running  $T_{\max}$  steps, we know that with probability

528 at least  $1 - \frac{T_{\max}}{\tau}$ ,  $\mathbb{1}_{\mathfrak{E}'_{T_{\max}}}$  holds. Therefore, by Lemma 7, we know

$$w_{T_{\max}} < w_0 - T_{\max}\eta b_+ - \sqrt{2T_{\max}\eta\nu} \log^{1/2}(2\tau) < 0 \quad \square$$

529 **Remark.** To make sure Theorem 6 is not vacuous, we need to make sure that  $T_{\min} \geq 1$ . If we want  
530 to make  $T_{\min}$ , say, at least 2, by Lemma 5, we have:

$$p_{\min} - 2\eta a_+ - 2\eta\nu \log^{1/2}(2\tau) \geq 0$$

531 Notice that  $p_{\min} > (c-1)\eta a_+ - 2\eta\nu$ , so we could solve the above inequality and get

$$\begin{aligned} & (c-1)\eta a_+ - 2\eta\nu - 2\eta a_+ - 2\eta\nu \log^{1/2}(2\tau) \geq 0 \\ \Rightarrow & \frac{(c-3)a_+ - 2\nu}{2\nu} \geq \log^{1/2}(2\tau) \\ \Rightarrow & \tau \leq \frac{e^{\left(\frac{(c-3)a_+ - 2\nu}{2\nu}\right)^2}}{2} \end{aligned}$$

532 Since we assume that  $c$  is a large constant and  $a_+ \geq \nu$ , so  $\tau$  can be fairly large in order to make sure  
533  $T_{\min} \geq 2$ . We also know that  $T_{\min} \leq \frac{-(a_- + a_+ + 2\nu)}{a_+} < c$ .

534 On the other hand, by simple calculation, we know  $T_{\max} \leq \frac{-(b_- - \nu)}{b_+} < c' < \frac{e^{c/3}}{6}$ . Therefore, we  
535 can always pick a  $\tau$  such that  $\frac{T_{\min} + T_{\max}}{\tau} \leq \frac{1}{2}$ . So finally, we are ready to prove Theorem 2.

536 *Proof of Theorem 2.* By Lemma 4 and Theorem 8,  $T_{\max}$  is an upper bound on the length of the  $i$ -th  
537 round. By Theorem 6, we know that SGD will stay at flat side for at least  $T_{\min}$  steps, and each step is  
538 lower bounded by  $w_t > w_0 - t\eta a_+ - \sqrt{2t\eta\nu} \log^{1/2}(2\tau)$ , therefore we know that with probability  
539  $1 - \frac{T_{\min} + T_{\max}}{\tau}$ :

$$\begin{aligned} \frac{1}{T_i} \sum_{j=0}^{T_i} w_j^i &\geq \frac{1}{T_{\max}} \left( \sum_{t=0}^{T_{\min}} [w_0 - t\eta a_+ - \sqrt{2t\eta\nu} \log^{1/2}(2\tau)] - \eta(a_+ + \nu) \right) \\ &\geq \frac{1}{T_{\max}} \left( \eta a_+ \frac{(T_{\min} + 1)T_{\min}}{2} + \sqrt{2T_{\min}\eta\nu} \log^{1/2}(2\tau) - \eta(a_+ + \nu) \right) \\ &\geq \frac{T_{\min}^2}{T_{\max}} \eta a_+ \end{aligned}$$

540 The above inequality discussed the scenario when Theorem 6 and Theorem 8 hold. If they do not hold,  
541 which happens with probability at most  $\frac{T_{\min} + T_{\max}}{\tau}$ , we need to get lower bound for  $\frac{1}{T_i} \sum_{j=0}^{T_i} w_j^i$ .

