



저작자표시-비영리-변경금지 2.0 대한민국

이용자는 아래의 조건을 따르는 경우에 한하여 자유롭게

- 이 저작물을 복제, 배포, 전송, 전시, 공연 및 방송할 수 있습니다.

다음과 같은 조건을 따라야 합니다:



저작자표시. 귀하는 원저작자를 표시하여야 합니다.



비영리. 귀하는 이 저작물을 영리 목적으로 이용할 수 없습니다.



변경금지. 귀하는 이 저작물을 개작, 변형 또는 가공할 수 없습니다.

- 귀하는, 이 저작물의 재이용이나 배포의 경우, 이 저작물에 적용된 이용허락조건을 명확하게 나타내어야 합니다.
- 저작권자로부터 별도의 허가를 받으면 이러한 조건들은 적용되지 않습니다.

저작권법에 따른 이용자의 권리는 위의 내용에 의하여 영향을 받지 않습니다.

이것은 [이용허락규약\(Legal Code\)](#)을 이해하기 쉽게 요약한 것입니다.

[Disclaimer](#)

이학석사 학위논문

Input Distribution Analysis for Learning Two Layer Neural Network

(2층 신경망의 학습을 위한 입력 분포 분석)

2019년 2월

서울대학교 대학원

수리과학부

박정명

Input Distribution Analysis for Learning Two Layer Neural Network

(2층 신경망의 학습을 위한 입력 분포 분석)

지도교수 강 명 주

이 논문을 이학석사 학위논문으로 제출함

2018년 10월

서울대학교 대학원

수 리 과 학 부

박 정 명

박 정 명의 이학석사 학위논문을 인준함

2018년 12월

위 원 장 _____ (인)

부 위 원 장 _____ (인)

위 원 _____ (인)

Input Distribution Analysis for Learning Two Layer Neural Network

by

Jeongmyeong Park

A DISSERTATION

Submitted to the faculty of the Graduate School
in partial fulfillment of the requirements
for the degree of Master of Science
in the Department of Mathematical Sciences
Seoul National University
February 2019

Abstract

Input Distribution Analysis for Learning Two Layer Neural Network

Jeongmyeong Park

Department of Mathematical Sciences
The Graduate School
Seoul National University

In recent years, Deep neural networks have achieved state-of-the-art performance in many tasks. Despite those empirical successes, it remains hard to explain why stochastic gradient descent can solve the highly non-convex optimization problem. In this paper, we analyze the convergence of stochastic gradient descent for weights in learning process. The model we used is two layer neural network with ReLU activation function. In particular, our theory used the notion of the smoothness of the input distribution so that we don't need any specific input distribution.

Key words: Optimization, Input distribution, Neural Network, Deep Learning

Student Number: 2016-20235

Contents

Abstract	i
1 Introduction	1
2 In case of $n=1$	3
3 Generalization for arbitrary n	6
4 Conclusion	9
The bibliography	10
5 appendix	11
5.1 Proof of Theorem 3.4	11
Abstract (in Korean)	15

Chapter 1

Introduction

In recent years, Deep neural networks have achieved state-of-the-art performance in many tasks such as computer vision [6], natural language processing [2], speech recognition [4]. Despite those empirical successes, it remains difficult to explain the performance theoretically. One of the theoretical problems is the optimization. To optimize the neural network, We use stochastic gradient descent(SGD). SGD is a simple algorithm using first order which is effective and easy to modify. Even though there are many different optimizers used in real works, all of them are based on SGD. Why such a simple algorithm can be so successful? Because of the simplicity of SGD, one may think that it won't be hard to explain it. But it is wrong. Deep neural networks are the highly non-convex optimization problems. Not only because it is deep, people usually adopt ReLU function as their activation function. So it is almost impossible to analyze all of the aspect in deep network. Although it is a difficult problem, researchers have tried many theoretical attempts to analyze it. Recently, there have been many papers which try to explain the convergence of shallow networks. Zhang et al.[7] analyzed two layer network on certain types of activation functions, but generalization on ReLU function remained unclear. Brutzkus et al. [1] showed optimization of convolutional network with ReLU function and gaussian input can be proved. Janzamin et al. [5] guarantees the reexery of parameters on 2-layered network with tensor decomposition. However, most of those attempts relies on the certain input distribution, like gaussian distribution. Instead of the specific form of input

