
An empirical analysis of the optimization of deep network loss surfaces

Daniel Jiwoong Im^{1,2*}

¹AIFounded Inc.

daniel.im@aifounded.com

Michael Tao

University of Toronto

mtao@dgp.toronto.edu

Kristin Branson

²Janelia Research Campus, HHMI

bransonk@janelia.hhmi.org

Abstract

The success of deep neural networks hinges on our ability to accurately and efficiently optimize high-dimensional, non-convex functions. In this paper, we empirically investigate the loss functions of state-of-the-art networks, and how commonly-used stochastic gradient descent variants optimize these loss functions. To do this, we visualize the loss function by projecting them down to low-dimensional spaces chosen based on the convergence points of different optimization algorithms. Our observations suggest that optimization algorithms encounter and choose different descent directions at many saddle points to find different final weights. Based on consistency we observe across re-runs of the same stochastic optimization algorithm, we hypothesize that each optimization algorithm makes characteristic choices at these saddle points.

1 Introduction

Deep neural networks are trained by optimizing extremely high-dimensional loss functions with respect to the networks' weights. These loss functions measure the error of the network's predictions based on these weights compared to training data. These loss functions are non-convex and are known to have many local minima. They are usually minimized using first-order gradient descent algorithms such as stochastic gradient descent (SGD) [3]. The success of deep learning critically depends on how well we can minimize these loss functions, both in terms of the quality of the convergence points and the time it takes to find them. Understanding the geometry of these loss functions and how optimization algorithms traverse them is thus of vital importance.

Several works have theoretically analyzed and characterized the shape of deep network loss functions. However, to make these analyses tractible, they have relied on simplifying assumptions. Some have characterized the critical points of deep *linear* neural networks [2, 28] and high-dimensional, random Gaussian error functions [4, 9, 11, 24]. Others [7, 8, 17] characterize the critical points of fully-connected, nonlinear deep networks, but make simplifying assumptions about the distributions and independence of various variables in the network. These works have used a variety of techniques to conclude that, given some simplifying assumptions, the loss functions of deep network have no or few bad local minima, but may instead have many saddle points. Many of these works have used empirical analyses to show that some properties of the simplified and real networks are similar.

If bad local minima are indeed rare, as suggested by theoretical analyses as well as the practical success of deep learning, then optimization algorithms do not need to take precautions to avoid them.

*Work done during an internship at Janelia Research Campus

Instead, they must only bypass saddle points and find *any* local minimum quickly. The speed of gradient-descent-based algorithms are generally measured on strictly convex, in particular, quadratic functions [5, 10, 22, 23, 25]. If different algorithms converge to the same local minimum from a common initialization, then such analyses might directly apply to the performance of these algorithms on deep networks. Conversely, if different optimization algorithms converge to different local minima that have different characteristics, then it may be necessary to evaluate the algorithms’ performance on real loss functions.

In this work, we empirically investigated the geometry of the real loss functions for state-of-the-art networks and data sets, and how commonly used optimization algorithms interact with these real loss surfaces. To our knowledge, this is the first work to jointly analyze deep network loss functions and optimization algorithms. To do this, we extended the methodology of Goodfellow et al. [13], and examined the loss function in a low-dimensional, projected space chosen to investigate properties of and the relationship between the convergence points of the different optimization algorithms. Our empirical results support the following novel conclusions:

- Different optimization algorithms find different solutions² within the projected space. This is true even when starting from the same initialization with the same mini-batch and dropout settings. Most surprisingly, this remained true when we switched from one optimization algorithm to another after the training error has nearly plateaued, suggesting that there are a plethora of saddle points even near the convergence points.
- Despite corresponding to different final points, the loss surfaces for the same algorithm from different initializations are remarkably consistent and characteristic of the optimization algorithm. The shapes of the loss functions near the final solution differ across algorithms, and we trace this back to different algorithms selecting weight vectors with a consistent norm. Switching from one optimization algorithm to another late in training results in a final point characteristic of the second optimization algorithm.
- Batch normalization is key to obtaining this consistency in the projected loss surface. Without it, we see much more variability across re-runs from different initializations.

2 Experimental setup

2.1 Network architectures and data sets

We conducted experiments on three different neural network architectures. They are all high-dimensional, deep networks and are currently used in many machine vision and learning tasks. Most importantly, their loss functions are highly non-convex.

The Network-in-Network (NIN) [21] and Visual Geometry Group (VGG) network [29] are feed-forward convolutional networks developed for image classification, and have excellent performance on the Imagenet [27] and CIFAR10 [19] data sets. Finally, we tested a two-layer fully-connected neural network (FC2).

In our experiments, we tested NIN and VGG on the CIFAR10 image classification data set. We tested FC2 on the MNIST digit recognition task [20].

Details of the network parameters and training data sets can be found in Section A.4.

2.2 Optimization methods

We analyzed the performance of five stochastic gradient-descent optimization methods commonly used for training deep neural networks: (vanilla) Stochastic Gradient Descent (SGD) [26], Stochastic Gradient Descent with Momentum (SGDM), RMSprop [32], Adadelta [33], and Adam [18]. These are all first-order gradient descent algorithms that estimate the gradients based on randomly-grouped minibatches of training examples. One of the major differences between these algorithms is the step-sizes chosen for each iteration. SGD and SGDM utilize fixed step-sizes, while RMSprop, Adadelta, and Adam use adaptive step-sizes based off of previous iterations. Details are provided in Section A.2.

²When we say solutions we specifically mean the result of running a reasonably-parameterized optimization method for a reasonable number of iterations

In addition to these five existing optimization methods, we compared to a new gradient descent method we developed based on the family of Runge-Kutta integrators [6]. In our experiments, we tested a second-order Runge-Kutta integrator in combination with SGD (RK2) and in combination with Adam (Adam&RK2). Details are provided in Section A.3).

2.3 Analysis methods

Several of our empirical analyses are based on the technique of Goodfellow et al. [13], in which properties of the loss function are examined by projecting it onto a single, carefully chosen dimension. The projection is chosen based on important weight configurations and they plot the value of the loss function along the line segment between two weight configurations. They perform two such analyses: one in which they interpolate between the initialization weights and the final learned weights, and one in which they interpolate between two sets of final weights³, each learned from different initializations. Based on their observations, they conclude that local minima do not need to be overcome by SGDM, which is in agreement with previous work suggesting there are few bad local minima.

In this work, we use a similar visualization technique, but choose different low-dimensional subspaces for the projection of the loss function. These subspaces are based on the initial weights as well as the final weights learned using the different optimization algorithms. They are chosen to answer a variety of questions about the loss function and how the different optimization algorithms interact with this loss function. In contrast, Goodfellow et al. only looked at SGDM. In addition, we explore the use of two-dimensional projections of the loss function, allowing us to better visualize the space between final parameter configurations. We do this via barycentric and bilinear interpolation for triplets and quartets of points, respectively (details in Section A.1). See Fig. 1 and Appendix Fig. 2(c) for examples.

We examine the loss surfaces around the final weight configurations learned by these variants of SGD. For all networks, we train until convergence, when fluctuations in training accuracy based on dropout are much larger than the trending increase over epochs. The total numbers of epochs of training for each experiment are shown in Table 2.

These final points are local minima in the projected space (Fig. 6(a,c)). If a weight vector is a local minimum in the intelligently-chosen projected space, it is suggestive, but *not* conclusive, that it is a local minimum in the original high-dimensional space. Hence, we will use the term solutions.

3 Experimental results

3.1 Different optimization methods find different solutions

We trained the neural networks described in Section 2.1 using each optimization method starting from the same initial weights and with the same minibatching. As shown in Fig. 2(a-b) for the VGG network, the quality of these final points were quite similar in terms of both training and test error. These results agree with previous work suggesting there are few bad solutions.

To investigate the shape of the loss function, we computed the value of the loss function for weight vectors interpolated between the initial weights, the final weights for one algorithm, and the final weights for a second algorithm for each pairing of algorithms (Fig. 2(c)) for the VGG network. For every pair of optimization algorithms, we observe that, within the projected space, the final points are

³Batch normalization parameters beta and gamma (from [16]) are interpolated, and batch mean and standard deviation were calculated using training data.

