

---

# FINITE-SUM OPTIMIZATION: A NEW PERSPECTIVE FOR CONVERGENCE TO A GLOBAL SOLUTION

---

**Lam M. Nguyen<sup>1\*</sup>, Trang H. Tran<sup>2\*</sup>, Marten van Dijk<sup>3</sup>**

<sup>1</sup> IBM Research, Thomas J. Watson Research Center, Yorktown Heights, NY, USA

<sup>2</sup> School of Operations Research and Information Engineering, Cornell University, Ithaca, NY, USA

<sup>3</sup> CWI Amsterdam, The Netherlands

LamNguyen.MLTD@ibm.com, htt27@cornell.edu, marten.van.dijk@cwi.nl

## ABSTRACT

Deep neural networks (DNNs) have shown great success in many machine learning tasks. Their training is challenging since the loss surface of the network architecture is generally non-convex, or even non-smooth. How and under what assumptions is guaranteed convergence to a *global* minimum possible? We propose a reformulation of the minimization problem allowing for a new recursive algorithmic framework. By using bounded style assumptions, we prove convergence to an  $\varepsilon$ -(global) minimum using  $\tilde{O}(1/\varepsilon^3)$  gradient computations. Our theoretical foundation motivates further study, implementation, and optimization of the new algorithmic framework and further investigation of its non-standard bounded style assumptions. This new direction broadens our understanding of why and under what circumstances training of a DNN converges to a global minimum.

## 1 Introduction

In recent years, deep neural networks (DNNs) have shown a great success in many machine learning tasks. However, training these neural networks is challenging since the loss surface of network architecture is generally non-convex, or even non-smooth. Thus, there have been a long-standing question on how optimization algorithms may converge to a global minimum. Many previous work have investigated Gradient Descent algorithm and its stochastic version for over-parameterized setting [Arora et al., 2018, Soudry et al., 2018, Allen-Zhu et al., 2019, Du et al., 2019a, Zou and Gu, 2019]. Although these works have shown promising convergence results under certain assumptions, there is still a lack of new efficient methods that can guarantee convergence to a global solution for machine learning optimization. In this paper, we address this problem using a different perspective. Instead of analyzing the traditional finite-sum formulation, we adopt a new *composite formulation* that exactly depicts the structure of machine learning where a data set is used to learn a common classifier.

**Representation.** Let  $\{(x^{(i)}, y^{(i)})\}_{i=1}^n$  be a given training set with  $x^{(i)} \in \mathbb{R}^m, y^{(i)} \in \mathbb{R}^c$ , we investigate the following novel representation for deep learning tasks:

$$\min_{w \in \mathbb{R}^d} \left\{ F(w) = \frac{1}{n} \sum_{i=1}^n \phi_i(h(w; i)) \right\}, \quad (1)$$

where  $h(\cdot; i) : \mathbb{R}^d \rightarrow \mathbb{R}^c, i \in [n] = \{1, \dots, n\}$ , is the classifier for each input data  $x^{(i)}$ ; and  $\phi_i : \mathbb{R}^c \rightarrow \mathbb{R}, i \in [n]$ , is the loss function corresponding to each output data  $y^{(i)}$ . Our *composite formulation* (1) is a special case of the finite-sum problem  $\min_{w \in \mathbb{R}^d} \{F(w) = \frac{1}{n} \sum_{i=1}^n f(w; i)\}$  where each individual function  $f(\cdot; i)$  is a composition of the loss function  $\phi_i$  and the classifier  $h(\cdot; i)$ . This problem covers various important applications in machine learning, including logistic regression and neural networks. The most common approach for the finite-sum problem is using

---

\* Equal contribution. Correspondence to: Lam M. Nguyen.

first-order methods such as (stochastic) gradient algorithms and making assumptions on the component functions  $f(\cdot; i)$ . As an alternative, we further investigate the structure of the loss function  $\phi_i$  and narrow our assumption on the classifier  $h(\cdot; i)$ . For the purpose of this work, we first consider convex and Lipschitz-smooth loss functions while the classifiers can be non-convex. Using this representation, we propose a new framework followed by two algorithms that guarantee convergence to a global solution for the minimization problem.

**Algorithmic Framework.** Representation (1) admits a new perspective. Our key insight is to (A) define  $z_i^{(t)} = h(w^{(t)}; i)$ , where  $t$  is an iteration count of the outer loop in our algorithmic framework. Next (B), we want to approximate the change  $z_i^{(t+1)} - z_i^{(t)}$  in terms of a step size times the gradient

$$\nabla\phi_i(z_i^{(t)}) = (\partial\phi_i(z)/\partial z_a)_{a \in [c]} \Big|_{z=z_i^{(t)}},$$

and (C) we approximate the change  $h(w^{(t+1)}; i) - h(w^{(t)}; i)$  in terms of the first order derivative

$$H_i^{(t)} = (\partial h_a(w; i)/\partial w_b)_{a \in [c], b \in [d]} \Big|_{w=w^{(t)}}.$$

Finally, we combine (A), (B), and (C) to equate the approximations of  $z_i^{(t+1)} - z_i^{(t)}$  and  $h(w^{(t+1)}; i) - h(w^{(t)}; i)$ . This leads to a recurrence on  $w^{(t)}$  of the form  $w^{(t+1)} = w^{(t)} - \eta^{(t)}v^{(t)}$ , where  $\eta^{(t)}$  is a step size and which involves computing  $v^{(t)}$  by solving a convex quadratic subproblem, see the details in Section 4. We explain two methods for approximating a solution for the derived subproblem. We show how to approximate the subproblem by transforming it into a strongly convex problem by adding a regularizer which can be solved in closed form. And we show how to use Gradient Descent (GD) on the subproblem to find an approximation  $v^{(t)}$  of its solution.

**Convergence Analysis.** Our analysis introduces non-standard bounded style assumptions. Intuitively, we assume that our convex and quadratic subproblem has a *bounded* solution. This allows us to prove a total complexity of  $\tilde{O}(\frac{1}{\varepsilon^3})$  to find an  $\varepsilon$ -(global) solution that satisfies  $F(\hat{w}) - F_* \leq \varepsilon$ , where  $F_*$  is the global minimizer of  $F$ . Our analysis applies to a wide range of applications in machine learning: Our results hold for squared loss and softmax cross-entropy loss and applicable for a range of activation functions in DNN as we only assume that the  $h(\cdot; i)$  are twice continuously differentiable and their Hessian matrices (second order derivatives) as well as their gradients (first order derivatives) are bounded.

**Contributions and Outline.** Our contributions in this paper can be summarized as follows.

- We propose a new representation (1) for analyzing the machine learning minimization problem. Our formulation utilizes the structure of machine learning tasks where a training data set of inputs and outputs is used to learn a common classifier. Related work in Section 2 shows how (1) is different from the classical finite-sum problem.
- Based on the new representation we propose a novel algorithm framework. The algorithmic framework approximates a solution to a subproblem for which we show two distinct approaches.
- For general DNNs and based on bounded style assumptions, we prove a total complexity of  $\tilde{O}(\frac{1}{\varepsilon^3})$  to find an  $\varepsilon$ -(global) solution that satisfies  $F(\hat{w}) - F_* \leq \varepsilon$ , where  $F_*$  is the global minimizer of  $F$ .

We emphasize that our focus is on developing a new theoretical foundation and that a translation to a practical implementation with empirical results is for future work. Our theoretical foundation motivates further study, implementation, and optimization of the new algorithmic framework and further investigation of its non-standard bounded style assumptions. This new direction broadens our understanding of why and under what circumstances training of a DNN converges to a global minimum.

The rest of this paper is organized as follows. Section 2 discusses related work. Section 3 describes our setting and deep learning representation. Section 4 explains our key insight and derives our Framework 1. Section 5 presents our algorithms and their convergence to a global solution. All technical proofs are deferred to the Appendix.

## 2 Related Work

**Formulation for Machine Learning Problems.** The finite-sum problem is one of the most important and fundamental problems in machine learning. Analyzing this model is the most popular approach in the machine learning literature and it has been studied intensively throughout the years [Bottou et al., 2018, Reddi et al., 2016, Duchi et al., 2011]. Our new formulation (1) is a special case of the finite-sum problem, however, it is much more complicated than the previous model since it involves the data index  $i$  both inside the classifiers  $h(\cdot; i)$  and the loss functions  $\phi_i$ . For a

comparison, previous works only consider a common loss function  $l(\hat{y}, y)$  for the predicted value  $\hat{y}$  and output data  $y$  [Zou et al., 2018, Soudry et al., 2018]. Our modified version of loss function  $\phi_i$  is a natural setting for machine learning. We note that when  $h(w; i)$  is the output produced by a model, our goal is to match this output with the corresponding target  $y^{(i)}$ . For that reason, the loss function for each output has a dependence on the output data  $y^{(i)}$ , and is denoted by  $\phi_i$ . This fact reflects the natural setting of machine learning where the outputs are designed to fit different targets, and the optimization process depends on both outer function  $\phi_i$  and inner functions  $h(\cdot; i)$ . This complication may potentially bring a challenge to theoretical analysis. However, with separate loss functions, we believe this model will help to exploit better the structure of machine learning problems and gain more insights on the neural network architecture.

Other related composite optimization models are also investigated thoroughly in [Lewis and Wright, 2016, Zhang and Xiao, 2019, Tran-Dinh et al., 2020]. Our model is different from these works as it does not have a common function wrapping outside the finite-sum term, as in [Lewis and Wright, 2016]. Note that a broad class of variance reduction algorithms (e.g. SAG [Le Roux et al., 2012], SAGA [Defazio et al., 2014], SVRG [Johnson and Zhang, 2013], SARAH [Nguyen et al., 2017]) is designed specifically for the finite-sum formulation and is known to have certain benefits over Gradient Descent. In addition, the multilevel composite problem considered in [Zhang and Xiao, 2021] also covers empirical risk minimization problem. However our formulation does not match their work since our inner function  $h(w; i)$  is not an independent expectation over some data distribution, but a specific function that depends on the current data.

**Global Convergence for Neural Networks.** A recent popular line of research is studying the dynamics of optimization methods on some specific neural network architectures. There are some early works that show the global convergence of Gradient Descent (GD) for simple linear network and two-layer network [Brutzkus et al., 2018, Soudry et al., 2018, Arora et al., 2019, Du et al., 2019b]. Some further works extend these results to deep learning architectures [Allen-Zhu et al., 2019, Du et al., 2019a, Zou and Gu, 2019]. These theoretical guarantees are generally proved for the case when the last output layer is fixed, which is not standard in practice. A recent work [Nguyen and Mondelli, 2020] prove the global convergence for GD when all layers are trained with some initial conditions. However, these results are for neural networks without bias neurons and it is unclear how these analyses can be extended to handle the bias terms of deep networks with different activations. Our novel framework and algorithms do not exclude learning bias layers as in [Nguyen and Mondelli, 2020].

Using a different algorithm, Brutzkus et al. [2018] investigate Stochastic Gradient Descent (SGD) for two-layer networks in a restricted linearly separable data setting. This line of research continues with the works from [Allen-Zhu et al., 2019, Zou et al., 2018] and later with [Zou and Gu, 2019]. They justify the global convergence of SGD for deep neural networks for some probability depending on the number of input data and the initialization process.

**Over-Parameterized Settings and other Assumptions for Machine Learning.** Most of the modern learning architectures are over-parameterized, which means that the number of parameters are very large and often far more than the number of input data. Some recent works prove the global convergence of Gradient Descent when the number of neurons are extensively large, e.g. [Zou and Gu, 2019] requires  $\Omega(n^8)$  neurons for every hidden layer, and [Nguyen and Mondelli, 2020] improves this number to  $\Omega(n^3)$ . If the initial point satisfies some special conditions, then they can show a better dependence of  $\Omega(n)$ . In [Allen-Zhu et al., 2019], the authors initialize the weights using a random Gaussian distribution where the variance depends on the dimension of the problem. In non-convex setting, they prove the convergence of SGD using the assumption that the dimension depends inversely on the tolerance  $\epsilon$ . We will discuss how these over-parameterized settings might be a necessary condition to develop our theory.

Other standard assumptions for machine learning include the bounded gradient assumption [Nemirovski et al., 2009, Shalev-Shwartz et al., 2007, Reddi et al., 2016, Tran et al., 2021]. It is also common to assume all the iterations of an algorithm stay in a bounded domain [Duchi et al., 2011, Levy et al., 2018, Gürbüzbalaban et al., 2019, Reddi et al., 2018, Vaswani et al., 2021]. Since we are analyzing a new *composite formulation*, it is understandable that our assumptions may also not be standard. However, we believe that there is a strong connection between our assumptions and the traditional setting of machine learning. We will discuss this point more clearly in Section 4.

### 3 Background

In this section, we discuss our formulation and notations in detail. Although this paper focuses on deep neural networks, our framework and theoretical analysis are general and applicable for other learning architectures.

**Deep Learning Representation.** Let  $\{(x^{(i)}, y^{(i)})\}_{i=1}^n$  be a training data set where  $x^{(i)} \in \mathbb{R}^m$  is a training input and  $y^{(i)} \in \mathbb{R}^c$  is a training output. We consider a fully-connected neural network with  $L$  layers, where the  $l$ -th layer,

$l \in \{0, 1, \dots, L\}$ , has  $n_l$  neurons. We represent layer 0-th and  $L$ -th layer as input and output layers, respectively, that is,  $n_0 = d$  and  $n_L = c$ . For  $l \in \{1, \dots, L\}$ , let  $W^{(l)} \in \mathbb{R}^{n_{l-1} \times n_l}$  and  $b^{(l)} \in \mathbb{R}^{n_l}$ , where  $\{(W^{(l)}, b^{(l)})_{l=1}^L\}$  represent the parameters of the neural network. A classifier  $h(w; i)$  is formulated as

$$h(w; i) = W^{(L)T} \sigma_{L-1}(W^{(L-1)T} \sigma_{L-2}(\dots \sigma_1(W^{(1)T} x^{(i)} + b^{(1)}) \dots) + b^{(L-1)}) + b^{(L)},$$

where  $w = \mathbf{vec}(\{W^{(1)}, b^{(1)}, \dots, W^{(L)}, b^{(L)}\}) \in \mathbb{R}^d$  is the vectorized weight and  $\{\sigma_l\}_{l=1}^{L-1}$  are some activation functions. The most common choices for machine learning are ReLU, sigmoid, hyperbolic tangent and softplus. For  $j \in [c]$ ,  $h_j(\cdot; i) : \mathbb{R}^d \rightarrow \mathbb{R}$  denotes the component function of the output  $h(\cdot; i)$ , for each data  $i \in [n]$  respectively. Moreover, we define  $h_i^* = \arg \min_{z \in \mathbb{R}^c} \phi_i(z)$ ,  $i \in [n]$ .

**Loss Functions.** The well-known loss functions in neural networks for solving classification and regression problems are *softmax cross-entropy loss* and *square loss*, respectively:

(*Softmax*) *Cross-Entropy Loss*:  $F(w) = \frac{1}{n} \sum_{i=1}^n f(w; i)$  with

$$f(w; i) = -y^{(i)T} \log(\text{softmax}(h(w; i))). \quad (2)$$

(*Squared Loss*):  $F(w) = \frac{1}{n} \sum_{i=1}^n f(w; i)$  with

$$f(w; i) = \frac{1}{2} \|h(w; i) - y^{(i)}\|^2. \quad (3)$$

We provide some basic definitions in optimization theory to support our theory.

**Definition 1** (*L-smooth*). *Function  $\phi : \mathbb{R}^c \rightarrow \mathbb{R}$  is  $L_\phi$ -smooth if there exists a constant  $L_\phi > 0$  such that,  $\forall x_1, x_2 \in \mathbb{R}^c$ ,*

$$\|\nabla\phi(x_1) - \nabla\phi(x_2)\| \leq L_\phi \|x_1 - x_2\|. \quad (4)$$

**Definition 2** (*Convex*). *Function  $\phi : \mathbb{R}^c \rightarrow \mathbb{R}$  is convex if  $\forall x_1, x_2 \in \mathbb{R}^c$ ,*

$$\phi(x_1) - \phi(x_2) \geq \langle \nabla\phi(x_2), x_1 - x_2 \rangle. \quad (5)$$

The following corollary shows the properties of softmax cross-entropy loss (2) and squared loss (3).

**Corollary 1.** *For softmax cross-entropy loss (2) and squared loss (3), there exist functions  $h(\cdot; i) : \mathbb{R}^d \rightarrow \mathbb{R}^c$  and  $\phi_i : \mathbb{R}^c \rightarrow \mathbb{R}$  such that, for  $i \in [n]$ ,  $\phi_i(z)$  is convex and  $L_\phi$ -smooth with  $L_\phi = 1$ , and*

$$f(w; i) = \phi_i(h(w; i)) = \phi_i(z) \Big|_{z=h(w; i)}. \quad (6)$$

The following lemma is a standard result in [Nesterov, 2004].

**Lemma 1** ([Nesterov, 2004]). *If  $\phi$  is  $L_\phi$ -smooth and convex, then for  $\forall z \in \mathbb{R}^c$ ,*

$$\|\nabla\phi(z)\|^2 \leq 2L_\phi(\phi(z) - \phi(z_*)), \quad (7)$$

where  $z_* = \arg \min_z \phi(z)$ .

The following useful derivations can be used later in our theoretical analysis. Since  $\phi_i$  is convex, by Definition 2 we have

$$\phi_i(h(w; i)) \geq \phi_i(h(w'; i)) + \left\langle \nabla_z \phi_i(z) \Big|_{z=h(w'; i)}, h(w; i) - h(w'; i) \right\rangle. \quad (8)$$

If  $\phi_i$  is convex and  $L_\phi$ -smooth, then by Lemma 1

$$\left\| \nabla_z \phi_i(z) \Big|_{z=h(w; i)} \right\|^2 \leq 2L_\phi [\phi_i(h(w; i)) - \phi_i(h_i^*)], \quad (9)$$

where  $h_i^* = \arg \min_{z \in \mathbb{R}^c} \phi_i(z)$ .

We compute gradients of  $f(w; i)$  in terms of  $\phi_i(h(w; i))$ .