542 Notice that by Lemma 4, we know that SGD stays at the sharp side for at most 1 iterate in each round,  
 543 and also the iterates on the flat sides are always positive with  $w_0 \geq p_{\min} > \eta(a_+ + \nu)$ . Therefore,  
 544 we have the following trivial bound:

$$\frac{1}{T_i} \sum_{j=0}^{T_i} w_j^i \geq \frac{-\eta(a_+ + \nu) + w_0}{2} > 0$$

545 Combining two cases together we get

$$\mathbb{E} \left[ \frac{1}{T_i} \sum_{j=0}^{T_i} w_j^i \right] \geq \left( 1 - \frac{T_{\min} + T_{\max}}{\tau} \right) \frac{T_{\min}^2}{T_{\max}} \eta a_+ + 0$$

546 Since we can pick  $\tau$  s.t.  $\frac{T_{\min} + T_{\max}}{\tau} \leq \frac{1}{2}$ , we have

$$\mathbb{E} \left[ \frac{1}{T_i} \sum_{j=0}^{T_i} w_j^i \right] \geq \frac{T_{\min}^2}{2T_{\max}} \eta a_+ \triangleq c_0 > 0$$

□

## 547 **E Additional Figures in Section 6.1: No Bumps Between SGD and SWA** 548 **Solutions**

- 549 Asymmetric valley of ResNet-56 on CIFAR-10,  $(\bar{r}, \underline{r}, p, c) = (3.7, 3.0, 0.016, 10)$ . See Figure 20.  
 550 Asymmetric valley of ResNet-110 on CIFAR-10,  $(\bar{r}, \underline{r}, p, c) = (5.3, 3.5, 0.0050, 11)$ . See Figure 21.  
 551 Asymmetric valley of ResNet-164 on CIFAR-10,  $(\bar{r}, \underline{r}, p, c) = (2.5, 2.0, 0.027, 4.3)$ . See Figure 22.  
 552 Asymmetric valley of VGG-16 on CIFAR-10,  $(\bar{r}, \underline{r}, p, c) = (5.6, 4.0, 0.0033, 30)$ . See Figure 23.  
 553 Asymmetric valley of DenseNet-100 on CIFAR-10,  $(\bar{r}, \underline{r}, p, c) = (13.0, 8.0, 0.0029, 7.4)$ . See Figure  
 554 24  
 555 Asymmetric valley of ResNet-56 on CIFAR-100,  $(\bar{r}, \underline{r}, p, c) = (11.0, 6.0, 0.034, 15)$ . See Figure 25.  
 556 Asymmetric valley of ResNet-110 on CIFAR-100,  $(\bar{r}, \underline{r}, p, c) = (7.5, 4.5, 0.053, 6.3)$ . See Figure 26.  
 557 Asymmetric valley of ResNet-164 on CIFAR-100,  $(\bar{r}, \underline{r}, p, c) = (11.0, 6.0, 0.012, 18)$ . See Figure  
 558 27.  
 559 Asymmetric valley of VGG-16 on CIFAR-100,  $(\bar{r}, \underline{r}, p, c) = (9.0, 6.0, 0.0084, 17)$ . See Figure 28.  
 560 Asymmetric valley of ResNet-56 on SVHN,  $(\bar{r}, \underline{r}, p, c) = (5.0, 4.0, 0.018, 15)$ . See Figure 29.  
 561 Asymmetric valley of ResNet-110 on SVHN,  $(\bar{r}, \underline{r}, p, c) = (4.5, 2.5, 0.010, 11)$ . See Figure 30.  
 562 Asymmetric valley of ResNet-164 on SVHN,  $(\bar{r}, \underline{r}, p, c) = (4.5, 2.5, 0.033, 7.0)$ . See Figure 31.  
 563 Asymmetric valley of VGG-16 on SVHN,  $(\bar{r}, \underline{r}, p, c) = (4.5, 2.5, 0.0043, 43)$ . See Figure 32.  
 564 Asymmetric valley of ResNet-56 on STL-10,  $(\bar{r}, \underline{r}, p, c) = (8.0, 5.0, 0.33, 2.4)$ . See Figure 33.  
 565 Asymmetric valley of ResNet-110 on STL-10,  $(\bar{r}, \underline{r}, p, c) = (11.0, 6.0, 0.51, 3.5)$ . See Figure 34.  
 566 Asymmetric valley of ResNet-164 on STL-10,  $(\bar{r}, \underline{r}, p, c) = (12.0, 7.0, 0.092, 16)$ . See Figure 35.  
 567 Asymmetric valley of VGG-16 on STL-10,  $(\bar{r}, \underline{r}, p, c) = (5.0, 3.0, 0.11, 12)$ . See Figure 36.

