
CS260 Project Report: About SGD Convergence Analysis of NN with Non-linear Activation

Song Jiang 605222804*

Department of Computer Science
University of California, Los Angeles
Los Angeles, CA 90095
songjiang@cs.ucla.edu

Yewen Wang 905229899*

Department of Computer Science
University of California, Los Angeles
Los Angeles, CA 90095
wyw10804@cs.ucla.edu

Zhiping Xiao 604775684*

Department of Computer Science
University of California, Los Angeles
Los Angeles, CA 90095
patricia.xiao@cs.ucla.edu

Abstract

1 Neural networks are increasingly popular during the past few years, with promising
2 performance on various tasks in diverse fields. Recently, SGD based method are
3 widely explored and became a standard training method for Neural Networks.
4 In this work, we tried to analyze the convergence of SGD on 1-hidden-layer
5 feedforward networks with Non-linear activation and "identity mapping" structure
6 under the constrain that the input obeys normal distribution. We also conduct
7 related experiments to give an intuitive support of our theorem.
8 Also, we made multiple attempts including exploring SGD's convergence on NN
9 with multiple layers or with special structures or with other activation functions.
10 Discussions about our attempts and corresponding problems we met will be pro-
11 vided in this report.
12 **Note:** Our experiment code will be submitted as supplementary materials on CCLE.
13 Also, there are still some unclear parts, and hopefully we can fix those parts later.
14 Most parts of our detailed proof could be provided(handwriting version) if needed.

15 1 Introduction

16 Neural networks are increasingly popular during the past few years, with promising performance on
17 various tasks in many fields. Stochastic gradient descent is applied in countless experiments, resulting
18 in satisfying outcome. However, with different network architectures, there are various landscapes
19 and some of them contains bad local minima and saddle points. It still remains not that clear when
20 can SGD guarantee a convergence to global minima and there still lacks solid complete theoretical
21 guarantees that SGD can have good performance when finding the desired weights for such neural
22 networks.

23 To bridge this gap, in this project, we conducted a comprehensive literature review about theoretical
24 guarantee for SGD convergence analysis of one-hidden-layer neural networks, summarized its related
25 theorems and proof methods, Also, followed the work of [13], we tried to prove the convergence of
26 stochastic gradient descent on feed-forward one-hidden-layer neural network with "identity mapping"
27 structure and Non-linear activation function under the constraint that the network's input satisfies
28 normal distribution. What is more, we took several other attempts including deepening the neural

29 network, changing the structure of neural networks(to Resnet[9, 8] and DenseNet[10]), or focusing
30 on other non-convex problems. Also, we conducted several auxiliary experiments to provide an
31 intuitive support of our theorem.

32 The rest of this report is developed as follows: Section 2 provides a thoroughly literature review about
33 SGD convergence analysis on NNs and NN with identity mapping structure. Section 3 illustrates
34 how we formulate our problem and gives our main theorem. Section 4 is an overview of proofs in
35 which we discuss the general proof roadmap and some key difficulties for the proof. Section 5 shows
36 some experiment results which could be an intuitive support for our theorems. Section 6 concludes
37 other attempts we made and discusses possible reasons why we failed in those attempts. Section 7
38 performs as a summary of our work.

39 2 Literature Review

40 Generally, the related works can be divided in to three types, the first type is about SGD convergence
41 analysis for one-hidden-layer neural networks, the second type is the work expanding the network to
42 multiple layers and analyze those deep neural networks, the last type is focusing on neural networks
43 with an advanced structure – identity mapping. In this literature review section, we tried to summarize
44 and elaborate these three types of related work.

45 2.1 Stochastic Gradient Descent on One-hidden-layer Neural Network

46 Recently, stochastic gradient descent has been widely leveraged in seeking the optimal parameters
47 of neural network. While proving its convergence to the global minima, problem might occurs
48 because even for the simplest setting where the neural network has only one hidden layer and is a
49 feed forward network, SGD will probably get stuck at local minima or saddle points. To deal with
50 such unsatisfactory situations, various solutions are brought by previous works.