Table 1: Table of notations

Notation	Meaning
$F_*$	Global minimization function of $F$ in (1) $F_* = \min_{w \in \mathbb{R}^d} F(w)$
$h_i^*$	$h_i^* = \arg \min_{z \in \mathbb{R}^c} \phi_i(z), i \in [n]$
$v_*^{(t)}$	Solution of the convex problem in (15) $\min_{v \in \mathbb{R}^d} \frac{1}{2} \frac{1}{n} \sum_{i=1}^n \ \eta^{(t)} H_i^{(t)} v - \alpha_i^{(t)} \nabla_z \phi_i(h(w^{(t)}; i))\ ^2$
$v^{(t)}$	An approximation of $v_*^{(t)}$ which is used as the search direction in Framework 1
$\hat{v}_{*\varepsilon}^{(t)}$	A vector that satisfies $\frac{1}{2} \frac{1}{n} \sum_{i=1}^n \ \eta^{(t)} H_i^{(t)} v - \alpha_i^{(t)} \nabla_z \phi_i(h(w^{(t)}; i))\ ^2 \leq \varepsilon^2$ for some $\varepsilon > 0$ and $\ \hat{v}_{*\varepsilon}^{(t)}\ ^2 \leq V$ , for some $V > 0$ .
$v_*^{(t) \text{ reg}}$	Solution of the strongly convex problem in (20) $\min_{v \in \mathbb{R}^d} \left\{ \frac{1}{2} \frac{1}{n} \sum_{i=1}^n \ \eta^{(t)} H_i^{(t)} v - \alpha_i^{(t)} \nabla_z \phi_i(h(w^{(t)}; i))\ ^2 + \frac{\varepsilon^2}{2} \ v\ ^2 \right\}$

- **Gradient of softmax cross-entropy loss:**

$$\nabla \phi_i(z) \Big|_{z=h(w;i)} = \left( \frac{\partial \phi_i(z)}{\partial z_1} \Big|_{z=h(w;i)}, \dots, \frac{\partial \phi_i(z)}{\partial z_c} \Big|_{z=h(w;i)} \right)^T,$$

where for  $j \in [c]$ ,  $\frac{\partial \phi_i(z)}{\partial z_j} \Big|_{z=h(w;i)}$  is

$$\begin{cases} \frac{\exp([h(w;i)]_j - [h(w;i)]_{I(y^{(i)})})}{\sum_{k=1}^c \exp([h(w;i)]_k - [h(w;i)]_{I(y^{(i)})})}, & j \neq I(y^{(i)}) \\ -\frac{\sum_{k \neq I(y^{(i)})} \exp([h(w;i)]_k - [h(w;i)]_{I(y^{(i)})})}{\sum_{k=1}^c \exp([h(w;i)]_k - [h(w;i)]_{I(y^{(i)})})}, & j = I(y^{(i)}) \end{cases}. \quad (10)$$

- **Gradient of squared loss:**

$$\nabla \phi_i(z) \Big|_{z=h(w;i)} = h(w; i) - y^{(i)}. \quad (11)$$

We introduce the notations that we use throughout the paper in Table 1.

## 4 New Algorithm Framework

### 4.1 Key Insight

We assume  $f(w; i) = \phi_i(h(w; i))$  with  $\phi_i$  convex and  $L_\phi$ -smooth. Our goal is to utilize the convexity of the outer function  $\phi_i$ . In order to simplify notation, we write  $\nabla_z \phi_i(h(w^{(t)}; i))$  instead of  $\nabla_z \phi_i(z) \Big|_{z=h(w^{(t)}; i)}$  and denote  $z_i^{(t)} = h(w^{(t)}; i)$ . Starting from the current weight  $w^{(t)}$ , we would like to find the next point  $w^{(t+1)}$  that satisfies the following approximation for all  $i \in [n]$ :

$$h(w^{(t+1)}; i) = z_i^{(t+1)} \approx z_i^{(t)} - \alpha_i^{(t)} \nabla_z \phi_i(z_i^{(t)}) = h(w^{(t)}; i) - \alpha_i^{(t)} \nabla_z \phi_i(h(w^{(t)}; i)). \quad (12)$$

We can see that this approximation is a ‘‘noisy’’ version of a gradient descent update for every function  $\phi_i$ , simultaneously for all  $i \in [n]$ . In order to do this, we use the following update

$$w^{(t+1)} = w^{(t)} - \eta^{(t)} v^{(t)}, \quad (13)$$

where  $\eta^{(t)} > 0$  is a learning rate and  $v^{(t)}$  is a search direction that helps us approximate equation (12). If the update term  $\eta^{(t)} v^{(t)}$  is small enough, and if  $h(\cdot; i)$  has some nice smooth properties, then from basic calculus we have the

following approximation:

$$h(w^{(t+1)}; i) = h(w^{(t)} - \eta^{(t)}v^{(t)}; i) \approx h(w^{(t)}; i) - H_i^{(t)}(\eta^{(t)}v^{(t)}), \quad (14)$$

where  $H_i^{(t)}$  is a matrix in  $\mathbb{R}^{c \times d}$  with first-order derivatives. Motivated by approximations (12) and (14), we consider the following optimization problem:

$$v_*^{(t)} = \arg \min_{v \in \mathbb{R}^d} \left\{ \frac{1}{2} \frac{1}{n} \sum_{i=1}^n \|H_i^{(t)}(\eta^{(t)}v) - \alpha_i^{(t)} \nabla_z \phi_i(h(w^{(t)}; i))\|^2 \right\}. \quad (15)$$

Hence, by solving for the solution  $v_*^{(t)}$  of problem (15) we are able to find a search direction for the key approximation (12). This yields our new algorithmic Framework 1, see below.

---

### Framework 1 New Algorithm Framework

---

**Initialization:** Choose an initial point  $w^{(0)} \in \mathbb{R}^d$ ;  
**for**  $t = 0, 1, \dots, T - 1$  **do**

Solve for an approximation  $v^{(t)}$  of the solution  $v_*^{(t)}$  of the problem in (15)

$$v^{(t)} = \arg \min_{v \in \mathbb{R}^d} \left\{ \frac{1}{2} \frac{1}{n} \sum_{i=1}^n \|\eta^{(t)} H_i^{(t)} v - \alpha_i^{(t)} \nabla_z \phi_i(h(w^{(t)}; i))\|^2 \right\}.$$

Update  $w^{(t+1)} = w^{(t)} - \eta^{(t)}v^{(t)}$

**end for**

---

## 4.2 Technical Assumptions

**Assumption 1.** *The loss function  $\phi_i$  is convex and  $L_\phi$ -smooth for  $i \in [n]$ . Moreover, we assume that it is lower bounded, i.e.  $\inf_{z \in \mathbb{R}^c} \phi_i(z) > -\infty$  for  $i \in [n]$ .*

We have shown the convexity and smoothness of squared loss and softmax cross-entropy loss in Section 3. The bounded property of  $\phi_i$  is required in any algorithm for the well-definedness of (1). Now, in order to use the Taylor series approximation, we need the following assumption on the neural network architecture  $h$ :

**Assumption 2.** *We assume that  $h(\cdot; i)$  is twice continuously differentiable for all  $i \in [n]$  (i.e. the second-order partial derivatives of all scalars  $h_j(\cdot; i)$  are continuous for all  $j \in [c]$  and  $i \in [n]$ ), and that their Hessian matrices are bounded, that is, there exists a  $G > 0$  such that for all  $w \in \mathbb{R}^d$ ,  $i \in [n]$  and  $j \in [c]$ ,*

$$\|M_{i,j}(w)\| = \|\mathbf{J}_w(\nabla_w h_j(w; i))\| \leq G, \quad (16)$$

where  $\mathbf{J}_w$  denotes the Jacobian<sup>1</sup>.

**Remark 1** (Relation to second-order methods). *Although our analysis requires an assumption on the Hessian matrices of  $h(w; i)$ , our algorithms do not use any second order information or try to approximate this information. Our theoretical analysis focused on the approximation of the classifier and the gradient information, therefore is not related to the second order type algorithms. It is currently unclear how to apply second order methods into our problem, however, this is an interesting research question to expand the scope of this work.*

Assumption 2 allows us to apply a Taylor approximation of each function  $h_j(\cdot; i)$  with which we prove the following Lemma that bounds the error in equation (14):

**Lemma 2.** *Suppose that Assumption 2 holds for the classifier  $h$ . Then for all  $i \in [n]$  and  $0 \leq t < T$ ,*

$$h(w^{(t+1)}; i) = h(w^{(t)} - \eta^{(t)}v^{(t)}; i) = h(w^{(t)}; i) - \eta^{(t)} H_i^{(t)} v^{(t)} + \epsilon_i^{(t)}, \quad (17)$$

where

$$H_i^{(t)} = \mathbf{J}_w(h(w; i))|_{w=w^{(t)}} \in \mathbb{R}^{c \times d} \quad (18)$$

is defined as the Jacobian matrix of  $h(w; i)$  at  $w^{(t)}$  and entries  $\epsilon_{i,j}^{(t)}$ ,  $j \in [c]$ , of vector  $\epsilon_i^{(t)}$  satisfy

$$|\epsilon_{i,j}^{(t)}| \leq \frac{1}{2} (\eta^{(t)})^2 \|v^{(t)}\|^2 G. \quad (19)$$

---

<sup>1</sup>For a continuously differentiable function  $g(w) : \mathbb{R}^d \rightarrow \mathbb{R}^c$  we define the Jacobian  $\mathbf{J}_w(g(w))$  as the matrix  $(\partial g_a(w) / \partial w_b)_{a \in [c], b \in [d]}$ .

In order to approximate (12) combined with (14), that is, to make sure the right hand sides of (12) and (14) are close to one another, we consider the optimization problem (15):

$$v_*^{(t)} = \arg \min_{v \in \mathbb{R}^d} \left\{ \frac{1}{2} \frac{1}{n} \sum_{i=1}^n \|\eta^{(t)} H_i^{(t)} v - \alpha_i^{(t)} \nabla_z \phi_i(h(w^{(t)}; i))\|^2 \right\}.$$

The optimal value of problem (15) is equal to 0 if there exists a vector  $v_*^{(t)}$  satisfying  $\eta^{(t)} H_i^{(t)} v_*^{(t)} = \alpha_i^{(t)} \nabla_z \phi_i(h(w^{(t)}; i))$  for every  $i \in [n]$ . Since the solution  $v_*^{(t)}$  is in  $\mathbb{R}^d$  and  $\nabla_z \phi_i(h(w^{(t)}; i))$  is in  $\mathbb{R}^c$ , this condition is equivalent to a linear system with  $n \cdot c$  constraints and  $d$  variables. In the over-parameterized setting where dimension  $d$  is sufficiently large ( $d \gg n \cdot c$ ) and there are no identical data, there exists almost surely a vector  $v_*^{(t)}$  that interpolates all the training set, see the Appendix for details.

Let us note that an approximation of  $v_*^{(t)}$  serves as the search direction for Framework 1. For this reason, the solution  $v_*^{(t)}$  of problem (15) plays a similar role as a gradient in the search direction of (stochastic) gradient descent method. It is standard to assume a bounded gradient in the machine learning literature [Nemirovski et al., 2009, Shalev-Shwartz et al., 2007, Reddi et al., 2016]. Motivated by these facts, we assume the following Assumption 3, which implies the existence of a near-optimal *bounded* solution of (15):

**Assumption 3.** *We consider an over-parameterized setting where dimension  $d$  is sufficiently large enough to interpolate all the data and the tolerance  $\varepsilon$ . We assume that there exists a bound  $V > 0$  such that for  $\varepsilon > 0$  and  $0 \leq t < T$  as in Framework 1, there exists a vector  $\hat{v}_{*\varepsilon}^{(t)}$  with  $\|\hat{v}_{*\varepsilon}^{(t)}\|^2 \leq V$  so that*

$$\frac{1}{2} \frac{1}{n} \sum_{i=1}^n \|\eta^{(t)} H_i^{(t)} \hat{v}_{*\varepsilon}^{(t)} - \alpha_i^{(t)} \nabla_z \phi_i(h(w^{(t)}; i))\|^2 \leq \varepsilon^2.$$

Our Assumption 3 requires a nice dependency on the tolerance  $\varepsilon$  for the gradient matrices  $H_i^{(t)}$  and  $\nabla_z \phi_i(h(w^{(t)}; i))$ . We note that at the starting point  $t = 0$ , these matrices may depend on  $\varepsilon$  due to the initialization process and the dependence of  $d$  on  $\varepsilon$ . This setting is similar to previous works, e.g. [Allen-Zhu et al., 2019]. In the Appendix, we show an example of neural network architecture where Assumption 3 is justified at the start of the training process.

## 5 New Algorithms and Convergence Results

### 5.1 Approximating the solution using regularizer

Since problem (15) is convex and quadratic, we consider the following regularized problem:

$$\min_{v \in \mathbb{R}^d} \left\{ \Psi^{(t)}(v) = \Phi^{(t)}(v) + \frac{\varepsilon^2}{2} \|v\|^2 \right\}, \quad (20)$$

where

$$\Phi^{(t)}(v) = \frac{1}{2} \frac{1}{n} \sum_{i=1}^n \|\eta^{(t)} H_i^{(t)} v - \alpha_i^{(t)} \nabla_z \phi_i(h(w^{(t)}; i))\|^2.$$

for some small  $\varepsilon > 0$  and  $t \geq 0$ . It is widely known that problem (20) is strongly convex, and has a unique minimizer  $v_{*\text{reg}}^{(t)}$ . The global minimizer satisfies  $\nabla_v \Psi^{(t)}(v_{*\text{reg}}^{(t)}) = 0$ . We have

$$\begin{aligned} \nabla_v \Psi^{(t)}(v) &= \frac{1}{n} \sum_{i=1}^n \left[ \eta^{(t)} H_i^{(t)T} H_i^{(t)} \eta^{(t)} v - \alpha_i^{(t)} \eta^{(t)} H_i^{(t)T} \nabla_z \phi_i(h(w^{(t)}; i)) \right] + \varepsilon^2 \cdot v \\ &= \left( \frac{1}{n} \sum_{i=1}^n \eta^{(t)} H_i^{(t)T} H_i^{(t)} \eta^{(t)} + \varepsilon^2 I \right) v - \left( \frac{1}{n} \sum_{i=1}^n \alpha_i^{(t)} \eta^{(t)} H_i^{(t)T} \nabla_z \phi_i(h(w^{(t)}; i)) \right). \end{aligned}$$

Therefore,

$$v_{*\text{reg}}^{(t)} = \left( \frac{1}{n} \sum_{i=1}^n \eta^{(t)} H_i^{(t)T} H_i^{(t)} \eta^{(t)} + \varepsilon^2 I \right)^{-1} \left( \frac{1}{n} \sum_{i=1}^n \alpha_i^{(t)} \eta^{(t)} H_i^{(t)T} \nabla_z \phi_i(h(w^{(t)}; i)) \right). \quad (21)$$

---

**Algorithm 1** Solve for the exact solution of the regularized problem
 

---

**Initialization:** Choose an initial point  $w^{(0)} \in \mathbb{R}^d$ , tolerance  $\varepsilon > 0$ ;  
**for**  $t = 0, 1, \dots, T - 1$  **do**

Update the search direction  $v^{(t)}$  as the solution  $v_{* \text{reg}}^{(t)}$  of problem in (20):

$$v^{(t)} = v_{* \text{reg}}^{(t)} = \left( \frac{1}{n} \sum_{i=1}^n \eta^{(t)} H_i^{(t)T} H_i^{(t)} \eta^{(t)} + \varepsilon^2 I \right)^{-1} \left( \frac{1}{n} \sum_{i=1}^n \alpha_i^{(t)} \eta^{(t)} H_i^{(t)T} \nabla_z \phi_i(h(w^{(t)}; i)) \right)$$

Update  $w^{(t+1)} = w^{(t)} - \eta^{(t)} v^{(t)}$

**end for**

---

If  $\varepsilon^2$  is small enough, then  $v_{* \text{reg}}^{(t)}$  is a close approximation of the solution  $v_*^{(t)}$  for problem (15). Our first algorithm updates Framework 1 based on this approximation.

The following Lemma shows the relation between the regularized solution  $v_{* \text{reg}}^{(t)}$  and the optimal solution of the original convex problem  $\hat{v}_{*\varepsilon}^{(t)}$ .

**Lemma 3.** For given  $\varepsilon > 0$ , suppose that Assumption 3 holds for bound  $V > 0$ . Then, for iteration  $0 \leq t < T$ , the optimal solution  $v_{* \text{reg}}^{(t)}$  of problem (20) satisfies  $\|v_{* \text{reg}}^{(t)}\|^2 \leq 2 + V$  and

$$\frac{1}{2} \frac{1}{n} \sum_{i=1}^n \|\eta^{(t)} H_i^{(t)} v_{* \text{reg}}^{(t)} - \alpha_i^{(t)} \nabla_z \phi_i(h(w^{(t)}; i))\|^2 \leq \left(1 + \frac{V}{2}\right) \varepsilon^2. \quad (22)$$

Based on Lemma 3, we guarantee the convergence to a global solution of Algorithm 1 and prove our first theorem. Since it is currently expensive to solve for the exact solution of problem (20), our algorithm serves as a theoretical method to obtain the convergence to a global solution for the finite-sum minimization.

**Theorem 1.** Let  $w^{(t)}$  be generated by Algorithm 1 where we use the closed form solution for the search direction. We execute Algorithm 1 for  $T = \frac{\beta}{\varepsilon}$  outer loops for some constant  $\beta > 0$ . We assume Assumption 1 holds. Suppose that Assumption 2 holds for  $G > 0$  and Assumption 3 holds for  $V > 0$ . We set the step size equal to  $\eta^{(t)} = D\sqrt{\varepsilon}$  for some  $D > 0$  and choose a learning rate  $\alpha_i^{(t)} = (1 + \varepsilon)\alpha_i^{(t-1)} = (1 + \varepsilon)^t \alpha_i^{(0)}$ . Based on  $\beta$ , we define  $\alpha_i^{(0)} = \frac{\alpha}{e^\beta L_\phi}$  with  $\alpha \in (0, \frac{1}{3})$ . Let  $F_*$  be the global minimizer of  $F$ , and  $h_i^* = \arg \min_{z \in \mathbb{R}^c} \phi_i(z)$ ,  $i \in [n]$ . Then

$$\begin{aligned} \frac{1}{T} \sum_{t=0}^{T-1} [F(w^{(t)}) - F_*] &\leq \frac{e^\beta L_\phi (1 + \varepsilon)}{2(1 - 3\alpha)\alpha\beta} \cdot \frac{1}{n} \sum_{i=1}^n \|h(w^{(0)}; i) - h_i^*\|^2 \cdot \varepsilon \\ &\quad + \frac{e^\beta L_\phi (3\varepsilon + 2)}{8\alpha(1 - 3\alpha)} [c(4 + (V + 2)GD^2)^2 + 8 + 4V] \cdot \varepsilon. \end{aligned} \quad (23)$$

We note that  $\beta$  is a constant for the purpose of choosing the number of iterations  $T$ . The analysis can be simplified by choosing  $\beta = 1$  with  $T = \frac{1}{\varepsilon}$ . Notice that the common convergence criteria for finding a stationary point for non-convex problems is  $\frac{1}{T} \sum_{t=1}^T \|\nabla F(w_t)\|^2 \leq O(\varepsilon)$ . This criteria has been widely used in the existing literature for non-convex optimization problems. Our convergence criteria  $\frac{1}{T} \sum_{t=1}^T [F(w_t) - F_*] \leq O(\varepsilon)$  is slightly different, in order to find a global solution for non-convex problems.

Our proof for Theorem 1 is novel and insightful. It is originally motivated by the Gradient Descent update (12) and the convexity of the loss functions  $\phi_i$ . For this reason it may not be a surprise that Algorithm 1 can find an  $\varepsilon$ -global solution after  $\mathcal{O}(\frac{1}{\varepsilon})$  iterations. However, computing the exact solution in every iteration might be extremely challenging, especially when the number of samples  $n$  is large. Therefore, we present a different approach to this problem in the following section.

## 5.2 Approximation using Gradient Descent

In this section, we use Gradient Descent (GD) algorithm to solve the strongly convex problem (20). It is well-known that if  $\psi(x) - \frac{\mu}{2}\|x\|^2$  is convex for  $\forall x \in \mathbb{R}^c$ , then  $\psi(x)$  is  $\mu$ -strongly convex (see e.g. [Nesterov, 2004]). Hence  $\Psi(\cdot)$



is  $\varepsilon^2$ -strongly convex. For each iteration  $t$ , we use GD to find a search direction  $v^{(t)}$  which is sufficiently close to the optimal solution  $v_{* \text{ reg}}^{(t)}$  in that

$$\|v^{(t)} - v_{* \text{ reg}}^{(t)}\| \leq \varepsilon. \quad (24)$$

Our Algorithm 2 is described as follows.