distribution, Du et al. [3] used the notion of smoothness of the input distribution. Their analysis shows that any kind of inputs can be optimized with some conditions. In this paper, we will follow the similar method. We build a shallow network and discuss about its convergence. More precisely, we are going to follow the Idea of Du et al. [3]. Their suggestion of smoothness will be used and there will be no assumption of input distribution except the smoothness. Since his idea is constrained on one layer convolutional network, we will adapt the notion of smoothness on two layer network with ReLU activation function. With the notion of smoothness, We are going to connect it to the convergence rate for weights recovery. The smoother input distribution is, the faster convergence happens. The content of this paper is organized as follows. We introduce our network and apply Du’s idea in cases that only one weight vector is used in **Chapter 2**. We generalize the idea in cases when multiple weight vectors are used in **Chapter 3**. Conclusions are in **Chapter 4**.

Chapter 2

In case of n=1

Now, we introduce our two layer neural network. We will analyze the optimization on this architecture. We call $z \in \mathbb{R}^k$ as an input vector, $\mathbf{W} = (w_1, w_2, \dots, w_n) \in \mathbb{R}^{k \times n}$ and $w_i \in \mathbb{R}^k$ as weights. $\sigma(x) = \max(x, 0)$ is the ReLU activation function. In this paper, we focus on the following two layer ReLU neural network.

$$f(\mathbf{W}, z) = \frac{1}{k} \sum_{i=1}^k \sigma(w_i^\top z) \quad (2.1)$$

After multiplication of input and weights, ReLU activation function follows. Lastly average pooling derives output. For Loss function,

$$\ell(w, z) = \frac{1}{2} (f(w, z) - f(w_*, z))^2 \quad (2.2)$$

When we apply stochastic gradient descent in real experiment, it works as follows

$$w_{t+1} = w^t - \eta^t \nabla \ell(w^t, z) \quad (2.3)$$

However, z is not a fixed input. It is a random variable of input distribution. So we use population gradient in analysis.

$$w_{t+1} = w^t - \eta^t \mathbb{E}[\nabla \ell(w^t, z)] \quad (2.4)$$

Before we start the discussion, we clarify some notations. $\lambda_{\max}(R)$ is the largest singular value of the matrix R and $\lambda_{\min}(R)$ is the smallest singular

value of the matrix \mathbf{R} . $\|\cdot\|_{oper}$ denotes the operator norm of a matrix. For the convenience in the calculation, we will use vector \mathbf{W}_{vec} instead of matrix \mathbf{W} . The definition is as follows. $\mathbf{W}_{vec} = [w_1^\top, w_2^\top, \dots, w_n^\top] \in \mathbb{R}^{kn}$. Also we make some assumptions. Assume that inputs, weights are bounded. As a result, gradient functions are uniformly bounded.

In our paper, activation function is ReLU function. ReLU function is hard to deal because of its vanishing properties. To solve the problem with minimum assumptions, we split regions so that we can verify the vanishment by ReLU function. Since our network is shallow network, it is affordable calculations. Also using following definitions, we will define the smoothness of input distribution.

Definition 2.1. (in Du[3]) We define two events and matrices. w can be any w_i in \mathbf{W}_{vec}

$$\begin{aligned} R(w, w_*) &= \{z : w^\top z \geq 0, w_*^\top z \geq 0\} \\ R(w, -w_*) &= \{z : w^\top z \geq 0, -w_*^\top z \geq 0\} \\ \mathbf{R}_{w, w_*} &= \mathbb{E} [zz^\top \mathbb{I}\{R(w, w_*)\}] \\ \mathbf{R}_{w, -w_*} &= \mathbb{E} [zz^\top \mathbb{I}\{R(w, -w_*)\}] \end{aligned} \tag{2.5}$$

Now it's time to introduce the smoothness in the case of $n=1$. We are going to generalize it for arbitrary n in Chapter 3.

Definition 2.2. (in Du[3]) For ϕ in $[0, \pi]$,

$$\begin{aligned} m_{min}^{w, w_*}(\phi) &= \min_{w: \angle w, w_* = \phi} \lambda_{\min}(\mathbf{R}_{w, w_*}) \\ m_{max}^{w, w_*}(\phi) &= \max_{w: \angle w, w_* = \phi} \lambda_{\max}(\mathbf{R}_{w, w_*}) \end{aligned} \tag{2.6}$$

Du defined the smoothness as the difference of the largest and smallest eigenvalue values of the \mathbf{R}_{w, w_*} . If the difference of the two eigenvalue is large, then we can consider it as the sign of the biased probability mass. When the input distribution is gaussian distribution or the rotationally invariant, Two eigenvalue will match by definition. So, closer two eigenvalues are, more even the probability mass is. Therefore the definition can be called the smoothness. This is the definition of smoothness when $n=1$. Du generalized the smoothness to use on the convolutional network. Since our network is not convolutional

network, we tried a different way of generalization. Before we discuss the two layer network for arbitrary n , we will analyze the case of $n=1$.