51 Some works may simplifies the model by ignoring the activation functions and turn the network into
52 a linear deep neural network. The work of Kenji in 2016 is a representative one of this type[11], it
53 first proves that for deep linear neural networks with any depth and any weights and squared loss
54 function, the loss is always non-convex and every local minima would be a global minima, then, this
55 work proves for deep nonlinear neural networks, we can always conduct a reduction to linear model
56 under the independence assumption, thus it can still have those properties they proved before for
57 linear models. With these properties, SGD for such kind of problem can achieve convergence to
58 global minima.

59 Some works may rely on some unrealistic assumptions with which local minima would have some nice
60 properties[6, 12, 14, 3]. For example, focusing on the simple non-trivial ReLU neural networks, [14]
61 proves that spurious local minima is guaranteed when k , the size of weight per layer, is fixed within
62 certain range. From this point of view, it provides an answer on why neural network are successfully
63 trained even if the associated optimization problem has non-convexity, and what assumptions could
64 mitigate this problem.

65 Some works design advanced initialization method to guarantee SGD’s convergence[16]. In the work
66 of [16] section 5.2, the author applied a Tensor Method to obtain a proper initialization point, with
67 which SGD can converge to global minima.

68 Some works guarantees SGD’s convergence on over-parameterization networks, that is, based on
69 the assumption that the network is wide enough. For example, [5] proved that even with random
70 initialization, first-order methods can achieve zero training loss even if the objective function is
71 neither smooth or convex as long as the hidden layer is wide enough.

72 Some works rely on specific network structures to ensure SGD can converge to global minima[13].
73 This is also a work that is most closely to our project. In this work, it proves that with a special
74 structure, “Identity Mapping”, stochastic gradient descent will converge to the global minimum of
75 two-layer neural network in polynomial number of steps. The “Identity Mapping” structure makes the
76 network asymmetric and thus guarantees a unique global minimum.

77 2.2 Convergence Analysis of Deep Neural Network

78 Despite some success on theoretical analysis for 2-layer neural network, recently, people starts to
79 tackle deeper neural networks and attempts to explore more about NN with multiple layers on the
80 theoretical side. Though the convergence of SGD for deep neural networks still remains an open
81 problem, there are already some existing works focus on such problem[2, 1, 17, 4].

82 [2] focuses on the convergence speed on Recurrent Neural Networks(RNN) and provides its related
83 theoretical understanding. It proved that as long as the number of neurons is sufficiently large, even
84 with random initialization, SGD can minimize regression loss in a linear convergence rate.

85 Another interesting work[4] gives theoretical analysis of gradient descent(instead of SGD) and proves
86 rgar GD can reach the global optima for a deep over-parameterized neural network with residual
87 connections. This work relies on such NN architecture to ensure the global optimality of the gradient
88 descent algorithm.

89 Also, there's one recent work[1] that studies how SGD can get to global minima of deep neural
90 networks(DNN). This work is also based on the assumption that the network is sufficiently wide,
91 and can be applied to fully-connected neural networks, convolutional neural networks (CNN), and
92 residual neural networks (ResNet).

93 2.3 Network with Identity Mapping

94 To improve the performance of neural networks, in some previous works, modifications of structures
95 are made to construct stronger networks. Identity mapping, which is presented in [9], is such kind of
96 modifications that gains widely application. It allows that signal could be directly propagated from
97 one neural unit to any others, in both forward and backward passe, which guarantees information to
98 flow unimpeded through the entire network.

99 Based on the idea of identity mapping, there are two representative neural network architecture that
100 outperform standard NN in most benchmark tasks. "ResNet"[9, 8] is a deep CNN architecture that has
101 impressive performance for tasks such as image classification, object detect, semantic segmentation.
102 Instead of simply stacking layers, ResNet added identity mapping, thus the deeper the network is, the
103 lower error rate it will get(since identity mapping could guarantee that the deeper on will have at least
104 same training error rate as the shallower one). Inspired by ResNet,[10] brought up "DenseNet", which
105 is an advanced convolutional network architecture that achieves the-state-of-art performance in image
106 classification task on various benchmark datasets. This architecture comprises several "dense blocks"
107 connected with convolution layers and pooling layers, each block is a convolutional feed-forward
108 network with skip-connection(which is identity mapping) to flow information between earlier layers
109 and later layers.