---

**Algorithm 2** Solve the regularized problem using Gradient Descent

---

**Initialization:** Choose an initial point  $w^{(0)} \in \mathbb{R}^d$ , tolerance  $\varepsilon > 0$ ;

**for**  $t = 0, 1, \dots, T - 1$  **do**

    Use Gradient Descent algorithm to solve Problem (20) and find a solution  $v^{(t)}$  that satisfies

$$\|v^{(t)} - v_{* \text{ reg}}^{(t)}\| \leq \varepsilon.$$

    Update  $w^{(t+1)} = w^{(t)} - \eta^{(t)}v^{(t)}$

**end for**

---

Since Algorithm 2 can only approximate a solution within some  $\varepsilon$ -preciseness, we need a supplemental assumption for the analysis of our next Theorem 2:

**Assumption 4.** Let  $H_i^{(t)}$  be the Jacobian matrix defined in Lemma 2. We assume that there exists some constant  $H > 0$  such that, for  $i \in [n]$ ,  $\varepsilon > 0$ , and  $0 \leq t < T$  as in Algorithm 2,

$$\|H_i^{(t)}\| \leq \frac{H}{\sqrt{\varepsilon}}. \quad (25)$$

Assumption 4 requires a mild condition on the bounded Jacobian of  $h(w; i)$ , and the upper bound may depend on  $\varepsilon$ . This flexibility allows us to accommodate a good dependence of  $\varepsilon$  for the theoretical analysis. We are now ready to present our convergence theorem for Algorithm 2.

**Theorem 2.** Let  $w^{(t)}$  be generated by Algorithm 2 where  $v^{(t)}$  satisfies (24). We execute Algorithm 2 for  $T = \frac{\beta}{\varepsilon}$  outer loops for some constant  $\beta > 0$ . We assume Assumption 1 holds. Suppose that Assumption 2 holds for  $G > \bar{0}$ , Assumption 3 holds for  $V > 0$  and Assumption 4 holds for  $H > 0$ . We set the step size equal to  $\eta^{(t)} = D\sqrt{\varepsilon}$  for some  $D > 0$  and choose a learning rate  $\alpha_i^{(t)} = (1 + \varepsilon)\alpha_i^{(t-1)} = (1 + \varepsilon)^t\alpha_i^{(0)}$ . Based on  $\beta$ , we define  $\alpha_i^{(0)} = \frac{\alpha}{e^\beta L_\phi}$  with  $\alpha \in (0, \frac{1}{4})$ . Let  $F_*$  be the global minimizer of  $F$ , and  $h_i^* = \arg \min_{z \in \mathbb{R}^c} \phi_i(z)$ ,  $i \in [n]$ . Then

$$\begin{aligned} \frac{1}{T} \sum_{t=0}^{T-1} [F(w^{(t)}) - F_*] &\leq \frac{e^\beta L_\phi (1 + \varepsilon)}{2(1 - 4\alpha)\alpha\beta} \cdot \frac{1}{n} \sum_{i=1}^n \|h(w^{(0)}; i) - h_i^*\|^2 \cdot \varepsilon \\ &\quad + \frac{e^\beta L_\phi (4\varepsilon + 3)}{2\alpha(1 - 4\alpha)} [D^2 H^2 + c(2 + (V + \varepsilon^2 + 2)GD^2)^2 + 2 + V] \cdot \varepsilon. \end{aligned}$$

Theorem 2 implies Corollary 2 which provides the computational complexity for Algorithm 2. Note that for (Stochastic) Gradient Descent, we derive the complexity in terms of component gradient calculations for the finite-sum problem (1). As an alternative, for Algorithm 2 we compare the number of component gradients in problem (20) where  $\Phi^{(t)}(v) = \frac{1}{n} \sum_{i=1}^n \psi_i^{(t)}(v)$ . Such individual gradient has the following form:

$$\nabla_v \psi_i^{(t)}(v) = \eta^{(t)} H_i^{(t)T} H_i^{(t)} \eta^{(t)} v - \alpha_i^{(t)} \eta^{(t)} H_i^{(t)T} \nabla_z \phi_i(h(w^{(t)}; i)).$$

In machine learning applications, the gradient of  $f(\cdot; i)$  is calculated using automatic differentiation (i.e. backpropagation). Since  $f(\cdot; i)$  is the composition of the network structure  $h(\cdot; i)$  and loss function  $\phi_i(\cdot)$ , this process also computes the Jacobian matrix  $H_i^{(t)}$  and the gradient  $\nabla_z \phi_i(h(w^{(t)}; i))$  at a specific weight  $w^{(t)}$ . Since matrix-vector multiplication computation is not expensive, the cost for computing the component gradient of problem (20) is similar to problem (1).

**Corollary 2.** Suppose that the conditions in Theorem 2 hold with  $\eta^{(t)} = \frac{D\sqrt{\varepsilon}}{\sqrt{N}}$  for some  $D > 0$  and  $0 < \hat{\varepsilon} \leq N$  (that is, we set  $\varepsilon = \hat{\varepsilon}/N$ ), where

$$N = \frac{e^\beta L_\phi \sum_{i=1}^n \|h(w^{(0)}; i) - h_i^*\|^2}{n(1 - 4\alpha)\alpha\beta} + \frac{7e^\beta L_\phi [D^2 H^2 + c(2 + (V + 3)GD^2)^2 + 2 + V]}{2\alpha(1 - 4\alpha)}.$$

Then, the total complexity to guarantee

$$\min_{0 \leq t \leq T-1} [F(w^{(t)}) - F_*] \leq \frac{1}{T} \sum_{t=0}^{T-1} [F(w^{(t)}) - F_*] \leq \hat{\varepsilon}$$

is  $\mathcal{O}\left(n \frac{N^3 \beta}{\varepsilon^3} (D^2 H^2 + (\hat{\varepsilon}^2/N)) \log\left(\frac{N}{\hat{\varepsilon}}\right)\right)$ .

**Remark 2.** Corollary 2 shows that  $\mathcal{O}(1/\hat{\varepsilon})$  outer loop iterations are needed in order to reach an  $\hat{\varepsilon}$ -global solution, and it proves that each iteration needs the equivalent of  $\mathcal{O}\left(\frac{n}{\varepsilon^2} \log\left(\frac{1}{\hat{\varepsilon}}\right)\right)$  gradient computations for computing an approximate solution. In total, Algorithm 2 has total complexity  $\mathcal{O}\left(\frac{n}{\varepsilon^3} \log\left(\frac{1}{\hat{\varepsilon}}\right)\right)$  for finding an  $\hat{\varepsilon}$ -global solution.

For a comparison, Stochastic Gradient Descent uses a total of  $\mathcal{O}\left(\frac{1}{\varepsilon^2}\right)$  gradient computations to find a stationary point satisfying  $\mathbb{E}[\|\nabla F(\hat{w})\|^2] \leq \varepsilon$  for non-convex problems [Ghadimi and Lan, 2013]. Gradient Descent has a better complexity in terms of  $\varepsilon$ , i.e.  $\mathcal{O}\left(\frac{n}{\varepsilon}\right)$  such that  $\|\nabla F(\hat{w})\|^2 \leq \varepsilon$  [Nesterov, 2004]. However, both methods may not be able to reach a global solution of (1). In order to guarantee global convergence for nonconvex settings, one may resort to use Polyak-Lojasiewicz (PL) inequality [Karimi et al., 2016, Gower et al., 2021]. This assumption is widely known to be strong, which implies that every stationary point is also a global minimizer.

## 6 Further Discussion and Conclusions

This paper presents an alternative *composite formulation* for solving the finite-sum optimization problem. Our formulation allows a new way of exploiting the structure of machine learning problems and the convexity of squared loss and softmax cross entropy loss, and leads to a novel algorithmic framework that guarantees convergence to a global solution (when the outer loss functions are convex and Lipschitz-smooth).

Our analysis is general and can be applied to various different learning architectures, in particular, our analysis and assumptions match practical neural networks; in recent years, there has been a great interest in the structure of deep learning architectures for over-parameterized settings [Arora et al., 2018, Allen-Zhu et al., 2019, Nguyen and Mondelli, 2020]. Algorithm 2 demonstrates a gradient method to solve the regularized problem, however, other methods can be applied to our framework (e.g. conjugate gradient descent).

Our theoretical foundation motivates further study, implementation, and optimization of the new algorithmic framework and further investigation of its non-standard bounded style assumptions. Possible research directions include more practical algorithm designs based on our Framework 1, and different related methods to solve the regularized problem and approximate the solution such as Stochastic Gradient Descent and its stochastic first-order variants (e.g. [Duchi et al., 2011, Kingma and Ba, 2014, Bottou et al., 2018, Nguyen et al., 2018, 2019, 2021]). This potentially leads to a new class of efficient algorithms for machine learning problems. This paper presents a new perspective to the research community.

## Appendix

### A Useful Results

The following lemmas provide key tools for our results.

**Lemma 4** (Squared loss). *Let  $b \in \mathbb{R}^c$  and define  $\phi(z) = \frac{1}{2}\|z - b\|^2$  for  $z \in \mathbb{R}^c$ . Then  $\phi$  is convex and  $L_\phi$ -smooth with  $L_\phi = 1$ .*

**Lemma 5** (Softmax cross-entropy loss). *Let index  $a \in [c]$  and define*

$$\phi(z) = \log \left[ \sum_{k=1}^c \exp(z_k - z_a) \right] = \log \left[ \sum_{k=1}^c \exp(w_k^T z) \right],$$

for  $z = (z_1, \dots, z_c)^T \in \mathbb{R}^c$ , where  $w_k = e_k - e_a$  with  $e_i$  representing the  $i$ -th unit vector (containing 1 at the  $i$ -th position and 0 elsewhere). Then  $\phi$  is convex and  $L_\phi$ -smooth with  $L_\phi = 1$ .

### B Additional Discussion

#### B.1 About Assumption 2

We make a formal assumption for the case  $h(\cdot; i)$  is closely approximated by  $k(\cdot; i)$ .

**Assumption 5.** We assume that for all  $i \in [n]$  there exists some approximations  $k(w; i) : \mathbb{R}^d \rightarrow \mathbb{R}^c$  such that

$$|k_j(w; i) - h_j(w; i)| \leq \varepsilon, \forall w \in \mathbb{R}^d, i \in [n] \text{ and } j \in [c], \quad (26)$$

where  $k(\cdot; i)$  are twice continuously differentiable (i.e. the second-order partial derivatives of all scalars  $k_j(\cdot; i)$  are continuous for all  $i \in [n]$ ), and that their Hessian matrices are bounded:

$$\|M_{i,j}(w)\| = \|\mathbf{J}_w (\nabla_w k_j(w; i))\| \leq G, \forall w \in \mathbb{R}^d, i \in [n] \text{ and } j \in [c]. \quad (27)$$

Assumption 5 allows us to prove the following Lemma that bound the error in equation (14):

**Lemma 6.** Suppose that Assumption 5 holds for the classifier  $h$ . Then for all  $i \in [n]$  and  $0 \leq t < T$ , we have:

$$h(w^{(t+1)}; i) = h(w^{(t)} - \eta^{(t)} v^{(t)}; i) = h(w^{(t)}; i) - \eta^{(t)} H_i^{(t)} v^{(t)} + \epsilon_i^{(t)}, \quad (28)$$

where  $H_i^{(t)}$  is defined to be the Jacobian matrix of the approximation  $k(w; i)$  at  $w^{(t)}$ :

$$H_i^{(t)} := \mathbf{J}_w k(w; i)|_{w=w^{(t)}} = \left[ \begin{array}{ccc} \frac{\partial k_1(w; i)}{\partial w_1} & \cdots & \frac{\partial k_1(w; i)}{\partial w_d} \\ \cdots & \cdots & \cdots \\ \frac{\partial k_c(w; i)}{\partial w_1} & \cdots & \frac{\partial k_c(w; i)}{\partial w_d} \end{array} \right] \Big|_{w=w^{(t)}} \in \mathbb{R}^{c \times d}. \quad (29)$$

Additionally we have,

$$|\epsilon_{i,j}^{(t)}| \leq \frac{1}{2} (\eta^{(t)})^2 \|v^{(t)}\|^2 G + 2\varepsilon, j \in [c]. \quad (30)$$

Note that these result recover the case when  $h(\cdot; i)$  is itself smooth. Hence we analyze our algorithms using the result of Lemma 6, which generalizes the result from Lemma 2.

## B.2 About Assumption 3

In this section, we justify the existence of the search direction in Assumption 3 (almost surely). We argue that there exists a vector  $\hat{v}_{*\varepsilon}^{(t)}$  satisfying

$$\frac{1}{2} \frac{1}{n} \sum_{i=1}^n \|\eta^{(t)} H_i^{(t)} \hat{v}_{*\varepsilon}^{(t)} - \alpha_i^{(t)} \nabla_z \phi_i(h(w^{(t)}; i))\|^2 \leq \varepsilon^2.$$

It is sufficient to find a vector  $v$  satisfying that

$$\eta^{(t)} H_i^{(t)} v = \alpha_i^{(t)} \nabla_z \phi_i(h(w^{(t)}; i)) \text{ for every } i \in [n].$$

Since the solution  $v$  is in  $\mathbb{R}^d$  and  $\nabla_z \phi_i(h(w^{(t)}; i))$  is in  $\mathbb{R}^c$ , this condition is equivalent to a linear system with  $n \cdot c$  constraints and  $d$  variables. Let  $A$  and  $b$  be the following stacked matrix and vector:

$$A = \begin{bmatrix} H_1^{(t)} \eta^{(t)} \\ \cdots \\ H_n^{(t)} \eta^{(t)} \end{bmatrix} \in \mathbb{R}^{n \cdot c \times d}, \text{ and } b = \begin{bmatrix} \alpha_1^{(t)} \nabla_z \phi_1(h(w^{(t)}; i)) \\ \cdots \\ \alpha_n^{(t)} \nabla_z \phi_n(h(w^{(t)}; i)) \end{bmatrix} \in \mathbb{R}^{n \cdot c},$$

then the problem reduce to finding the solution of the equation  $Av = b$ . In the over-parameterized setting where dimension  $d$  is sufficiently large ( $d \gg n \cdot c$ ), then  $\text{rank } A = n \cdot c$  almost surely and there exists almost surely a vector  $v$  that interpolates all the training set.

To demonstrate this fact easier, we consider a simple neural network where the classifier  $h(w; i)$  is formulated as

$$h(w; i) = W^{(2)T} \sigma(W^{(1)T} x^{(i)}),$$

where  $c = 1$ ,  $W^{(1)} \in \mathbb{R}^{m \times l}$  and  $W^{(2)} \in \mathbb{R}^{l \times 1}$ ,  $w = \mathbf{vec}(\{W^{(1)}, W^{(2)}\}) \in \mathbb{R}^d$  is the vectorized weight where  $d = l(m + 1)$  and  $\sigma$  is sigmoid activation function.

$H_i^{(t)}$  is defined to be the Jacobian matrix of  $h(w; i)$  at  $w^{(t)}$ :

$$H_i^{(t)} := \mathbf{J}_w h(w; i)|_{w=w^{(t)}} = \left[ \begin{array}{ccc} \frac{\partial h(w; i)}{\partial w_1} & \cdots & \frac{\partial h(w; i)}{\partial w_d} \end{array} \right] \Big|_{w=w^{(t)}} \in \mathbb{R}^{1 \times d},$$

then

$$A = \eta^{(t)} \begin{bmatrix} H_1^{(t)} \\ \dots \\ H_n^{(t)} \end{bmatrix} = \eta^{(t)} \begin{bmatrix} \frac{\partial h(w;1)}{\partial w_1} & \dots & \frac{\partial h(w;1)}{\partial w_d} \\ \dots & \dots & \dots \\ \frac{\partial h(w;n)}{\partial w_1} & \dots & \frac{\partial h(w;n)}{\partial w_d} \end{bmatrix} \in \mathbb{R}^{n \times d}.$$

We want to show that  $A$  has full rank, almost surely. We consider the over-parameterized setting where the last layer has at least  $n$  neuron (i.e.  $l = n$  and the simple version when  $c = 1$ ). We argue that rank of matrix  $A$  is greater than or equal to rank of the submatrix  $B$  created by the weights of the last layer  $W^{(2)} \in \mathbb{R}^n$ :

$$B = \begin{bmatrix} \frac{\partial h(w;1)}{\partial W_1^{(2)}} & \dots & \frac{\partial h(w;1)}{\partial W_n^{(2)}} \\ \dots & \dots & \dots \\ \frac{\partial h(w;n)}{\partial W_1^{(2)}} & \dots & \frac{\partial h(w;n)}{\partial W_n^{(2)}} \end{bmatrix} \in \mathbb{R}^{n \times n}.$$

Note that  $h(\cdot, i)$  is a linear function of the last weight layers (in this simple case  $W^{(2)} \in \mathbb{R}^n$  and  $\sigma(W^{(1)T}x^{(i)}) \in \mathbb{R}^n$ ), we can compute the partial derivatives as follows:

$$\frac{\partial h(w; i)}{\partial W^{(2)}} = \sigma(W^{(1)T}x^{(i)}); \quad i \in [n].$$

Hence

$$B = \begin{bmatrix} \sigma(W^{(1)T}x^{(1)}) \\ \dots \\ \sigma(W^{(1)T}x^{(n)}) \end{bmatrix} \in \mathbb{R}^{n \times n}.$$

Assuming that there are no identical data, and  $\sigma$  is the sigmoid activation, the set of weights  $W^{(1)}$  that make matrix  $B$  degenerate has measure zero. Hence  $B$  has full rank almost surely, and we have the same conclusion for  $A$ . Therefore we are able to prove the almost surely existence of a solution  $v$  of the linear equation  $Av = b$  for simple two layers network. Using the same argument, this result can be generalized for larger neural networks where the dimension  $d$  is sufficiently large ( $d \gg nc$ ).

### B.3 Initialization example

Our Assumption 3 requires a nice dependency on the tolerance  $\varepsilon$  for the gradient matrices  $H_i^{(0)}$  and  $\nabla_z \phi_i(h(w^{(0)}; i))$ . We note that at the starting point  $t = 0$ , these matrices may depend on  $\varepsilon$  due to the initialization process and the dependence of  $d$  on  $\varepsilon$ . In order to accommodate the choice of learning rate  $\eta^{(0)} = D\sqrt{\varepsilon}$  in our theorems, in this section we describe a network initialization that satisfies  $\|H_i^{(0)}\| = \Theta\left(\frac{1}{\sqrt{\varepsilon}}\right)$  where the gradient norm  $\|\nabla_z \phi_i(h(w^{(0)}; i))\|$  is at most constant order with respect to  $\varepsilon$ . To simplify the problem, we only consider small-dimension data and networks without activation.

**About the target vector:** We choose  $\phi_i$  to be the softmax cross-entropy loss. By Lemma 7 (see below), we have that the gradient norm is upper bounded by a constant  $c$ , where  $c$  is the output dimension of the problem and is not dependent on  $\varepsilon$ . Note that when we stack all gradients for  $n$  data points, then the size of new vector is still not dependent on  $\varepsilon$ .