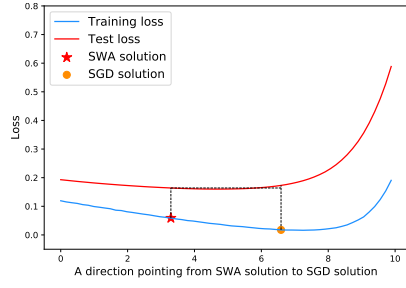


Figure 20: SWA and SGD interpolation (ResNet-56 on CIFAR-10)

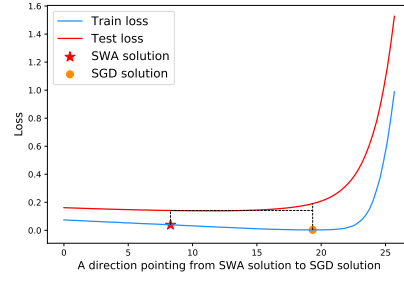


Figure 21: SWA and SGD interpolation (ResNet-110 on CIFAR-10)

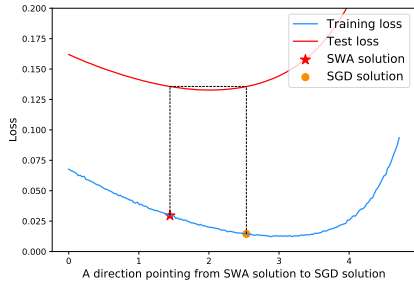


Figure 22: SWA and SGD interpolation (ResNet-164 on CIFAR-10)

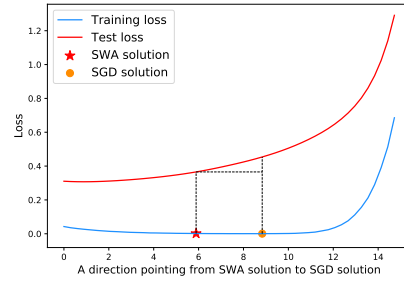


Figure 23: SWA and SGD interpolation (VGG-16 on CIFAR-10)

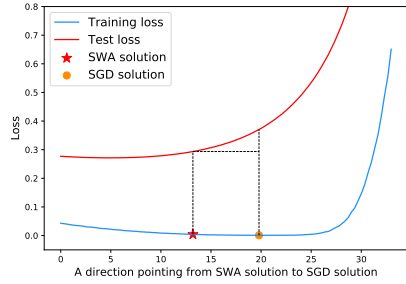


Figure 24: SWA and SGD interpolation (DenseNet-100 on CIFAR-10)

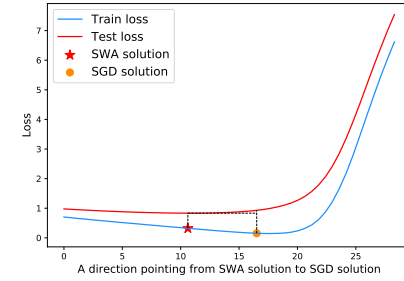


Figure 25: SWA and SGD interpolation (ResNet-56 on CIFAR-100)