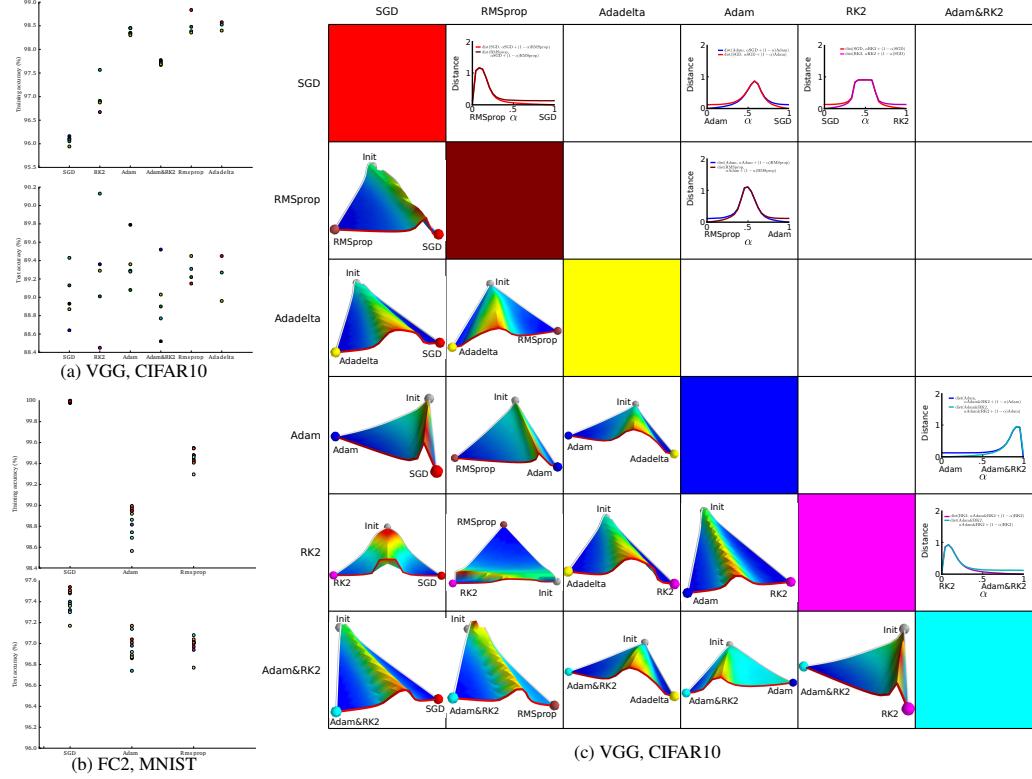
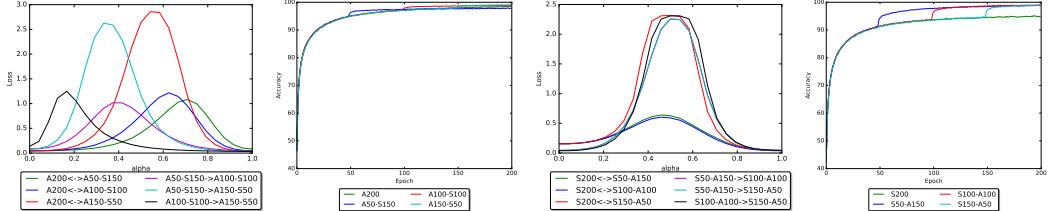


Figure 2: (a-b) Training and test accuracy for each of the optimization methods for (a) the VGG network on CIFAR10 and (b) the FC2 network on MNIST. Colors correspond to different initializations. (c) Visualization of the loss surface near and between the final points found by different optimization methods for the VGG network on CIFAR10. Each box corresponds to a pair of optimization methods. In the lower triangle, we plot the projection of the loss surface at weight vectors between the initial weight and the learned weights found by the two optimization methods. Color as well as height of the surface indicate the loss function value. In the upper triangle, we plot the functional difference between the network corresponding to the learned weights for the first algorithm and networks corresponding to weights linearly interpolated between the first and second algorithm's learned weights.

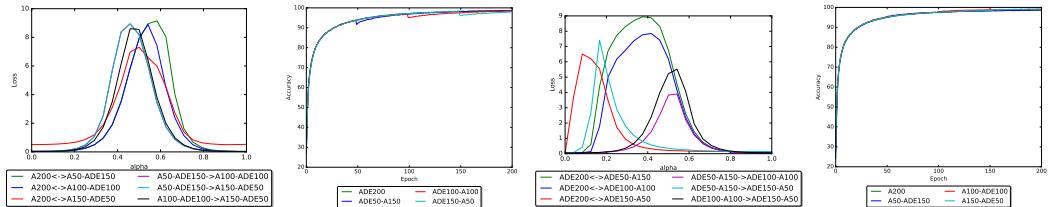
always separated by a high-loss region. With the caveat that we have only visualized a 2-dimensional projection of the loss surface, this suggests that each optimization algorithm found a solution within the basin of a different local minimum (in sliced low-dimensional space), despite starting at the same initialization. We observed similar phenomena for NIN on CIFAR10 and FC2 on MNIST (Fig. 5). We investigated the space between other triplets and quadruplets of weight vectors (Fig. 1), and even in these projections of the loss function, we still saw that the local minima returned by different algorithms are separated by high loss weight parameters.

Our observation that different optimization algorithms find very different solutions suggests that they choose different descent directions at saddle points encountered during optimization [9]. We next investigated whether these saddle points only occur in the early, *transient* [30] phase of optimization in which the loss is decreasing rapidly, or also occur in the later *minimization* [30] phase in which the loss decreases slowly. We investigated the effects of switching from one type of optimization method to another 25%, 50%, and 75% of the way through training. We emphasize that we are not switching methods to improve performance, but rather to investigate the shape of the loss function during the minimization phase of optimization.

In Figure 3, for a given pair of optimization algorithms, we plot the loss surface between final points found by switching algorithms at different points during training. For all pairs of algorithms and all switching points, the final points correspond to different critical points in the projected space. For example, different critical points in the projected space (left column) are found using Adam versus switching from Adam to SGD after 150 epochs of training, when the train-



(a) NIN: Switching from Adam (A, learning rate $\eta = .001$) to SGD (S, $\eta = .01$).
(b) NIN: Switching from SGD (S, $\eta = .1$) to Adam (A, $\eta = .0001$).



(c) VGG: Switching from ADAM (A, $\eta = .001$) to (d) VGG: Switching from Adadelta (ADE) to ADAM Adadelta (ADE, $\eta = .001$).

Figure 3: Effects of switching from one optimization method to another at epochs 50 and 100 and 150, or never. A50-S150 corresponds to training with Adam for 50 epochs, then switching to SGD and training for 150 epochs. A200 corresponds to training with Adam for the entire 200 epochs. *Left:* Loss at weight vectors interpolated between two final points. A200<->A50-S150 is the loss at weight interpolations between the final point corresponding to A200 ($\alpha = 0$) and that corresponding to A50-S150 ($\alpha = 1$). *Right:* Accuracy on training data as a function of training epoch. We see qualitatively similar results for FC2 on MNIST (Supp. Fig. 14).

ing accuracy has nearly plateaued (right column). This suggests that saddle points are encountered late in the optimization at which different descent directions are chosen by the different algorithms. This disagrees with the common lore that the local minimum has effectively been chosen in the transient phase, and instead suggests that which solution is found is still in flux late in optimization [30]. It appears that this switch from one critical point to another happens almost immediately after the optimization method switches, with the training accuracy jumping to the characteristic accuracy for the given method within a few epochs (Figure 3, right).

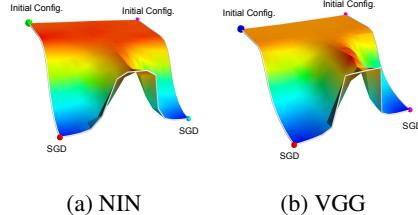


Figure 4: Visualization of the loss surface at weights interpolated between two initial configurations and the final weight vectors learned using SGD from these initializations, for the VGG and NIN networks on CIFAR10.

Our experiments thus far have suggested that deep network loss function have many similarly good solutions in terms of training and validation error (see Figure 2). However, deep networks are overparameterized. For example, if we switch all corresponding weights for a pair of nodes in our network, we will obtain effectively the same network, with both the original and permuted networks outputting the same prediction for a given input. To investigate whether the final points found by different algorithms corresponded to different parameterizations of equivalent networks, or to truly different networks, we compared the outputs of the networks on each example in a validation data set:

$$\text{dist}(\theta_1, \theta_2) = \sqrt{\frac{1}{N_{test}} \sum_{i=1}^{N_{test}} \|F(x_i, \theta_1) - F(x_i, \theta_2)\|^2},$$

where θ_1 and θ_2 are the weights learned by two different optimization algorithms, x_i is the input for a validation example, and $F(x, \theta)$ is the output of the network for weights θ on input x .

We found that, for all pairs of algorithms, the average distance between the outputs of the networks (Equation 3.1) was approximately 0.16, corresponding to a label disagreement of about 8% (upper

triangle of Fig. 2(c)). Given the generalization error of these networks (approximately 11%, Fig. 2(b)), the maximum disagreement we could see was 22%. Thus, these networks disagreed on a large fraction of these test examples – over one third, and the final points found by different algorithms appear to correspond to effectively different networks, not trivial reparameterizations of the same one.

3.2 Different optimization algorithms find different types of solutions

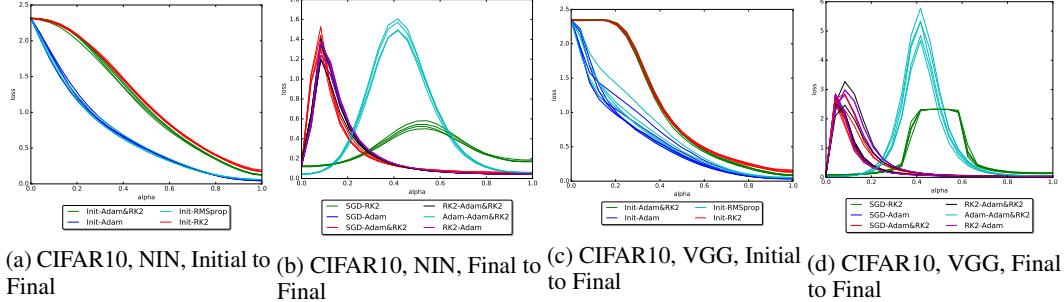


Figure 5: Loss function visualizations for multiple re-runs of each algorithm. Each re-run corresponds to a different initialization. We see that the loss function near the final point for a given algorithm has a characteristic geometry.