**About the network architecture:** For simplicity, we consider the following classification problem where

- The input data is in  $\mathbb{R}^2$ . There are only two data points  $\{x^{(1)}, x^{(2)}\}$ . Input data is bounded and non-degenerate (we will clarify this property later).
- The output data is (categorical) in  $\mathbb{R}^2$ :  $\{y^{(1)} = (1, 0), y^{(2)} = (0, 1)\}$ .

We want to have an over-parameterized setting where the dimension of weight vector is at least  $nc = 4$ . We consider a simple network with two layers, no biases and no activation functions. Let the number of neurons in the hidden layer be  $m$ . The flow of this network is (in)  $\mathbb{R}^2 \rightarrow \mathbb{R}^m \rightarrow \mathbb{R}^2$  (out). First, we consider the case where  $m = 1$ .

- The first layer has 2 parameters  $(w_1, w_2)$  and only 1 neuron that outputs  $z^{(i)} = w_1x_1^{(i)} + w_2x_2^{(i)}$  (the subscript is for the coordinate of input data  $x^{(i)}$ ).

- The second layer has 2 parameters  $(w_3, w_4)$ . The final output is

$$h(w, i) = [w_3(w_1x_1^{(i)} + w_2x_2^{(i)}), w_4(w_1x_1^{(i)} + w_2x_2^{(i)})]^T \in \mathbb{R}^2,$$

with  $w = [w_1, w_2, w_3, w_4]^T \in \mathbb{R}^4$ . This network satisfies that the Hessian matrices of  $h(w; i)$  are bounded. Let  $Q$  and  $b$  be the following stacked matrix and vector:

$$Q = \begin{bmatrix} H_1^{(0)} \\ H_2^{(0)} \end{bmatrix} \in \mathbb{R}^{4 \times 4}, \text{ and } b = \begin{bmatrix} \nabla_z \phi_1(h(w^{(0)}; 1)) \\ \nabla_z \phi_2(h(w^{(0)}; 2)) \end{bmatrix} \in \mathbb{R}^4,$$

Then we have the following:

$$\begin{aligned} Q = Q(w) &= \begin{bmatrix} H_1^{(0)} \\ H_2^{(0)} \end{bmatrix} = \begin{bmatrix} \nabla_w [w_3(w_1x_1^{(1)} + w_2x_2^{(1)})] \\ \nabla_w [w_4(w_1x_1^{(1)} + w_2x_2^{(1)})] \\ \nabla_w [w_3(w_1x_1^{(2)} + w_2x_2^{(2)})] \\ \nabla_w [w_4(w_1x_1^{(2)} + w_2x_2^{(2)})] \end{bmatrix} \\ &= \begin{bmatrix} w_3x_1^{(1)} & w_3x_2^{(1)} & w_1x_1^{(1)} + w_2x_2^{(1)} & 0 \\ w_4x_1^{(1)} & w_4x_2^{(1)} & 0 & w_1x_1^{(1)} + w_2x_2^{(1)} \\ w_3x_1^{(2)} & w_3x_2^{(2)} & w_1x_1^{(2)} + w_2x_2^{(2)} & 0 \\ w_4x_1^{(2)} & w_4x_2^{(2)} & 0 & w_1x_1^{(2)} + w_2x_2^{(2)} \end{bmatrix}. \end{aligned}$$

The determinant of this matrix is a polynomial of the weight  $w$  and the input data. Under some mild non-degenerate condition of the input data, we can choose some base point  $w'$  that made this matrix invertible (note that if this condition is not satisfied, we can rescale/add a very small noise to the data - which is the common procedure in machine learning).

Hence the system  $Qu = b$  always has a solution. Now we consider the following two initializations:

1. We choose to initialize the starting point at  $w^{(0)} = \frac{1}{\sqrt{\varepsilon}}w'$  and note that  $Q(w)$  is a linear function of  $w$  and  $Q(w')$  is independent of  $\varepsilon$ . Then the norm of matrix  $Q(w^{(0)})$  has the same scale with  $\frac{1}{\sqrt{\varepsilon}}$ .

2. Instead of choosing  $m = 1$ , we consider an over-parameterized network where  $m = \frac{1}{\varepsilon}$  (recall that  $m$  is the number of neurons in the hidden layer). The hidden layer in this case is:

$$z = \begin{cases} z_1^{(i)} &= w_{1,1}^{(1)}x_1^{(i)} + w_{2,1}^{(1)}x_2^{(i)} \\ &\dots \\ z_m^{(i)} &= w_{1,m}^{(1)}x_1^{(i)} + w_{2,m}^{(1)}x_2^{(i)} \end{cases}.$$

The output layer is:

$$\begin{cases} y_1^{(i)} &= z_1^{(i)}w_{1,1}^{(2)} + \dots + z_m^{(i)}w_{m,1}^{(2)} = (w_{1,1}^{(1)}x_1^{(i)} + w_{2,1}^{(1)}x_2^{(i)})w_{1,1}^{(2)} + \dots + (w_{1,m}^{(1)}x_1^{(i)} + w_{2,m}^{(1)}x_2^{(i)})w_{m,1}^{(2)} \\ y_2^{(i)} &= z_1^{(i)}w_{1,2}^{(2)} + \dots + z_m^{(i)}w_{m,2}^{(2)} = (w_{1,1}^{(1)}x_1^{(i)} + w_{2,1}^{(1)}x_2^{(i)})w_{1,2}^{(2)} + \dots + (w_{1,m}^{(1)}x_1^{(i)} + w_{2,m}^{(1)}x_2^{(i)})w_{m,2}^{(2)} \end{cases}$$

with  $w = [w_{1,1}^{(1)}, \dots, w_{1,m}^{(1)}, w_{2,1}^{(1)}, \dots, w_{2,m}^{(1)}, w_{1,1}^{(2)}, w_{1,2}^{(2)}, \dots, w_{m,1}^{(2)}, w_{m,2}^{(2)}]^T \in \mathbb{R}^{4m}$ .

Hence,

$$Q(w) = \begin{bmatrix} w_{1,1}^{(2)}x_1^{(1)} & \dots & w_{m,1}^{(2)}x_1^{(1)} & w_{1,1}^{(2)}x_2^{(1)} & \dots & w_{m,1}^{(2)}x_2^{(1)} & z_1^{(1)} & 0 & \dots & z_m^{(1)} & 0 \\ w_{1,2}^{(2)}x_1^{(1)} & \dots & w_{m,2}^{(2)}x_1^{(1)} & w_{1,2}^{(2)}x_2^{(1)} & \dots & w_{m,2}^{(2)}x_2^{(1)} & 0 & z_1^{(1)} & \dots & 0 & z_m^{(1)} \\ w_{1,1}^{(2)}x_1^{(2)} & \dots & w_{m,1}^{(2)}x_1^{(2)} & w_{1,1}^{(2)}x_2^{(2)} & \dots & w_{m,1}^{(2)}x_2^{(2)} & z_1^{(2)} & 0 & \dots & z_m^{(2)} & 0 \\ w_{1,2}^{(2)}x_1^{(2)} & \dots & w_{m,2}^{(2)}x_1^{(2)} & w_{1,2}^{(2)}x_2^{(2)} & \dots & w_{m,2}^{(2)}x_2^{(2)} & 0 & z_1^{(2)} & \dots & 0 & z_m^{(2)} \end{bmatrix}.$$

Hence, the number of (possibly) non-zero elements in each row is  $3m = \frac{3}{\varepsilon}$ .

For matrix  $A$  of rank  $r$ , we have  $\|A\|_2 \leq \|A\|_F \leq \sqrt{r}\|A\|_2$ . Since the rank of  $Q(w)$  is at most 4 ( $nc = 4$ , independent of  $\varepsilon$ ), we only need to find the Frobenius norm of  $Q(w)$ . We have

$$\|Q(w)\|_F = \sqrt{\sum_{i=1}^4 \sum_{j=1}^{4m} |q_{ij}|^2}.$$

Let  $q_{min}$  and  $q_{max}$  be the element with smallest/largest magnitude of  $Q(w)$ . Suppose that  $x^{(i)} \neq (0, 0)$  and choose  $w \neq 0$  such that  $z \neq 0$ ,  $q_{min} > 0$  and independent of  $\varepsilon$ . Hence,  $\frac{\sqrt{8}}{\sqrt{\varepsilon}}|q_{min}| \leq \|Q(w)\|_F \leq \frac{\sqrt{12}}{\sqrt{\varepsilon}}|q_{max}|$ .

Hence,  $\|Q(w)\| = \Theta\left(\frac{1}{\sqrt{\varepsilon}}\right)$ . Therefore this simple network initialization supports the dependence on  $\varepsilon$  for our Assumption 3. We note that a similar setting is found in [Allen-Zhu et al., 2019], where the authors initialize the weights using a random Gaussian distribution with a variance depending on the dimension of the problem. In non-convex setting, they prove the convergence of SGD using the assumption that the number of neurons  $m$  depends inversely on the tolerance  $\varepsilon$ .

**Lemma 7.** For softmax cross-entropy loss, and  $x = h(w; i) \in \mathbb{R}^c$ , for  $\forall w \in \mathbb{R}^d$  and  $i \in [n]$ , we have

$$\left\| \nabla_z \phi_i(x) \Big|_{x=h(w; i)} \right\|^2 \leq c. \quad (31)$$

*Proof.* By (10), we have for  $i = 1, \dots, n$ ,

- For  $j \neq I(y^{(i)})$ :

$$\begin{aligned} \left( \frac{\partial \phi_i(x)}{\partial x_j} \Big|_{x=h(w; i)} \right)^2 &= \left( \frac{\exp([h(w; i)]_j - [h(w; i)]_{I(y^{(i)})})}{\sum_{k=1}^c \exp([h(w; i)]_k - [h(w; i)]_{I(y^{(i)})})} \right)^2 \\ &= \left( \frac{\exp([h(w; i)]_j - [h(w; i)]_{I(y^{(i)})})}{1 + \sum_{k \neq I(y^{(i)})} \exp([h(w; i)]_k - [h(w; i)]_{I(y^{(i)})})} \right)^2 \leq 1. \end{aligned}$$

- For  $j = I(y^{(i)})$ :

$$\begin{aligned} \left( \frac{\partial \phi_i(x)}{\partial x_j} \Big|_{x=h(w; i)} \right)^2 &= \left( \frac{\sum_{k \neq I(y^{(i)})} \exp([h(w; i)]_k - [h(w; i)]_{I(y^{(i)})})}{\sum_{k=1}^c \exp([h(w; i)]_k - [h(w; i)]_{I(y^{(i)})})} \right)^2 \\ &= \left( \frac{\sum_{k \neq I(y^{(i)})} \exp([h(w; i)]_k - [h(w; i)]_{I(y^{(i)})})}{1 + \sum_{k \neq I(y^{(i)})} \exp([h(w; i)]_k - [h(w; i)]_{I(y^{(i)})})} \right)^2 \leq 1 \end{aligned}$$

Hence, for  $i = 1, \dots, n$ ,

$$\left\| \nabla_z \phi_i(x) \Big|_{x=h(w; i)} \right\|^2 = \sum_{j=1}^c \left( \frac{\partial \phi_i(x)}{\partial x_j} \Big|_{x=h(w; i)} \right)^2 \leq c.$$

This completes the proof.  $\square$

## C Proofs of Lemmas and Corollary 1

### Proof of Lemma 2

*Proof.* Since  $h(\cdot; i)$  are twice continuously differentiable for all  $i \in [n]$ , we have the following Taylor approximation for each component outputs  $h_j(\cdot; i)$  where  $j \in [c]$  and  $i \in [n]$ :

$$\begin{aligned} h_j(w^{(t+1)}; i) &= h_j(w^{(t)} - \eta^{(t)} v^{(t)}; i) \\ &= h_j(w^{(t)}; i) - \mathbf{J}_w h_j(w; i)|_{w=w^{(t)}} \eta^{(t)} v^{(t)} + \frac{1}{2} (\eta^{(t)} v^{(t)})^T M_{i,j}(\tilde{w}^{(t)}) (\eta^{(t)} v^{(t)}), \end{aligned} \quad (32)$$

where  $M_{i,j}(\tilde{w}^{(t)})$  is the Hessian matrices of  $h_j(\cdot; i)$  at  $\tilde{w}^{(t)}$  and  $\tilde{w}^{(t)} = \alpha w^{(t)} + (1 - \alpha) w^{(t+1)}$  for some  $\alpha \in [0, 1]$ . This leads to our desired statement:

$$h(w^{(t+1)}; i) = h(w^{(t)} - \eta^{(t)} v^{(t)}; i) = h(w^{(t)}; i) - \eta^{(t)} H_i^{(t)} v^{(t)} + \epsilon_i^{(t)},$$

where

$$\epsilon_{i,j}^{(t)} = \frac{1}{2} (\eta^{(t)} v^{(t)})^T M_{i,j}(\tilde{w}^{(t)}) (\eta^{(t)} v^{(t)}), \quad j \in [c],$$

Hence we get the final bound:

$$\begin{aligned}
|\epsilon_{i,j}^{(t)}| &\leq \frac{1}{2} \left| (\eta^{(t)} v^{(t)})^T M_{i,j}(\tilde{w}^{(t)}) (\eta^{(t)} v^{(t)}) \right| \\
&\leq \frac{1}{2} (\eta^{(t)})^2 \|v^{(t)}\|^2 \cdot \|M_{i,j}(\tilde{w}^{(t)})\| \\
&\stackrel{(16)}{\leq} \frac{1}{2} (\eta^{(t)})^2 \|v^{(t)}\|^2 G, \quad j \in [c].
\end{aligned}$$

□

### Proof of Lemma 3

*Proof.* From Assumption 3, we know that there exists  $\hat{v}_{*\varepsilon}^{(t)}$  so that

$$\frac{1}{2} \frac{1}{n} \sum_{i=1}^n \|\eta^{(t)} H_i^{(t)} \hat{v}_{*\varepsilon}^{(t)} - \alpha_i^{(t)} \nabla_z \phi_i(h(w^{(t)}; i))\|^2 \leq \varepsilon^2,$$

and  $\|\hat{v}_{*\varepsilon}^{(t)}\|^2 \leq V$ , for some  $V > 0$ . Hence,

$$\frac{1}{2} \frac{1}{n} \sum_{i=1}^n \|\eta^{(t)} H_i^{(t)} \hat{v}_{*\varepsilon}^{(t)} - \alpha_i^{(t)} \nabla_z \phi_i(h(w^{(t)}; i))\|^2 + \frac{\varepsilon^2}{2} \|\hat{v}_{*\varepsilon}^{(t)}\|^2 \leq \varepsilon^2 + \frac{\varepsilon^2}{2} V = (1 + \frac{V}{2}) \varepsilon^2.$$

Since  $v_{*\text{reg}}^{(t)}$  is the optimal solution of the problem in (20) for  $0 \leq t < T$ , we have

$$\frac{1}{2} \frac{1}{n} \sum_{i=1}^n \|\eta^{(t)} H_i^{(t)} v_{*\text{reg}}^{(t)} - \alpha_i^{(t)} \nabla_z \phi_i(h(w^{(t)}; i))\|^2 + \frac{\varepsilon^2}{2} \|v_{*\text{reg}}^{(t)}\|^2 \leq (1 + \frac{V}{2}) \varepsilon^2.$$

Therefore, we have (22) and  $\|v_{*\text{reg}}^{(t)}\|^2 \leq 2 + V$  for  $0 \leq t < T$ .

□

### Proof of Lemma 4

*Proof.* 1. We want to show that for any  $\alpha \in [0, 1]$

$$\phi(\alpha z_1 + (1 - \alpha) z_2) \leq \alpha \phi(z_1) + (1 - \alpha) \phi(z_2), \quad \forall z_1, z_2 \in \mathbb{R}^c, \quad (33)$$

in order to have the convexity of  $\phi$  with respect to  $z$  (see [Nesterov, 2004]).

For any  $\alpha \in [0, 1]$ , we have for  $\forall z_1, z_2 \in \mathbb{R}^c$ ,

$$\begin{aligned}
&\alpha \|z_1 - b\|^2 + (1 - \alpha) \|z_2 - b\|^2 - \|\alpha(z_1 - b) + (1 - \alpha)(z_2 - b)\|^2 \\
&= \alpha \|z_1 - b\|^2 + (1 - \alpha) \|z_2 - b\|^2 - \alpha^2 \|z_1 - b\|^2 - (1 - \alpha)^2 \|z_2 - b\|^2 \\
&\quad - 2\alpha(1 - \alpha) \langle z_1 - b, z_2 - b \rangle \\
&\geq \alpha(1 - \alpha) \|z_1 - b\|^2 + (1 - \alpha)\alpha \|z_2 - b\|^2 - 2\alpha(1 - \alpha) \|z_1 - b\| \cdot \|z_2 - b\| \\
&= \alpha(1 - \alpha) (\|z_1 - b\| - \|z_2 - b\|)^2 \geq 0,
\end{aligned}$$

where the first inequality follows according to Cauchy-Schwarz inequality  $\langle a, b \rangle \leq \|a\| \cdot \|b\|$ . Hence,

$$\frac{1}{2} \|\alpha z_1 + (1 - \alpha) z_2 - b\|^2 \leq \frac{\alpha}{2} \|z_1 - b\|^2 + \frac{(1 - \alpha)}{2} \|z_2 - b\|^2.$$

Therefore, (33) implies the convexity of  $\phi$  with respect to  $z$ .

2. We want to show that  $\exists L_\phi > 0$  such that

$$\|\nabla \phi(z_1) - \nabla \phi(z_2)\| \leq L_\phi \|z_1 - z_2\|, \quad \forall z_1, z_2 \in \mathbb{R}^c. \quad (34)$$

Notice that  $\nabla \phi(z) = z - b$ , then clearly  $\forall z_1, z_2 \in \mathbb{R}^c$ ,

$$\|\nabla \phi(z_1) - \nabla \phi(z_2)\| = \|z_1 - z_2\|.$$

Therefore, (34) implies the  $L_\phi$ -smoothness of  $\phi$  with respect to  $z$  with  $L_\phi = 1$ .

□

## Proof of Lemma 5

*Proof.* 1. For  $\forall z_1, z_2 \in \mathbb{R}^c$  and  $1 \leq k \leq c$ , denote  $u_{k,1} = \exp(w_k^T z_1)$  and  $u_{k,2} = \exp(w_k^T z_2)$  and using Holder inequality

$$\sum_{k=1}^c a_k \cdot b_k \leq \left( \sum_{k=1}^c |a_k|^p \right)^{\frac{1}{p}} \left( \sum_{k=1}^c |b_k|^q \right)^{\frac{1}{q}}, \text{ where } \frac{1}{p} + \frac{1}{q} = 1, \quad (35)$$

we have

$$\begin{aligned} \phi(\alpha z_1 + (1-\alpha)z_2) &= \log \left[ \sum_{k=1}^c \exp(w_k^T (\alpha z_1 + (1-\alpha)z_2)) \right] = \log \left[ \sum_{k=1}^c u_{k,1}^\alpha \cdot u_{k,2}^{(1-\alpha)} \right] \\ &\stackrel{(35)}{\leq} \log \left[ \left( \sum_{k=1}^c u_{k,1}^{\alpha \cdot \frac{1}{\alpha}} \right)^\alpha \left( \sum_{k=1}^c u_{k,2}^{(1-\alpha) \cdot \frac{1}{(1-\alpha)}} \right)^{1-\alpha} \right] \\ &= \alpha \log \left[ \sum_{k=1}^c \exp(w_k^T z_1) \right] + (1-\alpha) \log \left[ \sum_{k=1}^c \exp(w_k^T z_2) \right] \\ &= \alpha \phi(z_1) + (1-\alpha) \phi(z_2), \end{aligned}$$

where the first inequality since  $\log(x)$  is an increasing function for  $\forall x > 0$  and  $\exp(v) > 0$  for  $\forall v \in \mathbb{R}$ . Therefore, (33) implies the convexity of  $\phi$  with respect to  $z$ .