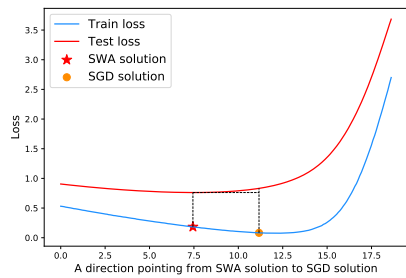


Figure 26: SWA and SGD interpolation (ResNet-110 on CIFAR-100)

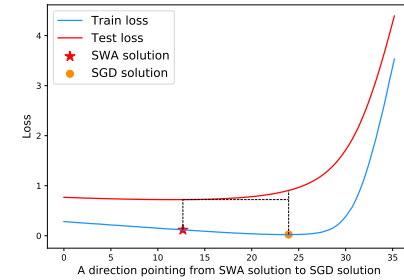


Figure 27: SWA and SGD interpolation (ResNet-164 on CIFAR-100)

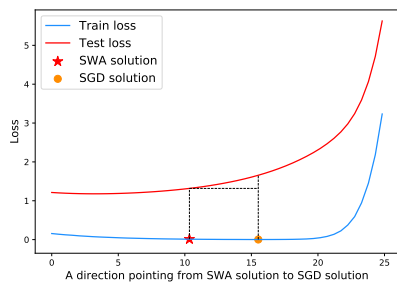


Figure 28: SWA and SGD interpolation (VGG-16 on CIFAR-100)

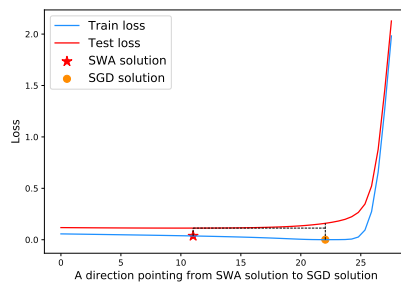


Figure 29: SWA and SGD interpolation (ResNet-56 on SVHN)

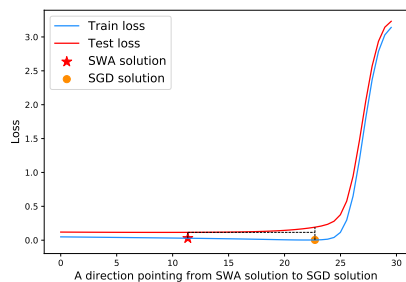


Figure 30: SWA and SGD interpolation (ResNet-110 on SVHN)

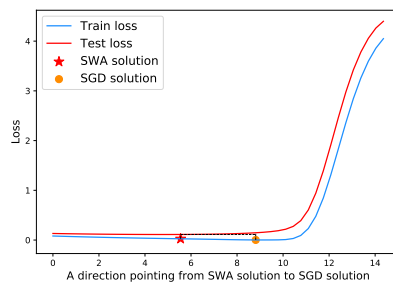


Figure 31: SWA and SGD interpolation (ResNet-164 on SVHN)

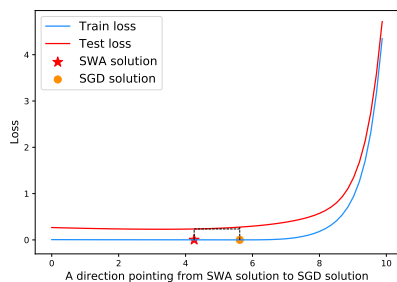


Figure 32: SWA and SGD interpolation (VGG-16 on SVHN)

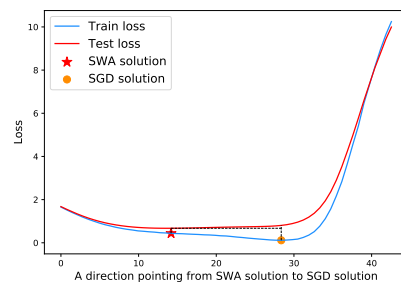


Figure 33: SWA and SGD interpolation (ResNet-56 on STL-10)