110 3 Main Theorem

111 In this section, we first puts how we formulate this problem and provide some preliminaries (including
112 notation and related definitions), then we develop our main theorems.

113 3.1 Problem Formulation

114 In our project, we tried to analyze SGD's convergence of single-hidden-layer neural network with
115 some non-linear activation. First, for the activation, a Parametric Rectified Linear Unit (PReLU)[7] is
116 proposed which generalizes the traditional rectified unit,which need almost no extra cost for compu-
117 tation and overfitting risk while improving model fitting. Also, inspired by Residual Network[8], we
118 add identity mapping. Thus, our network structure could be presented as Figure 1.

119 With this architecture, we can develop a formal form for this problem as follows:

120 Denote by $x \in \mathbb{R}^{n \times 1}$ the input of the network, denote by $W \in \mathbb{R}^{n \times n}$ the parameters forwarding
121 message from input layer to hidden layer, denote by σ the activation function(which is PReLU in our
122 case), the out put of our network architecture could be given as:

$$y(x, W) = \|\sigma(Wx + x)\|_1 \quad (1)$$

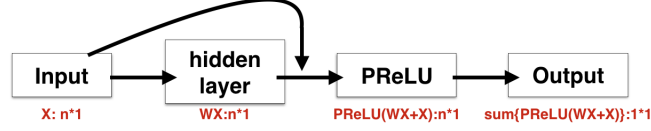


Figure 1: Network Architecture

123 where $\|\cdot\|_1$ denotes L1-norm which is the sum of all elements in “.”. Also, in our setting, input x
 124 should obey normal distribution, that is to say, $x \sim \mathcal{N}(0, I)$.

125 With the standard setting[15], we give the assumption that there exists a teaching network and a
 126 student network, both networks share the same structure which is a feed-forward 2 layer network
 127 with identity mapping and PReLU activation shown in Figure1. Denote by W^* the parameter(i.e.
 128 the weight) of teaching network, and W the parameter of student network. Our goal is to train the
 129 student network to learn the results given by teaching network. While training the student network,
 130 we apply quadratic loss, thus our objective function would be:

$$L(W) = \mathbb{E}_x[(y(x, W) - y(x, W^*))^2] \quad (2)$$

131 By applying equation 1, we will have:

$$L(W) = \mathbb{E}_x[(\|\sigma(Wx + x)\|_1 - \|\sigma(W^*x + x)\|_1)^2] \quad (3)$$

132 And with algebra skills, denote by $\mathbf{e} = e_1, \dots, e_n$ the base vectors, and $W = (w_1, \dots, w_n)$, $W^* =$
 133 (w_1^*, \dots, w_n^*) . Then the above equation3, would be written as:

$$L(W) = \mathbb{E}_x[(\sum_i(\sigma(\langle w_i + e_i, x \rangle) - \sigma(\langle w_i^* + e_i, x \rangle)))^2] \quad (4)$$

134 By formulating the problem into this way, we find it can perfectly fit into a two phase framework
 135 brought up by [13] which was designed to analyze the convergence of SGD. In the two phase
 136 framework, SGD could be regarded as a 2 stages process. In the first stage, W may head towards
 137 wrong direction but with shrinking probability and will definitely get into a specific region(i.e. one-
 138 point region, which we will define in section3.2), which lead to stage 2. Then in this stage 2, W will
 139 for sure move towards correct direction, and thus finally converge to the target W^* .

140 3.2 Preliminaries

141 In this section, we present some definitions and notations that will be used throughout the analysis.

142 **Basic Notations** Basic notations appeared in the above sections remain the same as before. Plus, we
 143 denote $\theta(v_1, v_2)$ as the angle between vector v_1 and vector v_2 , denote by \bar{M} the column-normalized
 144 version of matrix M .