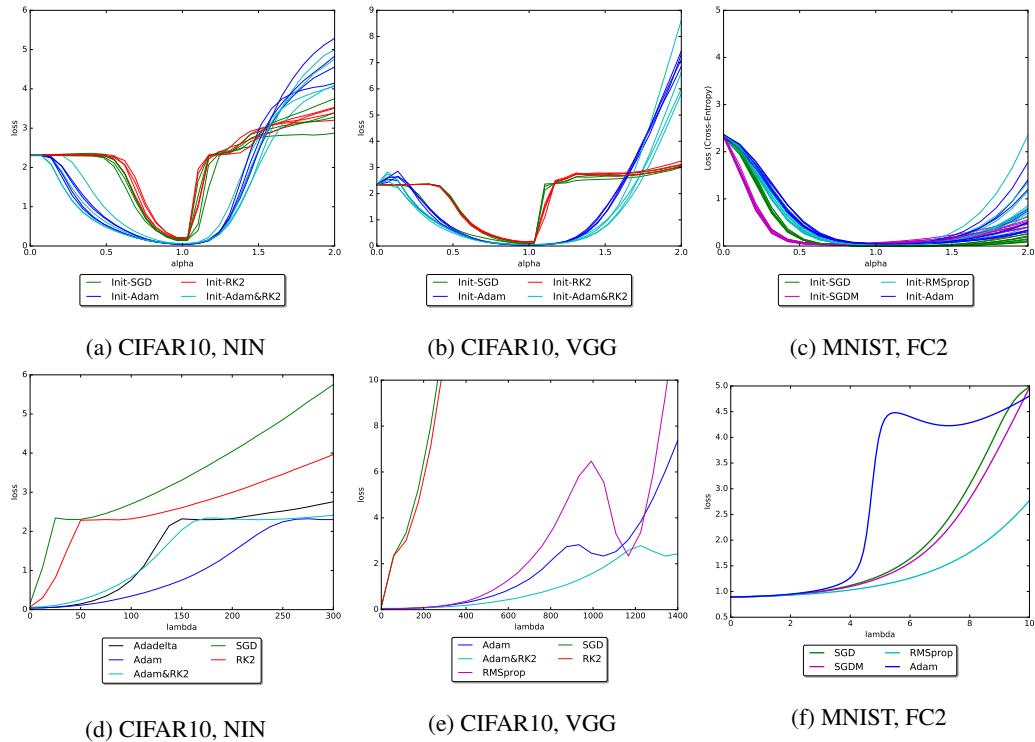


Figure 6: Comparison of the size of the projected space local minima basins for different algorithms for NIN (a-b) and VGG (c-d) on CIFAR10 and for FC2 on MNIST (e-f). In (a,c,e), the units of the x-axis (α) are relative to the scale of the final weight vector of the algorithm, with $\alpha=0$ corresponding to the initial weight vector, $\alpha=1$ corresponding to the final weight vector and $\alpha>1$ corresponding to points interpolated beyond the final weight vector. In (b,d,f), the units of the x-axis (λ) are the same for each algorithm – the units of the high-dimensional weight-vector space (Eq. 1). $\lambda=0$ indicates the final weight vector found by the algorithm. In (b,d), the kink near $\text{loss}=2.4$ corresponds to the initial weight vector.

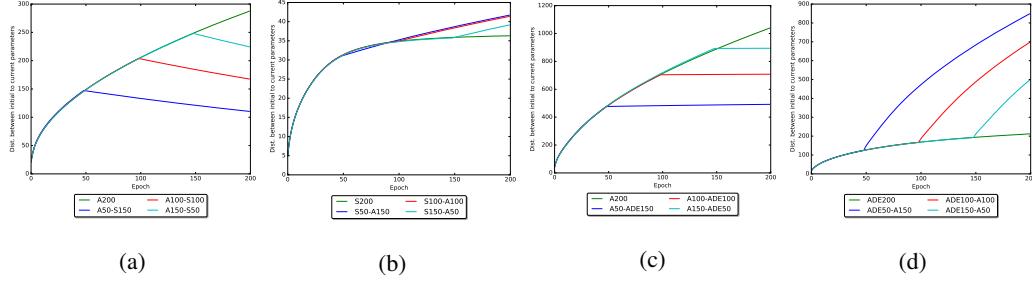


Figure 7: Distance from initial to current weights during training in which we switch from one optimization method to another at epochs 50 and 100 and 150, or never. Different plots correspond to different networks and pairs of algorithms, following Fig. 3. Different optimization algorithms have a characteristic distance. After switching, this distance quickly approaches the new characteristic distance. (a) NIN: Switching from Adam (A, learning rate $\eta = .001$) to SGD (S, $\eta = .01$). (b) NIN: Switching from SGD (S, $\eta = .1$) to Adam (A, $\eta = .0001$). (c) VGG: Switching from Adam (A, $\eta = .001$) to Adadelta (ADE). (d) VGG: Switching from Adadelta (ADE) to Adam (A, $\eta = .001$).

Next, we investigated whether the solutions found by the different optimization algorithms had distinguishing properties. To do this, we trained the networks with each optimization algorithm from different initializations. We then compared differences between runs of the same algorithm from different initializations to differences between different algorithms.

As shown in Figure 2(a), in terms of training accuracy, we see some stereotypy for the final points found by different algorithms, with SGD having the lowest training accuracy and ADAM, RMSprop, and Adadelta having the highest training accuracy. However, the generalization accuracy of these different final points on validation data was not different between algorithms (Figure 2(b)). We also did not see a relationship between the weight initialization and the validation accuracy. This supports the conclusions of previous work that most local minima are equally good. As found in Goodfellow et al. [13], these final points correspond to different local minima in the projected space 4.

We visualized the loss surface around each of the final points for the multiple runs from different initializations. To do this, we plotted the value of the loss function between the initial and final weights for each algorithm for several runs of the algorithm with different initializations (Figure 5(a,c)). In addition, we plotted the value of the loss function between the final weights for selected pairs of algorithms for each run (Figure 5(b,d)). We see that the surfaces look strikingly similar for different runs of the same algorithm, but characteristically different for different algorithms. Thus, we found evidence that the final weights found by different algorithms are of similar quality in terms of validation error, they are qualitatively different. This suggests that there is some stereotypy in the descent direction chosen at saddle points by each algorithm.

In particular, we see in Figure 6(a,c) that the size of the basins in the projected spaces around the local minima found by Adam and Adam&RK2 are larger than those found by SGD and RK2 for NIN and VGG on CIFAR10, i.e. that the training loss is small for a wider range of α values. The x-axis of Fig. 6(a,c,e), α , is a multiplier of the weight vector. If the norm of the weight vector found by one algorithm is larger than that found by another, then a change of $\Delta\alpha$ for this curve would correspond to a larger absolute change in the weight vector, $\Delta\alpha(\theta_1 - \theta_0)$, where θ_0 is the initial weight vector and θ_1 is the result found by a given optimization algorithm. In Figure 6(b,d,f) we again plot the basin, but normalize for the norm of the weight vector difference, and show the loss as a function of the absolute distance in parameter space:

$$\theta(\lambda) = \theta_1 + \lambda \frac{\theta_0 - \theta_1}{\|\theta_0 - \theta_1\|} \quad (1)$$

Even after normalization, we see the same result for NIN and VGG on CIFAR10: the sizes of the basins in the projected spaces around the local minima are bigger for Adam and Adam&RK2 (Fig. 6(b,d)). This trend does not repeat for FC2 on MNIST, however, with Adam corresponding to the smallest basin (Fig. 6(e,f)).

To try to understand why some algorithms find final points within bigger basins in the projected space, we looked for characteristic properties of the training process and the final weights found by different algorithms. We found that, for NIN on CIFAR10, Adam-based algorithms result in final weights that

are much farther from the initial weight vectors (Fig. 8). When we switch from one optimization algorithm to another, we see strong effects on the distance from the initial to current weight vectors (Fig. 7). Interestingly, when we switch from Adam to SGD (Fig. 7(a)), the weight vectors travel *back toward* the initialization. This implies that, for this pair of algorithms, it is not just that the step size is smaller for SGD, but that SGD has a characteristic range of distances from the initial weights. Because the norm of the initial weight vector is so small, this may instead correspond to a characteristic weight vector norm. Fig. 8(c) shows that the distance traveled and weight vector norm curves are nearly identical. Saddle points in deep *linear* networks, which have a scale ambiguity, have been shown to correspond to weight vector norm scale [28]. Batch normalization in deep nonlinear networks introduces such a scale ambiguity between the incoming weights and the normalization [1].

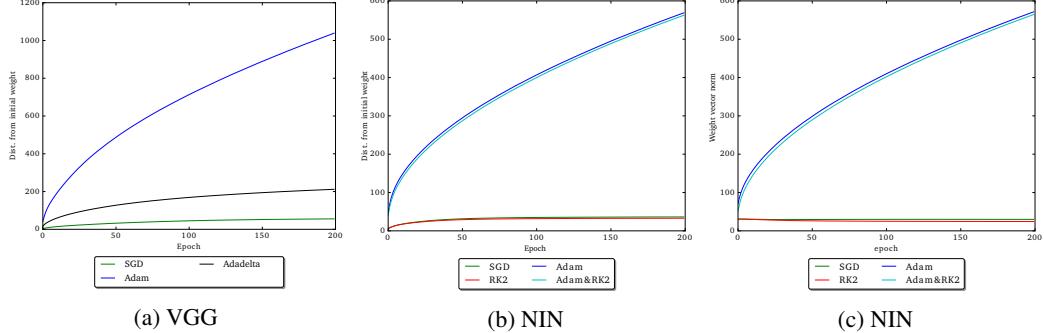


Figure 8: (a-b) Distance from initial to current weights during training for different optimization algorithms. Left: VGG network, right: NIN on CIFAR10. (c) Norm of weight vectors during training for different optimization algorithms for NIN on CIFAR10. (c) and (b) are nearly identical.

