

# Composite Optimization Algorithms for Sigmoid Networks

**Huixiong Chen**

*School of Mathematical Sciences  
South China Normal University  
Guangzhou 510631, China*

HXCHEM@M.SCNU.EDU.CN

**Qi Ye**

*School of Mathematical Sciences  
South China Normal University  
Guangzhou 510631, China*

YEQI@M.SCNU.EDU.CN

## Abstract

In this paper, we use composite optimization algorithms to solve sigmoid networks. We equivalently transfer the sigmoid networks to a convex composite optimization and propose the composite optimization algorithms based on the linearized proximal algorithms and the alternating direction method of multipliers. Under the assumptions of the weak sharp minima and the regularity condition, the algorithm is guaranteed to converge to a globally optimal solution of the objective function even in the case of non-convex and non-smooth problems. Furthermore, the convergence results can be directly related to the amount of training data and provide a general guide for setting the size of sigmoid networks. Numerical experiments on Franke's function fitting and handwritten digit recognition show that the proposed algorithms perform satisfactorily and robustly.

**Keywords:** sigmoid network, composite optimization, non-convex non-smooth algorithm, global convergence, adaptive network size

## 1. Introduction

The neural network is an important and popular branch of machine learning. People have already developed many useful and well-studied neural network models, such as artificial neural networks, convolutional neural networks, recurrent neural networks, and deep neural networks. Neural networks have been widely used in pattern recognition, image processing, computer vision, neuroinformatics, bioinformatics, and other various fields with great success (LeCun et al. 2015; Abiodun et al. 2018).

When the neural networks are used in practical tasks, they are commonly trained by the error BackPropagation (BP) algorithm which is the most distinguished and successful neural network learning algorithm up to now. The BP algorithm is based on the gradient descent strategy that updates the parameters to the negative gradient direction of the target. To accelerate the learning process, stochastic gradient descent (SGD) with momentum and adaptive methods including adaptive gradient (AdaGrad), root mean square prop (RMSProp), adaptive moment estimation (Adam), and so on have emerged one after another and made a huge impact. As we all know, most of these first-order methods can converge to the critical point only if the objective function is convex or smooth. But for non-convex and non-smooth functions, it remains ambiguous how to find the convergence to even first-

or second-order critical points (Burke et al. 2005). Typical cases are sigmoid networks with absolute or hinge loss functions. The BP algorithm can solve these non-convex and non-smooth problems as well, but they are not consistent with the convergence properties of the algorithm. Moreover, it is still non-trivial to find globally optimal solutions for traditional neural network algorithms. We take the state-of-the-art Adam as an example. Its theory is poorly understood in the literature, and it suffers from several deficiencies. For instance, Adam may miss globally optimal solutions (Wilson et al. 2017), and it can be shown that it does not converge on some simple test problems (Reddi et al. 2018).

In this paper, we use composite optimization algorithms to solve sigmoid networks; see Algorithms 1 and 2 for details. The algorithm is guaranteed to (even globally) converge to a globally optimal solution of the objective function even in the case of non-convex and non-smooth problems. That is the main contribution of this paper. The start of our work stems from the finding that sigmoid networks (2.1) can be equivalently transformed into a convex composite optimization (2.2), where the inner function is smooth and the outer function is convex. This provides a new perspective on sigmoid networks. In fact, composite optimization problems arise in many applications in engineering, such as compressed sensing, image processing, machine learning, and artificial intelligence (Boyd et al. 2011; Hong et al. 2017). The composite optimization is an area at the cutting edge of mathematical optimization, and how to efficiently solve composite optimization problems has been a popular subject. For the sigmoid networks with the structure (2.2), the traditional first-order methods do not take advantage of the convex property of the outer function, so sometimes they have certain limitations in practical applications. However, composite optimization methods can fully exploit the information in the structure for algorithm design. There are many iterative algorithms with theoretical foundations for the optimization (2.2), such as the famous Gauss-Newton method (GNM, Burke and Ferris 1995), the proximal descent algorithm (ProxDescent, Lewis and Wright 2016), and the linearized proximal algorithms (LPA, Hu et al. 2016). The basic idea of these algorithms is to transfer a complex optimization problem to a sequence of simple optimization problems whose optimal solutions are easy to compute or have explicit formulas. The LPA is one of the most advanced algorithms in convex composite optimization. It can transform a non-convex and possibly non-smooth problem into a series of unconstrained strongly convex optimization subproblems, which has an attractive computational advantage. The LPA has also been applied to sensor network localization, gene regulatory network inference, and other engineering problems with great success (Hu et al. 2016, 2020; Wang et al. 2017). Therefore, we use the LPA to solve sigmoid networks in this paper.

Under the assumptions of the weak sharp minima and the regularity condition, we establish the convergence behavior of the algorithms for sigmoid networks; see Theorems 3 and 5 for details. Furthermore, we prove that the weak sharp minima is often satisfied for sigmoid networks, and the full row rank of the Jacobian matrix of the inner function, namely  $\text{rank}(F'(\bar{\theta})) = m$ , where  $m$  is the amount of training data, is a sufficient condition of the regularity condition. Hence the convergence results can be directly related to the amount of training data; see Corollaries 4 and 6 for details. This conclusion is of great theoretical and applied significance, especially since it can provide a general guide for setting the size of sigmoid networks. By the full row rank of the Jacobian matrix, we obtain a lower bound on the network size in Corollary 8. We call this lower bound the “adaptive network

size”. In this paper, our numerical experiments verify that the adaptive network size is sufficient to construct an ideal sigmoid network that solves the problem effectively. Hence Corollary 8 does provide a good guide for setting the size of sigmoid networks. The essence is to guarantee that the number of parameters in neural networks is not smaller than the amount of training data, and that a sufficient number of parameters ensure the feasibility of the networks. It can also serve as a general guide for setting the size of neural networks. That is another contribution of this paper.

Our work is also motivated by the lack of convex composite optimization algorithms and related software packages for neural networks. To the best of our knowledge, the introduction of convex composite optimization into the area of neural networks has not been addressed in the literature before. This paper is the first piece of work combining neural networks and convex composite optimization.

We organize the paper as follows. In section 2, we introduce the three-layer sigmoid networks and transfer the problem to a convex composite optimization. In section 3, we use the LPA-type algorithms to solve sigmoid networks and employ the alternating direction method of multipliers (ADMM) to solve the non-smooth convex subproblems. In section 4, we prove some convergence properties of the proposed algorithms. In section 5, the numerical experiments are demonstrated including Franke’s function fitting and handwritten digit recognition. Finally, we conclude with an outlook in section 6.

## 2. Sigmoid Networks

To begin with, we introduce the two-layer real-output sigmoid network, which is known as ‘universal approximators’ (Anthony and Bartlett 1999). Using the standard sigmoid function  $\sigma : \mathbb{R} \rightarrow (0, 1)$  of the form

$$\sigma(a) = \frac{1}{1 + e^{-a}},$$

the sigmoid network computes a function  $f : \mathbb{R}^d \rightarrow \mathbb{R}$  of the form

$$f(\mathbf{x}) = \sum_{i=1}^q w_i \sigma(\mathbf{v}_i \cdot \mathbf{x} + u_i) + w_0,$$

where  $w_i \in \mathbb{R}$  ( $i = 0, 1, \dots, q$ ) are the output weights,  $\mathbf{v}_i \in \mathbb{R}^d$  and  $u_i \in \mathbb{R}$  ( $i = 1, 2, \dots, q$ ) are the input weights. We define these adjustable parameters by

$$\boldsymbol{\theta} = (\mathbf{w}, \mathbf{v}, \mathbf{u}, w_0)^T \in \mathbb{R}^n,$$

where  $\mathbf{w} = (w_1, w_2, \dots, w_q)$ ,  $\mathbf{v} = (\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_q)$ ,  $\mathbf{u} = (u_1, u_2, \dots, u_q)$ . In the following paragraphs, we replace  $f(\mathbf{x})$  with  $f(\mathbf{x}; \boldsymbol{\theta})$ . Given a training dataset

$$D = \{(\mathbf{x}_1, y_1), (\mathbf{x}_2, y_2), \dots, (\mathbf{x}_m, y_m) \mid \mathbf{x}_i \in \mathbb{R}^d, y_i \in \mathbb{R}\},$$

the goal of using this network for a supervised learning problem is to find parameters that minimize some measure of the error of the network output over the training dataset, that is,

$$\min_{\boldsymbol{\theta} \in \mathbb{R}^n} E(\boldsymbol{\theta}) := \frac{1}{m} \sum_{i=1}^m L(f(\mathbf{x}_i; \boldsymbol{\theta}), y_i), \quad (2.1)$$

where  $L : \mathbb{R} \times \mathbb{R} \rightarrow [0, \infty)$  is a loss function. To simplify the discussions, we focus on three convex loss functions including the quadratic loss function,  $(f(\mathbf{x}; \boldsymbol{\theta}), y) \mapsto (f(\mathbf{x}; \boldsymbol{\theta}) - y)^2$ , the absolute loss function,  $(f(\mathbf{x}; \boldsymbol{\theta}), y) \mapsto |f(\mathbf{x}; \boldsymbol{\theta}) - y|$ , and the hinge loss function,  $(f(\mathbf{x}; \boldsymbol{\theta}), y) \mapsto (1 - yf(\mathbf{x}; \boldsymbol{\theta}))_+$ .