145 **Definition1: One Point Strong Convexity**[13] A function $f(x)$ is called δ -one point strongly convex
 146 in domain D w.r.t. point x^* if $\forall x \in D, \langle -\nabla f(x), x^* - x \rangle > \delta \|x^* - x\|_2^2$.

147 **Definition2: Auxiliary Function** Denote $f_A = \sum_i(\|e_i + w_i^*\|_2 - \|e_i - w_i\|_2)$ the main auxiliary
 148 function, and denote $f_A(i) = f_A - (\|e_i + w_i^*\|_2 - \|e_i - w_i\|_2)$

149 **Definition3: Auxiliary Matrix** Denote $A = (W^* + I)\bar{W}^* + \bar{I}^T - (W + I)\bar{W} + \bar{I}^T$ the main
 150 auxiliary matrix, and denote $A(i) = A - ((e_i + w_i^*)(e_i + w_i^*)^T - (e_i + w_i)(e_i + w_i)^T)$

151 3.3 Theorems

152 **Main Theorem** While $x \sim \mathcal{N}(0, I)$, $\|w\|_2 \leq \gamma$ and $\|W^*\| \leq \gamma^*$ are both bounded with some small
 153 constant γ, γ^* , SGD with small learning rate η and initial W_0 (random/zero/standard all work) will
 154 converge to W^* within polynomial number of steps, in two phase.

155 This main theorem could be divided into two theorems, which are the theorem for phase 1 and phase
 156 2 respectively.

157 **Theorem for Phase 1** While $x \sim \mathcal{N}(0, I)$, $\|w\|_2 \leq \gamma$ and $\|W^*\| \leq \gamma^*$ are both bounded with some
 158 small constant γ, γ^* , then the auxiliary function f_A will keep decreasing after every step until $f_A \leq \epsilon$
 159 is small enough (which means we enter phase 2).

160 **Theorem for Phase 2** While $x \sim \mathcal{N}(0, I)$, $\|w\|_2 \leq \gamma$ and $\|W^*\| \leq \gamma^*$ are both bounded with some
 161 small constant γ, γ^* , and $f_A \leq \epsilon$, then $\langle -\nabla L(W), W^* - W \rangle > \delta \|W^* - W\|_2^2$ with constant δ .

162 4 Overview of Proofs

163 This section gives an overview of proofs of our main theorem. We first provide a flowchart as a clear
 164 roadmap to our proof, then we illustrate the key point of each part of the proofs. About detailed
 165 proofs, due to time limitation, we didn't type out all our proofs within this latex document. We can
 166 provide a handwriting version of detailed proof if needed.

167 4.1 Generally: A 2-stage Process

168 To proof our theorem, we divide our proof into 2 stages. In the first stage, we proof that our W
 169 will definitely enter the one-point convexity region. To check if its a one-point convexity region, we
 170 calculate $\langle -\nabla L(W), W^* - W \rangle$ and proof it will be less than or equal to $\delta \|W^* - W\|_2^2$. Here we
 171 have,

$$172 \quad \nabla L(W)_j = 2\mathbb{E}_x[(\sum_i \sigma(\langle e_i + w_i, x \rangle) - \sum_i \sigma(\langle e_i + w_i^*, x \rangle))x \mathbb{1}_{\langle e_j + w_j, x \rangle \geq 0}] \quad (5)$$

$$+ (\sum_i \sigma(\langle e_i + w_i, x \rangle) - \sum_i \sigma(\langle e_i + w_i^*, x \rangle))x \alpha \mathbb{1}_{\langle e_j + w_j, x \rangle < 0} \quad (6)$$

173 Then, after entering the one-point convexity region, we move to stage 2 and proof after each step, W
 174 will get closer to W^* with no exception.