2. Note that  $\|\nabla^2 \phi(z)\| \leq L_\phi$  if and only if  $\phi(z)$  is  $L_\phi$ -smooth (see [Nesterov, 2004]). First, we compute gradient of  $\phi(z)$ :

- For  $i \neq a$ :

$$\frac{\partial \phi(z)}{\partial z_i} = \frac{\exp(z_i - z_a)}{\sum_{k=1}^c \exp(z_k - z_a)}.$$

- For  $i = a$ :

$$\begin{aligned} \frac{\partial \phi(z)}{\partial z_i} &= \frac{-\sum_{k \neq a} \exp(z_k - z_a)}{\sum_{k=1}^c \exp(z_k - z_a)} = \frac{-\sum_{k=1}^c \exp(z_k - z_a) + 1}{\sum_{k=1}^c \exp(z_k - z_a)} \\ &= -1 + \frac{1}{\sum_{k=1}^c \exp(z_k - z_a)} = -1 + \frac{\exp(z_i - z_a)}{\sum_{k=1}^c \exp(z_k - z_a)}. \end{aligned}$$

We then calculate  $\frac{\partial^2 \phi(z)}{\partial z_j \partial z_i} = \frac{\partial}{\partial z_j} \left( \frac{\partial \phi(z)}{\partial z_i} \right)$

- For  $i = j$ :

$$\begin{aligned} \frac{\partial^2 \phi(z)}{\partial z_j \partial z_i} &= \frac{\exp(z_i - z_a) [\sum_{k=1}^c \exp(z_k - z_a)] - \exp(z_i - z_a) \exp(z_i - z_a)}{[\sum_{k=1}^c \exp(z_k - z_a)]^2} \\ &= \frac{\exp(z_i - z_a) [\sum_{k=1}^c \exp(z_k - z_a) - \exp(z_i - z_a)]}{[\sum_{k=1}^c \exp(z_k - z_a)]^2}. \end{aligned}$$

- For  $i \neq j$ :

$$\frac{\partial^2 \phi(z)}{\partial z_j \partial z_i} = \frac{-\exp(z_j - z_a) \exp(z_i - z_a)}{[\sum_{k=1}^c \exp(z_k - z_a)]^2}.$$

Denote that  $y_i = \exp(z_i - z_a) \geq 0$ ,  $i \in [c]$ , we have:

- For  $i = j$ :

$$\left| \frac{\partial^2 \phi(z)}{\partial z_j \partial z_i} \right| = \left| \frac{y_i (\sum_{k=1}^c y_k - y_i)}{(\sum_{k=1}^c y_k)^2} \right|.$$



- For  $i \neq j$ :

$$\left| \frac{\partial^2 \phi(z)}{\partial z_j \partial z_i} \right| = \frac{|y_i y_j|}{\left( \sum_{k=1}^c y_k \right)^2}.$$

Recall that for matrix  $A = (a_{ij}) \in \mathbb{R}^{c \times c}$ :  $\|A\|^2 \leq \|A\|_F^2 = \sum_{i=1}^c \sum_{j=1}^c |a_{ij}|^2$ . We have:

$$\begin{aligned} \sum_{j=1}^c \left| \frac{\partial^2 \phi(z)}{\partial z_j \partial z_i} \right|^2 &\leq \frac{1}{\left( \sum_{k=1}^c y_k \right)^4} \left[ y_i^2 \left( \sum_{k=1}^c y_k - y_i \right)^2 + \sum_{j \neq i} (y_i y_j)^2 \right] \\ &= \frac{1}{\left( \sum_{k=1}^c y_k \right)^4} \left[ y_i^2 \left( \sum_{k=1}^c y_k \right)^2 - 2y_i^2 \sum_{k=1}^c y_k \cdot y_i + y_i^4 + \sum_{j \neq i} (y_i y_j)^2 \right] \\ &= \frac{1}{\left( \sum_{k=1}^c y_k \right)^4} \left[ y_i^2 \left( \sum_{k=1}^c y_k \right)^2 - 2y_i^3 \sum_{k=1}^c y_k + y_i^2 \sum_{k=1}^c y_k^2 \right] \end{aligned}$$

Therefore,

$$\begin{aligned} \|\nabla^2 \phi(z)\|^2 &\leq \sum_{i=1}^c \sum_{j=1}^c \left| \frac{\partial^2 \phi(z)}{\partial z_j \partial z_i} \right|^2 \\ &\leq \frac{1}{\left( \sum_{k=1}^c y_k \right)^4} \left[ \left( \sum_{i=1}^c y_i^2 \right) \left( \sum_{k=1}^c y_k \right)^2 - 2 \left( \sum_{i=1}^c y_i^3 \right) \left( \sum_{k=1}^c y_k \right) + \left( \sum_{i=1}^c y_i^2 \right) \left( \sum_{k=1}^c y_k^2 \right) \right] \\ &\leq \frac{\left( \sum_{i=1}^c y_i^2 \right) \left( \sum_{k=1}^c y_k \right)^2}{\left( \sum_{k=1}^c y_k \right)^4} \leq \frac{\left( \sum_{k=1}^c y_k \right)^4}{\left( \sum_{k=1}^c y_k \right)^4} = 1, \end{aligned}$$

where the last inequality holds since

$$\left( \sum_{i=1}^c y_i^2 \right) \left( \sum_{k=1}^c y_k^2 \right) \leq \left( \sum_{i=1}^c y_i^3 \right) \left( \sum_{k=1}^c y_k \right) \Leftrightarrow \left( \sum_{k=1}^c y_k^2 \right) \leq \sqrt{\left( \sum_{i=1}^c y_i^3 \right) \left( \sum_{k=1}^c y_k \right)},$$

which follows by the application of Holder inequality (35) with  $p = 2$ ,  $q = 2$ ,  $a_k = y_k^{3/2}$ , and  $b_k = y_k^{1/2}$  (Note that  $y_k \geq 0$ ,  $k \in [c]$ ). Hence,  $\|\nabla^2 \phi(z)\| \leq L_\phi$  with  $L_\phi = 1$  which is equivalent to  $L_\phi$ -smoothness of  $\phi$ .  $\square$

### Proof of Lemma 6

*Proof.* Since  $k(\cdot; i)$  are twice continuously differentiable for all  $i \in [n]$ , we have the following Taylor approximation for each component outputs  $k_j(\cdot; i)$  where  $j \in [c]$  and  $i \in [n]$ :

$$\begin{aligned} k_j(w^{(t+1)}; i) &= k_j(w^{(t)} - \eta^{(t)} v^{(t)}; i) \\ &= k_j(w^{(t)}; i) - \mathbf{J}_w k_j(w; i)|_{w=w^{(t)}} \eta^{(t)} v^{(t)} + \frac{1}{2} (\eta^{(t)} v^{(t)})^T M_{i,j}(\tilde{w}^{(t)}) (\eta^{(t)} v^{(t)}), \end{aligned} \quad (36)$$

where  $M_{i,j}(\tilde{w}^{(t)})$  is the Hessian matrices of  $k_j(\cdot; i)$  at  $\tilde{w}^{(t)}$  and  $\tilde{w}^{(t)} = \alpha w^{(t)} + (1 - \alpha) w^{(t+1)}$  for some  $\alpha \in [0, 1]$ .

Shifting this back to the original function  $h_j(\cdot; i)$  we have:

$$\begin{aligned} h_j(w^{(t+1)}; i) &= k_j(w^{(t+1)}; i) + (h_j(w^{(t+1)}; i) - k_j(w^{(t+1)}; i)) \\ &\stackrel{(36)}{=} k_j(w^{(t)}; i) - \mathbf{J}_w k_j(w; i)|_{w=w^{(t)}} \eta^{(t)} v^{(t)} + \frac{1}{2} (\eta^{(t)} v^{(t)})^T M_{i,j}(\tilde{w}^{(t)}) (\eta^{(t)} v^{(t)}) \\ &\quad + (h_j(w^{(t+1)}; i) - k_j(w^{(t+1)}; i)), \\ &= h_j(w^{(t)}; i) - \mathbf{J}_w k_j(w; i)|_{w=w^{(t)}} \eta^{(t)} v^{(t)} + \frac{1}{2} (\eta^{(t)} v^{(t)})^T M_{i,j}(\tilde{w}^{(t)}) (\eta^{(t)} v^{(t)}) \\ &\quad + (h_j(w^{(t+1)}; i) - k_j(w^{(t+1)}; i)) + (k_j(w^{(t)}; i) - h_j(w^{(t)}; i)), \end{aligned}$$

which leads to our desired statement:

$$h(w^{(t+1)}; i) = h(w^{(t)} - \eta^{(t)} v^{(t)}; i) = h(w^{(t)}; i) - \eta^{(t)} H_i^{(t)} v^{(t)} + \epsilon_i^{(t)},$$

where

$$\begin{aligned} \epsilon_{i,j}^{(t)} &= \frac{1}{2} (\eta^{(t)} v^{(t)})^T M_{i,j}(\tilde{w}^{(t)}) (\eta^{(t)} v^{(t)}) \\ &\quad + (h_j(w^{(t+1)}; i) - k_j(w^{(t+1)}; i)) + (k_j(w^{(t)}; i) - h_j(w^{(t)}; i)), \quad j \in [c], \end{aligned}$$

Hence we get the final bound:

$$\begin{aligned} |\epsilon_{i,j}^{(t)}| &\leq \frac{1}{2} \left| (\eta^{(t)} v^{(t)})^T M_{i,j}(\tilde{w}^{(t)}) (\eta^{(t)} v^{(t)}) \right| \\ &\quad + |h_j(w^{(t+1)}; i) - k_j(w^{(t+1)}; i)| + |k_j(w^{(t)}; i) - h_j(w^{(t)}; i)| \\ &\stackrel{(26)}{\leq} \frac{1}{2} \left| (\eta^{(t)} v^{(t)})^T M_{i,j}(\tilde{w}^{(t)}) (\eta^{(t)} v^{(t)}) \right| + 2\varepsilon, \\ &\leq \frac{1}{2} (\eta^{(t)})^2 \|v^{(t)}\|^2 \cdot \|M_{i,j}(\tilde{w}^{(t)})\| + 2\varepsilon \\ &\stackrel{(16)}{\leq} \frac{1}{2} (\eta^{(t)})^2 \|v^{(t)}\|^2 G + 2\varepsilon, \quad j \in [c]. \end{aligned}$$

□

### Proof of Corollary 1

*Proof.* The proof of this corollary follows directly by the applications of Lemmas 4 and 5. □

## D Technical Proofs for Theorem 1

**Lemma 8.** *Suppose that Assumption 2 holds for  $G > 0$  and Assumption 3 holds for  $V > 0$ , and  $v^{(t)} = v_{* \text{reg}}^{(t)}$ . Consider  $\eta^{(t)} = D\sqrt{\varepsilon}$  for some  $D > 0$  and  $\varepsilon > 0$ . For  $i \in [n]$  and  $0 \leq t < T$ , we have*

$$\|\epsilon_i^{(t)}\|^2 \leq \frac{1}{4} c(4 + (V + 2)GD^2)^2 \varepsilon^2. \quad (37)$$

*Proof.* From (19), for  $i \in [n]$ ,  $j \in [c]$ , and for  $0 \leq t < T$ , by Lemma 2 and Lemma 6 we have

$$|\epsilon_{i,j}^{(t)}| \leq \frac{1}{2} (\eta^{(t)})^2 \|v^{(t)}\|^2 G + 2\varepsilon \leq \frac{1}{2} (V + 2)GD^2 \varepsilon + 2\varepsilon = \frac{1}{2} \varepsilon (4 + (V + 2)GD^2),$$

where the last inequality follows by the fact  $\|v^{(t)}\|^2 = \|v_{* \text{reg}}^{(t)}\|^2 \leq 2 + V$  of Lemma 3 and  $\eta^{(t)} = D\sqrt{\varepsilon}$ . Hence,

$$\|\epsilon_i^{(t)}\|^2 = \sum_{j=1}^c |\epsilon_{i,j}^{(t)}|^2 \leq \frac{1}{4} c(4 + (V + 2)GD^2)^2 \varepsilon^2.$$

□

**Lemma 9.** *Let  $w^{(t)}$  be generated by Algorithm 1 where we use the closed form solution for the search direction. We execute Algorithm 1 for  $T = \frac{\beta}{\varepsilon}$  outer loops for some constant  $\beta > 0$ . We assume Assumption 1 holds. Suppose that Assumption 2 holds for  $G > 0$  and Assumption 3 holds for  $V > 0$ . We set the step size equal to  $\eta^{(t)} = D\sqrt{\varepsilon}$  for some  $D > 0$  and choose a learning rate  $\alpha_i^{(t)} \leq \frac{\alpha}{L_\phi}$ , for some  $\alpha \in (0, \frac{1}{3})$ . For  $i \in [n]$  and  $0 \leq t < T$ , we have*

$$\begin{aligned} \|h(w^{(t+1)}; i) - h_i^*\|^2 &\leq (1 + \varepsilon) \|h(w^{(t)}; i) - h_i^*\|^2 - 2(1 - 3\alpha)\alpha_i^{(t)} [\phi_i(h(w^{(t)}; i)) - \phi_i(h_i^*)] \\ &\quad + \frac{(3\varepsilon + 2)}{4} c(4 + (V + 2)GD^2)^2 \cdot \varepsilon \\ &\quad + \frac{3\varepsilon + 2}{\varepsilon} \|\eta^{(t)} H_i^{(t)} v_{* \text{reg}}^{(t)} - \alpha_i^{(t)} \nabla_z \phi_i(h(w^{(t)}; i))\|^2 \end{aligned} \quad (38)$$

*Proof.* Note that we have the optimal solution  $v_{* \text{reg}}^{(t)}$  for the optimization problem (20) for  $0 \leq t < T$ . From (17), we have, for  $i \in [n]$ ,

$$\begin{aligned} h(w^{(t+1)}; i) &= h(w^{(t)} - \eta^{(t)} v_{* \text{reg}}^{(t)}; i) \\ &= h(w^{(t)}; i) - \eta^{(t)} H_i^{(t)} v_{* \text{reg}}^{(t)} + \epsilon_i^{(t)} \\ &= h(w^{(t)}; i) - \alpha_i^{(t)} \nabla_z \phi_i(h(w^{(t)}; i)) + \epsilon_i^{(t)} - [\eta^{(t)} H_i^{(t)} v_{* \text{reg}}^{(t)} - \alpha_i^{(t)} \nabla_z \phi_i(h(w^{(t)}; i))]. \end{aligned}$$

Hence, we have

$$\begin{aligned} & \|h(w^{(t+1)}; i) - h_i^*\|^2 \\ &= \|h(w^{(t)}; i) - h_i^* - \alpha_i^{(t)} \nabla_z \phi_i(h(w^{(t)}; i)) + \epsilon_i^{(t)} - [\eta^{(t)} H_i^{(t)} v_{* \text{reg}}^{(t)} - \alpha_i^{(t)} \nabla_z \phi_i(h(w^{(t)}; i))]\|^2 \\ &= \|h(w^{(t)}; i) - h_i^*\|^2 + (\alpha_i^{(t)})^2 \|\nabla_z \phi_i(h(w^{(t)}; i))\|^2 \\ &\quad + \|\epsilon_i^{(t)}\|^2 + \|\eta^{(t)} H_i^{(t)} v_{* \text{reg}}^{(t)} - \alpha_i^{(t)} \nabla_z \phi_i(h(w^{(t)}; i))\|^2 \\ &\quad - 2 \cdot \langle h(w^{(t)}; i) - h_i^*, \alpha_i^{(t)} \nabla_z \phi_i(h(w^{(t)}; i)) \rangle \\ &\quad + 2 \cdot \langle h(w^{(t)}; i) - h_i^*, \epsilon_i^{(t)} \rangle \\ &\quad - 2 \cdot \langle h(w^{(t)}; i) - h_i^*, \eta^{(t)} H_i^{(t)} v_{* \text{reg}}^{(t)} - \alpha_i^{(t)} \nabla_z \phi_i(h(w^{(t)}; i)) \rangle \\ &\quad - 2 \cdot \langle \alpha_i^{(t)} \nabla_z \phi_i(h(w^{(t)}; i)), \epsilon_i^{(t)} \rangle \\ &\quad + 2 \cdot \langle \alpha_i^{(t)} \nabla_z \phi_i(h(w^{(t)}; i)), \eta^{(t)} H_i^{(t)} v_{* \text{reg}}^{(t)} - \alpha_i^{(t)} \nabla_z \phi_i(h(w^{(t)}; i)) \rangle \\ &\quad - 2 \cdot \langle \epsilon_i^{(t)}, \eta^{(t)} H_i^{(t)} v_{* \text{reg}}^{(t)} - \alpha_i^{(t)} \nabla_z \phi_i(h(w^{(t)}; i)) \rangle, \end{aligned}$$

where we expand the square term. Now applying Young's inequalities:  $2|\langle u, v \rangle| \leq \frac{\|u\|^2}{\varepsilon/2} + (\varepsilon/2)\|v\|^2$  for  $\varepsilon > 0$  and  $2|\langle u, v \rangle| \leq \|u\|^2 + \|v\|^2$  we have:

$$\begin{aligned} & \|h(w^{(t+1)}; i) - h_i^*\|^2 \\ &= \|h(w^{(t)}; i) - h_i^*\|^2 + (\alpha_i^{(t)})^2 \|\nabla_z \phi_i(h(w^{(t)}; i))\|^2 \\ &\quad + \|\epsilon_i^{(t)}\|^2 + \|\eta^{(t)} H_i^{(t)} v_{* \text{reg}}^{(t)} - \alpha_i^{(t)} \nabla_z \phi_i(h(w^{(t)}; i))\|^2 \\ &\quad - 2\alpha_i^{(t)} \langle h(w^{(t)}; i) - h_i^*, \nabla_z \phi_i(h(w^{(t)}; i)) \rangle \\ &\quad + \frac{\varepsilon}{2} \|h(w^{(t)}; i) - h_i^*\|^2 + \frac{2}{\varepsilon} \|\epsilon_i^{(t)}\|^2 \\ &\quad + \frac{\varepsilon}{2} \|h(w^{(t)}; i) - h_i^*\|^2 + \frac{2}{\varepsilon} \|\eta^{(t)} H_i^{(t)} v_{* \text{reg}}^{(t)} - \alpha_i^{(t)} \nabla_z \phi_i(h(w^{(t)}; i))\|^2 \\ &\quad + 2(\alpha_i^{(t)})^2 \|\nabla_z \phi_i(h(w^{(t)}; i))\|^2 + 2\|\epsilon_i^{(t)}\|^2 \\ &\quad + 2\|\eta^{(t)} H_i^{(t)} v_{* \text{reg}}^{(t)} - \alpha_i^{(t)} \nabla_z \phi_i(h(w^{(t)}; i))\|^2 \\ &\stackrel{(8)}{\leq} (1 + \varepsilon) \|h(w^{(t)}; i) - h_i^*\|^2 + 3(\alpha_i^{(t)})^2 \|\nabla_z \phi_i(h(w^{(t)}; i))\|^2 \\ &\quad + \left(3 + \frac{2}{\varepsilon}\right) \|\epsilon_i^{(t)}\|^2 + \left(3 + \frac{2}{\varepsilon}\right) \|\eta^{(t)} H_i^{(t)} v_{* \text{reg}}^{(t)} - \alpha_i^{(t)} \nabla_z \phi_i(h(w^{(t)}; i))\|^2 \\ &\quad - 2\alpha_i^{(t)} [\phi_i(h(w^{(t)}; i)) - \phi_i(h_i^*)]. \end{aligned}$$