The model of the sigmoid network is usually non-convex and non-smooth. Interestingly, we discover that this problem can be seen as a convex composite optimization problem of the form

$$\min_{\boldsymbol{\theta} \in \mathbb{R}^n} E(\boldsymbol{\theta}) = \mathbb{L}(F(\boldsymbol{\theta})), \quad (2.2)$$

where the inner function  $F : \mathbb{R}^n \rightarrow \mathbb{R}^m$  is smooth, and the outer function  $\mathbb{L} : \mathbb{R}^m \rightarrow \mathbb{R}$  is convex. Specifically, for the absolute or quadratic loss functions, we can set

$$F(\boldsymbol{\theta}) = \begin{pmatrix} f(\mathbf{x}_1; \boldsymbol{\theta}) - y_1 \\ f(\mathbf{x}_2; \boldsymbol{\theta}) - y_2 \\ \vdots \\ f(\mathbf{x}_m; \boldsymbol{\theta}) - y_m \end{pmatrix}, \quad \mathbb{L}(\mathbf{z}) = \frac{1}{m} \|\mathbf{z}\|_p^p, \quad \text{where } p = 1 \text{ or } 2.$$

In general, we replace  $\|\cdot\|_2$  with  $\|\cdot\|$ . For the hinge loss function, we can set

$$F(\boldsymbol{\theta}) = \begin{pmatrix} y_1 f(\mathbf{x}_1; \boldsymbol{\theta}) \\ y_2 f(\mathbf{x}_2; \boldsymbol{\theta}) \\ \vdots \\ y_m f(\mathbf{x}_m; \boldsymbol{\theta}) \end{pmatrix}, \quad \mathbb{L}(\mathbf{z}) = \frac{1}{m} \sum_{i=1}^m (1 - z_i)_+ = \frac{1}{m} \|(\mathbf{1} - \mathbf{z})_+\|_1,$$

where  $\mathbf{z}_+$  denotes the componentwise non-negative part of  $\mathbf{z}$ . As we can see, all the outer functions are separable and have the form

$$\mathbb{L}(\mathbf{z}) = \frac{1}{m} \sum_{i=1}^m \mathbb{L}(z_i),$$

where  $\mathbb{L} : \mathbb{R} \rightarrow [0, \infty)$  is a convex function. It is the special property of  $\mathbb{L}$  in sigmoid networks.

### 3. Composite Optimization Algorithms for Sigmoid Networks

In this section, we show how to solve the sigmoid networks based on the composite optimization algorithms including the linearized proximal algorithms (LPA) and the alternating direction method of multipliers (ADMM).

**3.1 LPA for Sigmoid Networks.** The LPA is one of the most advanced algorithms in convex composite optimization. It is proposed under the inspiration of the GNM and the proximal point algorithm (PPA), and maintains the same convergence rate as that but also overcomes some of their disadvantages. Each subproblem of the LPA is constructed from a linearized approximation to the composite function and a regularization term at the current iterate. Since the subproblem is an unconstrained strongly convex optimization problem whose optimal solution is global and unique, it is easier to solve than that of the GNM. Consequently, the LPA has an attractive computational advantage, although it is generally

not a descent algorithm. Moreover, there are some connections of the LPA with other algorithms mentioned in this paper. The ProxDescent for solving (2.2) is a special case of the LPA. As the descent directions are used, the ProxDescent is a descent algorithm. The case when the inner function is simply identity mapping has a long history. The iteration  $\boldsymbol{\theta}_{k+1} = \boldsymbol{\theta}_k + \Delta\boldsymbol{\theta}_k$ , where  $\Delta\boldsymbol{\theta}_k$  minimizes the function  $\Delta\boldsymbol{\theta} \mapsto h(\boldsymbol{\theta}_k + \Delta\boldsymbol{\theta}) + \frac{1}{2t}\|\Delta\boldsymbol{\theta}\|^2$ , is the well-known PPA.

Applying the LPA directly to (2.2), we get the following algorithm for sigmoid networks.

---

**Algorithm 1.** LPA for sigmoid networks.

---

**Input:** Model  $f$ , training dataset  $D = \{(\mathbf{x}_i, y_i)\}_{i=1}^m$ , outer function  $\mathbb{L}$ , inner function  $F$ .

- 1: **Initialization:**  $t > 0$ ,  $\boldsymbol{\theta}_0 \in \mathbb{R}^n$ ,  $k \leftarrow 0$ , accept  $\leftarrow$  false;
- 2: **while** not accept **do**
- 3:   calculate the search direction

$$\Delta\boldsymbol{\theta}_k := \arg \min_{\Delta\boldsymbol{\theta} \in \mathbb{R}^n} \left\{ \mathbb{L}(F(\boldsymbol{\theta}_k) + F'(\boldsymbol{\theta}_k)\Delta\boldsymbol{\theta}) + \frac{1}{2t}\|\Delta\boldsymbol{\theta}\|^2 \right\}, \quad (3.1)$$

where  $F'(\boldsymbol{\theta}) = (\nabla_{\boldsymbol{\theta}}^T f(\mathbf{x}_i; \boldsymbol{\theta}))_{i=1}^m \in \mathbb{R}^{m \times n}$  is the Jacobian matrix of  $F(\boldsymbol{\theta})$ ;

- 4:   **if**  $\Delta\boldsymbol{\theta}_k = \mathbf{0}$  **then**
- 5:     accept  $\leftarrow$  true;
- 6:   **end if**
- 7:    $\boldsymbol{\theta}_{k+1} \leftarrow \boldsymbol{\theta}_k + \Delta\boldsymbol{\theta}_k$ ;
- 8:    $k \leftarrow k + 1$ ;
- 9: **end while**
- 10:  $\boldsymbol{\theta}^* \leftarrow \boldsymbol{\theta}_k$ .

**Output:**  $f(\mathbf{x}; \boldsymbol{\theta}^*)$ .

---

The focus of Algorithm 1 is how to solve the subproblem (3.1) accurately and efficiently. Now, we discuss some numerical algorithms for the special loss functions.

For the **quadratic loss function**, Algorithm 1 is reduced to the well-known Levenberg-Marquardt method for solving the following nonlinear least squares problem of the form

$$\min_{\boldsymbol{\theta} \in \mathbb{R}^n} E(\boldsymbol{\theta}) = \frac{1}{m}\|F(\boldsymbol{\theta})\|^2.$$

The **smooth convex subproblem** can be written as

$$\min_{\Delta\boldsymbol{\theta} \in \mathbb{R}^n} \frac{1}{m}\|F(\boldsymbol{\theta}_k) + F'(\boldsymbol{\theta}_k)\Delta\boldsymbol{\theta}\|^2 + \frac{1}{2t}\|\Delta\boldsymbol{\theta}\|^2,$$

and its necessary and sufficient optimality conditions imply that

$$\frac{2}{m}F'(\boldsymbol{\theta}_k)^T (F(\boldsymbol{\theta}_k) + F'(\boldsymbol{\theta}_k)\Delta\boldsymbol{\theta}_k) + \frac{1}{t}\Delta\boldsymbol{\theta}_k = \mathbf{0}.$$

Hence the closed formula of the search direction is given by

$$\Delta\boldsymbol{\theta}_k = - \left( \frac{2}{m}F'(\boldsymbol{\theta}_k)^T F'(\boldsymbol{\theta}_k) + \frac{1}{t}\mathbf{I} \right)^{-1} \left( \frac{2}{m}F'(\boldsymbol{\theta}_k)^T F(\boldsymbol{\theta}_k) \right) \triangleq -B_k^{-1}\nabla E(\boldsymbol{\theta}_k),$$

where  $B_k$  is always a positive-definite and invertible matrix, and  $\nabla E(\boldsymbol{\theta})$  is the gradient of the objective function in the original problem. Thus, the iteration  $\boldsymbol{\theta}_{k+1} = \boldsymbol{\theta}_k - B_k^{-1} \nabla E(\boldsymbol{\theta}_k)$  can be regarded as a variant of gradient descent algorithm, where  $B_k^{-1}$  is an adaptive learning rate. Moreover, the stopping criterion  $\Delta\boldsymbol{\theta}_k = \mathbf{0}$  shows that  $\nabla E(\boldsymbol{\theta}_k) = \mathbf{0}$ , which is the first-order necessary condition of the original problem. In section 4, we will give a first-order sufficient condition of the original problem in Theorem 7.

**3.2 ADMM for Non-Smooth Convex Subproblems.** For the non-smooth convex loss functions, the subproblem of Algorithm 1 is more complex, but luckily it is convex. There are many widely used convex optimization methods and heuristic algorithms to solve it, such as gradient or subgradient methods, approximation or composite optimization methods (Bertsekas 2015), and simulated annealing algorithms. Moreover, there are many related software packages to implement these algorithms, such as CVXPY in Python, qpOASES in C++, and CVX toolbox in Matlab. So it is not difficult to calculate the search direction from the subproblem. We use a mapping  $A$  to represent a specific algorithm to solve the subproblem, then the search direction can be presented in

$$\Delta\boldsymbol{\theta}_k = A(\mathbb{L}, F, F', t, \boldsymbol{\theta}_k).$$

Here we use the ADMM to solve the subproblem. The ADMM is a simple scheme that often works well and has a good reliability with a wide range of applications, especially for convex problems. It is also easy to understand and implement for many composite optimization problems with complex structures (Boyd et al. 2011).