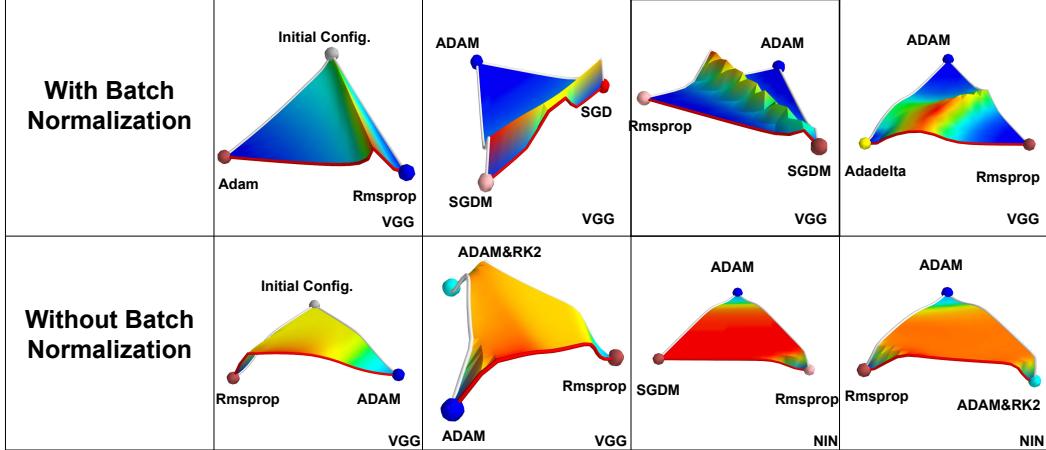


Figure 9: Visualization of the loss surface with and without batch-normalization.

3.3 Effects of batch-normalization

To understand how batch normalization affects the types of solutions found, we performed a set of experiments comparing loss surfaces near solutions found with and without batch normalization for each of the optimization methods. We visualized the surface near these solutions by interpolating between the initial weights and the final weights as well as between pairs of final weights found with different algorithms (Fig. 9).

We observed clear qualitative differences between optimization with (Figure 5) and without (Figure 10) batch normalization. We see that, without batch normalization, the quality of the solutions found by a given algorithm is much more dependent on the initialization. In addition, the surfaces between different solutions are more complex in appearance: with batch normalization we see sharp unimodal jumps in performance but without batch normalization we obtain wide bumpy shapes that aren't necessarily unimodal.

The neural networks are typically initialized with very small parameter values [12, 15]. Instead, we trained NIN with exotic intializations such as initial parameters drawn from $\mathcal{N}(-10.0, 0.01)$ or $\mathcal{N}(-1.0, 1.0)$ and observe the loss surface behaviours (We used same initialization distribution as [31].) The details of results are discussed in Appendix A.7.

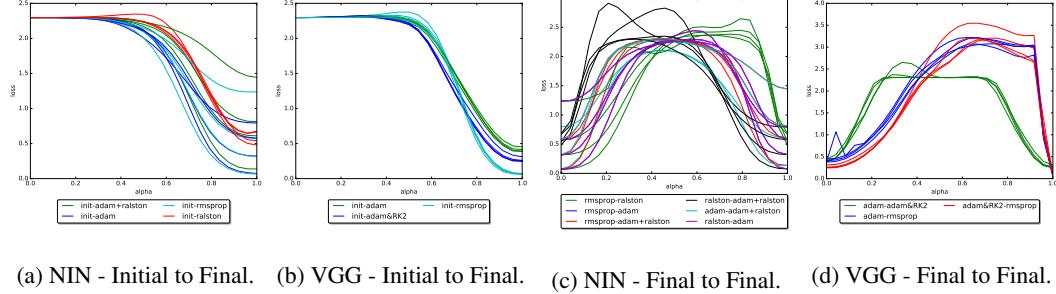


Figure 10: Loss function visualizations for multiple re-runs of each algorithm without batch normalization. Each re-run corresponds to a different initialization. Without batch normalization, re-runs are much less consistent.

4 Conclusions

In this work, we have presented and analyzed a series of experiments with neural networks. The intent of these experiments was to elaborate our understanding on the nature of neural network loss functions, specifically with respect to the geometry of these loss functions and the sensitivities of different optimization methods for minimizing them.

We found that different optimization techniques find different final weights, even when they were seeded by the same initialization point. Furthermore, we saw that when batch normalization causes the validation error for these final networks are similar. In the projected spaces we chose for visualizing the loss function, any pair of final weights was always separated by a high loss bump. This is suggestive that the different algorithms find different solutions, and thus diverge at saddle points encountered during optimization. We also found that switching from one optimization algorithm to another late in training, when the training loss is decreasing slowly, still resulted in different local minima within the projected space. This is suggestive that there are saddle points even at this late stage in training, at which different optimization algorithms diverge. Our observations suggest that the loss surfaces of deep networks have a large number of saddle points that are encountered during optimization.

While our observations are consistent with those of Goodfellow et al. [13], they offer a more complex picture. Goodfellow et al. showed that the projected loss between initial and final weights monotonically decreases. They concluded that the loss function shape was simple, and that optimization could proceed without overcoming local minima or saddle points. Instead, we see that, during optimization, many saddle points are encountered even late in optimization, at which points different optimization algorithms choose different descent directions. Our observations are thus more consistent with previous theoretical analyses on simplified networks, which have suggested that deep network loss surfaces have many saddle points and local minima are similar in terms of training and validation error (Sec. 1). Most surprisingly, we see a characteristic shape to the loss function projected around the final weights found by different algorithms, which suggests that different algorithms make consistently different choices at saddle points, which, to our knowledge, is a novel hypothesis.

While we found evidence that these different final weights are not simply different parameterizations of the same network, we do not yet know if the qualitative differences in final weights found by different algorithms matter in a practical sense, i.e. in some quality of the predictions of the network, and we hope this is a fruitful direction of future research. Because of the many saddle points encountered during optimization, and the qualitative differences in the shape of the loss function around the final points found by different algorithms, we conclude that measuring the efficiency of different optimization algorithms on the same strictly convex loss function is insufficient for estimating the efficiency of these algorithms on deep network loss surfaces.

References

- [1] Jimmy Lei Ba, Jamie Ryan Kiros, and Geoffrey E Hinton. Layer normalization. *arXiv preprint arXiv:1607.06450*, 2016.
- [2] Pierre Baldi and K. Hornik. Neural networks and principal component analysis: Learning from examples without local minima. *Neural Networks*, 2:53—58, 1989.
- [3] Leon Bottou. Stochastic gradient learning in neural networks. In *Proceedings of Nuero-Nimes*, 1991.
- [4] Alan J. Bray and David S. Dean. The statistics of critical points of gaussian fields on large-dimensional spaces. In *Physics Review Letter*, 2007.
- [5] C. G Broyden. The convergence of a class of double-rank minimization algorithms 1. general considerations. *Journal of Applied Mathematics*, 6:76—90, 1970.
- [6] John C. Butcher. Coefficients for the study of runge-kutta integration processes. *Society for Industrial and Applied Mathematics*, 3:185—201, 1963.
- [7] Anna Choromanska, Mikael Henaf, Michael Mathieu, Gerard Ben Arous, and Yann LeCun. The loss surfaces of multilayer networks. *arXiv preprint arXiv:1406.2572*, 2015.
- [8] Anna Choromanska, Yann LeCun, and Gerard Ben Arous. Open problem: The landscape of the loss surfaces of multilayer networks. *JMLR: Workshop and Conference Proceedings*, 40:1—5, 2015.
- [9] 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. *arXiv preprint arXiv:1406.2572*, 2014.
- [10] Murat A. Erdogdu and Andrea Montanari. Convergence rates of sub-sampled newton methods. In *Proceedings of the Neural Information Processing Systems (NIPS)*, 2015.
- [11] Yan V. Fyodorov and Ian Williams. Replica symmetry breaking condition exposed by random matrix calculation of landscape complexity. *Journal of Statistical Physics*, 129:1081–1161, 2007.
- [12] Xavier Glorot and Y. Bengio. Understanding the difficulty of training deep feedforward neural networks. In *International conference on artificial intelligence and statistics*, 2010.
- [13] Ian J. Goodfellow, Oriol Vinyals, and Andrew M. Saxe. Qualitatively characterizing neural network optimization problems. In *Proceedings of the International Conference on Learning Representations (ICLR)*, 2015.
- [14] Ernst Hairer, Nørsett Syvert P., and Gerhard Wanner. *Solving Ordinary Differential Equations I – Nonstiff*. Springer, 1987.
- [15] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Delving deep into rectifiers: Surpassing human-level performance on imagenet classification. In *arXiv preprint arXiv:1502.01852*, 2015.
- [16] Sergey Ioffe and Christian Szegedy. Batch normalization: Accelerating deep network training by reducing internal covariate shift. In *arXiv preprint arXiv:1502.03167*, 2015.
- [17] Kenji Kawaguchi. Deep learning without poor local minima. *arXiv preprint arXiv:1605.07110*, 2016.
- [18] Diederik Kingma and Jimmy Ba. Adam: A method for stochastic optimization. In *Proceedings of the International Conference on Learning Representations (ICLR)*, 2014.
- [19] Alex Krizhevsky. Learning multiple layers of features from tiny images. In *MSc thesis, Univesity of Toronto*, 2009.
- [20] Yann LeCun, Corinna Cortes, and Christopher JC Burges. The mnist database of handwritten digits, 1998.