Theorem 2.3. *(in Du[3]) Assume that there exists positive d s.t. $m_{max}^{w,-w_*}(\phi) = \max_{w: \angle w, w_* = \phi} \lambda_{\max}(\mathbf{R}_{w,-w_*}) < d\phi$. and $\|w_0 - w_*\|_2 < \|w_*\|_2$ for the initialization w_0 . Let $\phi^t = \arcsin(\|w_0 - w_*\|_2 / \|w_*\|_2)$. If $0 \leq \eta_t \leq \min_{0 \leq \phi \leq \phi^t} \frac{m_{min}(\phi)}{2(m_{max}(\phi) + 2\alpha)^2}$, then*

$$\|w_{t+1} - w_*\|_2^2 \leq \left(1 - \frac{\eta_t m_{min}(\phi_t)}{2}\right) \|w_t - w_*\|_2^2 \quad (2.7)$$

Now, by this theorem, we can sure that network will converge in a proper time. This theorem show the convergence rate of the $n=1$ network. It suggests the relationship between smoothness and converging time. With small α , small m_{max} and Large m_{min} , we can choose large learning rate. Also the bigger m_{min} is, the faster convergence comes.

Chapter 3

Generalization for arbitrary n

It is time to generalize for our network. Our network has two layers. In the hidden layer, each w_i 's multiplies z and pass through the ReLU function. So when we investigate the gradient of loss function, we have to face the intersection of ReLU functions. Therefore, we need more complex split of regions.

Definition 3.1.

$$b_{i,j}^{k,l,m} = z z^\top R(w_i, (-1)^k w_{i,*}) R((-1)^l w_j, (-1)^m w_{j,*})$$

$$\mathbf{B}_{k,l,m} = \frac{1}{n^2} \mathbb{E} \begin{bmatrix} b_{1,1}^{k,l,m} & \dots & b_{1,n}^{k,l,m} \\ \vdots & \ddots & \vdots \\ b_{n,1}^{k,l,m} & \dots & b_{n,n}^{k,l,m} \end{bmatrix} \quad (3.1)$$

Definition 3.2.

$$m_{\max}(\phi, w_{i,*}) = \max_{\forall i, \angle w_i, w_{i,*} = \phi} \lambda_{\max}(\mathbf{B}_{0,0,0})$$

$$m_{\min}(\phi, w_{vec,*}) = \min_{\forall i, \angle w_i, w_i = \phi} \lambda_{\min}(\mathbf{B}_{0,0,0}) \quad (3.2)$$

Assumption 3.3. *We assume that exists positive α and m_{\max}^{ex} such that*

$$\max_{(k,l,m)=(1,0,0),(0,0,1),(0,1,0),(1,1,0)} \left\{ \max_{\forall i, \angle w_i, w_{i,*} = \phi} \lambda_{\max}(\mathbf{B}_{k,l,m}) \right\} \leq m_{\max}^{ex} \phi \quad (3.3)$$

$$\max_{\forall i, \angle w_i, w_{i,*} = \phi} \lambda_{\max}(\mathbf{B}_{1,0,1}) \leq \alpha \phi$$

By converging n weights in one vector, we derived the very similar detail with Du's generalization on Convolutional Neural Network. So the following theorems are adaptation of his work on two layer neural network. By following theorem, we show the convergence guarantee of the two neural network in the following theorem.