The subproblem (3.1) can be seen as an equivalent problem of the form

$$\begin{aligned} \min_{\boldsymbol{\mu} \in \mathbb{R}^m, \Delta\boldsymbol{\theta} \in \mathbb{R}^n} \quad & \mathbb{L}(\boldsymbol{\mu}) + \frac{1}{2t} \|\Delta\boldsymbol{\theta}\|^2, \\ \text{s.t.} \quad & \boldsymbol{\mu} - F(\boldsymbol{\theta}_k) - F'(\boldsymbol{\theta}_k)\Delta\boldsymbol{\theta} = \mathbf{0}. \end{aligned}$$

The augmented Lagrangian function of the above problem is

$$\mathcal{L}_\rho(\boldsymbol{\mu}, \Delta\boldsymbol{\theta}, \boldsymbol{\lambda}) = \mathbb{L}(\boldsymbol{\mu}) + \frac{1}{2t} \|\Delta\boldsymbol{\theta}\|^2 + \boldsymbol{\lambda}^T (\boldsymbol{\mu} - F(\boldsymbol{\theta}_k) - F'(\boldsymbol{\theta}_k)\Delta\boldsymbol{\theta}) + \frac{\rho}{2} \|\boldsymbol{\mu} - F(\boldsymbol{\theta}_k) - F'(\boldsymbol{\theta}_k)\Delta\boldsymbol{\theta}\|^2,$$

where  $\rho > 0$  is the penalty parameter. The ADMM consists of the iterations

$$\begin{aligned} \boldsymbol{\mu}^{i+1} &:= \arg \min_{\boldsymbol{\mu} \in \mathbb{R}^m} \mathcal{L}_\rho(\boldsymbol{\mu}, \Delta\boldsymbol{\theta}^i, \boldsymbol{\lambda}^i), \\ \Delta\boldsymbol{\theta}^{i+1} &:= \arg \min_{\Delta\boldsymbol{\theta} \in \mathbb{R}^n} \mathcal{L}_\rho(\boldsymbol{\mu}^{i+1}, \Delta\boldsymbol{\theta}, \boldsymbol{\lambda}^i), \\ \boldsymbol{\lambda}^{i+1} &:= \boldsymbol{\lambda}^i + \rho (\boldsymbol{\mu}^{i+1} - F(\boldsymbol{\theta}_k) - F'(\boldsymbol{\theta}_k)\Delta\boldsymbol{\theta}^{i+1}). \end{aligned} \tag{3.2}$$

The calculation for  $\boldsymbol{\mu}^{i+1}$  is as follows.

$$\begin{aligned}
\boldsymbol{\mu}^{i+1} &= \arg \min_{\boldsymbol{\mu} \in \mathbb{R}^m} \left\{ \mathbb{L}(\boldsymbol{\mu}) + \frac{\rho}{2} \|\boldsymbol{\mu} - F(\boldsymbol{\theta}_k) - F'(\boldsymbol{\theta}_k)\Delta\boldsymbol{\theta}^i + \frac{1}{\rho}\boldsymbol{\lambda}^i\|^2 \right\}, \\
&= \arg \min_{\boldsymbol{\mu} \in \mathbb{R}^m} \left\{ \mathbb{L}(\boldsymbol{\mu}) + \frac{\rho}{2} \|\boldsymbol{\mu} - \mathbf{a}^i\|^2 \right\}, \\
&= \arg \min_{\boldsymbol{\mu} \in \mathbb{R}^m} \left\{ \sum_{j=1}^m \left( \frac{1}{m} \mathbb{L}(\mu_j) + \frac{\rho}{2} (\mu_j - a_j^i)^2 \right) \right\}, \\
&= \left( \arg \min_{\mu_j \in \mathbb{R}} \left\{ \frac{1}{m} \mathbb{L}(\mu_j) + \frac{\rho}{2} (\mu_j - a_j^i)^2 \right\} \right)_{j=1}^m, \\
&= (\Phi_{1/m\rho}(a_j^i))_{j=1}^m, \tag{3.3}
\end{aligned}$$

where  $\mathbf{a}^i = F(\boldsymbol{\theta}_k) + F'(\boldsymbol{\theta}_k)\Delta\boldsymbol{\theta}^i - \frac{1}{\rho}\boldsymbol{\lambda}^i$ , and  $\Phi_\kappa$  is the proximity operator of  $\mathbb{L}$  with the penalty  $\frac{1}{\kappa}$  (Boyd et al. 2011). Specifically, for the **absolute loss function**, the proximity operator  $\Phi$ , also called the soft thresholding operator, is defined as

$$\Phi_\kappa(a) = \begin{cases} a - \kappa, & a > \kappa, \\ 0, & |a| \leq \kappa, \\ a + \kappa, & a < -\kappa. \end{cases}$$

For the **hinge loss function**, the proximity operator  $\Phi$  is defined as

$$\Phi_\kappa(a) = \begin{cases} a, & a > 1, \\ 1, & 1 - \kappa \leq a \leq 1, \\ a + \kappa, & a < 1 - \kappa. \end{cases}$$

The calculation for  $\Delta\boldsymbol{\theta}^{i+1}$  is as follows. Since

$$\Delta\boldsymbol{\theta}^{i+1} = \arg \min_{\Delta\boldsymbol{\theta} \in \mathbb{R}^n} \left\{ \frac{1}{2t} \|\Delta\boldsymbol{\theta}\|^2 + \frac{\rho}{2} \|\boldsymbol{\mu}^{i+1} - F(\boldsymbol{\theta}_k) - F'(\boldsymbol{\theta}_k)\Delta\boldsymbol{\theta} + \frac{1}{\rho}\boldsymbol{\lambda}^i\|^2 \right\},$$

by its necessary and sufficient optimality conditions, we obtain that

$$\Delta\boldsymbol{\theta}^{i+1} = \left( \rho F'(\boldsymbol{\theta}_k)^T F'(\boldsymbol{\theta}_k) + \frac{1}{t} \mathbf{I} \right)^{-1} \left( \rho F'(\boldsymbol{\theta}_k)^T \left( \boldsymbol{\mu}^{i+1} - F(\boldsymbol{\theta}_k) + \frac{1}{\rho}\boldsymbol{\lambda}^i \right) \right). \tag{3.4}$$

As we can see, the iterations of the ADMM for the non-smooth convex subproblems have explicit formulas, which is one of the advantages of the ADMM. Defining the primal residual of the optimality conditions at iteration  $i$  as

$$\mathbf{r}^i = \boldsymbol{\mu}^i - F(\boldsymbol{\theta}_k) - F'(\boldsymbol{\theta}_k)\Delta\boldsymbol{\theta}^i, \tag{3.5}$$

and the dual residual at iteration  $i$  as

$$\mathbf{s}^i = \rho F'(\boldsymbol{\theta}_k)(\Delta\boldsymbol{\theta}^i - \Delta\boldsymbol{\theta}^{i-1}), \tag{3.6}$$

we set the stopping criterion as  $\|\mathbf{r}^i\| \approx 0$  and  $\|\mathbf{s}^i\| \approx 0$ .

---

**Algorithm A\***. ADMM for non-smooth convex subproblems.

---

**Input:** Numbers  $m$  and  $t$ , matrices  $F(\boldsymbol{\theta}_k)$  and  $F'(\boldsymbol{\theta}_k)$ , non-smooth convex function  $\mathbb{L}$ .

- 1: **Initialization:**  $\rho > 0$ ,  $\Delta\boldsymbol{\theta}^0 \in \mathbb{R}^n$ ,  $\boldsymbol{\lambda}^0 \in \mathbb{R}^m$ ,  $\epsilon > 0$ ,  $i \leftarrow 0$ ;
- 2: **repeat**
- 3:    $i \leftarrow i + 1$ ;
- 4:   calculate  $\boldsymbol{\mu}^i$  from (3.3);
- 5:   calculate  $\Delta\boldsymbol{\theta}^i$  from (3.4);
- 6:   calculate  $\boldsymbol{\lambda}^i$  from (3.2);
- 7:   calculate  $\boldsymbol{r}^i$  and  $\boldsymbol{s}^i$  from (3.5) and (3.6), respectively;
- 8: **until**  $\|\boldsymbol{r}^i\| < \epsilon$  and  $\|\boldsymbol{s}^i\| < \epsilon$ ;
- 9:  $\Delta\boldsymbol{\theta}_k \leftarrow \Delta\boldsymbol{\theta}^i$ .

**Output:**  $\Delta\boldsymbol{\theta}_k$ .

---

**3.3 A Globalization Strategy for Algorithm 1.** Moreover, we show the following algorithm by employing the globalized LPA (GLPA) that adopts a backtracking line-search as a globalization strategy. The choice of the stepsize is based on the virtue of the backtracking line-search, which guarantees the monotone decrease of the objective function at each iteration. As a result, it ensures that the GLPA is a descent algorithm. In the algorithm implementation, the backtracking strategy finds the first point satisfying the inequality (3.7) by continuously decreasing the trial stepsize in an exponential way. That makes the stepsize with the descent property as large as possible.

---

**Algorithm 2.** GLPA for sigmoid networks.

---

**Input:** Model  $f$ , training dataset  $D = \{(\boldsymbol{x}_i, y_i)\}_{i=1}^m$ , outer function  $\mathbb{L}$ , inner function  $F$ .