Note that from (9) we get that  $\|\nabla_z \phi_i(h(w^{(t)}; i))\|^2 \leq 2L_\phi [\phi_i(h(w^{(t)}; i)) - \phi_i(h_i^*)]$ . Applying this and using the fact that  $\alpha_i^{(t)} \leq \frac{\alpha}{L_\phi}$ , for some  $\alpha \in (0, \frac{1}{3})$ , we are able to derive:

$$\begin{aligned} & \|h(w^{(t+1)}; i) - h_i^*\|^2 \\ &\leq (1 + \varepsilon) \|h(w^{(t)}; i) - h_i^*\|^2 - 2(1 - 3\alpha)\alpha_i^{(t)} [\phi_i(h(w^{(t)}; i)) - \phi_i(h_i^*)] \\ &\quad + \frac{3\varepsilon + 2}{\varepsilon} \|\epsilon_i^{(t)}\|^2 + \frac{3\varepsilon + 2}{\varepsilon} \|\eta^{(t)} H_i^{(t)} v_{* \text{reg}}^{(t)} - \alpha_i^{(t)} \nabla_z \phi_i(h(w^{(t)}; i))\|^2 \\ &\leq (1 + \varepsilon) \|h(w^{(t)}; i) - h_i^*\|^2 - 2(1 - 3\alpha)\alpha_i^{(t)} [\phi_i(h(w^{(t)}; i)) - \phi_i(h_i^*)] \end{aligned}$$

$$+ \frac{3\varepsilon + 2}{\varepsilon} \frac{1}{4} c(4 + (V + 2)GD^2)^2 \varepsilon^2 + \frac{3\varepsilon + 2}{\varepsilon} \|\eta^{(t)} H_i^{(t)} v_{* \text{reg}}^{(t)} - \alpha_i^{(t)} \nabla_z \phi_i(h(w^{(t)}; i))\|^2$$

where the last inequality follows by Lemma 8.  $\square$

**Lemma 10.** *Let  $w^{(t)}$  be generated by Algorithm 1 where we use the closed form solution for the search direction. We execute Algorithm 1 for  $T = \frac{\beta}{\varepsilon}$  outer loops for some constant  $\beta > 0$ . We assume Assumption 1 holds. Suppose that Assumption 2 holds for  $G > 0$  and Assumption 3 holds for  $V > 0$ . We set the step size equal to  $\eta^{(t)} = D\sqrt{\varepsilon}$  for some  $D > 0$  and choose a learning rate  $\alpha_i^{(t)} = (1 + \varepsilon)\alpha_i^{(t-1)} = (1 + \varepsilon)^t \alpha_i^{(0)}$ . Based on  $\beta$ , we define  $\alpha_i^{(0)} = \frac{\alpha}{e^{\beta} L_\phi}$  with  $\alpha \in (0, \frac{1}{3})$ . We have*

$$\begin{aligned} \frac{1}{T} \sum_{t=0}^{T-1} \frac{1}{n} \sum_{i=1}^n [f(w^{(t)}; i) - \phi_i(h_i^*)] &\leq \frac{e^\beta L_\phi (1 + \varepsilon)}{2(1 - 3\alpha)\alpha\beta} \cdot \frac{1}{n} \sum_{i=1}^n \|h(w^{(0)}; i) - h_i^*\|^2 \cdot \varepsilon \\ &+ \frac{e^\beta L_\phi}{8\alpha(1 - 3\alpha)} (3\varepsilon + 2) [c(4 + (V + 2)GD^2)^2 + 8 + 4V] \cdot \varepsilon. \end{aligned} \quad (39)$$

*Proof.* Rearranging the terms in Lemma 9, we have

$$\begin{aligned} \phi_i(h(w^{(t)}; i)) - \phi_i(h_i^*) &\leq \frac{1}{2(1 - 3\alpha)} \left( \frac{(1 + \varepsilon)}{\alpha_i^{(t)}} \|h(w^{(t)}; i) - h_i^*\|^2 - \frac{1}{\alpha_i^{(t)}} \|h(w^{(t+1)}; i) - h_i^*\|^2 \right) \\ &+ \frac{1}{8(1 - 3\alpha)} \cdot \frac{1}{\alpha_i^{(t)}} \cdot \varepsilon(3\varepsilon + 2)c(4 + (V + 2)GD^2)^2 \\ &+ \frac{1}{2(1 - 3\alpha)} \cdot \frac{1}{\alpha_i^{(t)}} \cdot \frac{3\varepsilon + 2}{\varepsilon} \|\eta^{(t)} H_i^{(t)} v_{* \text{reg}}^{(t)} - \alpha_i^{(t)} \nabla_z \phi_i(h(w^{(t)}; i))\|^2 \\ &\leq \frac{1}{2(1 - 3\alpha)} \left( \frac{(1 + \varepsilon)}{\alpha_i^{(t)}} \|h(w^{(t)}; i) - h_i^*\|^2 - \frac{(1 + \varepsilon)}{\alpha_i^{(t+1)}} \|h(w^{(t+1)}; i) - h_i^*\|^2 \right) \\ &+ \frac{e^\beta L_\phi}{8\alpha(1 - 3\alpha)} \cdot \varepsilon(3\varepsilon + 2)c(4 + (V + 2)GD^2)^2 \\ &+ \frac{e^\beta L_\phi}{2\alpha(1 - 3\alpha)} \cdot \frac{3\varepsilon + 2}{\varepsilon} \|\eta^{(t)} H_i^{(t)} v_{* \text{reg}}^{(t)} - \alpha_i^{(t)} \nabla_z \phi_i(h(w^{(t)}; i))\|^2. \end{aligned}$$

The last inequality follows because the learning rate satisfies  $\alpha_i^{(0)} = \frac{\alpha}{e^{\beta} L_\phi} \leq \frac{\alpha}{L_\phi}$  and for  $t = 1, \dots, T = \frac{\beta}{\varepsilon}$  for some  $\beta > 0$

$$\alpha_i^{(t)} = (1 + \varepsilon)\alpha_i^{(t-1)} = (1 + \varepsilon)^t \alpha_i^{(0)} \leq (1 + \varepsilon)^T \alpha_i^{(0)} = (1 + \varepsilon)^{\beta/\varepsilon} \frac{\alpha}{e^{\beta} L_\phi} \leq \frac{\alpha}{L_\phi},$$

since  $(1 + x)^{1/x} \leq e$ ,  $x > 0$ . Moreover, we have  $\frac{1}{\alpha_i^{(t)}} \leq \frac{1}{\alpha_i^{(0)}} = \frac{e^{\beta} L_\phi}{\alpha}$ ,  $t = 0, \dots, T - 1$ .

Taking the average sum from  $t = 0, \dots, T - 1$ , we have

$$\begin{aligned} \frac{1}{T} \sum_{t=0}^{T-1} [\phi_i(h(w^{(t)}; i)) - \phi_i(h_i^*)] &\leq \frac{1}{2(1 - 3\alpha)T} \cdot \frac{(1 + \varepsilon)}{\alpha_i^{(0)}} \|h(w^{(0)}; i) - h_i^*\|^2 \\ &+ \frac{e^\beta L_\phi}{8\alpha(1 - 3\alpha)} \cdot \varepsilon(3\varepsilon + 2)c(4 + (V + 2)GD^2)^2 \\ &+ \frac{e^\beta L_\phi}{2\alpha(1 - 3\alpha)} \cdot \frac{3\varepsilon + 2}{\varepsilon} \frac{1}{T} \sum_{t=0}^{T-1} \|\eta^{(t)} H_i^{(t)} v_{* \text{reg}}^{(t)} - \alpha_i^{(t)} \nabla_z \phi_i(h(w^{(t)}; i))\|^2 \\ &= \frac{e^\beta L_\phi (1 + \varepsilon)}{2(1 - 3\alpha)\alpha\beta} \varepsilon \cdot \|h(w^{(0)}; i) - h_i^*\|^2 \\ &+ \frac{e^\beta L_\phi}{8\alpha(1 - 3\alpha)} \cdot \varepsilon(3\varepsilon + 2)c(4 + (V + 2)GD^2)^2 \end{aligned}$$

$$+ \frac{e^\beta L_\phi}{2\alpha(1-3\alpha)} \cdot \frac{3\varepsilon+2}{\varepsilon} \frac{1}{T} \sum_{t=0}^{T-1} \|\eta^{(t)} H_i^{(t)} v_{* \text{reg}}^{(t)} - \alpha_i^{(t)} \nabla_z \phi_i(h(w^{(t)}; i))\|^2.$$

Taking the average sum from  $i = 1, \dots, n$ , we have

$$\begin{aligned} & \frac{1}{T} \sum_{t=0}^{T-1} \frac{1}{n} \sum_{i=1}^n [\phi_i(h(w^{(t)}; i)) - \phi_i(h_i^*)] \\ & \leq \frac{e^\beta L_\phi(1+\varepsilon)}{2(1-3\alpha)\alpha\beta} \varepsilon \cdot \frac{1}{n} \sum_{i=1}^n \|h(w^{(0)}; i) - h_i^*\|^2 \\ & \quad + \frac{e^\beta L_\phi}{8\alpha(1-3\alpha)} \cdot \varepsilon(3\varepsilon+2)c(4+(V+2)GD^2)^2 \\ & \quad + \frac{e^\beta L_\phi}{2\alpha(1-3\alpha)} \cdot \frac{3\varepsilon+2}{\varepsilon} \frac{1}{T} \sum_{t=0}^{T-1} \frac{1}{n} \sum_{i=1}^n \|\eta^{(t)} H_i^{(t)} v_{* \text{reg}}^{(t)} - \alpha_i^{(t)} \nabla_z \phi_i(h(w^{(t)}; i))\|^2 \\ & \stackrel{(22)}{\leq} \frac{e^\beta L_\phi(1+\varepsilon)}{2(1-3\alpha)\alpha\beta} \varepsilon \cdot \frac{1}{n} \sum_{i=1}^n \|h(w^{(0)}; i) - h_i^*\|^2 \\ & \quad + \frac{e^\beta L_\phi}{8\alpha(1-3\alpha)} \cdot \varepsilon(3\varepsilon+2)c(4+(V+2)GD^2)^2 \\ & \quad + \frac{e^\beta L_\phi}{2\alpha(1-3\alpha)} \cdot \frac{3\varepsilon+2}{\varepsilon} (2+V)\varepsilon^2. \end{aligned} \tag{40}$$

Note that

$$\frac{1}{T} \sum_{t=0}^{T-1} \frac{1}{n} \sum_{i=1}^n [\phi_i(h(w^{(t)}; i)) - \phi_i(h_i^*)] = \frac{1}{T} \sum_{t=0}^{T-1} \frac{1}{n} \sum_{i=1}^n [f(w^{(t)}; i) - \phi_i(h_i^*)]. \tag{41}$$

Therefore, applying (41) to (40), we have

$$\begin{aligned} \frac{1}{T} \sum_{t=0}^{T-1} \frac{1}{n} \sum_{i=1}^n [f(w^{(t)}; i) - \phi_i(h_i^*)] & \leq \frac{e^\beta L_\phi(1+\varepsilon)}{2(1-3\alpha)\alpha\beta} \cdot \frac{1}{n} \sum_{i=1}^n \|h(w^{(0)}; i) - h_i^*\|^2 \cdot \varepsilon \\ & \quad + \frac{e^\beta L_\phi}{8\alpha(1-3\alpha)} (3\varepsilon+2) [c(4+(V+2)GD^2)^2 + 8 + 4V] \cdot \varepsilon. \end{aligned}$$

which is our desired result.  $\square$

### Proof of Theorem 1

*Proof.* We have

$$\begin{aligned} F_* & = \min_{w \in \mathbb{R}^d} F(w) = \min_{w \in \mathbb{R}^d} \left( \frac{1}{n} \sum_{i=1}^n f_i(w) \right) = \frac{1}{n} \min_{w \in \mathbb{R}^d} \left( \sum_{i=1}^n f_i(w) \right) \\ & \geq \frac{1}{n} \sum_{i=1}^n \min_{w \in \mathbb{R}^d} (f_i(w)) = \frac{1}{n} \sum_{i=1}^n f_i^* \geq \frac{1}{n} \sum_{i=1}^n \phi_i(h_i^*). \end{aligned} \tag{42}$$

Hence  $F_* - \frac{1}{n} \sum_{i=1}^n \phi_i(h_i^*) \geq 0$ . Therefore

$$\begin{aligned} \frac{1}{T} \sum_{t=0}^{T-1} [F(w^{(t)}) - F_*] & = \frac{1}{T} \sum_{t=0}^{T-1} \left( \frac{1}{n} \sum_{i=1}^n [f(w^{(t)}; i) - \phi_i(h_i^*)] - \left[ F_* - \frac{1}{n} \sum_{i=1}^n \phi_i(h_i^*) \right] \right) \\ & \leq \frac{1}{T} \sum_{t=0}^{T-1} \frac{1}{n} \sum_{i=1}^n [f(w^{(t)}; i) - \phi_i(h_i^*)] \end{aligned}$$

$$\begin{aligned}
&\stackrel{(39)}{\leq} \frac{e^\beta L_\phi(1+\varepsilon)}{2(1-3\alpha)\alpha\beta} \cdot \frac{1}{n} \sum_{i=1}^n \|h(w^{(0)}; i) - h_i^*\|^2 \cdot \varepsilon \\
&\quad + \frac{e^\beta L_\phi(3\varepsilon+2)}{8\alpha(1-3\alpha)} [c(4+(V+2)GD^2)^2 + 8 + 4V] \cdot \varepsilon.
\end{aligned}$$

□

## E Technical Proofs for Theorem 2

**Lemma 11.** For  $0 \leq t < T$ , suppose that Assumption 3 holds for  $V \geq 0$  and  $v^{(t)}$  satisfies (24). Then

$$\|v^{(t)}\|^2 \leq 2(\varepsilon^2 + V + 2).$$

*Proof.* From  $\|v^{(t)} - v_{* \text{reg}}^{(t)}\| \leq \varepsilon$ . Using  $\|a\|^2 \leq 2\|a-b\|^2 + 2\|b\|^2$ , we have

$$\|v^{(t)}\|^2 \leq 2\|v^{(t)} - v_{* \text{reg}}^{(t)}\|^2 + 2\|v_{* \text{reg}}^{(t)}\|^2 \stackrel{(24)}{\leq} 2\varepsilon^2 + 4 + 2V.$$

where the last inequality follows since  $\|v_{* \text{reg}}^{(t)}\|^2 \leq 2 + V$  for some  $V > 0$  in Lemma 3.

□

**Lemma 12.** Suppose that Assumption 2 holds for  $G > 0$  and Assumption 3 holds for  $V > 0$ . Consider  $\eta^{(t)} = D\sqrt{\varepsilon}$  for some  $D > 0$  and  $\varepsilon > 0$ . For  $i \in [n]$  and  $0 \leq t < T$ , we have

$$\|\epsilon_i^{(t)}\|^2 \leq c(2 + (V + \varepsilon^2 + 2)GD^2)^2 \varepsilon^2. \quad (43)$$

*Proof.* From (19), for  $i \in [n]$ ,  $j \in [c]$ , and for  $0 \leq t < T$ , by Lemma 2 and Lemma 6 we have

$$|\epsilon_{i,j}^{(t)}| \leq \frac{1}{2}(\eta^{(t)})^2 \|v^{(t)}\|^2 G + 2\varepsilon \leq \frac{1}{2}2(\varepsilon^2 + V + 2)GD^2\varepsilon + 2\varepsilon = \varepsilon(2 + (V + \varepsilon^2 + 2)GD^2),$$

where the last inequality follows by the application of Lemma 11 and  $\eta^{(t)} = D\sqrt{\varepsilon}$ . Hence,

$$\|\epsilon_i^{(t)}\|^2 = \sum_{j=1}^c |\epsilon_{i,j}^{(t)}|^2 \leq c(2 + (V + \varepsilon^2 + 2)GD^2)^2 \varepsilon^2.$$

□

**Lemma 13.** Let  $w^{(t)}$  be generated by Algorithm 2 where  $v^{(t)}$  satisfies (24). We execute Algorithm 2 for  $T = \frac{\beta}{\varepsilon}$  outer loops for some constant  $\beta > 0$ . We assume Assumption 1 holds. Suppose that Assumption 2 holds for  $G > 0$ , Assumption 3 holds for  $V > 0$  and Assumption 4 holds for  $H > 0$ . We set the step size equal to  $\eta^{(t)} = D\sqrt{\varepsilon}$  for some  $D > 0$  and choose a learning rate  $\alpha_i^{(t)} \leq \frac{\alpha}{L_\phi}$ , for some  $\alpha \in (0, \frac{1}{4})$ . For  $i \in [n]$  and  $0 \leq t < T$ , we have

$$\begin{aligned}
\|h(w^{(t+1)}; i) - h_i^*\|^2 &\leq (1+\varepsilon)\|h(w^{(t)}; i) - h_i^*\|^2 - 2(1-4\alpha)\alpha_i^{(t)}[\phi_i(h(w^{(t)}; i)) - \phi_i(h_i^*)] \\
&\quad + \varepsilon(4\varepsilon+3) [D^2H^2 + c(2+(V+\varepsilon^2+2)GD^2)^2] \\
&\quad + \frac{4\varepsilon+3}{\varepsilon} \|\eta^{(t)}H_i^{(t)}v_{* \text{reg}}^{(t)} - \alpha_i^{(t)}\nabla_z\phi_i(h(w^{(t)}; i))\|^2
\end{aligned} \quad (44)$$

*Proof.* Note that  $v^{(t)}$  is obtained from the optimization problem (20) for  $0 \leq t < T$ . From (14), we have, for  $i \in [n]$ ,

$$\begin{aligned}
h(w^{(t+1)}; i) &= h(w^{(t)} - \eta^{(t)}v^{(t)}; i) \\
&= h(w^{(t)}; i) - \eta^{(t)}H_i^{(t)}v^{(t)} + \epsilon_i^{(t)} \\
&= h(w^{(t)}; i) - \eta^{(t)}H_i^{(t)}(v^{(t)} - v_{* \text{reg}}^{(t)}) - \alpha_i^{(t)}\nabla_z\phi_i(h(w^{(t)}; i)) + \epsilon_i^{(t)} \\
&\quad - [\eta^{(t)}H_i^{(t)}v_{* \text{reg}}^{(t)} - \alpha_i^{(t)}\nabla_z\phi_i(h(w^{(t)}; i))].
\end{aligned}$$