175 4.2 For Phase I

176 Phase I aims at proving that, $\exists \gamma_0 \in (0, \gamma)$ s.t. If $\|\mathbf{W}_0\|_2, \|\mathbf{W}^*\|_2 \leq \gamma_0$, and d has a constant lower
 177 bounded, η has a upper bound determined by γ and G the gradient, ϵ upper bounded by a term
 178 depending solely on γ , then the auxiliary function f_A will keep decreasing until it reaches $O(1)$,
 179 which is the condition of starting *Phase 2*.

180 The decreasing factor for each step is influenced by η and d , the total number of iterations needed is
 181 determined by η , and the final value of f_A after Phase I ends is decided by γ .

182 In order to prove the above claims and to give a clear bound of all terms, we use approximations.

183 The way we prove the decreasing trend of f_A is by first introducing an auxiliary variable $s =$
 184 $(\mathbf{W}^* - \mathbf{W})u$, where u is the all-one vector, and then express $s^{(t)}$ and $f_A^{(t)}$, using $s^{(t-1)}$ and $f_A^{(t-1)}$.
 185 Then we could solve the dynamics from the previous step to show that g , a potential function that is
 186 expected to depend on f_A , approaches to and stays around $O(\gamma)$.

187 The second task of Phase I is to use the conclusion that f_A decreases to prove that $\|\mathbf{W}\|$ remains
 188 small. This limitation will guarantee that once we move on to Phase II, there is no possibility of
 189 coming back to Phase I.

190 A and f_A are useful tools to help with constructing a matrix P to approximate $-\nabla L(W)$. The proof
 191 of that P is an appropriate approximation should also be provided in Phase I.

192 4.3 For Phase II

193 The goal of Phase II is to prove that SGD can obtain optimal parameters in the small region derived
 194 from Phase I. i.e., Prove that $\exists \gamma$ with a small enough f_A , s.t. $L(W)$ is a δ one point strongly convexity.
 195 The formal is shown as follow.

$$\langle -\nabla L(W), W^* - W \rangle = \sum_{j=1}^d \langle -\nabla L(W)_j, w_j^* - w_j \rangle > \delta \|W^* - W\|_F^2 \quad (7)$$

196 Here we use Taylor expansion and control the higher order term, which is shown in figure 2 Then,
 197 lower bound each part of Taylor expansion separately. Note that when $W \approx W^*$, we will use

$$\langle \text{constant} + \text{1 order} + \text{higher order}, W^* - W \rangle$$

Figure 2: Taylor expansion

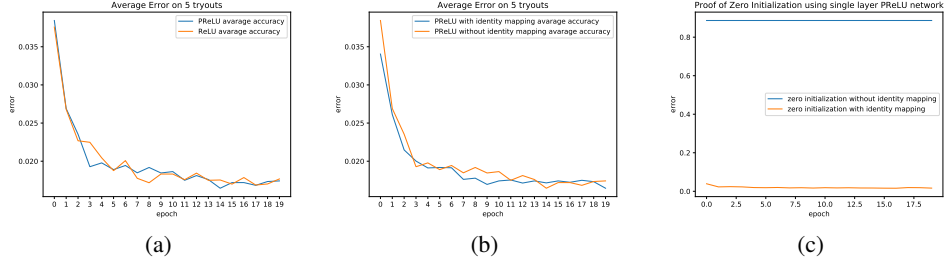


Figure 3: (a).Accuracy for NN with ReLU activation and PReLU activation, (b).Accuracy Curve for NN architecture with and without identity mapping structure, (c).Performance of NN with or without Identity Mapping while given Zero Initialization

198 joint Taylor expansion to overcome a large higher term. Then, we lower bound each part of Taylor
 199 expansion by $c * ||W^* - W||_F^2$, where c is a small constant related to f_A . Note that the value of
 200 f_A is derived via Phase I.

201 5 Experiments

202 We designed a 4-stages experiment, each stage of experiment will perform as an auxiliary support for
 203 our main theorem.