- 1: **Initialization:**  $t > 0$ ,  $c, \tau \in (0, 1)$ ,  $\boldsymbol{\theta}_0 \in \mathbb{R}^n$ ,  $k \leftarrow 0$ ,  $\text{accept} \leftarrow \text{false}$ ;
- 2: **while** not  $\text{accept}$  **do**
- 3:   calculate the search direction

$$\Delta\boldsymbol{\theta}_k = \arg \min_{\Delta\boldsymbol{\theta} \in \mathbb{R}^n} \left\{ \mathbb{L}(F(\boldsymbol{\theta}_k) + F'(\boldsymbol{\theta}_k)\Delta\boldsymbol{\theta}) + \frac{1}{2t} \|\Delta\boldsymbol{\theta}\|^2 \right\};$$

- 4:   **if**  $\Delta\boldsymbol{\theta}_k = \mathbf{0}$  **then**
- 5:      $\text{accept} \leftarrow \text{true}$ ;
- 6:   **end if**
- 7:    $\eta \leftarrow 1/\tau$ ;
- 8:   **repeat**
- 9:      $\eta \leftarrow \tau\eta$
- 10: **until**

$$\mathbb{L}(F(\boldsymbol{\theta}_k + \eta\Delta\boldsymbol{\theta}_k)) - \mathbb{L}(F(\boldsymbol{\theta}_k)) \leq c\eta \left( \mathbb{L}(F(\boldsymbol{\theta}_k) + F'(\boldsymbol{\theta}_k)\Delta\boldsymbol{\theta}_k) + \frac{1}{2t} \|\Delta\boldsymbol{\theta}_k\|^2 - \mathbb{L}(F(\boldsymbol{\theta}_k)) \right); \quad (3.7)$$

- 11:    $\eta_k \leftarrow \eta$ ;
- 12:    $\boldsymbol{\theta}_{k+1} \leftarrow \boldsymbol{\theta}_k + \eta_k\Delta\boldsymbol{\theta}_k$ ;
- 13:    $k \leftarrow k + 1$ ;

14: **end while**  
 15:  $\boldsymbol{\theta}^* \leftarrow \boldsymbol{\theta}_k$ .  
**Output:**  $f(\boldsymbol{x}; \boldsymbol{\theta}^*)$ .

---

#### 4. Convergence Analysis

In this section, we prove some convergence properties of the proposed algorithms under the assumptions of the weak sharp minima and the regularity condition or full row rank of the Jacobian matrix, a stronger condition. Before giving the main results, we introduce the following useful definitions and lemmas.

**4.1 Theoretical Foundations of LPA-type Algorithms.** Here we consider the convex composite optimization of the form

$$\min_{\boldsymbol{\omega} \in \mathbb{R}^b} h(G(\boldsymbol{\omega})), \quad (4.1)$$

where the inner function  $G : \mathbb{R}^b \rightarrow \mathbb{R}^l$  is continuously differentiable, and the outer function  $h : \mathbb{R}^l \rightarrow \mathbb{R}$  is convex. It is a more general mathematical form of the problem (2.2).

First, we introduce the concept of the Lipschitz continuous gradient, which has played an important role in investigating the convergence behavior of many optimization algorithms. For a differentiable function  $G$  and  $\Omega \subseteq \mathbb{R}^b$ , if there exists an  $K > 0$  such that

$$\|G'(\tilde{\boldsymbol{\omega}}_1) - G'(\tilde{\boldsymbol{\omega}}_2)\| \leq K \|\tilde{\boldsymbol{\omega}}_1 - \tilde{\boldsymbol{\omega}}_2\| \text{ for each } \tilde{\boldsymbol{\omega}}_1, \tilde{\boldsymbol{\omega}}_2 \in \Omega,$$

we say that  $G$  is  $K$ -smooth or has a *Lipschitz continuous gradient* with modulus  $K$  on  $\Omega$ .

Next, we give the notion of the weak sharp minima introduced in (Burke and Ferris 1993), which has far-reaching consequences for the convergence analysis of many iterative procedures. For a function  $h$ , the minimum value and the set of minima for  $h$ , denoted by  $h_{\min}$  and  $C_h$ , are defined by

$$h_{\min} := \min_{\boldsymbol{z} \in \mathbb{R}^l} h(\boldsymbol{z}) \quad \text{and} \quad C_h := \arg \min_{\boldsymbol{z} \in \mathbb{R}^l} h(\boldsymbol{z}).$$

Let  $C_h \subseteq S \subseteq \mathbb{R}^l$ , if there exist  $\alpha > 0$  and  $\beta \geq 1$  such that

$$h(\boldsymbol{z}) \geq h_{\min} + \alpha \text{dist}^\beta(\boldsymbol{z}, C_h) \text{ for each } \boldsymbol{z} \in S,$$

where  $\text{dist}(\boldsymbol{z}, C) := \inf_{\boldsymbol{c} \in C} \|\boldsymbol{z} - \boldsymbol{c}\|$ , then we say that  $C_h$  is the set of *weak sharp minima* of order  $\beta$  for  $h$  on  $S$  with modulus  $\alpha$ .

We now introduce the regularity condition proposed in (Burke and Ferris 1995), which is a crucial assumption applied to establish the convergence of several convex composite optimization algorithms. Let  $h$  and  $G$  be defined by (4.1), then a point  $\bar{\boldsymbol{\omega}} \in \mathbb{R}^b$  is said to be a *regular point* of the inclusion  $G(\boldsymbol{\omega}) \in C_h$  if

$$\ker(G'(\bar{\boldsymbol{\omega}})^T) \cap (C_h - G(\bar{\boldsymbol{\omega}}))^\ominus = \{\mathbf{0}\},$$

where  $\ker(W) := \{\boldsymbol{y} : W\boldsymbol{y} = \mathbf{0}\}$  is the nullspace of  $W$ , and  $Z^\ominus := \{\boldsymbol{y} : \langle \boldsymbol{y}, \boldsymbol{z} \rangle \leq 0, \forall \boldsymbol{z} \in Z\}$  is the negative polar of  $Z$ .

In the following lemmas, we give the local convergence of the LPA and the global convergence of the GLPA for solving optimization (4.1). They are based on three main conditions including Lipschitz continuous gradient, weak sharp minima and quasi-regularity or regularity condition. Note that the definition of quasi-regularity condition will only be described in the proof of Theorem 3. Since this condition is hard to verify in practice, we replace it with the regularity condition in the related theorem for sigmoid networks.

**Lemma 1.** (Hu et al. 2016, Corollary 14) *Let  $\bar{\omega} \in \mathbb{R}^b$  satisfy  $G(\bar{\omega}) \in C_h$ , and let  $C_h$  be the set of weak sharp minima of order  $\beta$  for  $h$  near  $G(\bar{\omega})$  with constant  $\alpha$ . Suppose that  $G$  is continuously differentiable with a Lipschitz continuous gradient  $G'$  near  $\bar{\omega}$ , and that  $\bar{\omega}$  is a quasi-regular point of the inclusion with constant  $\delta$ . Suppose further that  $\beta \in [1, 2)$  or the stepsize  $t > \frac{2\delta^2}{\alpha}$  (if  $\beta = 2$ ). Then there exists a neighborhood  $N(\bar{\omega})$  of  $\bar{\omega}$  such that, for any  $\omega_0 \in N(\bar{\omega})$ , the sequence  $\{\omega_k\}$  generated by the LPA with initial point  $\omega_0$  converges at a rate of  $\frac{2}{\beta}$  to a solution  $\omega^*$  satisfying  $G(\omega^*) \in C_h$ .*

**Lemma 2.** (Hu et al. 2016, Theorem 18) *Let  $\{\omega_k\}$  be a sequence generated by the GLPA and assume that  $\{\omega_k\}$  has a cluster point  $\omega^*$ . Suppose that  $\beta \in [1, 2)$  and that  $C_h$  be the set of weak sharp minima of order  $\beta$  for  $h$  near  $G(\omega^*)$ . Suppose further that  $G$  is continuously differentiable with a Lipschitz continuous gradient  $G'$  near  $\omega^*$ , and that  $\omega^*$  is a regular point of the inclusion. Then  $G(\omega^*) \in C_h$ , and  $\{\omega_k\}$  converges to  $\omega^*$  at a rate of  $\frac{2}{\beta}$ .*

Note that  $\beta \in [1, 2)$  in Lemma 2 is lightly different from  $\beta \in [1, 2]$  in Lemma 1, but both of them can find a globally optimal solution to optimization (4.1) since that  $G(\omega^*) \in C_h$ , equivalently,  $h(G(\omega^*)) = h_{\min} = (h \circ G)_{\min}$ .

**4.2 Convergence Analysis for Sigmoid Networks.** Let  $B(z, r)$  denote an open ball of radius  $r$  centered at  $z$ , then we establish the local convergence of Algorithm 1 by virtue of Lemma 1.

**Theorem 3. (Local Convergence).** *Let  $\beta \in [1, 2]$  and  $r > 0$ . Let  $\{\theta_k\}$  be a sequence generated by Algorithm 1, and  $\bar{\theta} \in \mathbb{R}^n$  be such that  $F(\bar{\theta}) \in C_{\mathbb{L}}$  and  $C_{\mathbb{L}}$  is the set of weak sharp minima of order  $\beta$  for  $\mathbb{L}$  on  $B(F(\bar{\theta}), r)$ . If  $\bar{\theta}$  is a regular point of the inclusion, then there exist  $t_0 \geq 0$  and  $\bar{r} > 0$  such that for any  $t > t_0$  and initial point  $\theta_0 \in B(\bar{\theta}, \bar{r})$ , the sequence  $\{\theta_k\}$  converges at a rate of  $\frac{2}{\beta}$  to a globally optimal solution  $\theta^*$  and  $F(\theta^*) \in C_{\mathbb{L}}$ .*

*Proof.* According to the assumptions of Lemma 1, we need to verify the following four conditions.