Theorem 3.4. (adapted from Du[3]) Assume that $\|\mathbf{W}_{vec}^0 - \mathbf{W}_{vec,*}\|_2 < \|\mathbf{W}_{vec,*}\|_2$ for the initialization \mathbf{W}_{vec}^0 . Let $\phi^t = \arcsin(\|\mathbf{W}_{vec}^t - \mathbf{W}_{vec,*}\|_2 / \|\mathbf{W}_{vec,*}\|_2)$. If $\eta_t \leq \min_{0 \leq \phi \leq \phi^t} \frac{m_{min}(\phi) - 10m_{max}^{ex}}{2(L(\phi) + 14L_{ex} + 4\alpha)^2}$, then

$$\|\mathbf{W}_{vec}^{t+1} - \mathbf{W}_{vec,*}\|_2^2 \leq \|\mathbf{W}_{vec}^t - \mathbf{W}_{*,vec}\|_2^2 \left(1 - \frac{\eta_t(m_{min} - 10m_{max}^{ex})}{2}\right)$$

Proof. In Appendix □

Now, we have our result. Population gradient leads the network to the convergence. For the faster convergence, we have to get a smaller m_{max}^{ex} , larger m_{min} , smaller L . Then we can have larger learning rate and better convergence rate. But we haven't show whether $m_{min} - 10m_{max}^{ex}$ is positive or not. Therefore, we introduce the next theorem.

Theorem 3.5. (adapted from Du[3]) Suppose Z has unit norm. ϕ is angle between $w_i, w_{i,*}$. Assume that there exists β such that for all i

$$\mathbb{P}[R(w_i, -w_{i,*})], \mathbb{P}[R(-w_i, w_{i,*})] \leq \beta\phi$$

then we have $m_{max}^{ex} \leq \beta\phi$

Proof. Assume $[x_1, x_2, \dots, x_n]$ has unit norm.

$$\begin{aligned}
& [x_1, x_2, \dots, x_n] \mathbf{B}_{0,0,1} [x_1, x_2, \dots, x_n]^\top \\
&= \frac{1}{n^2} \sum E [x_i^\top Z Z^\top x_j R(w_i, w_{i,*}) R(w_j, -w_{j,*})] \\
&\leq \frac{1}{n^2} \sum \|x_i\| \|x_j\| E [R(w_i, w_{i,*}) R(w_j, -w_{j,*})] \\
&\leq \frac{1}{n^2} \sum \frac{1}{2} (\|x_i\|^2 + \|x_j\|^2) E [R(w_i, w_{i,*}) R(w_j, -w_{j,*})] \\
&\leq \frac{1}{n^2} \sum \frac{1}{2} (\|x_i\|^2 + \|x_j\|^2) E [R(w_j, -w_{j,*})] \\
&\leq \frac{1}{n^2} \sum \frac{1}{2} (\|x_i\|^2 + \|x_j\|^2) \beta \phi \\
&= \frac{\beta \phi}{n}
\end{aligned}$$

For $\mathbf{B}_{1,0,0}$, $\mathbf{B}_{0,1,0}$, $\mathbf{B}_{1,1,0}$, the same as above. \square

Now we have upper bound of m_{max}^{ex} . If input distribution has low probability in $R(w_i, -w_{i,*})$, $R(-w_i, w_{i,*})$, then the assumptions above can be fulfilled.

Chapter 4

Conclusion

In this paper, we have discussed about the conditions on the convergence of two layer neural network. First, we discussed the case of $n=1$. We defined the smoothness and established theorem for convergence. Secondly, We generalized the smoothness for arbitrary n . Using the notion, we provided the conditions to guarantee the convergence for two layer neural network. While input distribution satisfies the conditions, it doesn't heed gaussian distribution. Also, with the convergence rate we calculated, we can check the two layer neural network can learn weights in a proper time. For the future works, we have to think about how to expand the smoothness for deeper network. Basically, our method is based on the dividing the regions. If a network has more layers, there are too many regions to consider. So it is not likely to work on a deeper network.

Bibliography

- [1] A. BRUTZKUS AND A. GLOBERSON, *Globally optimal gradient descent for a convnet with gaussian inputs*, arXiv preprint arXiv:1702.07966, (2017).
- [2] Y. N. DAUPHIN, A. FAN, M. AULI, AND D. GRANGIER, *Language modeling with gated convolutional networks*, arXiv preprint arXiv:1612.08083, (2016).
- [3] S. S. DU, J. D. LEE, AND Y. TIAN, *When is a convolutional filter easy to learn?*, arXiv preprint arXiv:1709.06129, (2017).
- [4] G. HINTON, L. DENG, D. YU, G. E. DAHL, A.-R. MOHAMED, N. JAITLEY, A. SENIOR, V. VANHOUCKE, P. NGUYEN, T. N. SAINATH, ET AL., *Deep neural networks for acoustic modeling in speech recognition: The shared views of four research groups*, IEEE Signal processing magazine, 29 (2012), pp. 82–97.
- [5] M. JANZAMIN, H. SEDGHI, AND A. ANANDKUMAR, *Beating the perils of non-convexity: Guaranteed training of neural networks using tensor methods*, arXiv preprint arXiv:1506.08473, (2015).
- [6] A. KRIZHEVSKY, I. SUTSKEVER, AND G. E. HINTON, *Imagenet classification with deep convolutional neural networks*, in Advances in neural information processing systems, 2012, pp. 1097–1105.
- [7] X. ZHANG, Y. YU, L. WANG, AND Q. GU, *Learning one-hidden-layer relu networks via gradient descent*, arXiv preprint arXiv:1806.07808, (2018).