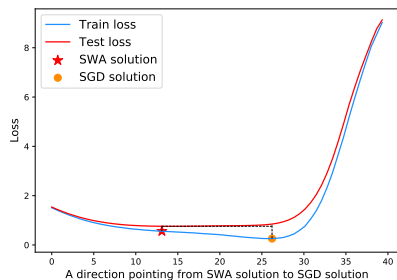


Figure 34: SWA and SGD interpolation (ResNet-110 on STL-10)

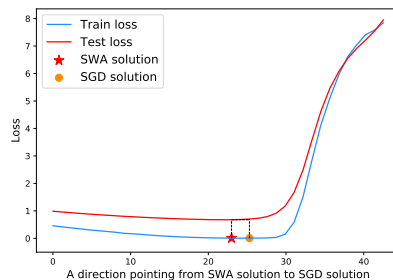


Figure 35: SWA and SGD interpolation (ResNet-164 on STL-10)





Figure 36: SWA and SGD interpolation (VGG-16 on STL-10)

## 568 **F Additional Figures in Section 6.1: SGD Averaging Generates Good Bias**

Examples for asymmetric directions of ResNet-110 on CIFAR-100 in Figure 37.

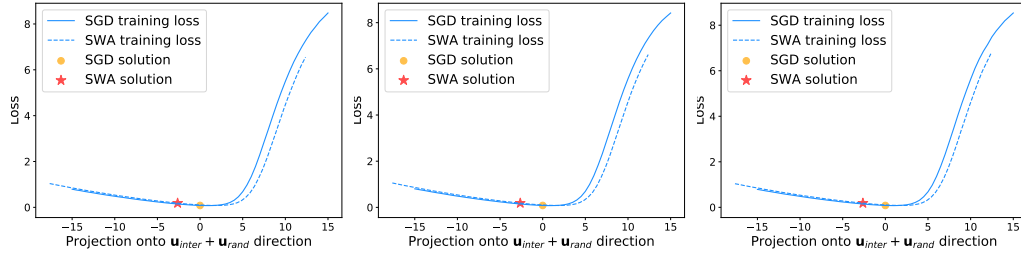


Figure 37: The average of SGD has a bias on flat side (ResNet-110 on CIFAR-100).

569

570 Examples for asymmetric directions of ResNet-164 on CIFAR-100 in Figure 38,

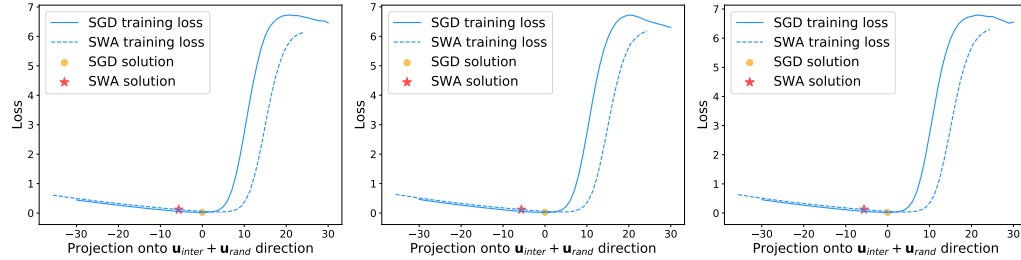


Figure 38: The average of SGD has a bias on flat side (ResNet-164 on CIFAR-100).

571 Examples for asymmetric directions of ResNet-110 on CIFAR-10 in Figure 39.

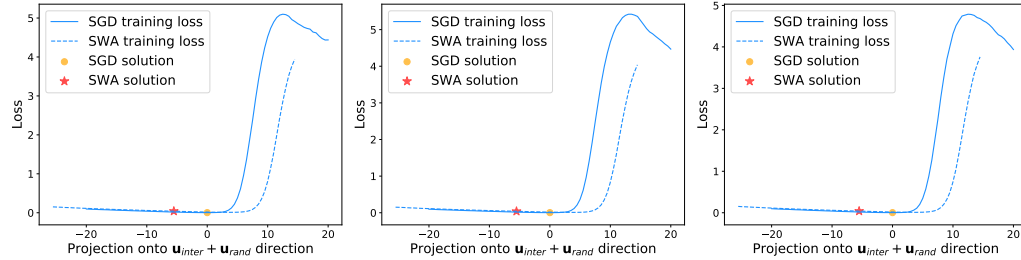


Figure 39: The average of SGD has a bias on flat side (ResNet-110 on CIFAR-10).