- (i) *Quasi-regularity condition.* By Proposition 3.3 in (Burke and Ferris 1995), we know that any regular point of the inclusion  $F(\theta) \in C_{\mathbb{L}}$  is also a quasi-regular point. Since  $\bar{\theta}$  is a regular point,  $\bar{\theta}$  is also a quasi-regular point of the inclusion  $F(\theta) \in C_{\mathbb{L}}$ , namely there exist  $\delta > 0$  and  $r_0 > 0$  such that

$$\Pi(\theta) \neq \emptyset \text{ and } \text{dist}(\mathbf{0}, \Pi(\theta)) \leq \delta \text{dist}(F(\theta), C_{\mathbb{L}}) \text{ for each } \theta \in B(\bar{\theta}, r_0),$$

where  $\Pi(\theta) := \{\Delta\theta \in \mathbb{R}^n : F(\theta) + F'(\theta)\Delta\theta \in C_{\mathbb{L}}\}$  is the solution set of the linearized inclusion  $F(\theta) + F'(\theta)\Delta\theta \in C_{\mathbb{L}}$ .

- (ii) *Weak sharp minima.* In particular, we set  $r_0 \in (0, r)$ . Naturally,  $C_{\mathbb{L}}$  is the set of local weak sharp minima of order  $\beta$  for  $\mathbb{L}$  on  $B(F(\bar{\theta}), r_0)$  with constant  $\alpha$  for some  $\alpha > 0$ , due to the assumption and definition of the weak sharp minima.

- (iii) *Lipschitz continuous gradient.* Note that a differentiable function with a Lipschitz continuous gradient is second-order differentiable almost everywhere on  $\Omega$ . If  $G$  is a second-order differentiable function, by the differential mean value theorem, it is obvious that the  $K$ -smoothness of  $G$  is equivalent to the boundedness of  $G''$ , that is,  $\|G''(\boldsymbol{\omega})\| \leq K$  for each  $\boldsymbol{\omega} \in \Omega$ . On the other hand, since  $F$  defined by (2.2) is smooth on  $\mathbb{R}^n$ ,  $F''$  is continuous on  $\mathbb{R}^n$ . Naturally,  $F''$  is bounded on the bounded subset  $B(\bar{\boldsymbol{\theta}}, r_0)$ . Therefore,  $F$  is continuously differentiable with a Lipschitz continuous gradient  $F'$  on  $B(\bar{\boldsymbol{\theta}}, r_0)$ .
- (iv) *Large stepsize.* If  $\beta = 2$ , we set  $t_0 = \frac{2\delta^2}{\alpha}$ ; otherwise, set  $t_0 = 0$ .

Hence, Lemma 1 is applicable and the conclusion follows.  $\square$

Furthermore, we analyze the convergence properties of Algorithm 1 for the three common sigmoid networks.

**Corollary 4.** *Let  $\{\boldsymbol{\theta}_k\}$  be a sequence generated by Algorithm 1, and  $\bar{\boldsymbol{\theta}} \in \mathbb{R}^n$  be such that  $F(\bar{\boldsymbol{\theta}}) \in C_{\mathbb{L}}$ . If  $F'(\bar{\boldsymbol{\theta}})$  has full row rank, then there exists an  $\bar{r} > 0$  such that for any initial point  $\boldsymbol{\theta}_0 \in B(\bar{\boldsymbol{\theta}}, \bar{r})$ , we have*

- (i) *for the sigmoid networks with the quadratic loss function, the sequence  $\{\boldsymbol{\theta}_k\}$  linearly converges to a globally optimal solution  $\boldsymbol{\theta}^*$  and  $F(\boldsymbol{\theta}^*) = \mathbf{0}$ , if  $t$  is sufficiently large.*
- (ii) *for the sigmoid networks with the absolute loss function, the sequence  $\{\boldsymbol{\theta}_k\}$  quadratically converges to a globally optimal solution  $\boldsymbol{\theta}^*$  and  $F(\boldsymbol{\theta}^*) = \mathbf{0}$ .*
- (iii) *for the sigmoid networks with the hinge loss function, the sequence  $\{\boldsymbol{\theta}_k\}$  quadratically converges to a globally optimal solution  $\boldsymbol{\theta}^*$  and  $F(\boldsymbol{\theta}^*) \geq \mathbf{1}$ .*

*Proof.* According to the assumptions of Theorem 3, we need to verify the following two conditions.

- (a) *Regularity condition.* Since the system of linear equations  $W\mathbf{y} = \mathbf{0}$  has only zero solution if and only if the matrix  $W$  has full column rank,  $F'(\bar{\boldsymbol{\theta}})$  with full row rank is equivalent to  $\ker(F'(\bar{\boldsymbol{\theta}})^T) = \{\mathbf{0}\}$ . Then, it follows that

$$\ker(F'(\bar{\boldsymbol{\theta}})^T) \cap (C_{\mathbb{L}} - F(\bar{\boldsymbol{\theta}}))^{\ominus} = \{\mathbf{0}\}.$$

Therefore, the regularity condition is satisfied.

- (b) *Weak sharp minima.* Note that  $\mathbb{L}_{\min} = 0$ ;  $C_{\mathbb{L}} = \{\mathbf{0}\}$  for the quadratic or absolute loss functions, and  $C_{\mathbb{L}} \geq \mathbf{1}$  for the hinge loss function.
  - (i) In the case when  $\mathbb{L}(\mathbf{z}) = \frac{1}{m}\|\mathbf{z}\|^2$ ,  $\mathbb{L}(\mathbf{z}) = \mathbb{L}_{\min} + \frac{1}{m}\text{dist}^2(\mathbf{z}, C_{\mathbb{L}})$  for each  $\mathbf{z} \in \mathbb{R}^m$ . By the definition of weak sharp minima, we know that  $C_{\mathbb{L}} = \{\mathbf{0}\}$  is the set of weak sharp minima of order 2 for  $\mathbb{L}$  on  $\mathbb{R}^m$  with modulus  $\frac{1}{m}$ .
  - (ii) In the case when  $\mathbb{L}(\mathbf{z}) = \frac{1}{m}\|\mathbf{z}\|_1$ ,  $\mathbb{L}(\mathbf{z}) \geq \frac{1}{m}\|\mathbf{z}\| = \mathbb{L}_{\min} + \frac{1}{m}\text{dist}(\mathbf{z}, C_{\mathbb{L}})$  for each  $\mathbf{z} \in \mathbb{R}^m$ . In the same way, it shows that  $C_{\mathbb{L}} = \{\mathbf{0}\}$  is the set of weak sharp minima of order 1 for  $\mathbb{L}$  on  $\mathbb{R}^m$  with modulus  $\frac{1}{m}$ .
  - (iii) In the case when  $\mathbb{L}(\mathbf{z}) = \frac{1}{m}\|(\mathbf{1} - \mathbf{z})_+\|_1$ ,  $\mathbb{L}(\mathbf{z}) \geq \frac{1}{m}\|(\mathbf{1} - \mathbf{z})_+\| = \mathbb{L}_{\min} + \frac{1}{m}\text{dist}(\mathbf{z}, C_{\mathbb{L}})$  for each  $\mathbf{z} \in \mathbb{R}^m$ , which implies that  $C_{\mathbb{L}} \geq \mathbf{1}$  is the set of weak sharp minima of order 1 for  $\mathbb{L}$  on  $\mathbb{R}^m$  with modulus  $\frac{1}{m}$ . Therefore, the local weak sharp minima is satisfied for the three common sigmoid networks.

Hence, Theorem 3 is applicable and the conclusion follows.  $\square$

As we have seen, the weak sharp minima is often satisfied for sigmoid networks, and its order determines the convergence rate of the algorithm. To our surprise, a first-order algorithm even has a second-order convergence rate. In the following paragraphs, we establish the global convergence of Algorithm 2 by virtue of Lemma 2.

**Theorem 5. (Global Convergence).** *Let  $\beta \in [1, 2)$  and  $r > 0$ . Let  $\{\boldsymbol{\theta}_k\}$  be a sequence generated by Algorithm 2, and  $\{\boldsymbol{\theta}_k\}$  have a cluster point  $\boldsymbol{\theta}^*$  such that  $C_{\mathbb{L}}$  be the set of weak sharp minima of order  $\beta$  for  $\mathbb{L}$  on  $B(F(\boldsymbol{\theta}^*), r)$ . If  $\boldsymbol{\theta}^*$  is a regular point of the inclusion, then  $\{\boldsymbol{\theta}_k\}$  converges at a rate of  $\frac{2}{\beta}$  to a globally optimal solution  $\boldsymbol{\theta}^*$  and  $F(\boldsymbol{\theta}^*) \in C_{\mathbb{L}}$ .*

*Proof.* According to the assumptions of Lemma 2, we need to verify the following three conditions.

- (i) *Regularity condition.* Since the cluster point  $\boldsymbol{\theta}^*$  is a regular point of the inclusion  $F(\boldsymbol{\theta}) \in C_{\mathbb{L}}$ , the regularity condition is satisfied.
- (ii) *Weak sharp minima.* Since  $C_{\mathbb{L}}$  is the set of weak sharp minima of order  $\beta$  for  $\mathbb{L}$  on  $B(F(\boldsymbol{\theta}^*), r)$  for some  $r > 0$  and  $\beta \in [1, 2)$ , the local weak sharp minima is satisfied.
- (iii) *Lipschitz continuous gradient.* By (iii) in the proof of Theorem 3, we know that  $F$  is continuously differentiable with a Lipschitz continuous gradient  $F'$  on  $B(\boldsymbol{\theta}^*, r)$ .

Hence, Lemma 2 is applicable and the conclusion follows.  $\square$

We can see that Algorithm 2 has the same conclusion and convergence rate as Algorithm 1 under the same assumptions. Next, we show the global convergence of Algorithm 2 for two non-convex and non-smooth sigmoid networks.