Chapter 5

appendix

5.1 Proof of Theorem 3.4

Assume that $\|\mathbf{W}_{vec}^0 - \mathbf{W}_{vec,*}\|_2 < \|\mathbf{W}_{vec,*}\|_2$ for the initialization \mathbf{W}_{vec}^0 . Let $\phi^t = \arcsin(\|\mathbf{W}_{vec}^t - \mathbf{W}_{vec,*}\|_2 / \|\mathbf{W}_{vec,*}\|_2)$. If $\eta_t \leq \min_{0 \leq \phi \leq \phi^t} \frac{m_{min}(\phi) - 10m_{max}^{ex}}{2(L(\phi) + 14L_{ex} + 4\alpha)^2}$, then

$$\|\mathbf{W}_{vec}^{t+1} - \mathbf{W}_{vec,*}\|_2^2 \leq \|\mathbf{W}_{vec}^t - \mathbf{W}_{*,vec}\|_2^2 \left(1 - \frac{\eta_t(m_{min} - 10m_{max}^{ex})}{2}\right)$$

Proof. For w_i without time mark, consider it as weights of time t.

$$\begin{aligned} \mathbb{E}[\nabla_{w_i} \ell(\mathbf{W}^t, z)] &= \\ \frac{1}{n^2} \sum_{j=1}^n \mathbb{E} [zz^\top \mathbb{I} \{R(w_i, w_{i,*}) R(w_j, w_{j,*}) + R(w_i, -w_{i,*}) R(w_j, w_{j,*})\}] (w_j - w_{j,*}) \\ &+ \frac{1}{n^2} \sum_{j=1}^n \mathbb{E} [zz^\top \mathbb{I} \{R(w_i, w_{i,*}) R(w_j, -w_{j,*}) + R(w_i, -w_{i,*}) R(w_j, -w_{j,*})\}] w_j \\ &+ \frac{1}{n^2} \sum_{j=1}^n \mathbb{E} [zz^\top \mathbb{I} \{R(w_i, w_{i,*}) R(-w_j, w_{j,*}) + R(w_i, -w_{i,*}) R(-w_j, w_{j,*})\}] (-w_{j,*}) \end{aligned}$$

$$\begin{aligned}
& \sum_{i=1}^n \mathbb{E}[\langle \nabla_{w_i} \ell(\mathbf{W}^t, z), w_i - w_{i,*} \rangle] = \\
& \frac{1}{n^2} \sum_{(i,j)=(1,1)}^{(n,n)} (w_i - w_{i,*})^\top \mathbb{E} [zz^\top \mathbb{I} \{R(w_i, w_{i,*}) R(w_j, w_{j,*}) + \\
& R(w_i, -w_{i,*}) R(w_j, w_{j,*}) (w_j - w_{j,*}) \\
& + \frac{1}{n^2} \sum_{(i,j)=(1,1)}^{(n,n)} (w_i - w_{i,*})^\top \mathbb{E} [zz^\top \mathbb{I} \{R(w_i, w_{i,*}) R(w_j, -w_{j,*}) + \\
& R(w_i, -w_{i,*}) R(w_j, -w_{j,*}) w_j \\
& + \frac{1}{n^2} \sum_{(i,j)=(1,1)}^{(n,n)} (w_i - w_{i,*})^\top \mathbb{E} [zz^\top \mathbb{I} \{R(w_i, w_{i,*}) R(-w_j, w_{j,*}) + \\
& R(w_i, -w_{i,*}) R(-w_j, w_{j,*}) (-w_{j,*})
\end{aligned}$$

$$\begin{aligned}
& = (\mathbf{W}_{vec} - \mathbf{W}_{*,vec})^\top (\mathbf{B}_{0,0,0} + \mathbf{B}_{1,0,0}) (\mathbf{W}_{vec} - \mathbf{W}_{*,vec}) \\
& + (\mathbf{W}_{vec} - \mathbf{W}_{*,vec})^\top (\mathbf{B}_{0,0,1} + \mathbf{B}_{1,0,1}) (\mathbf{W}_{vec}) \\
& + (\mathbf{W}_{vec} - \mathbf{W}_{*,vec})^\top (\mathbf{B}_{0,1,0} + \mathbf{B}_{1,1,0}) (-\mathbf{W}_{*,vec})
\end{aligned}$$

