051

052

053

# Statistical and Computational Complexities of BFGS Quasi-Newton Method for Generalized Linear Models

**Anonymous Authors**<sup>1</sup>

### Abstract

The gradient descent (GD) method has been used widely to solve parameter estimation in generalized linear models (GLMs), a generalization of linear models when the link function can be nonlinear. While GD has optimal statistical and computational complexities for estimating the true parameter under the high signal-to-noise ratio (SNR) regime of the GLMs, it has sub-optimal complexities when the SNR is low, namely, the iterates of GD require polynomial number of iterations to reach the final statistical radius due to the local convexity of the least-square loss functions of the GLMs in this case. Even though Newton's method can be used to resolve the flat curvature of the loss functions in the low SNR case, its computational cost is prohibitive in high-dimensional settings as it is  $\mathcal{O}(d^3)$ . To address the shortcomings of GD and Newton's method, we propose the use of BFGS quasi-Newton method to solve parameter estimation of the GLMs, which has a per iteration cost of  $\mathcal{O}(d^2)$ . On the optimization side, when the SNR is low, we demonstrate that iterates of BFGS converge linearly to the optimal solution of the population least-square loss function, and the contraction coefficient of the BFGS algorithm is comparable to that of Newton's method. On the statistical side, we prove that the iterates of BFGS reach the final statistical radius of the low SNR GLMs after a logarithmic number of iterations, which is much lower than the polynomial number of iterations of GD.

#### 1. Introduction

In supervised machine learning, we are given a set of n independent samples denoted by  $X_1, \ldots, X_n$  with corre-

sponding labels  $Y_1, \ldots, Y_n$ , that are drawn from some unknown distribution and our goal is to train a model that maps the feature vectors to their corresponding labels. We assume that the data is generated according to distribution  $\mathcal{P}_{\theta^*}$  which is parameterized by a ground truth parameter  $\theta^*$ . Our goal as the learner is to find  $\theta^*$  by solving the empirical risk minimization (ERM) problem defined as

$$\min_{\theta \in \mathbb{R}^d} \mathcal{L}_n(\theta) := \frac{1}{n} \sum_{i=1}^n \ell(\theta; (X_i, Y_i)), \tag{1}$$

where  $\ell(\theta; (X_i, Y_i))$  is the loss function that measures the error between the predicted output of  $X_i$  using parameter  $\theta$  and its true label  $Y_i$ . If we define  $\theta_n^*$  as an optimal solution of the above optimization problem, i.e.,  $\theta_n^* \in \arg\min_{\theta \in \mathbb{R}^d} \mathcal{L}_n(\theta)$ , it can be considered as an approximation of the ground-truth solution  $\theta^*$ , where  $\theta^*$  is also a minimizer of the population loss defined as

$$\min_{\theta \in \mathbb{R}^d} \mathcal{L}(\theta) := \mathbb{E}\left[\ell(\theta; (X, Y))\right].$$
 (2)

If one can solve the empirical risk efficiently, the output model could be close to  $\theta^*$ , when *n* is sufficiently large. Several works have studied the complexity of iterative methods for solving ERM or directly the population loss, for the case that the objective function is convex or strongly convex with respect to  $\theta$  (Balakrishnan et al., 2017; Ho et al., 2020; Loh & Wainwright, 2015; Agarwal et al., 2012; Yuan & Zhang, 2013; Dwivedi et al., 2020b; Hardt et al., 2016; Candes et al., 2011). However, when we move beyond linear models, the underlying loss becomes non-convex and therefore the behavior of iterative methods could substantially change, and it is not even clear if they can reach a neighborhood of a global minimizer of the ERM problem.

The focus of this paper is on the generalized linear model (GLM) (Carroll et al., 1997; Netrapalli et al., 2015; Fienup, 1982; Shechtman et al., 2015; Feiyan Tian, 2021) where the labels and features are generated according to a polynomial link function and we have  $Y_i = (X_i^{\top}\theta^*)^p + \zeta_i$ , where  $\zeta_i$  is an additive noise and  $p \ge 2$  is an integer. Due to nonlinear structure of the generative model, even if we select a convex loss function  $\ell$ , the ERM problem denoted to the considered GLM could be non-convex with respect to  $\theta$ . Interestingly, depending on the norm of  $\theta^*$ , the curvature of the ERM

<sup>&</sup>lt;sup>1</sup>Anonymous Institution, Anonymous City, Anonymous Region, Anonymous Country. Correspondence to: Anonymous Author <anon.email@domain.com>.