**Corollary 6.** *Let  $\{\boldsymbol{\theta}_k\}$  be a sequence generated by Algorithm 2 for the sigmoid networks with absolute or hinge loss functions, and  $\{\boldsymbol{\theta}_k\}$  have a cluster point  $\boldsymbol{\theta}^*$ . If  $F'(\boldsymbol{\theta}^*)$  has full row rank, then  $\{\boldsymbol{\theta}_k\}$  quadratically converges to a globally optimal solution  $\boldsymbol{\theta}^*$  and  $F(\boldsymbol{\theta}^*) \in C_{\mathbb{L}}$ .*

*Proof.* According to the assumptions of Theorem 5, we need to verify the following two conditions.

- (a) *Regularity condition.* By (a) in the proof of Corollary 4, the full row rank of  $F'(\boldsymbol{\theta}^*)$  implies that the cluster point  $\boldsymbol{\theta}^*$  is a regular point of the inclusion. Therefore, the regularity condition is satisfied.
- (b) *Weak sharp minima.* By (b) in the proof of Corollary 4, we know that  $C_{\mathbb{L}}$  is the set of weak sharp minima of order 1 for  $\mathbb{L}$  on  $\mathbb{R}^m$  with modulus  $\frac{1}{m}$ . Therefore, the local weak sharp minima is satisfied for the two sigmoid networks.

Hence, Theorem 5 is applicable and the conclusion follows.  $\square$

Note that  $F'(\bar{\boldsymbol{\theta}})$  with full row rank, namely  $\text{rank}(F'(\bar{\boldsymbol{\theta}})) = m$ , where  $m$  is the amount of training data, is the sufficient condition of the regularity condition; and it is also the necessary condition when  $C_{\mathbb{L}}$  is a singleton set and  $F'(\bar{\boldsymbol{\theta}}) \in C_{\mathbb{L}}$ . Hence the convergence results can be directly related to the amount of training data. Next, we show the following convergence property of the LPA-type algorithms in a finite number of iterations.

**Theorem 7. (Sufficient Condition).** *If the LPA-type algorithm stops at the  $k$ th iteration with  $\text{rank}(G'(\boldsymbol{\omega}_k)) = l$ , then  $\boldsymbol{\omega}_k$  is a globally optimal solution to the convex composite optimization (4.1) and  $G(\boldsymbol{\omega}_k) \in C_h$ .*

*Proof.* Since the subproblem of the LPA-type algorithms is an unconstrained convex optimization problem, its necessary and sufficient optimality conditions imply that

$$\mathbf{0} \in G'(\boldsymbol{\omega}_k)^T \partial h(G(\boldsymbol{\omega}_k) + G'(\boldsymbol{\omega}_k)\Delta\boldsymbol{\omega}_k) + \frac{1}{t}\Delta\boldsymbol{\omega}_k \text{ for each } k,$$

where  $\partial h(\mathbf{z})$  is the subdifferential of the convex function  $h(\mathbf{z})$ . The stopping criterion  $\Delta\boldsymbol{\omega}_k = \mathbf{0}$  of the algorithms shows that

$$\mathbf{0} \in G'(\boldsymbol{\omega}_k)^T \partial h(G(\boldsymbol{\omega}_k)).$$

By  $\text{rank}(G'(\boldsymbol{\omega}_k)) = l$ , equivalently, the full column rank of  $G'(\boldsymbol{\omega}_k)^T$ , it follows that

$$\mathbf{0} \in \partial h(G(\boldsymbol{\omega}_k)).$$

By the necessary and sufficient optimality conditions of the convex optimization, it shows that  $G(\boldsymbol{\omega}_k)$  is a globally optimal solution to  $h$ , equivalently,  $G(\boldsymbol{\omega}_k) \in C_h$ . Hence the proof is complete.  $\square$

Theorem 7 also shows that  $\text{rank}(F'(\boldsymbol{\theta}_k)) = m$  is the first-order sufficient condition of sigmoid networks when the LPA-type algorithm stops at the  $k$ th iteration. It is no surprise that there is a unified conclusion on the non-convex and possibly non-smooth sigmoid networks, thanks to the unified composite optimization framework and the convex subproblem.

We have seen that the full row rank is a critical condition for the convergence analysis of sigmoid networks. This condition is of great theoretical and applied significance, especially since it can provide a general guide for setting the network size. In order to guarantee the reliability of the algorithm, we can ensure that  $F'(\boldsymbol{\theta}) \in \mathbb{R}^{m \times n}$  is of full row rank, which implies that  $n = (d + 2)q + 1 \geq m$ , where  $d$  is the dimension of the input, and  $q$  is the number of hidden neurons. So we have the following corollary.

**Corollary 8.** *If  $\text{rank}(F'(\bar{\boldsymbol{\theta}})) = m$ , then we have a lower bound on the network size given by*

$$q \geq \left\lceil \frac{m - 1}{d + 2} \right\rceil. \quad (4.2)$$

Clearly, the lower bound on the network size is directly proportional to the amount of training data and inversely proportional to the dimension of the input. That is, the lower bound on the network size is adapted to the problem size, so we call this lower bound the ‘‘adaptive network size’’. Moreover, each row of the Jacobian matrix  $F'(\boldsymbol{\theta})$  is the gradient of the fitting function  $f(\mathbf{x}; \boldsymbol{\theta})$  at the corresponding data point. In a general sense, as the number of hidden neurons increases, the information contained in the gradient increases. As a result, the rank of the Jacobian matrix will also increase or be equal to  $m$ . Thus, the full row rank of  $F'(\bar{\boldsymbol{\theta}})$  can be satisfied in a theoretical sense by choosing the network size sufficiently large. In conclusion, the LPA-type algorithms are almost always reliable.

## 5. Numerical Experiment

Sigmoid networks are often used to solve regression and classification tasks, so we shall use our algorithms for both tasks. We train the sigmoid networks on the training dataset and demonstrate the performance on the test dataset. Note that we will use the adaptive network size, namely the lower bound on the network size given by Corollary 8, to build the sigmoid networks, which is sufficient to solve problems effectively.

**5.1 Regression on Scattered Data.** Franke's function is a standard test function for 2D scattered data fitting of the form

$$g(x_1, x_2) = \frac{3}{4}e^{-1/4((9x_1-2)^2+(9x_2-2)^2)} + \frac{3}{4}e^{-(1/49)(9x_1+1)^2-(1/10)(9x_2+1)^2} \\ + \frac{1}{2}e^{-1/4((9x_1-7)^2+(9x_2-3)^2)} - \frac{1}{5}e^{-(9x_1-4)^2-(9x_2-7)^2},$$

and its graph in the unit square in  $\mathbb{R}^2$  is shown on the left of Figure 1. One can see that Franke's function is a complex function with two Gaussian peaks and a small trough. We generate 289 training data points and 121 test data points using the Halton sequence. The points are uniformly distributed in the unit square in  $\mathbb{R}^2$ , and the result is shown on the right of Figure 1.

Considering the observational errors, we also add small white Gaussian noise to the training data to reflect the real case, that is,  $y_i = g(x_i^1, x_i^2) + |\xi_i|$ , and  $\xi_i \sim N(0, \tilde{\sigma}^2)$ , where  $N(0, \tilde{\sigma}^2)$  is a Gaussian distribution with a mean of 0 and a standard deviation of  $\tilde{\sigma}$ . All numerical experiments are implemented in Python 3.9. We generate the positive Gaussian noise using  $\frac{1}{\sqrt{2\pi}\tilde{\sigma}} \cdot \text{uniform}(0, 1)$ . The performance measure we choose for the regression task is the root mean squared error (RMS-error):

$$\text{RMS-error} = \frac{1}{\sqrt{M}} \left( \sum_{i=1}^M (\tilde{y}_i - y_i)^2 \right)^{\frac{1}{2}},$$

where  $\tilde{y}_i$  is the predicted value and  $y_i$  is the actual value.

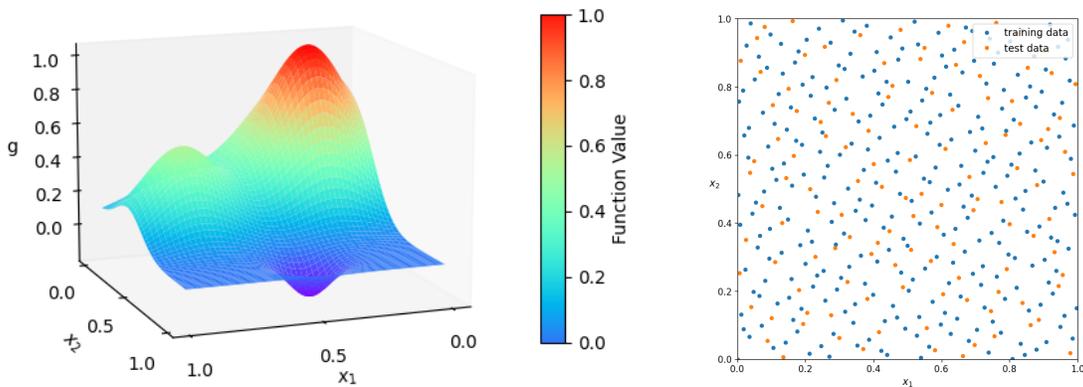


Figure 1: The graph of Franke's test function (left) and a set of 289 training data points and 121 test data points in the unit square in  $\mathbb{R}^2$  (right).