## G Batch size effect

Keskar et al. [32] observed that training with small batch size using SGD algorithm generalizes better than training with large batch size. They argue that it is because large batch SGD tends to converge to sharp minima, while small batch SGD generally converges to flat minima. Here we present a slightly different view that batch size has an influence on choosing sides of an asymmetric valley.

We use a PreResNet-164 trained on CIFAR-100 as an example. We first running SGD with a batch size of 128 for 200 epochs to find a solution (denoted as *Large batch solution*), and then continue the training with batch size 32 for another 80 epoch to find a nearby solution (denoted as *Small batch solution*). The reason for fine-tune is that we hope the two solutions are not far from each other, and we want to show that small batch size ensures a bias towards flat side.

From the results shown in Figure 40, it is clear that the small batch solution has worse training accuracy but better test accuracy. Meanwhile, there is no 'bump' between these solutions which suggests they are in the same basin. Therefore, small batch SGD generalizes better because it could find a better biased solution in the asymmetric valley under our training scheme, not because it finds a different wider or flatter minimum.

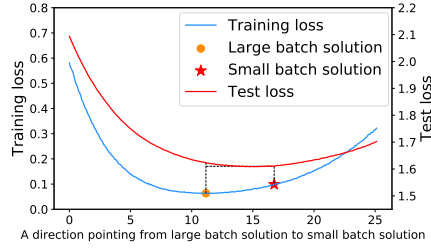


Figure 40: Large and small minibatch interpolation(batch size 128 to 32 of PreResNet-164 on CIFAR-100)

## H Batch Norm and Asymmetric Valleys

In this section, we present empirical evidences that the Batch Normalization (BN) [24] adopted by modern neural networks seems to be a major cause for asymmetric valleys.

**Directions on BN parameters are more asymmetric.** For a given SGD solution, if we take a random direction where only the BN parameters have non-zero entries, and compare it with a random direction where only the non-BN parameters have non-zero entries, we observe that those BN-related directions are usually more asymmetric. The result with ResNet-110 on CIFAR-10 is shown in Figure 41, . As we can see, the Non-BN direction is sharp on both sides, but BN direction is flat on one side, and sharp on the other side. We also conducted trials with different networks and datasets, and obtained similar results (see Figure 42, 43 and 44).

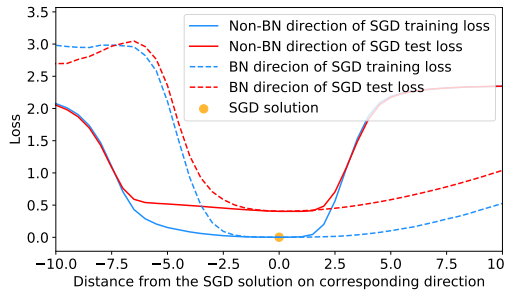


Figure 41: BN and Non-BN directions through a local minimum of ResNet-110 on CIFAR-10.

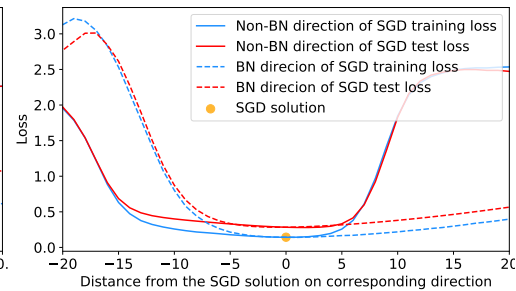


Figure 42: BN and Non-BN directions through a local minimum of of ResNet-164 on CIFAR-10.