204 5.1 Stage1: Show the Advancement of PReLU Activation

205 In the first stage of our experiment, we tried to show the advancement of PReLU activation even with
 206 one-hidden-layer neural network. Here we used MNIST handwritten digits dataset, trained with two
 207 layer NN without identity mapping, and with PReLU activation and ReLU activation respectively(so
 208 there are two different settings). We conducted each setting of experiments for 5 times and took the
 209 average of those tryouts. The curve of error rate we obtained is presented in Figure3(a).

210 For ReLU without identity mapping, random initialization, we get an average error rate of 0.01768.
 211 And for PReLU, without identity mapping, random initialization, we get an average error rate of
 212 0.01742. The results indicates PReLU's advancement successfully.

213 5.2 Stage2: Validate Identity Mapping Helps in Improving Accuracy

214 In the second stage of our experiment, we tried to show the advancement of identity mapping with
 215 one-hidden-layer neural network. The dataset we applied remains the same as in stage 1, trained with
 216 two layer NN with and without identity mapping respectively, and with PReLU activation(so there
 217 are still two different settings). We conducted each setting of experiments for 5 times and took the
 218 average of those tryouts. The curve of error rate we obtained is presented in Figure3(b).

219 For PReLU without identity mapping, random initialization, we get an average error rate of 0.01742.
 220 And for PReLU with identity mapping, random initialization, we get an error rate of 0.01646. The
 221 result indicates identity mapping's advancement successfully.

222 5.3 Stage3: Validate Zero Initialization Works

223 In the third stage of our experiment, we tried to show the important role identity mapping plays while
 224 with zero initialization. Again, most the experiment settings remains the same as in stage 2, but the

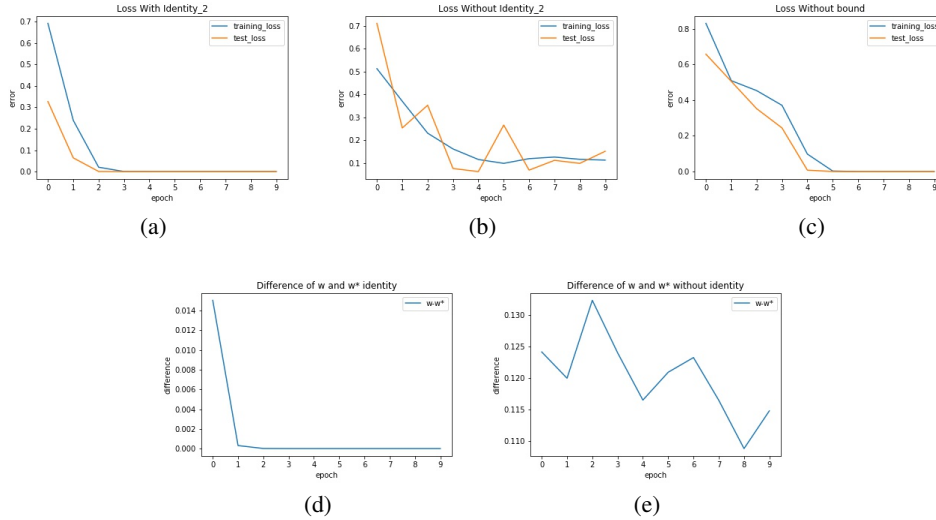


Figure 4: (a).Loss with identity mapping and bound, (b).Loss without Identity mapping, (c).Loss without bound, (d) $\|W - W^*\|_2$ with identity mapping, (e) $\|W - W^*\|_2$ without identity mapping

225 for the initialization we use zeros instead of some random values. The curve of error rate we obtained
 226 is presented in Figure3(c).

227 For PReLU without identity mapping, zero initialization, we get an average error rate of 0.8865. And
 228 for PReLU with identity mapping, zero initialization, we get an error rate of 0.0165. This result
 229 sufficiently illustrates that with identity mapping, the zero initialization will still work(and almost
 230 have the same performance as random initialization).

231 5.4 Stage4: Global Minimum Convergence

232 Following the work of[15, 13], we utilized one teacher network and at least one student network,
 233 where teacher network knows the ground truth optima parameters W^* , while the students will learn
 234 W via l_2 loss.