Preliminary work. Under review by the International Conference on Machine Learning (ICML). Do not distribute.

problem and its corresponding population risk minimization problem could change substantially. More precisely, in the 057 case that  $\|\theta^*\|$  is sufficiently large, which we refer to this 058 case as the high signal-to-noise ratio (SNR) regime, the 059 underlying population loss of the problem of interest is 060 locally strongly convex and smooth. On the other hand, in 061 the regime that  $\|\theta^*\|$  is close to zero, denoted by the low 062 SNR regime, then the underlying problem is neither strongly 063 convex nor smooth, and in fact, it is ill-conditioned.

064 These observations lead to the conclusion that in the high 065 SNR setting, due to strong convexity and smoothness of 066 the underlying problem, gradient descent (GD) reaches the 067 final statistical radius exponentially fast and overall it only 068 requires logarithmic number of iterations. However, in the 069 low SNR case, as the problem becomes locally convex, GD 070 converges at a sublinear rate to the final statistical radius and thus requires polynomial number of iterations in terms 072 of the sample size. To resolve this issue, (Ren et al., 2022a) 073 recommended the use of GD with Polyak step size to ac-074 celerate the convergence of GD in the low SNR case and 075 showed that the number of iterations becomes logarithmic 076 function of the sample size. However, as this method is still 077 a first-order method, its overall complexity scales linearly 078 by the condition number of the problem which depends on 079 the condition number of the feature vectors covariance and the norm  $\|\theta^*\|$ . Moreover, implementation of Polyak step 081 size requires access to the optimal objective function value. 082 Even in the cases that we can approximate the optimal value 083 sufficiently well, the Polyak iterates still have instability during training due to the influence of noise in the models. 085

086 An alternative approach is using Newton's method to han-087 dle the ill-conditioning issue in the low SNR case as well 088 as eliminating the need to estimate the optimal function 089 value. As we show in this paper, this idea indeed addresses the issue of poor curvature of the problem and leads to an 090 exponentially fast rate with contraction factor  $\frac{2p-2}{2p-1}$ , when the sample size is infinite. Moreover, in the high SNR 091 092 093 setting, Newton's method converges at a quadratic rate as 094 the problem is strongly convex and smooth. Alas, these 095 improvements come at the expense of increasing the com-096 putational complexity of each iteration to  $\mathcal{O}(d^3)$  which is 097 indeed more than the per iteration computational cost of GD 098 that is  $\mathcal{O}(d)$ . These points lead to this question:

## Is there a computationally-efficient method that performs well in both high and low SNR settings at a reasonable per iteration computational cost?

099

100

104 105 106 106 106 106 107 108 109 **Contributions.** In this paper, we address this question and show that the BFGS method is capable of achieving these goals. BFGS is a quasi-Newton method that approximates the objective function curvature using gradient information and its per iteration cost is  $\mathcal{O}(d^2)$ . It is well-known that it enjoys a superlinear convergence rate that is independent of condition number in strongly convex and smooth settings, and hence, in the high SNR setting it outperforms GD. In the low SNR setting, where the Hessian at the optimal solution could be singular, we show that the BFGS method converges linearly and outperforms the sublinear convergence rate of GD. Next, we formally summarize our contributions.

- Infinite sample, low SNR: For the infinite sample case where we minimize the population loss, we show that in the low SNR case the iterates of BFGS converge to the ground truth  $\theta^*$  at an exponentially fast rate that is independent of all problem parameters except the power of link function p. We further show that the linear convergence contraction coefficient of BFGS is comparable to that of Newton's method. This convergence result of BFGS is also of general interest as it provides the first global linear convergence of BFGS without line-search for a setting that is neither strictly nor strongly convex.
- Finite sample, low SNR: By leveraging the results developed for the population loss of the low SNR regime, we show that in the finite sample case, the BFGS iterates converge to the final statistical radius  $\mathcal{O}(1/n^{1/(2p+2)})$  within the true parameter after a logarithmic number of iterations  $\mathcal{O}(\log(n))$ . It is substantially lower than the required number of iterations for fixed-step size GD, which is  $\mathcal{O}(n^{(p-1)/p})$ , to reach a similar statistical radius. This improvement is the direct outcome of the linear convergence of BFGS versus the sublinear convergence rate of GD in the low SNR case. Further, while the iteration complexity of BFGS is comparable to the logarithmic number of iterations of GD with Polyak step size, we show that BFGS removes the dependency of the overall complexity on the condition number of the problem as well as the need to estimate the optimal function value.
- Experiments: We conduct numerical experiments for both infinite and finite sample cases to compare the performance of GD (with constant stepsize and Polyak stepsize), Newton's method and BFGS. The provided empirical results are consistent with our theoretical findings and show the advantages of BFGS in the low SNR regime.