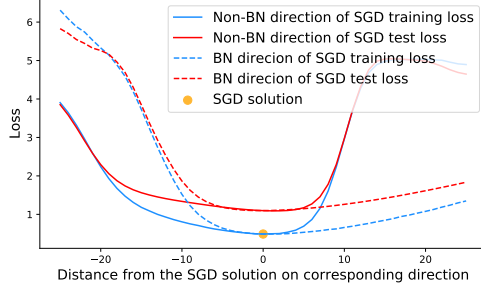


Figure 43: BN and Non-BN directions comparison of ResNet-110 on CIFAR-100

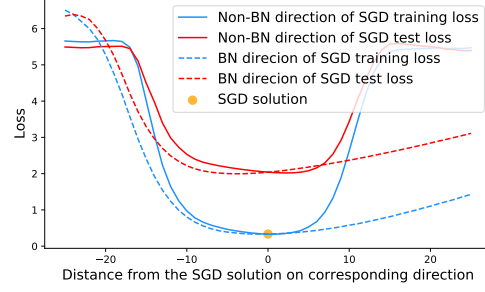


Figure 44: BN and Non-BN directions comparison of DenseNet-100 on CIFAR-100

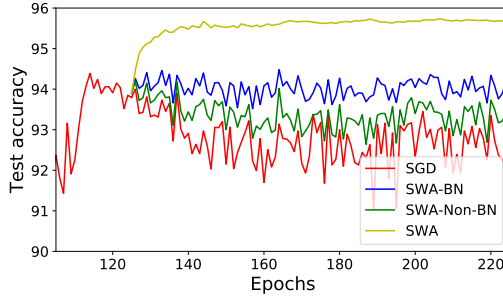


Figure 45: SGD averaging on BN parameters give better test accuracy compared with SGD averaging on non-BN parameters.

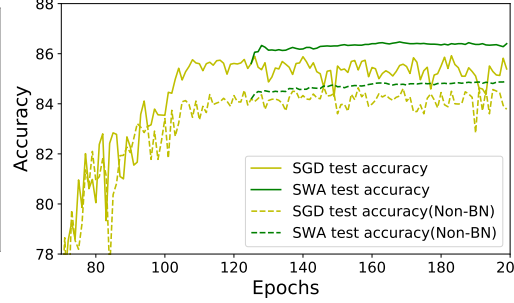


Figure 46: Test accuracy of ResNet-8 with and without BN layers, after running weight averaging (SWA).

**SGD averaging is more effective on BN parameters.** By Theorem 1 and 2, we know that SGD averaging could lead to biased solutions on asymmetric directions with better generalization. If BN indeed creates many asymmetric directions, can we improve the model performance by only averaging the weights of BN layers?

Note that BN parameters only constitute a small fraction of the total model parameters, e.g., 1.41% in a ResNet-110. In the follow experiment on ResNet-110 for CIFAR-10, we perform SGD averaging only on BN parameters, denoted as SWA-BN; and also averaging randomly selected non-BN parameters of the same amount (1.41% of the total parameters), denoted as SWA-Non-BN. The results are shown in Figure 45. It can be observed that averaging only BN parameters (blue curve) is more effective than averaging non-BN parameters (green curve), although there is still a gap comparing to averaging all the weights (yellow curve).

Moreover, we also conduct experiments with two 8-layer ResNets on CIFAR-10, one with BN layers and one without. We choose shallow networks here as deeper models without BN can not be effectively trained.

As shown in figure 46, we start weight averaging at the 126-th epoch. Although in both networks, we observe an improvement in test accuracy after averaging, it is clear that the network with BN layers have larger improvement compared with the network without BN layers. This again indicates that SGD averaging is more effective on BN parameters.