235 In this stage, we tried to show that, while $\|W\|_2$ and $\|W^*\|_2$ are bounded by some small value, when
 236 applying identity mapping, SGD can converge to global minima with zero error rate. Here we used
 237 random generalized input dataset of dimension $d = 10$ and size of 5000, both our teaching network
 238 and student network share the same structure of 2-layer with PReLU activation, for the training
 239 process of student network, we set learning rate as 0.1 and epoch as 10. We observed the loss with
 240 and without identity mapping, and with and without bound of $\|W^*\|_2$. For the bound, we set $\|w^*\|_2$
 241 = 0.1 and $\|w^*\|_F = 1$. To avoid coincidence, we conducted each setting of experiments for 5 times
 242 and took the average of those tryouts. The curves we obtained is presented in Figure4.

243 This result sufficiently presents that both the bound of $\|W^*\|_2$ and the identity mapping are crucial
 244 to SGD's convergence.

245 6 Other Attempts

246 We also took several steps in the SGD convergence analysis for neural networks, trying to see if this
 247 nice 2-phase framework is still applicable. Those attempts include: extending the proof framework
 248 to multiple layers neural networks, varying the network structures, and formulating several other
 249 non-convex problems and applying the proof framework to them. Though failed to achieve complete
 250 and solid proofs (or sometimes even failed to find a proper way to formulate the main theorem), we
 251 still want to discuss about those attempts and analyze why we failed in those attempts.

252 **6.1 Deepen the Network**

253 Our first step is to deepening the network by stacking multiple layers with identity mapping and
254 trying to extend the 2 phase framework[13] to deep neural networks. Then, by doing so, the problem
255 would be formulated as(take N-hidden-layers as example):

$$y(x, W) = \|\sigma((W_N + I) \dots \sigma((W_2 + I) \sigma((W_1 + I)x))\|_1 \quad (8)$$

256 Then, while dealing with $L(W)$, we find it hard to formulate it in an elegant way like equation4 does,
257 thus, by referring to the work[11], we tried to applied the way brought up by this work to turn our
258 non-linear deep neural network into a linear form via reduction, thus it would be equipped with the
259 following formula:

$$y(x, W) = \|\sigma((W_N \dots W_2 W_1 + I)x)\|_1 \quad (9)$$

260 which could be the exactly the same as the problem setting the 2-phase framework has. However,
261 this would cause some problems. Because to simplify the network into a linear one, we ignored the
262 activation functions by adding some constrains, which makes this problem less flexible, and thus
263 made sure that all the global minima would be local minima, so that the proof of SGD's convergence
264 to global minima no longer meaningful(since one of the key contribution to the convergence analysis
265 is to proof it can avoid sticking at local minima).

266 According to the above statement, we found applying the two-stage framework on deep neural
267 networks seems not to be an applicable idea.

268 **6.2 Vary the Network Structures**

269 We also tried to explore more about the convergence of SGD in diverse network structures. Based
270 on the nice properties identity mapping has, naturally we came to ResNet and DenseNet, which are
271 mentioned in section 2.3.

272 However, with such network structures, it would then turn to deep neural network, thus we encountered
273 the same problems in section 6.1, where the 2-stage proof framework no longer works for this
274 situation.

275 **6.3 For Other Non-convex Problems**

276 Since in previous sections of this report, we've already formulate this problem in a nice way in which
277 network function and loss function could be represent with the same equation14. By simply changing
278 the activation function σ , our problem will be different(not slightly) while still remains a non-convex
279 problem. However, even though, and thus we encountered lots of difficulties while figuring out the
280 corresponding auxiliary function and auxiliary matrix. Since for each different activation function,
281 there's no uniform form of auxiliary function and auxiliary matrix(which means we need a different
282 form of auxiliary function and auxiliary matrix), this attempt would turn to different direction and
283 need diverse proof techniques corresponding with different activation. And we do not have enough
284 time go down every branch of this path.