Use definition so fm_{min} and remove positive terms

$$\begin{aligned}
& \geq \|\mathbf{W}_{vec} - \mathbf{W}_{*,vec}\|_2^2 m_{min} \\
& + (\mathbf{W}_{vec} - \mathbf{W}_{*,vec})^\top (\mathbf{B}_{1,0,0}) (-\mathbf{W}_{*,vec}) \\
& + (\mathbf{W}_{vec} - \mathbf{W}_{*,vec})^\top (\mathbf{B}_{0,0,1}) (\mathbf{W}_{vec}) \\
& + (\mathbf{W}_{vec} - \mathbf{W}_{*,vec})^\top (\mathbf{B}_{0,1,0} + \mathbf{B}_{1,1,0}) (-\mathbf{W}_{*,vec})
\end{aligned}$$

Now use norm.

$$\begin{aligned}
&\geq \|\mathbf{W}_{vec} - \mathbf{W}_{*,vec}\|_2^2 m_{min} \\
&- \|\mathbf{W}_{*,vec}\|_2 \|\mathbf{W}_{vec} - \mathbf{W}_{*,vec}\|_2 \|\mathbf{B}_{1,0,0}\|_{oper} \\
&- \|\mathbf{W}_{vec}\|_2 \|\mathbf{W}_{vec} - \mathbf{W}_{*,vec}\|_2 \|\mathbf{B}_{0,0,1}\|_{oper} \\
&- \|\mathbf{W}_{*,vec}\|_2 \|\mathbf{W}_{vec} - \mathbf{W}_{*,vec}\|_2 \|\mathbf{B}_{0,1,0}\|_{oper} \\
&- \|\mathbf{W}_{*,vec}\|_2 \|\mathbf{W}_{vec} - \mathbf{W}_{*,vec}\|_2 \|\mathbf{B}_{1,1,0}\|_{oper}
\end{aligned}$$

With properly assigned initial weights, we can assume $\|\mathbf{W}_{vec} - \mathbf{W}_{*,vec}\|_2 \leq \|\mathbf{W}_{*,vec}\|_2$ when $t=0$. If this theorem satisfies, then we can use it inductively for all t . Therefore, $\|\mathbf{W}_{vec}\|_2 \leq 2\|\mathbf{W}_{*,vec}\|_2$

$$\begin{aligned}
&\geq \|\mathbf{W}_{vec} - \mathbf{W}_{*,vec}\|_2^2 m_{min} \\
&- \|\mathbf{W}_{*,vec}\|_2 \|\mathbf{W}_{vec} - \mathbf{W}_{*,vec}\|_2 \|\mathbf{B}_{1,0,0}\|_{oper} \\
&- 2\|\mathbf{W}_{*,vec}\|_2 \|\mathbf{W}_{vec} - \mathbf{W}_{*,vec}\|_2 \|\mathbf{B}_{0,0,1}\|_{oper} \\
&- \|\mathbf{W}_{*,vec}\|_2 \|\mathbf{W}_{vec} - \mathbf{W}_{*,vec}\|_2 \|\mathbf{B}_{0,1,0}\|_{oper} \\
&- \|\mathbf{W}_{*,vec}\|_2 \|\mathbf{W}_{vec} - \mathbf{W}_{*,vec}\|_2 \|\mathbf{B}_{1,1,0}\|_{oper}
\end{aligned}$$

By assumption, operation norm of \mathbf{B} can be replaced.

$$\geq \|\mathbf{W}_{vec} - \mathbf{W}_{*,vec}\|_2^2 m_{min} - 5\|\mathbf{W}_{*,vec}\|_2 \|\mathbf{W}_{vec} - \mathbf{W}_{*,vec}\|_2 m_{max}^{ex} \phi$$