Hence, we have

$$\begin{aligned}
& \|h(w^{(t+1)}; i) - h_i^*\|^2 \\
&= \|h(w^{(t)}; i) - h_i^* - \eta^{(t)} H_i^{(t)}(v^{(t)} - v_{* \text{reg}}^{(t)}) - \alpha_i^{(t)} \nabla_z \phi_i(h(w^{(t)}; i)) \\
&\quad + \epsilon_i^{(t)} - [\eta^{(t)} H_i^{(t)} v_{* \text{reg}}^{(t)} - \alpha_i^{(t)} \nabla_z \phi_i(h(w^{(t)}; i))]\|^2 \\
&= \|h(w^{(t)}; i) - h_i^*\|^2 + \|\eta^{(t)} H_i^{(t)}(v^{(t)} - v_{* \text{reg}}^{(t)})\|^2 + (\alpha_i^{(t)})^2 \|\nabla_z \phi_i(h(w^{(t)}; i))\|^2 \\
&\quad + \|\epsilon_i^{(t)}\|^2 + \|\eta^{(t)} H_i^{(t)} v_{* \text{reg}}^{(t)} - \alpha_i^{(t)} \nabla_z \phi_i(h(w^{(t)}; i))\|^2 \\
&\quad - 2 \cdot \langle h(w^{(t)}; i) - h_i^*, \eta^{(t)} H_i^{(t)}(v^{(t)} - v_{* \text{reg}}^{(t)}) \rangle \\
&\quad - 2 \cdot \langle h(w^{(t)}; i) - h_i^*, \alpha_i^{(t)} \nabla_z \phi_i(h(w^{(t)}; i)) \rangle \\
&\quad + 2 \cdot \langle h(w^{(t)}; i) - h_i^*, \epsilon_i^{(t)} \rangle \\
&\quad - 2 \cdot \langle h(w^{(t)}; i) - h_i^*, \eta^{(t)} H_i^{(t)} v_{* \text{reg}}^{(t)} - \alpha_i^{(t)} \nabla_z \phi_i(h(w^{(t)}; i)) \rangle \\
&\quad + 2 \cdot \langle \eta^{(t)} H_i^{(t)}(v^{(t)} - v_{* \text{reg}}^{(t)}), \alpha_i^{(t)} \nabla_z \phi_i(h(w^{(t)}; i)) \rangle \\
&\quad - 2 \cdot \langle \eta^{(t)} H_i^{(t)}(v^{(t)} - v_{* \text{reg}}^{(t)}), \epsilon_i^{(t)} \rangle \\
&\quad + 2 \cdot \langle \eta^{(t)} H_i^{(t)}(v^{(t)} - v_{* \text{reg}}^{(t)}), \eta^{(t)} H_i^{(t)} v_{* \text{reg}}^{(t)} - \alpha_i^{(t)} \nabla_z \phi_i(h(w^{(t)}; i)) \rangle \\
&\quad - 2 \cdot \langle \alpha_i^{(t)} \nabla_z \phi_i(h(w^{(t)}; i)), \epsilon_i^{(t)} \rangle \\
&\quad + 2 \cdot \langle \alpha_i^{(t)} \nabla_z \phi_i(h(w^{(t)}; i)), \eta^{(t)} H_i^{(t)} v_{* \text{reg}}^{(t)} - \alpha_i^{(t)} \nabla_z \phi_i(h(w^{(t)}; i)) \rangle \\
&\quad - 2 \cdot \langle \epsilon_i^{(t)}, \eta^{(t)} H_i^{(t)} v_{* \text{reg}}^{(t)} - \alpha_i^{(t)} \nabla_z \phi_i(h(w^{(t)}; i)) \rangle,
\end{aligned}$$

where we expand the square term. Now applying Young's inequalities:  $2|\langle u, v \rangle| \leq \frac{\|u\|^2}{\varepsilon/3} + (\varepsilon/3)\|v\|^2$  for  $\varepsilon > 0$  and  $2|\langle u, v \rangle| \leq \|u\|^2 + \|v\|^2$  we have:

$$\begin{aligned}
& \|h(w^{(t+1)}; i) - h_i^*\|^2 \\
&= \|h(w^{(t)}; i) - h_i^*\|^2 + \|\eta^{(t)} H_i^{(t)}(v^{(t)} - v_{* \text{reg}}^{(t)})\|^2 + (\alpha_i^{(t)})^2 \|\nabla_z \phi_i(h(w^{(t)}; i))\|^2 \\
&\quad + \|\epsilon_i^{(t)}\|^2 + \|\eta^{(t)} H_i^{(t)} v_{* \text{reg}}^{(t)} - \alpha_i^{(t)} \nabla_z \phi_i(h(w^{(t)}; i))\|^2 \\
&\quad + \frac{\varepsilon}{3} \|h(w^{(t)}; i) - h_i^*\|^2 + \frac{3}{\varepsilon} \|\eta^{(t)} H_i^{(t)}(v^{(t)} - v_{* \text{reg}}^{(t)})\|^2 \\
&\quad - 2\alpha_i^{(t)} \langle h(w^{(t)}; i) - h_i^*, \nabla_z \phi_i(h(w^{(t)}; i)) \rangle \\
&\quad + \frac{\varepsilon}{3} \|h(w^{(t)}; i) - h_i^*\|^2 + \frac{3}{\varepsilon} \|\epsilon_i^{(t)}\|^2 \\
&\quad + \frac{\varepsilon}{3} \|h(w^{(t)}; i) - h_i^*\|^2 + \frac{3}{\varepsilon} \|\eta^{(t)} H_i^{(t)} v_{* \text{reg}}^{(t)} - \alpha_i^{(t)} \nabla_z \phi_i(h(w^{(t)}; i))\|^2 \\
&\quad + 3(\eta^{(t)})^2 \|H_i^{(t)}(v^{(t)} - v_{* \text{reg}}^{(t)})\|^2 + 3(\alpha_i^{(t)})^2 \|\nabla_z \phi_i(h(w^{(t)}; i))\|^2 + 3\|\epsilon_i^{(t)}\|^2 \\
&\quad + 3\|\eta^{(t)} H_i^{(t)} v_{* \text{reg}}^{(t)} - \alpha_i^{(t)} \nabla_z \phi_i(h(w^{(t)}; i))\|^2 \\
&\stackrel{(8)}{\leq} (1 + \varepsilon) \|h(w^{(t)}; i) - h_i^*\|^2 + 4(\alpha_i^{(t)})^2 \|\nabla_z \phi_i(h(w^{(t)}; i))\|^2 \\
&\quad + \left(4 + \frac{3}{\varepsilon}\right) \|\eta^{(t)} H_i^{(t)}(v^{(t)} - v_{* \text{reg}}^{(t)})\|^2 + \left(4 + \frac{3}{\varepsilon}\right) \|\epsilon_i^{(t)}\|^2 \\
&\quad + \left(4 + \frac{3}{\varepsilon}\right) \|\eta^{(t)} H_i^{(t)} v_{* \text{reg}}^{(t)} - \alpha_i^{(t)} \nabla_z \phi_i(h(w^{(t)}; i))\|^2 \\
&\quad - 2\alpha_i^{(t)} [\phi_i(h(w^{(t)}; i)) - \phi_i(h_i^*)]
\end{aligned}$$

Note that from (9) we get that  $\|\nabla_z \phi_i(h(w^{(t)}; i))\|^2 \leq 2L_\phi [\phi_i(h(w^{(t)}; i)) - \phi_i(h_i^*)]$ . Applying this and using the fact that  $\alpha_i^{(t)} \leq \frac{\alpha}{L_\phi}$ , for some  $\alpha \in (0, \frac{1}{4})$ , we are able to derive:

$$\|h(w^{(t+1)}; i) - h_i^*\|^2$$

$$\begin{aligned}
&\leq (1 + \varepsilon)\|h(w^{(t)}; i) - h_i^*\|^2 - 2(1 - 4\alpha)\alpha_i^{(t)}[\phi_i(h(w^{(t)}; i)) - \phi_i(h_i^*)] \\
&\quad + \frac{4\varepsilon + 3}{\varepsilon}\|\eta^{(t)}H_i^{(t)}(v^{(t)} - v_{* \text{ reg}}^{(t)})\|^2 + \frac{4\varepsilon + 3}{\varepsilon}\|\epsilon_i^{(t)}\|^2 \\
&\quad + \frac{4\varepsilon + 3}{\varepsilon}\|\eta^{(t)}H_i^{(t)}v_{* \text{ reg}}^{(t)} - \alpha_i^{(t)}\nabla_z\phi_i(h(w^{(t)}; i))\|^2 \\
&\stackrel{(a)}{\leq} (1 + \varepsilon)\|h(w^{(t)}; i) - h_i^*\|^2 - 2(1 - 4\alpha)\alpha_i^{(t)}[\phi_i(h(w^{(t)}; i)) - \phi_i(h_i^*)] \\
&\quad + \frac{4\varepsilon + 3}{\varepsilon}D^2\varepsilon\frac{H^2}{\varepsilon}\|v^{(t)} - v_{* \text{ reg}}^{(t)}\|^2 + \frac{4\varepsilon + 3}{\varepsilon}\|\epsilon_i^{(t)}\|^2 \\
&\quad + \frac{4\varepsilon + 3}{\varepsilon}\|\eta^{(t)}H_i^{(t)}v_{* \text{ reg}}^{(t)} - \alpha_i^{(t)}\nabla_z\phi_i(h(w^{(t)}; i))\|^2 \\
&\stackrel{(b)}{\leq} (1 + \varepsilon)\|h(w^{(t)}; i) - h_i^*\|^2 - 2(1 - 4\alpha)\alpha_i^{(t)}[\phi_i(h(w^{(t)}; i)) - \phi_i(h_i^*)] \\
&\quad + \frac{4\varepsilon + 3}{\varepsilon}D^2H^2 \cdot \varepsilon^2 + \frac{4\varepsilon + 3}{\varepsilon} \cdot c(2 + (V + \varepsilon^2 + 2)GD^2)^2\varepsilon^2 \\
&\quad + \frac{4\varepsilon + 3}{\varepsilon}\|\eta^{(t)}H_i^{(t)}v_{* \text{ reg}}^{(t)} - \alpha_i^{(t)}\nabla_z\phi_i(h(w^{(t)}; i))\|^2 \\
&= (1 + \varepsilon)\|h(w^{(t)}; i) - h_i^*\|^2 - 2(1 - 4\alpha)\alpha_i^{(t)}[\phi_i(h(w^{(t)}; i)) - \phi_i(h_i^*)] \\
&\quad + \varepsilon(4\varepsilon + 3) [D^2H^2 + c(2 + (V + \varepsilon^2 + 2)GD^2)^2] \\
&\quad + \frac{4\varepsilon + 3}{\varepsilon}\|\eta^{(t)}H_i^{(t)}v_{* \text{ reg}}^{(t)} - \alpha_i^{(t)}\nabla_z\phi_i(h(w^{(t)}; i))\|^2
\end{aligned}$$

where (a) follows by using matrix vector inequality  $\|Hv\| \leq \|H\|\|v\|$ , where  $H \in \mathbb{R}^{c \times d}$  and  $v \in \mathbb{R}^d$  and Assumption 4 in (25) and  $\eta^{(t)} = D\sqrt{\varepsilon}$  for some  $D > 0$  and  $\varepsilon > 0$ ; (b) follows by the fact that  $\|v^{(t)} - v_{* \text{ reg}}^{(t)}\|^2 \leq \varepsilon^2$  in (24) and Lemma 12.  $\square$

**Lemma 14.** *Let  $w^{(t)}$  be generated by Algorithm 2 where  $v^{(t)}$  satisfies (24). We execute Algorithm 2 for  $T = \frac{\beta}{\varepsilon}$  outer loops for some constant  $\beta > 0$ . We assume Assumption 1 holds. Suppose that Assumption 2 holds for  $G > 0$ , Assumption 3 holds for  $V > 0$  and Assumption 4 holds for  $H > 0$ . We set the step size equal to  $\eta^{(t)} = D\sqrt{\varepsilon}$  for some  $D > 0$  and choose a learning rate  $\alpha_i^{(t)} = (1 + \varepsilon)\alpha_i^{(t-1)} = (1 + \varepsilon)^t\alpha_i^{(0)}$ . Based on  $\beta$ , we define  $\alpha_i^{(0)} = \frac{\alpha}{e^\beta L_\phi}$  with  $\alpha \in (0, \frac{1}{4})$ . We have*

$$\begin{aligned}
\frac{1}{T} \sum_{t=0}^{T-1} \frac{1}{n} \sum_{i=1}^n [f(w^{(t)}; i) - \phi_i(h_i^*)] &\leq \frac{e^\beta L_\phi (1 + \varepsilon)}{2(1 - 4\alpha)\alpha\beta} \cdot \frac{1}{n} \sum_{i=1}^n \|h(w^{(0)}; i) - h_i^*\|^2 \cdot \varepsilon \\
&\quad + \frac{e^\beta L_\phi (4\varepsilon + 3)}{2\alpha(1 - 4\alpha)} [D^2H^2 + c(2 + (V + \varepsilon^2 + 2)GD^2)^2 + 2 + V] \cdot \varepsilon. \quad (45)
\end{aligned}$$

*Proof.* Rearranging the terms in Lemma 13, we have

$$\begin{aligned}
\phi_i(h(w^{(t)}; i)) - \phi_i(h_i^*) &\leq \frac{1}{2(1 - 4\alpha)} \left( \frac{(1 + \varepsilon)}{\alpha_i^{(t)}} \|h(w^{(t)}; i) - h_i^*\|^2 - \frac{1}{\alpha_i^{(t)}} \|h(w^{(t+1)}; i) - h_i^*\|^2 \right) \\
&\quad + \frac{1}{2(1 - 4\alpha)} \cdot \frac{1}{\alpha_i^{(t)}} \cdot \varepsilon(4\varepsilon + 3) [D^2H^2 + c(2 + (V + \varepsilon^2 + 2)GD^2)^2] \\
&\quad + \frac{1}{2(1 - 4\alpha)} \cdot \frac{1}{\alpha_i^{(t)}} \cdot \frac{4\varepsilon + 3}{\varepsilon} \|\eta^{(t)}H_i^{(t)}v_{* \text{ reg}}^{(t)} - \alpha_i^{(t)}\nabla_z\phi_i(h(w^{(t)}; i))\|^2 \\
&\leq \frac{1}{2(1 - 4\alpha)} \left( \frac{(1 + \varepsilon)}{\alpha_i^{(t)}} \|h(w^{(t)}; i) - h_i^*\|^2 - \frac{(1 + \varepsilon)}{\alpha_i^{(t+1)}} \|h(w^{(t+1)}; i) - h_i^*\|^2 \right) \\
&\quad + \frac{e^\beta L_\phi}{2\alpha(1 - 4\alpha)} \cdot \varepsilon(4\varepsilon + 3) [D^2H^2 + c(2 + (V + \varepsilon^2 + 2)GD^2)^2] \\
&\quad + \frac{e^\beta L_\phi}{2\alpha(1 - 4\alpha)} \cdot \frac{4\varepsilon + 3}{\varepsilon} \|\eta^{(t)}H_i^{(t)}v_{* \text{ reg}}^{(t)} - \alpha_i^{(t)}\nabla_z\phi_i(h(w^{(t)}; i))\|^2. \quad (46)
\end{aligned}$$



The last inequality follows because the learning rate satisfies  $\alpha_i^{(0)} = \frac{\alpha}{e^\beta L_\phi} \leq \frac{\alpha}{L_\phi}$  and for  $t = 1, \dots, T = \frac{\beta}{\varepsilon}$  for some  $\beta > 0$

$$\alpha_i^{(t)} = (1 + \varepsilon)\alpha_i^{(t-1)} = (1 + \varepsilon)^t \alpha_i^{(0)} \leq (1 + \varepsilon)^T \alpha_i^{(0)} = (1 + \varepsilon)^{\beta/\varepsilon} \frac{\alpha}{e^\beta L_\phi} \leq \frac{\alpha}{L_\phi},$$

since  $(1 + x)^{1/x} \leq e, x > 0$ . Moreover, we have  $\frac{1}{\alpha_i^{(t)}} \leq \frac{1}{\alpha_i^{(0)}} = \frac{e^\beta L_\phi}{\alpha}, t = 0, \dots, T - 1$ .

Taking the average sum from  $t = 0, \dots, T - 1$ , we have

$$\begin{aligned} \frac{1}{T} \sum_{t=0}^{T-1} [\phi_i(h(w^{(t)}; i)) - \phi_i(h_i^*)] &\leq \frac{1}{2(1 - 4\alpha)T} \cdot \frac{(1 + \varepsilon)}{\alpha_i^{(0)}} \|h(w^{(0)}; i) - h_i^*\|^2 \\ &\quad + \frac{e^\beta L_\phi}{2\alpha(1 - 4\alpha)} \cdot \varepsilon(4\varepsilon + 3) [D^2 H^2 + c(2 + (V + \varepsilon^2 + 2)GD^2)^2] \\ &\quad + \frac{e^\beta L_\phi}{2\alpha(1 - 4\alpha)} \cdot \frac{4\varepsilon + 3}{\varepsilon} \frac{1}{T} \sum_{t=0}^{T-1} \|\eta^{(t)} H_i^{(t)} v_{* \text{reg}}^{(t)} - \alpha_i^{(t)} \nabla_z \phi_i(h(w^{(t)}; i))\|^2 \\ &= \frac{e^\beta L_\phi(1 + \varepsilon)}{2(1 - 4\alpha)\alpha\beta} \varepsilon \cdot \|h(w^{(0)}; i) - h_i^*\|^2 \\ &\quad + \frac{e^\beta L_\phi}{2\alpha(1 - 4\alpha)} \cdot \varepsilon(4\varepsilon + 3) [D^2 H^2 + c(2 + (V + \varepsilon^2 + 2)GD^2)^2] \\ &\quad + \frac{e^\beta L_\phi}{2\alpha(1 - 4\alpha)} \cdot \frac{4\varepsilon + 3}{\varepsilon} \frac{1}{T} \sum_{t=0}^{T-1} \|\eta^{(t)} H_i^{(t)} v_{* \text{reg}}^{(t)} - \alpha_i^{(t)} \nabla_z \phi_i(h(w^{(t)}; i))\|^2. \end{aligned}$$