- [21] Min Lin, Qiang Chen, and Shuicheng Yan. Network in network. In *Proceedings of the International Conference on Learning Representations (ICLR)*, 2014.
- [22] James Martens. Deep learning via hessian-free optimization. In *Proceedings of the International Conference of Machine Learning (ICML)*, 2010.
- [23] Yurii Nesterov. A method of solving a convex programming problem with convergence rate $o(1/\sqrt{k})$. *Soviet Mathematics Doklady*, 27:372—376, 1983.
- [24] Giorgio Parisi. Probabilistic line searches for stochastic optimization. *arXiv preprint arXiv:0706.0094*, 2016.
- [25] B.T Polyak. Some methods of speeding up the convergence of iteration methods. *USSR Computational Mathematics and Mathematical Physics*, 4(5):1–17, 1964.
- [26] Herbert Robbins and Sutton Monro. A stochastic approximation method. *Annals of Mathematical Statistics*, 22(3):400—407, 1951.
- [27] O. Russakovsky, H. Deng, J. Su, J. Krause, S. Satheesh, S. Ma, Z. Huang, A. Karpathy, A. Khosla, M. Bernstein, A. C. Berg, and L. Fei-Fei. Imagenet large scale visual recognition challenge. In *arXiv preprint arXiv:1409.0575*, 2014.
- [28] Andrew M Saxe, James L McClelland, and Surya Ganguli. Exact solutions to the nonlinear dynamics of learning in deep linear neural networks. In *In International Conference on Learning Representations.*, 2014.
- [29] Karen Simonyan and Andrew Zisserman. Very deep convolutional networks for large-scale image recognition. In *Proceedings of the International Conference on Learning Representations (ICLR)*, 2015.
- [30] Ilya Sutskever, James Martens, George Dahl, and Geoffery Hinton. On the importance of momentum and initialization in deep learning. In *Proceedings of the International Conference of Machine Learning (ICML)*, 2013.
- [31] Grzegorz Swirszcz, Wojciech Marian Czarnecki, and Razvan Pascanu. Local minima in training of deep networks. In *International conference on artificial intelligence and statistics*, 2016.
- [32] Tijmen Tieleman and Geoffery Hinton. Rmsprop gradient optimization. In *Neural Networks for Machine Learning slide: http://www.cs.toronto.edu/~tijmen/csc321/slides/lecture_slides_lec6.pdf*, 2012.
- [33] Matthew D. Zeiler, Graham W. Taylor, and Rob Fergus. Adaptive deconvolutional networks for mid and high level feature learning. In *International Conference on Computer Visio*, 2011.

A Supplementary Materials

A.1 3D Visualization

[13] introduced the idea of visualizing 1D subspace of the loss surface between the parameters. Here, we propose to visualize loss surface in 3D space through interpolating over three and four vertices.

Linear Interpolation Given two parameters θ_1 and θ_2 ,

$$\theta_i = \alpha\theta_1 + (1 - \alpha)\theta_2, \quad \forall \alpha \in [0, 1]. \quad (2)$$

Bilinear Interpolation Given four parameters θ_0 , θ_1 , θ_2 , and θ_3 ,

$$\phi_i = \alpha\theta_1 + (1 - \alpha)\theta_2 \quad (3)$$

$$\varphi_i = \alpha\theta_3 + (1 - \alpha)\theta_4 \quad (4)$$

$$\theta_j = \beta\phi_i + (1 - \beta)\varphi_i \quad (5)$$

for all $\alpha \in [0, 1]$ and $\beta \in [0, 1]$.

Barycentric Interpolation Given four parameters θ_0 , θ_1 , and θ_2 , let $d_1 = \theta_1 - \theta_0$ and $d_2 = \theta_2 - \theta_0$. Then, the formulation of the interpolation is

$$\phi_i = \alpha d_1 + \theta_0 \quad (6)$$

$$\varphi_i = \alpha d_2 + \theta_0 \quad (7)$$

$$\theta_j = \beta\phi_i + (1 - \beta)\varphi_i \quad (8)$$

for all $\alpha \in [0, 1]$ and $\beta \in [0, 1]$.

A.2 Optimization Methods

A.2.1 Stochastic Gradient Descent

In many deep learning applications both the number of parameters and quantity of input data points can be quite large. This makes the full evaluation of $U(\theta)$ be prohibitively expensive. A standard technique for alleviating computational loadis to apply an stochastic approximation to the gradient [26]. More precisely, one approximates U by a subset of n data points, denoted by $\{\sigma_j\}_{j=1}^N$ at each timestep:

$$U^n(\theta) = \frac{1}{n} \sum_{j=1}^n \ell(\theta, \mathbf{x}_{\sigma_j}) \simeq \frac{1}{N} \sum_{i=1}^N \ell(\theta, \mathbf{x}_i) = U(\theta) \quad (9)$$

Of course this approximation also carries over to the gradient, which is of vital importance to optimization techniques:

$$\nabla U^n(\theta) = \frac{1}{n} \sum_{j=1}^n \nabla \ell(\theta, \mathbf{x}_{\sigma_j}) \simeq \nabla U(\theta) \quad (10)$$

This method is what is commonly called *Stochastic Gradient Descent* or *SGD*. So long as the data is distributed nicely the approximation error of U^n should be sufficiently small such that not only will SGD still behave like normal GD , but it's wall clock time for to converge should be significantly lower as well.

Usually one uses the stochastic gradient rather than the true gradient, but the inherent noisiness must be kept in mind. In what follows we will always mean the stochastic gradient.

A.2.2 Momentum

In order to alleviate both noise in the input data as well as noise from stochasticity used in computing quantities one often maintains history of previous evaluations. In order to only require one extra variable one usually stores variables of the form

$$\mathbb{E}[F]_t = \alpha F_t + \beta \mathbb{E}[F]_{t-1}. \quad (11)$$

where F_t is some value changing over time and $\mathbb{E}[F]_t$ is the averaged quantity.

An easy scheme to apply this method to is to compute a rolling weighted average of gradients such as

$$\mathbb{E}[\mathbf{g}]_t = (1 - \alpha)\mathbf{g}_t + \alpha\mathbb{E}[\mathbf{g}]_{t-1}$$

but there will be other uses in the future.

A.2.3 Pertinent Methods

With the aforementioned tools there are a variety of methods that can be constructed. We choose to view these algorithms as implementations of Explicit Euler on a variety of different vector fields to remove the ambiguity between η and \mathbf{g}_t . We therefore can define a method by the vector field \mathcal{X}_t that explicit Euler is applied to with a single η that is never changed.

SGD with Momentum (SGDM) By simply applying momentum to \mathbf{g}_t one obtains this stabilized stochastic version of gradient descent:

$$\mathcal{X}_t = -\mathbb{E}[\mathbf{g}]_t. \quad (12)$$

This is the most fundamental method that is used in practice and the basis for everything that follows.

Adagrad Adagrad rescales \mathcal{X}_t by summing up the squares of all previous gradients in a coefficient-wise fashion:

$$\mathcal{X}_t = -\frac{\mathbf{g}_t}{\sqrt{\sum_{i=1}^t \mathbf{g}_i^2 + \epsilon}}. \quad (13)$$

Here ϵ is simply set to some small positive value to prevent division-by-zero. In the future we will neglect this term in denominators because it is always necessary.

The concept is to accentuate variations in \mathbf{g}_t , but because the denominator is monotonically nondecreasing over time this method is doomed to retard its own progress over time. The denominator can also be seen as a form of momentum where α and β are both set to 1.

Rmsprop A simple generalization of ADAGrad is to simply allow for α and β to be changed from 1. In particular one usually chooses a β less than 1, and presumably $\alpha = 1 - \beta$. Thus one arrives at a method where the effects of the distance history are diminished:

$$\mathcal{X}_t = -\frac{\mathbf{g}_t}{\sqrt{\mathbb{E}[\mathbf{g}^2]_t}}. \quad (14)$$

Adadelta Adadelta adds another term to RMSprop in order to guarantee that the magnitude of \mathcal{X} is balanced with \mathbf{g}_t [33]. More precisely it maintains

$$\frac{\mathcal{X}_t}{\sqrt{\mathbb{E}[\mathcal{X}_t^2]}} = -\frac{\mathbf{g}_t}{\sqrt{\mathbb{E}[\mathbf{g}_t^2]}} \quad (15)$$

which results in the following vector field:

$$\mathcal{X}_t = -\frac{\sqrt{\mathbb{E}[\mathcal{X}_t^2]}}{\sqrt{\mathbb{E}[\mathbf{g}_t^2]}} \mathbf{g}_t. \quad (16)$$

and η is set to 1.

ADAM By applying momentum to both \mathbf{g}_t and \mathbf{g}_t^2 one arrives at what is called ADAM. This is often considered a combination of SGDM + RMSprop,

$$\mathcal{X}_t = c_t \frac{\mathbb{E}[\mathbf{g}]_t}{\sqrt{\mathbb{E}[\mathbf{g}^2]_t}}. \quad (17)$$

$c_t = \frac{\sqrt{1-\beta_2^t}}{1-\beta_1^t}$ is the *initialization bias correction term* with $\beta_1, \beta_2 \in [0, 1]$ being the β parameters used in momentum for \mathbf{g} and \mathbf{g}^2 respectively. Initialization bias is caused by the history of the momentum variable being initially set to zero.

A.3 Runge Kutta

Runge-Kutta methods [6] are a broad class of numerical integrators categorized by their truncation error. Because the ordinary differential equations Runge-Kutta methods solve generalize gradient descent, our augmentation is quite straightforward. Although our method applies to all explicit Runge-Kutta methods we will only describe second order methods for simplicity.

The general form of second-order explicit Runge-Kutta on a time-independent vector field is

$$\theta_{t+1} = \theta_t + (a_1 k_1 + a_2 k_2)h \quad (18)$$

$$k_1 = \mathcal{X}(\theta_t) \quad (19)$$

$$k_2 = \mathcal{X}(\theta_t + q_1 h k_1) \quad (20)$$

where a_1, a_2 , and q_1 are parameters that define a given Runge-Kutta method. Table 1 refers to the parameters used for the different Runge-Kutta variants we use in our experiments.