When implementing the LPA-type algorithms for the sigmoid networks with a **quadratic loss function**, we set  $\tilde{\sigma} = 100$ ,  $\theta_0 = \mathbf{0}$ , and the stopping criterion as  $\|\Delta\theta_k\| < 1e-2$ . For the inequality (3.7) in Algorithm 2, we set  $\tau = 0.5$ ,  $c = 1e-3$ , and the maximum number of iterations for the backtracking line-search as 10 (indeed, one iteration is enough in most cases, that is,  $\eta_k = 1$  is often used). According to (4.2), we can set  $q \geq 72$  to guarantee the reliability of the algorithms. For the case when  $q = 72$  and  $t = 1e5$ , the performance of the algorithms is shown in Table 1 and Figure 2.

Table 1: The performance of regression on Franke’s function (using quadratic loss).

	LPA		GLPA	
	RMS-error	Max-error	RMS-error	Max-error
No noise	2.9525e-3	1.4736e-2	2.7790e-3	1.1547e-2
Gaussian noise	3.4364e-3	1.2678e-2	3.7613e-3	1.5765e-2

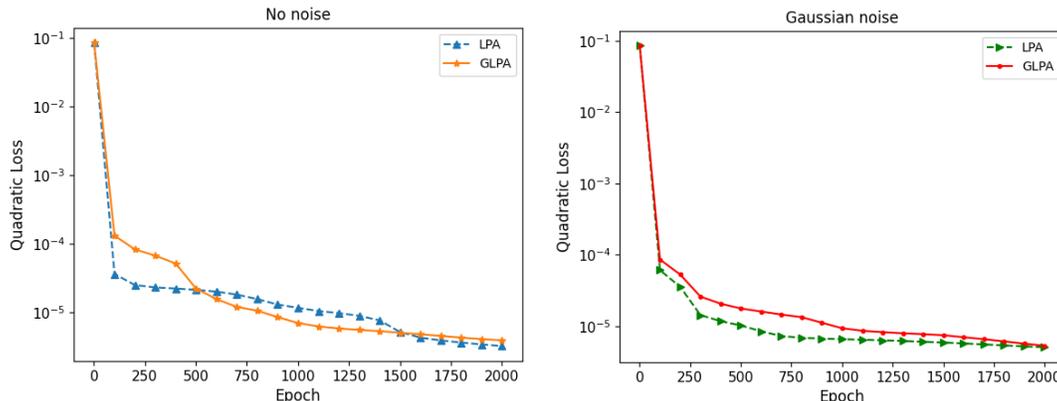


Figure 2: The variation of the quadratic loss during training. The training loss of the four experiments: (i) No noise: 3.1935e-6 and 5.0306e-6; (ii) Gaussian noise: 3.8656e-6 and 5.2940e-6.

As we can see, the LPA-type algorithms solve the regression tasks well, and they are robust even when the data is perturbed by the noise with a mean of 2.0094e-3 and a maximum of 3.9894e-3. The results show that the training loss is less than 5.2940e-6 for all test cases. In other words, our algorithms can obtain an ideal solution for this task. We find that the monotonic decrease of the objective function occurs at almost every iteration of the LPA. It is almost a descent algorithm. Through multiple experiments, we also find that the performance of the LPA depends on the choice of the initial point, but the GLPA is not affected by this. Thus, we conjecture that the GLPA for sigmoid networks with the quadratic loss function can converge globally under certain conditions. This will be explored in our future work.

Indeed, the LPA-type algorithms using small-scale networks can solve the problem as well. The illustration is shown on the left of Figure 3. Moreover, the performance of the algorithms is also affected by the stepsize of the subproblem. This is shown on the right of Figure 3.

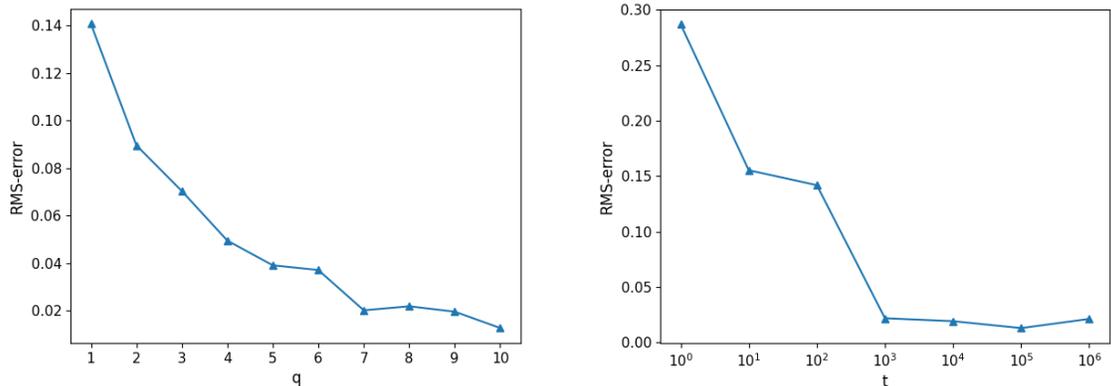


Figure 3: Execute Algorithm 2 by varying the number  $q$  of hidden neurons when  $t = 1e5$  (left) and the stepsize  $t$  of the subproblem when  $q = 10$  (right).

Corollary 6 shows that Algorithm 2 using absolute or hinge loss functions can converge globally. For simplicity, the rest of this section is devoted to demonstrating the performance of Algorithm 2. When implementing the GLPA for the sigmoid networks with an **absolute loss function**, we still use the same parameter values as in the previous experiments. For Algorithm A\*, we set  $\epsilon = \rho = 1e-2$ ,  $\Delta\theta^0 = \mathbf{0}$ ,  $\lambda^0 = \mathbf{0}$ , and the maximum number of ADMM iterations as 20. For the case when  $q = 72$  and  $t = 1e5$ , the performance of the algorithm is shown in Table 2 and Figure 4.

Table 2: The performance of regression on Franke’s function (using absolute loss).

GLPA		
	RMS-error	max-error
No noise	2.2093e-4	8.4516e-4
Gaussian noise	8.4138e-4	4.3988e-3

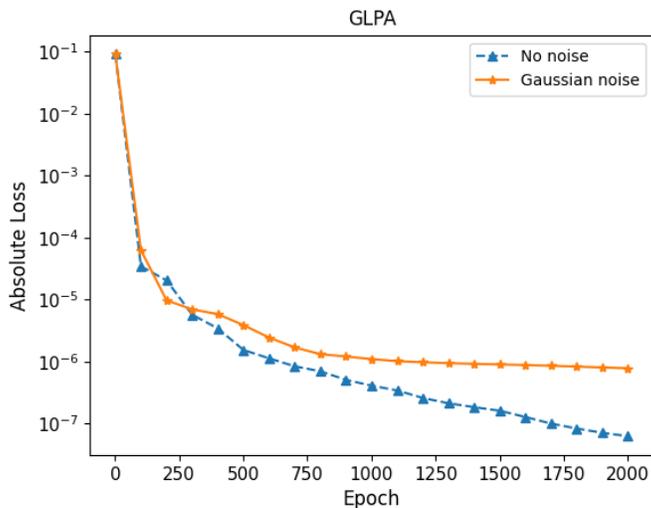


Figure 4: The variation of the absolute loss during training. The training loss in both experiments is  $6.2393e-8$  and  $7.7930e-7$ .

The training loss in both experiments is less than  $7.7930e-7$ , which shows that the GLPA obtains a better solution for sigmoid networks. Obviously, this result is more in line with the actual needs of regression tasks.

**5.2 Classification on Handwritten Digits.** The digits dataset from scikit-learn contains 1797 samples, each with 64 elements corresponding to an image of  $8 \times 8$  pixels, and with target attribute  $0, 1, \dots, 9$ . Some of the samples are shown in Figure 5.

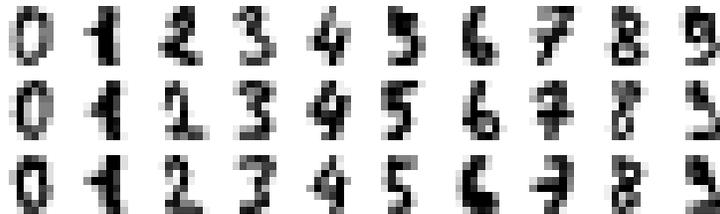


Figure 5: The first 30 samples of the digits dataset from scikit-learn.

We create four binary classification tasks, each to classify two digits: 0 and 1; 2 and 5; 3 and 7; 6 and 9. For each task, we take 70% of the selected samples as the training data and the rest as the test data. Here we run four algorithms on these tasks, including the GLPA and three other popular and practical tools in the machine learning community, SGDM, RMSProp and Adam. We also use the same parameter settings for the GLPA as the previous experiments. The only difference is that we set  $q = 4$  by Corollary 8 and the number of ADMM iterations does not exceed 10. For the other algorithms, implemented with PyTorch, we set the learning rate as  $1e-3$ , the momentum as 0.9, and the number of iterations as 1000. For the case when  $q = 4$ , the running results of the four algorithms are shown in Table 3 and Figure 6.

Three observations are indicated by the running results: (i) The small training loss shows that the GLPA can obtain excellent solutions to classification problems, and the training loss of the GLPA is generally smaller than the other algorithms. (ii) The GLPA has a much smaller number of iterations, thanks to its quadratic convergence rate in this case. It is striking that a first-order algorithm (GLPA) even has a second-order convergence rate. (iii) The adaptive network size given by Corollary 8 is sufficient to construct an ideal sigmoid network that solves the problem effectively. Hence Corollary 8 does provide a good guide for setting the size of sigmoid networks.