Taking the average sum from  $i = 1, \dots, n$ , we have

$$\begin{aligned} &\frac{1}{T} \sum_{t=0}^{T-1} \frac{1}{n} \sum_{i=1}^n [\phi_i(h(w^{(t)}; i)) - \phi_i(h_i^*)] \\ &\leq \frac{e^\beta L_\phi(1 + \varepsilon)}{2(1 - 4\alpha)\alpha\beta} \varepsilon \cdot \frac{1}{n} \sum_{i=1}^n \|h(w^{(0)}; i) - h_i^*\|^2 \\ &\quad + \frac{e^\beta L_\phi}{2\alpha(1 - 4\alpha)} \cdot \varepsilon(4\varepsilon + 3) [D^2 H^2 + c(2 + (V + \varepsilon^2 + 2)GD^2)^2] \\ &\quad + \frac{e^\beta L_\phi}{2\alpha(1 - 4\alpha)} \cdot \frac{4\varepsilon + 3}{\varepsilon} \frac{1}{T} \sum_{t=0}^{T-1} \frac{1}{n} \sum_{i=1}^n \|\eta^{(t)} H_i^{(t)} v_{* \text{reg}}^{(t)} - \alpha_i^{(t)} \nabla_z \phi_i(h(w^{(t)}; i))\|^2 \\ &\stackrel{(22)}{\leq} \frac{e^\beta L_\phi(1 + \varepsilon)}{2(1 - 4\alpha)\alpha\beta} \varepsilon \cdot \frac{1}{n} \sum_{i=1}^n \|h(w^{(0)}; i) - h_i^*\|^2 \\ &\quad + \frac{e^\beta L_\phi}{2\alpha(1 - 4\alpha)} \cdot \varepsilon(4\varepsilon + 3) [D^2 H^2 + c(2 + (V + \varepsilon^2 + 2)GD^2)^2] \\ &\quad + \frac{e^\beta L_\phi}{2\alpha(1 - 4\alpha)} \cdot \frac{4\varepsilon + 3}{\varepsilon} (2 + V)\varepsilon^2. \end{aligned} \tag{47}$$

Note that

$$\frac{1}{T} \sum_{t=0}^{T-1} \frac{1}{n} \sum_{i=1}^n [\phi_i(h(w^{(t)}; i)) - \phi_i(h_i^*)] = \frac{1}{T} \sum_{t=0}^{T-1} \frac{1}{n} \sum_{i=1}^n [f(w^{(t)}; i) - \phi_i(h_i^*)]. \tag{48}$$

Therefore, applying (48) to (47), we have

$$\frac{1}{T} \sum_{t=0}^{T-1} \frac{1}{n} \sum_{i=1}^n [f(w^{(t)}; i) - \phi_i(h_i^*)] \leq \frac{e^\beta L_\phi(1 + \varepsilon)}{2(1 - 4\alpha)\alpha\beta} \cdot \frac{1}{n} \sum_{i=1}^n \|h(w^{(0)}; i) - h_i^*\|^2 \cdot \varepsilon$$

$$+ \frac{e^\beta L_\phi}{2\alpha(1-4\alpha)} (4\varepsilon + 3) [D^2 H^2 + c(2 + (V + \varepsilon^2 + 2)GD^2)^2 + 2 + V] \cdot \varepsilon.$$

□

### Proof of Theorem 2

*Proof.* From (42) we have  $F_* - \frac{1}{n} \sum_{i=1}^n \phi_i(h_i^*) \geq 0$ . This leads to

$$\begin{aligned} \frac{1}{T} \sum_{t=0}^{T-1} [F(w^{(t)}) - F_*] &= \frac{1}{T} \sum_{t=0}^{T-1} \left( \frac{1}{n} \sum_{i=1}^n [f(w^{(t)}; i) - \phi_i(h_i^*)] - \left[ F_* - \frac{1}{n} \sum_{i=1}^n \phi_i(h_i^*) \right] \right) \\ &\stackrel{(42)}{\leq} \frac{1}{T} \sum_{t=0}^{T-1} \frac{1}{n} \sum_{i=1}^n [f(w^{(t)}; i) - \phi_i(h_i^*)] \\ &\stackrel{(45)}{\leq} \frac{e^\beta L_\phi (1 + \varepsilon)}{2(1-4\alpha)\alpha\beta} \cdot \frac{1}{n} \sum_{i=1}^n \|h(w^{(0)}; i) - h_i^*\|^2 \cdot \varepsilon \\ &\quad + \frac{e^\beta L_\phi (4\varepsilon + 3)}{2\alpha(1-4\alpha)} [D^2 H^2 + c(2 + (V + \varepsilon^2 + 2)GD^2)^2 + 2 + V] \cdot \varepsilon. \end{aligned} \quad (49)$$

□

### Proof of Corollary 2

*Proof.* For each iteration  $0 \leq t < T$ , we need to find  $v^{(t)}$  satisfying the following criteria:

$$\|v^{(t)} - v_{\text{reg}}^{(t)}\|^2 \leq \varepsilon^2,$$

for some  $\varepsilon > 0$ . Using Gradient Descent we need  $\mathcal{O}(n \frac{L}{\mu} \log(\frac{1}{\varepsilon})) = \mathcal{O}(2n \frac{L}{\mu} \log(\frac{1}{\varepsilon}))$  number of gradient evaluations [Nesterov, 2004], where  $L$  and  $\mu = \varepsilon^2$  are the smooth and strongly convex constants, respectively, of  $\Psi$ . Let

$$\psi_i^{(t)}(v) = \frac{1}{2} \|\eta^{(t)} H_i^{(t)} v - \alpha_i^{(t)} \nabla_z \phi_i(h(w^{(t)}; i))\|^2, \quad i \in [n]. \quad (50)$$

Then, for any  $v \in \mathbb{R}^c$

$$\nabla_v \psi_i^{(t)}(v) = \eta^{(t)} H_i^{(t)T} [\eta^{(t)} H_i^{(t)} v - \alpha_i^{(t)} \nabla_z \phi_i(h(w^{(t)}; i))], \quad i \in [n]. \quad (51)$$

Consider  $\eta^{(t)} = D\sqrt{\varepsilon}$  for some  $D > 0$  and  $\varepsilon > 0$ , we have for  $i \in [n]$  and  $0 \leq t < T$

$$\|\nabla_v \psi_i^{(t)}(v)\| = (\eta^{(t)})^2 \|H_i^{(t)T} H_i^{(t)}\| \leq (\eta^{(t)})^2 \|H_i^{(t)}\| \cdot \|H_i^{(t)}\| \stackrel{(25)}{\leq} D^2 H^2.$$

Hence,  $\|\nabla_v \Psi^{(t)}(v)\| \leq D^2 H^2 + \varepsilon^2$  for any  $v \in \mathbb{R}^c$  which implies that  $L = D^2 H^2 + \varepsilon^2$  (Nesterov [2004]) and  $\frac{L}{\mu} = \frac{D^2 H^2 + \varepsilon^2}{\varepsilon^2}$ . Therefore, the complexity to find  $v^{(t)}$  for each iteration  $t$  is  $\mathcal{O}(2n \frac{D^2 H^2 + \varepsilon^2}{\varepsilon^2} \log(\frac{1}{\varepsilon}))$ .

Let us choose  $0 < \varepsilon \leq 1$ . From (49), we have

$$\begin{aligned} \frac{1}{T} \sum_{t=0}^{T-1} [F(w^{(t)}) - F_*] &\leq \frac{e^\beta L_\phi}{(1-4\alpha)\alpha\beta} \cdot \frac{1}{n} \sum_{i=1}^n \|h(w^{(0)}; i) - h_i^*\|^2 \cdot \varepsilon \\ &\quad + \frac{7e^\beta L_\phi}{2\alpha(1-4\alpha)} [D^2 H^2 + c(2 + (V + 3)GD^2)^2 + 2 + V] \cdot \varepsilon = N\varepsilon, \end{aligned}$$

where

$$N = \frac{e^\beta L_\phi}{(1-4\alpha)\alpha\beta} \frac{1}{n} \sum_{i=1}^n \|h(w^{(0)}; i) - h_i^*\|^2 + \frac{7e^\beta L_\phi}{2\alpha(1-4\alpha)} [D^2 H^2 + c(2 + (V + 3)GD^2)^2 + 2 + V].$$

Let  $\hat{\varepsilon} = N\varepsilon$  with  $0 < \hat{\varepsilon} \leq N$ . Then, we need  $T = \frac{N\beta}{\hat{\varepsilon}}$  for some  $\beta > 0$  to guarantee  $\min_{0 \leq t \leq T-1} [F(w^{(t)}) - F_*] \leq \frac{1}{T} \sum_{t=0}^{T-1} [F(w^{(t)}) - F_*] \leq \hat{\varepsilon}$ . Hence, the total complexity is  $\mathcal{O}\left(n \frac{N^3 \beta}{\hat{\varepsilon}^3} (D^2 H^2 + (\hat{\varepsilon}^2/N)) \log(\frac{N}{\hat{\varepsilon}})\right)$ . □

## References

- Zeyuan Allen-Zhu, Yuanzhi Li, and Zhao Song. A convergence theory for deep learning via over-parameterization. In Kamalika Chaudhuri and Ruslan Salakhutdinov, editors, *Proceedings of the 36th International Conference on Machine Learning*, volume 97 of *Proceedings of Machine Learning Research*, pages 242–252. PMLR, 09–15 Jun 2019. URL <http://proceedings.mlr.press/v97/allen-zhu19a.html>.
- Sanjeev Arora, Nadav Cohen, and Elad Hazan. On the optimization of deep networks: Implicit acceleration by overparameterization. In Jennifer Dy and Andreas Krause, editors, *Proceedings of the 35th International Conference on Machine Learning*, volume 80 of *Proceedings of Machine Learning Research*, pages 244–253. PMLR, 10–15 Jul 2018. URL <http://proceedings.mlr.press/v80/arora18a.html>.
- Sanjeev Arora, Simon Du, Wei Hu, Zhiyuan Li, and Ruosong Wang. Fine-grained analysis of optimization and generalization for overparameterized two-layer neural networks. In Kamalika Chaudhuri and Ruslan Salakhutdinov, editors, *Proceedings of the 36th International Conference on Machine Learning*, volume 97 of *Proceedings of Machine Learning Research*, pages 322–332. PMLR, 09–15 Jun 2019. URL <http://proceedings.mlr.press/v97/arora19a.html>.
- Léon Bottou, Frank E. Curtis, and Jorge Nocedal. Optimization methods for large-scale machine learning. *SIAM Review*, 60(2):223–311, 2018. doi: 10.1137/16M1080173.
- Alon Brutzkus, Amir Globerson, Eran Malach, and Shai Shalev-Shwartz. SGD learns over-parameterized networks that provably generalize on linearly separable data. In *International Conference on Learning Representations*, 2018. URL <https://openreview.net/forum?id=rJ33wwxBb>.
- Aaron Defazio, Francis Bach, and Simon Lacoste-Julien. Saga: A fast incremental gradient method with support for non-strongly convex composite objectives. In *Advances in Neural Information Processing Systems*, pages 1646–1654, 2014.
- Simon Du, Jason Lee, Haochuan Li, Liwei Wang, and Xiyu Zhai. Gradient descent finds global minima of deep neural networks. In Kamalika Chaudhuri and Ruslan Salakhutdinov, editors, *Proceedings of the 36th International Conference on Machine Learning*, volume 97 of *Proceedings of Machine Learning Research*, pages 1675–1685. PMLR, 09–15 Jun 2019a. URL <http://proceedings.mlr.press/v97/du19c.html>.
- Simon S. Du, Xiyu Zhai, Barnabas Poczos, and Aarti Singh. Gradient descent provably optimizes over-parameterized neural networks. In *International Conference on Learning Representations*, 2019b. URL <https://openreview.net/forum?id=S1eK3i09YQ>.
- John Duchi, Elad Hazan, and Yoram Singer. Adaptive subgradient methods for online learning and stochastic optimization. *Journal of Machine Learning Research*, 12(61):2121–2159, 2011. URL <http://jmlr.org/papers/v12/duchi11a.html>.
- S. Ghadimi and G. Lan. Stochastic first-and zeroth-order methods for nonconvex stochastic programming. *SIAM J. Optim.*, 23(4):2341–2368, 2013.
- Robert Gower, Othmane Sebbouh, and Nicolas Loizou. Sgd for structured nonconvex functions: Learning rates, minibatching and interpolation. In Arindam Banerjee and Kenji Fukumizu, editors, *Proceedings of The 24th International Conference on Artificial Intelligence and Statistics*, volume 130 of *Proceedings of Machine Learning Research*, pages 1315–1323. PMLR, 13–15 Apr 2021. URL <https://proceedings.mlr.press/v130/gower21a.html>.
- M. Gürbüzbalaban, A. Ozdaglar, and P. A. Parrilo. Convergence rate of incremental gradient and incremental newton methods. *SIAM Journal on Optimization*, 29(4):2542–2565, 2019. doi: 10.1137/17M1147846. URL <https://doi.org/10.1137/17M1147846>.
- Rie Johnson and Tong Zhang. Accelerating stochastic gradient descent using predictive variance reduction. In *Advances in Neural Information Processing Systems*, pages 315–323, 2013.
- Hamed Karimi, Julie Nutini, and Mark Schmidt. Linear convergence of gradient and proximal-gradient methods under the Polyak-Łojasiewicz condition. In Paolo Frasconi, Niels Landwehr, Giuseppe Manco, and Jilles Vreeken, editors, *Machine Learning and Knowledge Discovery in Databases*, pages 795–811, Cham, 2016. Springer International Publishing.
- Diederik P. Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *CoRR*, abs/1412.6980, 2014.
- Nicolas Le Roux, Mark Schmidt, and Francis Bach. A stochastic gradient method with an exponential convergence rate for finite training sets. In *NIPS*, pages 2663–2671, 2012.

- Kfir Y. Levy, Alp Yurtsever, and Volkan Cevher. Online adaptive methods, universality and acceleration. In S. Bengio, H. Wallach, H. Larochelle, K. Grauman, N. Cesa-Bianchi, and R. Garnett, editors, *Advances in Neural Information Processing Systems*, volume 31. Curran Associates, Inc., 2018. URL <https://proceedings.neurips.cc/paper/2018/file/b0169350cd35566c47ba83c6ec1d6f82-Paper.pdf>.
- Adrian S. Lewis and Stephen J. Wright. A proximal method for composite minimization. *Mathematical Programming*, 158:501–546, 2016.
- A. Nemirovski, A. Juditsky, G. Lan, and A. Shapiro. Robust stochastic approximation approach to stochastic programming. *SIAM J. on Optimization*, 19(4):1574–1609, 2009.
- Yurii Nesterov. *Introductory lectures on convex optimization : a basic course*. Applied optimization. Kluwer Academic Publ., Boston, Dordrecht, London, 2004. ISBN 1-4020-7553-7.
- Lam Nguyen, Phuong Ha Nguyen, Marten van Dijk, Peter Richtarik, Katya Scheinberg, and Martin Takac. SGD and Hogwild! convergence without the bounded gradients assumption. In *Proceedings of the 35th International Conference on Machine Learning-Volume 80*, pages 3747–3755, 2018.
- Lam M Nguyen, Jie Liu, Katya Scheinberg, and Martin Takáč. Sarah: A novel method for machine learning problems using stochastic recursive gradient. In *Proceedings of the 34th International Conference on Machine Learning-Volume 70*, pages 2613–2621. JMLR. org, 2017.
- Lam M. Nguyen, Phuong Ha Nguyen, Peter Richtárik, Katya Scheinberg, Martin Takáč, and Marten van Dijk. New convergence aspects of stochastic gradient algorithms. *Journal of Machine Learning Research*, 20(176):1–49, 2019. URL <http://jmlr.org/papers/v20/18-759.html>.
- Lam M. Nguyen, Quoc Tran-Dinh, Dzung T. Phan, Phuong Ha Nguyen, and Marten van Dijk. A unified convergence analysis for shuffling-type gradient methods. *Journal of Machine Learning Research*, 22(207):1–44, 2021.
- Quynh N Nguyen and Marco Mondelli. Global convergence of deep networks with one wide layer followed by pyramidal topology. In H. Larochelle, M. Ranzato, R. Hadsell, M. F. Balcan, and H. Lin, editors, *Advances in Neural Information Processing Systems*, volume 33, pages 11961–11972. Curran Associates, Inc., 2020. URL <https://proceedings.neurips.cc/paper/2020/file/8abfe8ac9ec214d68541fcb888c0b4c3-Paper.pdf>.
- Sashank J. Reddi, Ahmed Hefny, Suvrit Sra, Barnabas Poczos, and Alex Smola. Stochastic variance reduction for nonconvex optimization. In Maria Florina Balcan and Kilian Q. Weinberger, editors, *Proceedings of The 33rd International Conference on Machine Learning*, volume 48 of *Proceedings of Machine Learning Research*, pages 314–323, New York, New York, USA, 20–22 Jun 2016. PMLR. URL <https://proceedings.mlr.press/v48/reddi16.html>.
- Sashank J. Reddi, Satyen Kale, and Sanjiv Kumar. On the convergence of adam and beyond. In *International Conference on Learning Representations*, 2018. URL <https://openreview.net/forum?id=ryQu7f-RZ>.
- Shai Shalev-Shwartz, Yoram Singer, and Nathan Srebro. Pegasos: Primal estimated sub-gradient solver for svm. *Association for Computing Machinery*, 2007. doi: 10.1145/1273496.1273598. URL <https://doi.org/10.1145/1273496.1273598>.
- Daniel Soudry, Elad Hoffer, Mor Shpigel Nacson, Suriya Gunasekar, and Nathan Srebro. The implicit bias of gradient descent on separable data. *J. Mach. Learn. Res.*, 19(1):2822–2878, January 2018. ISSN 1532-4435.
- Trang H Tran, Lam M Nguyen, and Quoc Tran-Dinh. SMG: A shuffling gradient-based method with momentum. In Marina Meila and Tong Zhang, editors, *Proceedings of the 38th International Conference on Machine Learning*, volume 139 of *Proceedings of Machine Learning Research*, pages 10379–10389. PMLR, 18–24 Jul 2021. URL <https://proceedings.mlr.press/v139/tran21b.html>.
- Quoc Tran-Dinh, Nhan Pham, and Lam Nguyen. Stochastic Gauss-Newton algorithms for nonconvex compositional optimization. In Hal Daumé III and Aarti Singh, editors, *Proceedings of the 37th International Conference on Machine Learning*, volume 119 of *Proceedings of Machine Learning Research*, pages 9572–9582. PMLR, 13–18 Jul 2020. URL <https://proceedings.mlr.press/v119/tran-dinh20a.html>.
- Sharan Vaswani, Issam Laradji, Frederik Kunstner, Si Yi Meng, Mark Schmidt, and Simon Lacoste-Julien. Adaptive gradient methods converge faster with over-parameterization (but you should do a line-search), 2021.
- Junyu Zhang and Lin Xiao. A stochastic composite gradient method with incremental variance reduction. In H. Wallach, H. Larochelle, A. Beygelzimer, F. d'Alché-Buc, E. Fox, and R. Garnett, editors, *Advances in Neural Information Processing Systems*, volume 32. Curran Associates, Inc., 2019. URL <https://proceedings.neurips.cc/paper/2019/file/a68259547f3d25ab3c0a5c0adb4e3498-Paper.pdf>.
- Junyu Zhang and Lin Xiao. Multilevel composite stochastic optimization via nested variance reduction. *SIAM Journal on Optimization*, 31(2):1131–1157, 2021. doi: 10.1137/19M1285457. URL <https://doi.org/10.1137/19M1285457>.

- Difan Zou and Quanquan Gu. An improved analysis of training over-parameterized deep neural networks. In Hanna M. Wallach, Hugo Larochelle, Alina Beygelzimer, Florence d'Alché-Buc, Emily B. Fox, and Roman Garnett, editors, *Advances in Neural Information Processing Systems 32: Annual Conference on Neural Information Processing Systems 2019, NeurIPS 2019, December 8-14, 2019, Vancouver, BC, Canada*, pages 2053–2062, 2019.
- Difan Zou, Yuan Cao, Dongruo Zhou, and Quanquan Gu. Stochastic gradient descent optimizes over-parameterized deep relu networks, 2018.