Table 1: The coefficients of various second order Runge-Kutta methods [14]

| Method Name | a_1 | a_2 | q_1 |
|-------------|---------------|---------------|---------------|
| Midpoint | 0 | 1 | $\frac{1}{2}$ |
| Heun | $\frac{1}{2}$ | $\frac{1}{2}$ | 1 |
| Ralston | $\frac{1}{3}$ | $\frac{2}{3}$ | $\frac{3}{4}$ |

A.3.1 Augmenting Optimization with Runge Kutta

For a given timestep, explicit integrators can be seen as a morphism over vector fields $\mathcal{X} \rightarrow \bar{\mathcal{X}}^h$. For a gradient $\mathbf{g}_t = \nabla_{\theta} U$ we can solve a modified RK2 gradient $\bar{\mathbf{g}}_t$ in the following fashion:

$$\theta_{t+1} = \theta_t + \bar{\mathbf{g}}_t h = \text{Advect}_{\mathbf{g}}^{rk2}(\theta, h) \quad (21)$$

rearranged with respect to $\bar{\mathbf{g}}_t$

$$\bar{\mathbf{g}}_t = \frac{\text{Advect}_{\mathbf{g}}^{rk2}(\theta, h) - \theta_t}{h} \quad (22)$$

$$= \frac{\theta_t + (a_1 k_1 + a_2 k_2)h - \theta_t}{h} \quad (23)$$

$$= (a_1 k_1 + a_2 k_2). \quad (24)$$

If we simply substitute the gradient \mathbf{g}_t with $\bar{\mathbf{g}}_t$ one obtains an RK2-augmented optimization technique.

A.4 Network architecture and data set details

We used the VGG and NIN implementations from <https://github.com/szagoruyko/cifar.torch.git>. As well, we used two layer full connected neural network (FC) with 50 hidden units and batch normalization on each layer. Table 2 presents experiment status for each figure.

The batch size was set to 128 and the number of epochs was set to 100 for FC and 200 for NIN and VGG. The FC parameters were initialized using Xavier initialization, $\mathcal{U}\left(\frac{-\sqrt{6}}{N_{in}+N_{out}}, \frac{\sqrt{6}}{N_{in}+N_{out}}\right)$ where N_{in} and N_{out} are input and output dimension of each layer. NIN and VGG were initialized from $\mathcal{N}(0, 0.05)$. NIN had 50% dropouts and VGG had 40% dropouts on convolutional layer and 50% dropouts for fully connected layers during the training. The learning rate was chosen from the discrete range between $[0.2, 0.1, 0.05, 0.01]$ for SGD and $[0.002, 0.001, 0.0005, 0.0001]$ for adaptive learning methods⁴. (Table 3 presents the learning rate of the optimizers that was used for different experiments.) We doubled the learning rates when we ran our augmented versions with Runge-Kutta because they required two stochastic gradient computations per epoch. We used batch-normalization and dropout to regularize our networks. For all parts of the optimization algorithms that require randomness, we use the same random choices: we used the same initializations, dropout masks, and sequences of mini-batch samples across algorithms. We'll make this explicit in our Experimental setup section. All experiments were run on a 6-core Intel(R) Xeon(R) CPU @ 2.40GHz with a TITAN X.

⁴Learning rate was fixed for different optimization methods. When we switched one method to another, we switched the learning rate accordingly. We also experimented the learning rates that were applied after switching. For example, we tried with 0.001, 0.0001, 0.00001 for ADAM. We'll make this more explicit in our Experimental results section.

Table 2: The Neural Network architecture status for different experiments.

| Figure Ref. | Network Type | Batch-Norm. | data | # of Epoch |
|-------------------|------------------|-------------|---------|-------------------------|
| Fig. 13 & 14 | 2-layer NN | X | MNIST | 100 |
| Fig. 4, 1, 5 & 11 | NIN & VGG | ✓ | CIFAR10 | 200 |
| Fig. 2, 3, 6 & 16 | NIN | ✓ | CIFAR10 | 200 |
| Table 9 | NIN & VGG | Both | CIFAR10 | 200 w. BN & 600 w.o. BN |
| Fig. 10 | NIN, VGG, & LSTM | X | CIFAR10 | 600 |
| Fig. 15 | NIN | Both | CIFAR10 | 200 w. BN & 600 w.o. BN |
| Fig. 18, 19 & 17 | VGG | ✓ | CIFAR10 | 200 |

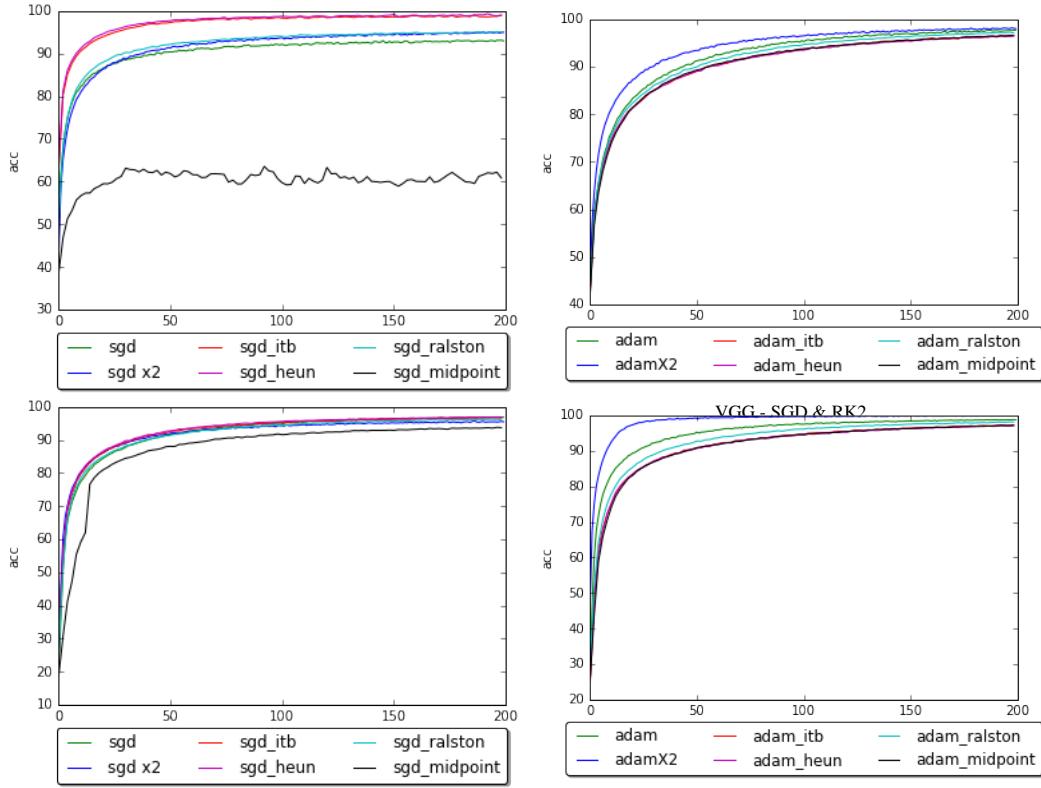
Table 3: Learning rate of the optimizers that is used for various experiments.

| Learning rate (η) | Optimizer | Figures |
|--------------------------|---------------|---|
| - | Adadelta | Fig. 2, 6, 19, 17 |
| 0.1 | SGD, SGDM | Fig. 4, 2, 5, 6, 3(b), 3(c), 10, 11, 13, 14, 6, 16, 18(a), ?? |
| 0.05 | SGD, SGDM | Fig. 16, 18(b), ?? |
| 0.01 | SGD, SGDM | Fig. 3(a), 3(d), 15, 16, 18(c), ??, |
| 0.001 | ADAM, RMSProp | Fig. 3(a), 3(c), 10, 14, 16, 18(a), 17 |
| 0.0005 | ADAM, RMSProp | Fig. 16, 18(b) |
| 0.0001 | ADAM, RMSProp | Fig. 4, 2, 5, 6, 3(b), 3(d), 11, 13, 6, 16, 18(c) |

A.5 Experiments with Runge-Kutta integrator

The results in Figure 11 illustrates that, with the exception of the Midpoint method, stochastic Runge-Kutta methods outperform SGD. “SGD x2” is the stochastic gradient descent with twice of the learning rate of “SGD”. From the figure, we observe that the Runge-Kutta methods perform even better with half the number of gradient computed by SGD. The reason is because SGD has the accumulated truncated error of $O(h)$ while second-order Runge-Kutta methods have the accumulated truncated error of $O(h^2)$.

Unfortunately, ADAM outperforms ADAM+RK2 methods. We speculate that this is because the way how ADAM’s renormalization of input gradients in conjunction with momentum eliminates the value added by using our RK-based descent directions.



NIN - ADAM & ADAM+ VGG - ADAM & ADAM+RK2
Figure 11: Training accuracy curve

A.6 The results on Two layer Fully Connected Neural Network on MNIST dataset

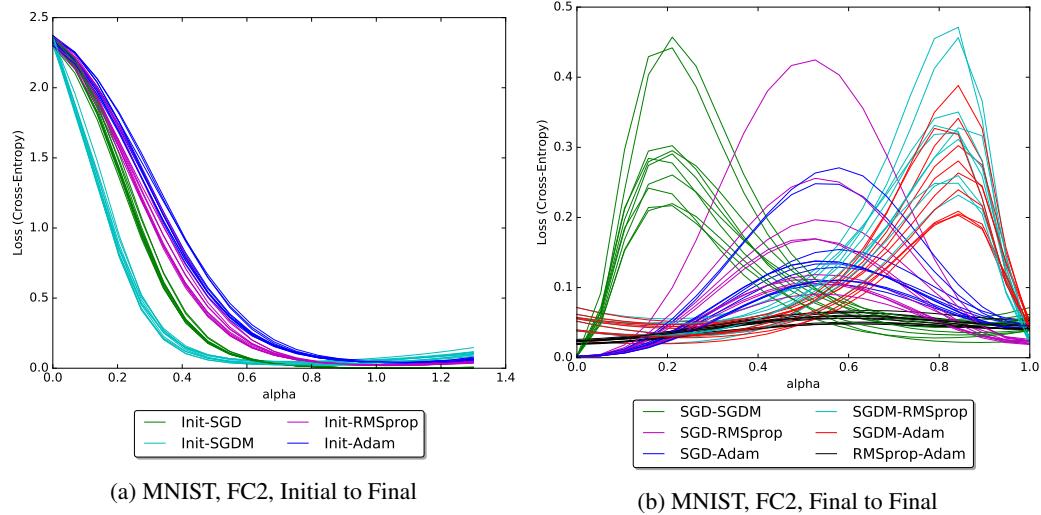


Figure 12: Loss function visualizations for multiple re-runs of each algorithm. Each re-run corresponds to a different initialization. We see that the loss function near the final point for a given algorithm has a characteristic geometry.