Since $\phi \leq 2\sin\phi$ for $0 \leq \phi \leq 2/\pi$, $\sin\phi = \|\mathbf{W}_{vec} - \mathbf{W}_{*,vec}\| / \|\mathbf{W}_{*,vec}\|_2$

$$\begin{aligned}
&\geq \|\mathbf{W}_{vec} - \mathbf{W}_{*,vec}\|_2^2 m_{min} - 10\|\mathbf{W}_{*,vec}\|_2 \|\mathbf{W}_{vec} - \mathbf{W}_{*,vec}\|_2 m_{max}^{ex} \sin(\phi) \\
&\geq \|\mathbf{W}_{vec} - \mathbf{W}_{*,vec}\|_2^2 m_{min} - 10\|\mathbf{W}_{vec} - \mathbf{W}_{*,vec}\|_2^2 m_{max}^{ex} \\
&= \|\mathbf{W}_{vec} - \mathbf{W}_{*,vec}\|_2^2 (m_{min} - 10m_{max}^{ex})
\end{aligned}$$

In the same way,

$$\begin{aligned}
& \left(\sum \left\| \mathbb{E}[\nabla_{w_i} \ell(\mathbf{W}^t, z)] \right\|_2^2 \right)^{0.5} \\
& \leq \|\mathbf{B}_{0,0,0}\|_{oper} \|\mathbf{W}_{vec} - \mathbf{W}_{*,vec}\|_2 + \|\mathbf{B}_{1,0,0}\|_{oper} \|\mathbf{W}_{vec} - \mathbf{W}_{*,vec}\|_2 \\
& + \|\mathbf{B}_{0,0,1}\|_{oper} \|\mathbf{W}_{vec}\|_2 + \|\mathbf{B}_{1,0,1}\|_{oper} \|\mathbf{W}_{vec}\|_2 \\
& + \|\mathbf{B}_{0,1,0}\|_{oper} \|\mathbf{W}_{vec,*}\|_2 + \|\mathbf{B}_{1,1,0}\|_{oper} \|\mathbf{W}_{vec,*}\|_2 \\
& \leq L \|\mathbf{W}_{vec} - \mathbf{W}_{*,vec}\|_2 + \|\mathbf{B}_{1,0,0}\|_{oper} 3 \|\mathbf{W}_{*,vec}\|_2 \\
& + \|\mathbf{B}_{0,0,1}\|_{oper} 2 \|\mathbf{W}_{vec,*}\|_2 + \|\mathbf{B}_{1,0,1}\|_{oper} 2 \|\mathbf{W}_{vec,*}\|_2 \\
& + \|\mathbf{B}_{0,1,0}\|_{oper} \|\mathbf{W}_{vec,*}\|_2 + \|\mathbf{B}_{1,1,0}\|_{oper} \|\mathbf{W}_{vec,*}\|_2 \\
& \leq \|\mathbf{W}_{vec} - \mathbf{W}_{*,vec}\|_2 (L + 4\alpha + 14m_{max}^{ex})
\end{aligned}$$

If η_t is small enough,

$$\begin{aligned}
\|\mathbf{W}_{vec}^{t+1} - \mathbf{W}_{vec,*}\|_2^2 & \leq \|\mathbf{W}_{vec}^t - \mathbf{W}_{*,vec}\|_2^2 (1 - \eta_t(m_{min} - 10m_{max}^{ex}) + \eta_t^2(L + 4\alpha + 14m_{max}^{ex})) \\
& \leq \|\mathbf{W}_{vec}^t - \mathbf{W}_{*,vec}\|_2^2 \left(1 - \frac{\eta_t(m_{min} - 10m_{max}^{ex})}{2}\right)
\end{aligned}$$

□

국문초록

최근 몇년간 딥러닝은 여러 분야에서 최고의 성능을 보여줬다. 그러한 경험적인 성공에도 불구하고 확률적 경사하강법이 왜 최적화 문제에서 높은 성능을 보여주는지에 대해 이론적인 설명이 아직 충분하지 않다. 이 논문에서는 확률적 경사하강법이 각 파라미터들을 학습과정에서 수렴시키는 과정에 대하여 분석하였다. 이 논문에서 사용된 딥러닝 모델은 ReLU 함수를 사용한 2층 신경망이다. 특히 입력 분포의 smoothness 개념을 사용하여 특정 분포가 아닌 일반적인 입력분포에 대해 적용 가능한 분석을 이끌어낼 수 있었다.

주요어휘: 최적화, 입력분포, 신경망, 딥러닝

학번: 2016-20235