285 **7 Conclusion**

286 To sum up, in this project, our main work could be summarized into 4 parts, first, we conducted a
287 comprehensive literature review including 3 types of related work(details are in section 2), second,
288 we proved the convergence of SGD with one-hidden-layer neural network with PReLU activation and
289 without constrains on initialization, third, though failed, we made several attempts on applying the
290 nice 2-phase framework on NN with different layers, structures, or for different non-convex problems
291 and analyzed why we can't get expected theoretical results, last, we conducted a 4-stage experiment
292 to support our theory intuitively.

293 **References**

294 [1] Zeyuan Allen-Zhu, Yuanzhi Li, and Zhao Song. A convergence theory for deep learning via
295 over-parameterization. *arXiv preprint arXiv:1811.03962*, 2018.

- 296 [2] Zeyuan Allen-Zhu, Yuanzhi Li, and Zhao Song. On the convergence rate of training recurrent
297 neural networks. *arXiv preprint arXiv:1810.12065*, 2018.
- 298 [3] Alon Brutzkus, Amir Globerson, Eran Malach, and Shai Shalev-Shwartz. Sgd learns over-
299 parameterized networks that provably generalize on linearly separable data. *arXiv preprint*
300 *arXiv:1710.10174*, 2017.
- 301 [4] Simon S Du, Jason D Lee, Haochuan Li, Liwei Wang, and Xiyu Zhai. Gradient descent finds
302 global minima of deep neural networks. *arXiv preprint arXiv:1811.03804*, 2018.
- 303 [5] Simon S Du, Xiyu Zhai, Barnabas Poczos, and Aarti Singh. Gradient descent provably optimizes
304 over-parameterized neural networks. *arXiv preprint arXiv:1810.02054*, 2018.
- 305 [6] Moritz Hardt, Tengyu Ma, and Benjamin Recht. Gradient descent learns linear dynamical
306 systems. *arXiv preprint arXiv:1609.05191*, 2016.
- 307 [7] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Delving deep into rectifiers:
308 Surpassing human-level performance on imagenet classification. In *Proceedings of the IEEE*
309 *international conference on computer vision*, pages 1026–1034, 2015.
- 310 [8] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image
311 recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*,
312 pages 770–778, 2016.
- 313 [9] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Identity mappings in deep residual
314 networks. In *European conference on computer vision*, pages 630–645. Springer, 2016.
- 315 [10] Gao Huang, Zhuang Liu, Laurens Van Der Maaten, and Kilian Q Weinberger. Densely connected
316 convolutional networks. In *CVPR*, volume 1, page 3, 2017.
- 317 [11] Kenji Kawaguchi. Deep learning without poor local minima. In *Advances in Neural Information*
318 *Processing Systems*, pages 586–594, 2016.
- 319 [12] Ilja Kuzborskij and Christoph H Lampert. Data-dependent stability of stochastic gradient
320 descent. *arXiv preprint arXiv:1703.01678*, 2017.
- 321 [13] Yuanzhi Li and Yang Yuan. Convergence analysis of two-layer neural networks with relu
322 activation. In *Advances in Neural Information Processing Systems*, pages 597–607, 2017.
- 323 [14] Itay Safran and Ohad Shamir. Spurious local minima are common in two-layer relu neural
324 networks. *arXiv preprint arXiv:1712.08968*, 2017.
- 325 [15] Yuandong Tian. Symmetry-breaking convergence analysis of certain two-layered neural net-
326 works with relu nonlinearity. 2017.
- 327 [16] Kai Zhong, Zhao Song, Prateek Jain, Peter L Bartlett, and Inderjit S Dhillon. Recovery
328 guarantees for one-hidden-layer neural networks. *arXiv preprint arXiv:1706.03175*, 2017.
- 329 [17] Difan Zou, Yuan Cao, Dongruo Zhou, and Quanquan Gu. Stochastic gradient descent optimizes
330 over-parameterized deep relu networks. *arXiv preprint arXiv:1811.08888*, 2018.