Table 3: The performance of classification on handwritten digit (using hinge loss).

Classified Digits	GLPA		SGDM (RMSProp, Adam)	
	Training errors	Test errors	Training errors	Test errors
0 - 1	0 / 252	0 / 108	0 / 252	0 / 108
2 - 5	0 / 251	0 / 108	0 / 251	0 / 108
3 - 7	0 / 253	0 / 109	0 / 253	0 / 109
6 - 9	0 / 252	1 / 109	0 / 252	1 / 109

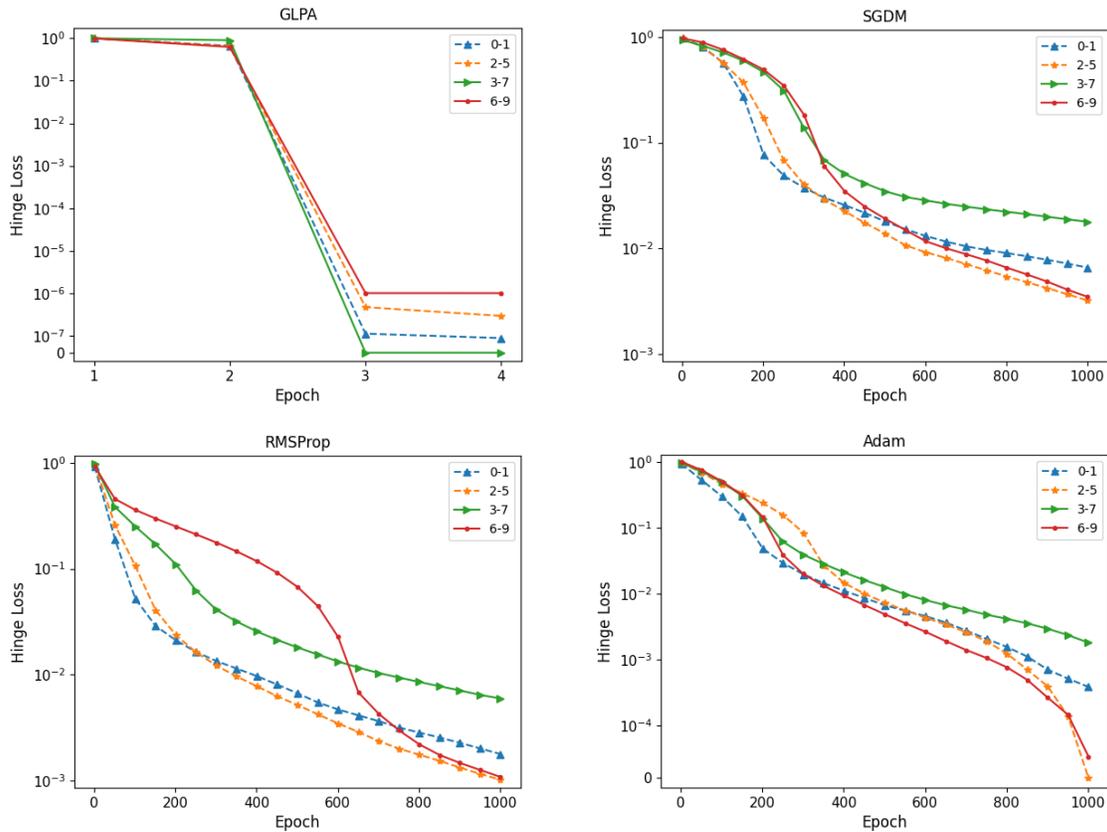


Figure 6: The variation of the hinge loss during training. The training loss of the four binary classification tasks: (i) GLPA:  $8.8007e-8$ ,  $2.9285e-7$ ,  $0.0$  and  $1.0065e-6$ ; (ii) SGDM:  $6.5701e-3$ ,  $3.2118e-3$ ,  $1.7844e-2$  and  $3.4731e-3$ ; (iii) RMSProp:  $1.7785e-3$ ,  $1.0178e-3$ ,  $5.9491e-3$  and  $1.0813e-3$ ; (iv) Adam:  $3.8339e-4$ ,  $0.0$ ,  $1.8042e-3$  and  $3.3567e-5$ .

The essence of Corollary 8 is to guarantee that the number of parameters in neural networks is not smaller than the amount of training data, and that a sufficient number of parameters ensure the feasibility of the networks. In our view, it is as if the information of a data point could be extracted by a single parameter in the model. Inspired by this, we think it can also serve as a general guide for setting the size of neural networks. It is well known that how to set the number of hidden neurons in neural networks is still an open problem, and it is usually adjusted by trial and error in practice. As stated above, we suggest that the number of hidden neurons can be specified by trial and error starting from the adaptive network size, which can avoid certain blindness at the beginning of the trial. This general rule deserves to be tried and further verified in practice.

## 6. Future Work

Although we only show the composite optimization algorithms for the three-layer sigmoid networks, our algorithms are also applicable to the more complex sigmoid networks, such as

the sigmoid networks with multiple hidden layers, with multiple outputs, and with output layer neurons that are processed with sigmoid functions. In the design of model (2.2), the convexity of the outer function  $\mathbb{L}$  is due to the convex loss function  $L$ , and the smoothness of the inner function  $F$  is due to the smooth fitting function  $f$ . So the algorithms can be used to solve the sigmoid networks whenever we maintain the convexity of  $L$  and the smoothness of  $f$  (note that  $f$  is always smooth in sigmoid networks). It is not difficult to solve the general sigmoid networks with convex loss functions using our algorithms by setting the same form of  $\mathbb{L}$  and  $F$  as the case of one hidden layer. As a matter of fact, the composite structure (2.2) can provide a unified framework for the development and analysis of sigmoid networks, especially for the non-convex and non-smooth optimization problems. Moreover, the various composite structures in neural networks pose more challenges for the study of composite optimization algorithms. The breakthrough of composite optimization algorithms will also drive the development of neural network learning algorithms. Last but not least, the convergence results of convex composite optimization (4.1) in the literature all seem to be established on  $G(\omega^*) \in C_h$ . While the more general convergence theorems should be established possibly on  $G(\omega^*) \notin C_h$ , which is still an open problem in the area of composite optimization. In view of this, we will explore this issue further.

## Acknowledgments

The research was supported in part by the National Natural Science Foundation of China under grants 12071157 and 12026602, and the Natural Science Foundation of Guangdong 2020B1515310013. Qi Ye is the corresponding author.

## References

- Oludare I. Abiodun, Aman Jantan, Abiodun E. Omolara, Kemi V. Dada, Nachaat A. Mohamed, and Humaira Arshad. State-of-the-art in artificial neural network applications: A survey. *Heliyon*, 4(11):e00938, 2018.
- Martin Anthony and Peter L. Bartlett. *Neural Network Learning: Theoretical Foundations*. Cambridge University Press, Cambridge, 1999.
- Dimitri P. Bertsekas. *Convex Optimization Algorithms*. Athena Scientific, Belmont, MA, 2015.
- Stephen Boyd, Neal Parikh, Eric Chu, Borja Peleato, and Jonathan Eckstein. Distributed optimization and statistical learning via the alternating direction method of multipliers. *Foundations and Trends in Machine learning*, 3(1):1–122, 2011.
- James V. Burke and Michael C. Ferris. Weak sharp minima in mathematical programming. *SIAM Journal on Control and Optimization*, 31(5):1340–1359, 1993.
- James V. Burke and Michael C. Ferris. A gauss-newton method for convex composite optimization. *Mathematical Programming*, 71(2):179–194, 1995.

- James V. Burke, Adrian S. Lewis, and Michael L. Overton. A robust gradient sampling algorithm for nonsmooth, nonconvex optimization. *SIAM Journal on Optimization*, 15(3):751–779, 2005.
- Byung-Woo Hong, Ja-Keoung Koo, Hendrik Dirks, and Martin Burger. Adaptive regularization in convex composite optimization for variational imaging problems. In *Pattern Recognition: 39th German Conference, GCPR 2017, Basel, Switzerland, September 12–15, 2017, Proceedings 39*, pages 268–280. Springer, 2017.
- Xinlin Hu, Yaohua Hu, Fanjie Wu, Ricky Wai Tak Leung, and Jing Qin. Integration of single-cell multi-omics for gene regulatory network inference. *Computational and Structural Biotechnology Journal*, 18:1925–1938, 2020.
- Yaohua Hu, Chong Li, and Xiaoqi Yang. On convergence rates of linearized proximal algorithms for convex composite optimization with applications. *SIAM Journal on Optimization*, 26(2):1207–1235, 2016.
- Yann LeCun, Yoshua Bengio, and Geoffrey Hinton. Deep learning. *Nature*, 521(7553):436–444, 2015.
- Adrian S. Lewis and Stephen J. Wright. A proximal method for composite minimization. *Mathematical Programming*, 158(1):501–546, 2016.
- Sashank J. Reddi, Satyen Kale, and Sanjiv Kumar. On the convergence of adam and beyond. In *International Conference on Learning Representations*, 2018.
- Jinhua Wang, Yaohua Hu, Chong Li, and Jen-Chih Yao. Linear convergence of CQ algorithms and applications in gene regulatory network inference. *Inverse Problems*, 33(5):055017, 2017.
- Ashia C. Wilson, Rebecca Roelofs, Mitchell Stern, Nati Srebro, and Benjamin Recht. The marginal value of adaptive gradient methods in machine learning. *Advances in Neural Information Processing Systems*, 30, 2017.