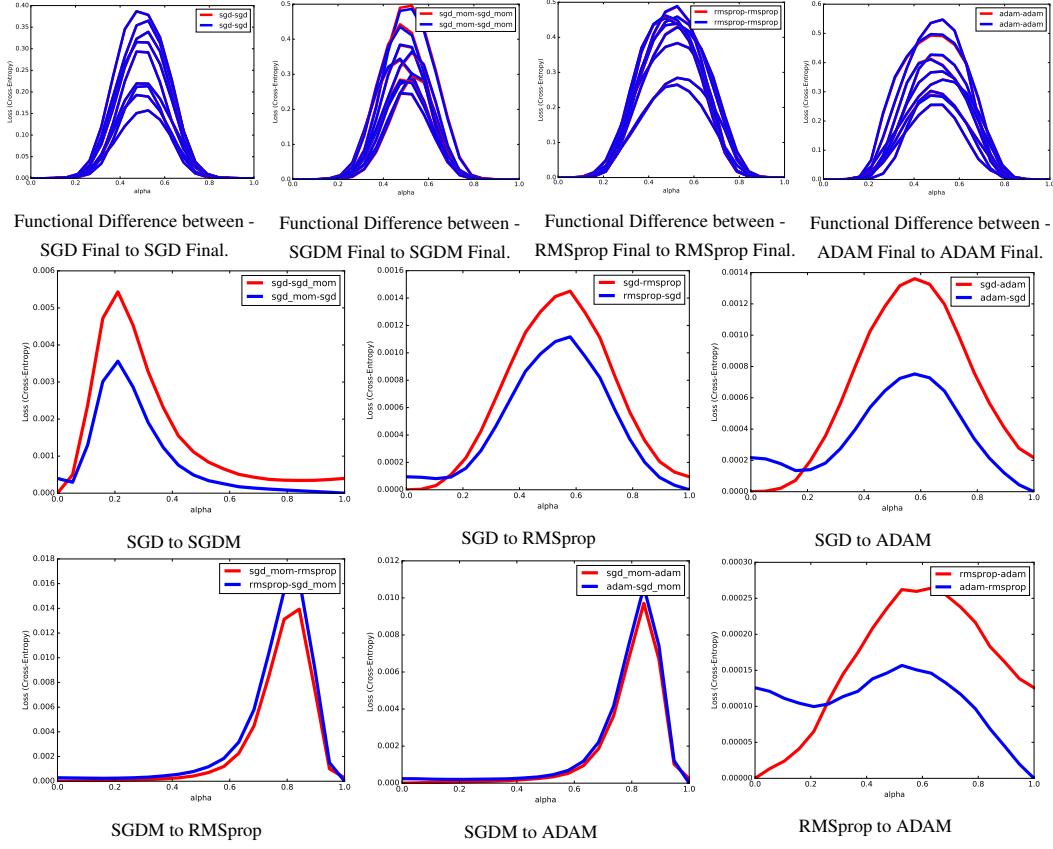


Figure 13: The projection of the loss surface at weight vectors between the initial weight and the learned weights found by the two optimization methods. Color as well as height of the surface indicate the loss function value. In the upper triangle, we plot the functional difference between the network corresponding to the learned weights for the first algorithm and networks corresponding to weights linearly interpolated between the first and second algorithm's learned weights.

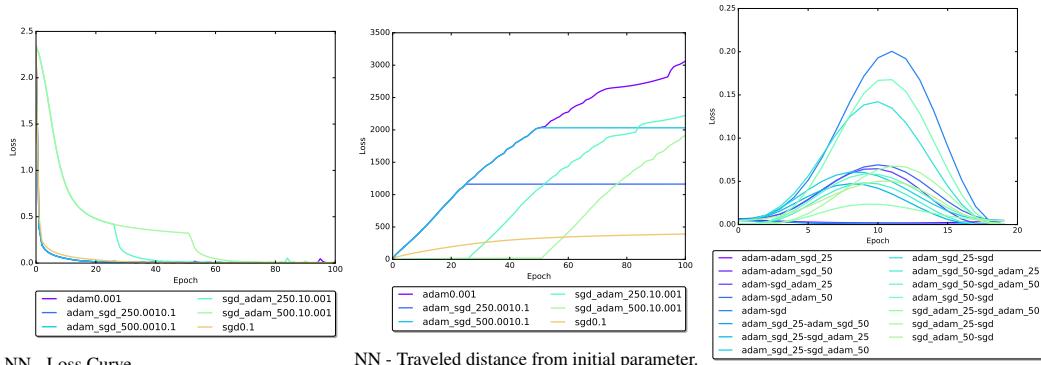


Figure 14: Switch algorithms

A.7 Effects of Batch-Normalization and Extreme Initializations

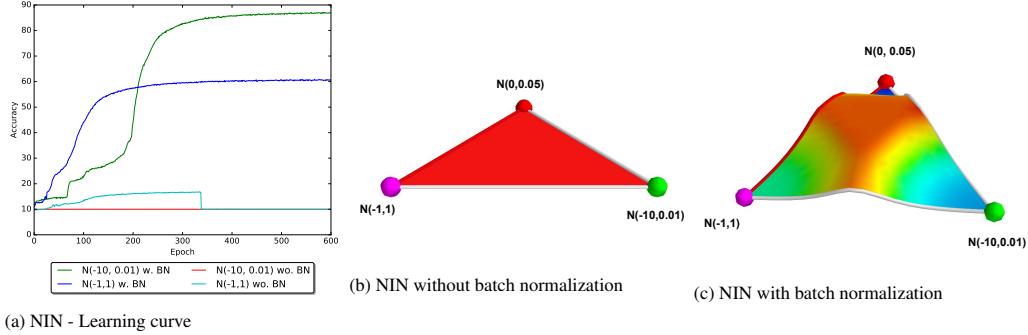


Figure 15: NIN trained from different initializations.

The neural networks are typically initialized with very small parameter values [12, 15]. Instead, we trained NIN with exotic intializations such as initial parameters drawn from $\mathcal{N}(-10.0, 0.01)$ or $\mathcal{N}(-1.0, 1.0)$ and observe the loss surface behaviours. The results are shown in Figure 15. We can see that NIN without BN does not train at all with any of these initializations. Swirszc z 2016 *et al.* mentioned that bad performance of neural networks trained with these initializations are due to finding a bad local minima. However, we see that loss surface region around these initializations are plateau ⁵ rather than a bad local minima as shown in Figure 15(b). On the other hand, NIN with BN does train slowly over time but finds a local minima. This implies that BN redeems the ill-posed loss surface (plateau region). Nevertheless, the local minima it found was not good as when the parameters were initialized with small values. However, it is not totally clear whether this is due to difficulty of training or due to falling in a bad local minima.

A.8 Switching optimization methods

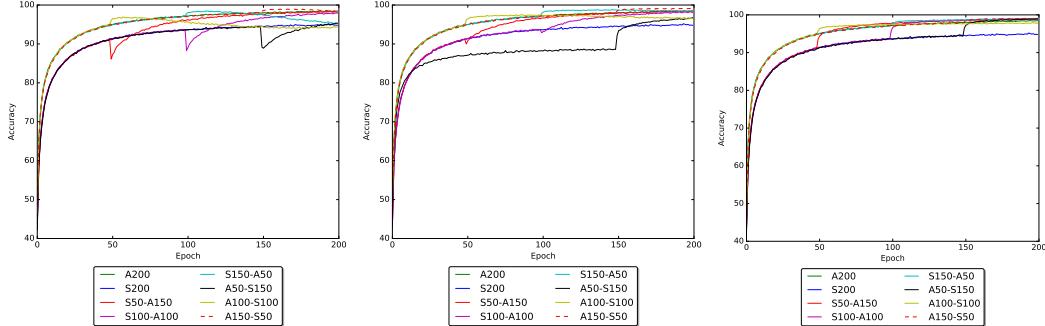


Figure 16: NIN - Learning curve when switching methods from SGD to ADAM and visa versa at epoch 50 and 100. Learning rate switched from SGD (ADAM) to ADAM (SGD) at (left) 0.001 (0.1) to 0.1 (0.001), (middle) 0.001 (0.1) to 0.05 (0.001), and (right) 0.001 (0.1) to 0.01 (0.001).

⁵We used same initializations as [31] but we trained different neural networks with SGD on a different dataset. We used NIN and CIFAR10 and Swirszc z 2016 *et al.* used smaller neural network and MNIST.

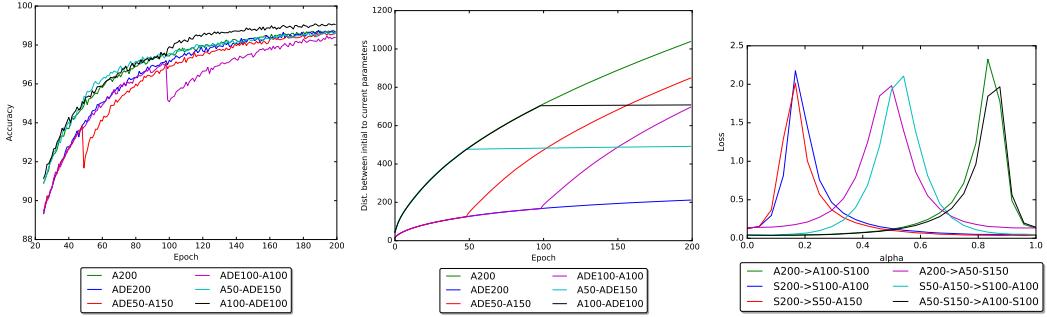
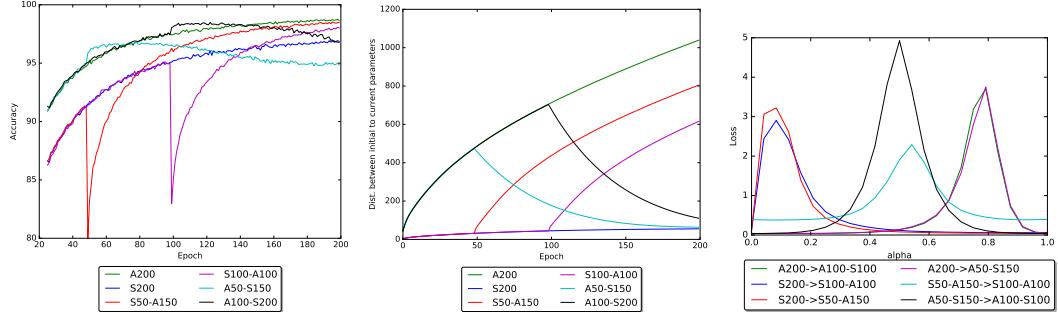
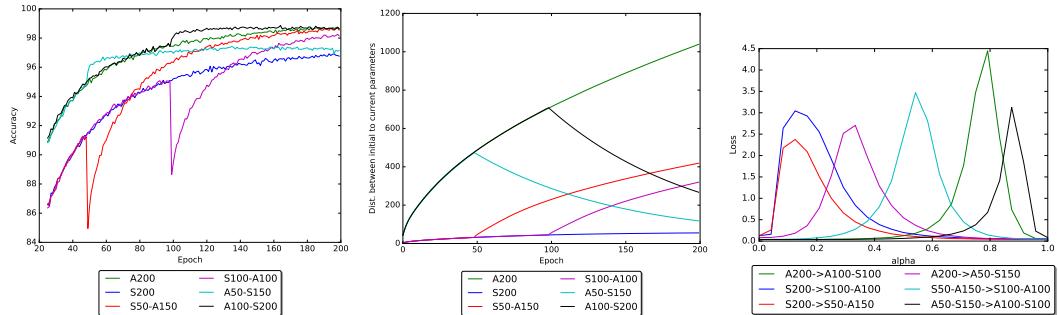


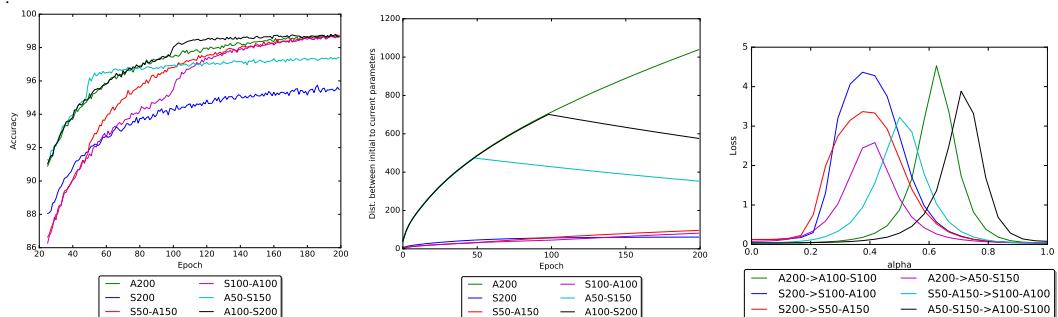
Figure 17: VGG - Switching methods from ADAM to Adadelta and Adadelta to ADAM at epoch 50 and 100. Zoomed in version (Left). Distance between initial weights to weights at each epoch (Middle). The interpolation between different convergence parameters (Right). Each figure shows the results of switching methods at different learning rate. We label the switch of methods in terms of ratio. For instance, ADE50-A50 as trained with ADAM in the first 100 epoch and switched to Adadelta for the rest of the epoch.



(a) The learning rates is set to 0.001 and 0.05 for ADAM and SGD in the beginning, and then switched it to 0.05, 0.001 for SGD and ADAM.

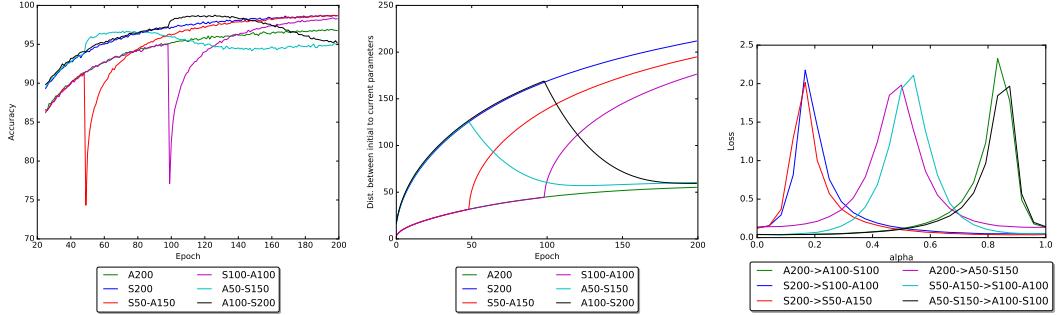


(b) The learning rates is set to 0.001 and 0.05 for ADAM and SGD in the beginning, and then switched it to 0.05, 0.0005 for SGD and ADAM

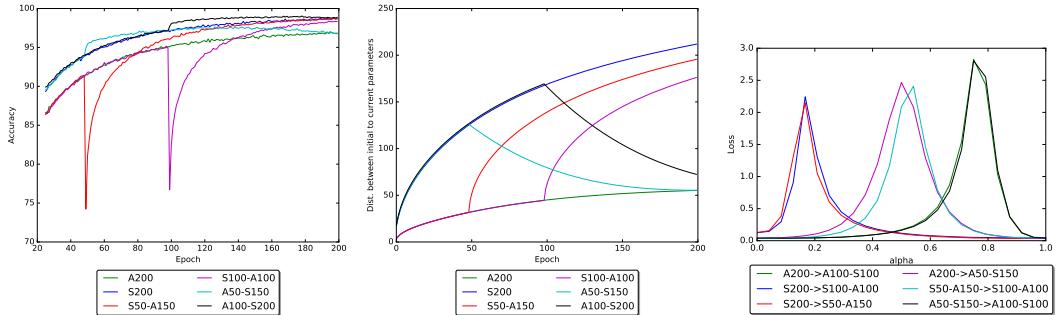


(c) The learning rates is set to 0.001 and 0.05 for ADAM and SGD in the beginning, and then switched it to 0.01, 0.0001 for SGD and ADAM.

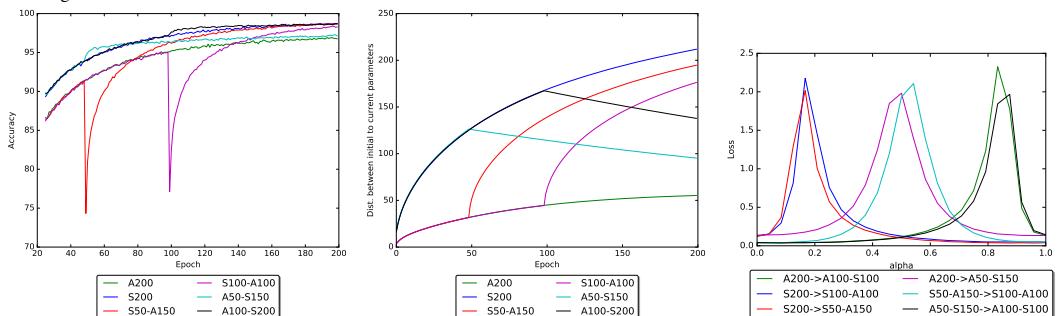
Figure 18: VGG - Switching methods from SGD to ADAM and ADAM to SGD at epoch 50 and 100. Zoomed in version (Left). Distance between initial weights to weights at each epoch (Middle). The interpolation between different convergence parameters (Right). Each figure shows the results of switching methods at different learning rate. We label the switch of methods in terms of ratio. For instance, S100-A100 as trained with SGD in the first 100 epoch and swithced to ADAM for the rest of the epoch.



Learning rate is not required for Adadelta. Learning rate is set to 0.05 for SGD in the beginning, and then switched it to 0.1.



Learning rate is set to 0.05 for SGD.



Learning rate is set to 0.05 for SGD in the beginning, and then switched it to 0.01.

Figure 19: VGG - Switching methods from SGD to Adadelta and Adadelta to SGD at epoch 50 and 100. Zoomed in version (Left). Distance between initial weights to weights at each epoch (Middle). The interpolation between different convergence parameters (Right). Each figure shows the results of switching methods at different learning rate. We label the switch of methods in terms of ratio. For instance, S50-A50 as trained with SGD in the first 100 epoch and switched to Adadelta for the rest of the epoch.