## 2. BFGS Algorithm

In this section, we review the basics of the BFGS quasi-Newton method, which is the main algorithm we analyze. Consider the case that we aim to minimize a differentiable convex function  $f : \mathbb{R}^d \to \mathbb{R}$ . The BFGS update is given by

$$\theta_{k+1} = \theta_k - \eta_k H_k \nabla f(\theta_k), \qquad \forall k \ge 0, \qquad (3)$$

where  $\eta_k$  is the step size and  $H_k \in \mathbb{R}^{d \times d}$  is a positive definite matrix that aims to approximate the true Hessian 111 inverse  $\nabla^2 f(\theta_k)^{-1}$ . There are several approaches for ap-112 113 proximating  $H_k$  leading to different quasi-Newton methods, 114 (Conn et al., 1991; Broyden, 1965; Broyden et al., 1973; 115 Gay, 1979; Davidon, 1959; Fletcher & Powell, 1963; Broy-116 den, 1970; Fletcher, 1970; Goldfarb, 1970; Shanno, 1970; 117 Nocedal, 1980; Liu & Nocedal, 1989), but in this paper, 118 we focus on the celebrated BFGS method, in which  $H_k$  is 119 updated as 120

$$H_{k} = \left(I - \frac{s_{k-1}u_{k-1}^{\top}}{s_{k-1}^{\top}u_{k-1}}\right)H_{k-1}\left(I - \frac{u_{k-1}s_{k-1}^{\top}}{s_{k-1}^{\top}u_{k-1}}\right) + \frac{s_{k-1}s_{k-1}^{\top}}{s_{k-1}^{\top}u_{k-1}}, \quad \forall k \ge 1,$$
(4)

121

122

124

125

126

127

128

129 130

131

159

160

161

162

163

164

where the variable variation  $s_{k-1}$  and gradient displacement  $u_{k-1}$  are defined as

$$s_{k-1} := \theta_k - \theta_{k-1},$$
  

$$u_{k-1} := \nabla f(\theta_k) - \nabla f(\theta_{k-1}), \qquad \forall k \ge 1,$$
(5)

132 133 respectively. The logic behind the update in (4) is to en-134 sure that the Hessian inverse approximation  $H_k$  satisfies the 135 secant condition  $H_k u_{k-1} = s_{k-1}$ , while it stays close to 136 the previous approximation matrix  $H_{k-1}$ . The update in (4) 137 only requires matrix-vector multiplications, and hence, the 138 computational cost per iteration of BFGS is  $O(d^2)$ .

139 The main advantage of BFGS is its fast superlinear conver-140 gence rate under the assumption of strict convexity, i.e., 141  $\lim_{k\to\infty} \|\theta_k - \theta_{opt}\| / \|\theta_{k-1} - \theta_{opt}\| = 0$ , where  $\theta_{opt}$  is 142 the optimal solution. Previous results on the superlinear 143 convergence of quasi-Newton methods were all asymptotic, 144 until the recent advancements on the non-asymptotic anal-145 ysis of quasi-Newton methods (Rodomanov & Nesterov, 146 2021a;b;c; Jin & Mokhtari, 2020; Jin et al., 2022; Ye et al., 147 2021; Lin et al., 2021a;b). For instance, (Jin & Mokhtari, 148 2020) established a local superlinear convergence rate of 149  $(1/\sqrt{k})^k$  for BFGS. However, all these superlinear conver-150 gence analyses require the objective function to be smooth 151 and strictly or strongly convex. Alas, these conditions are 152 not satisfied in the low SNR setting, since the Hessian at 153 the optimal solution could be singular, and hence the loss 154 function is neither strongly convex nor strictly convex; we 155 further discuss this point in Section 3. This observation 156 implies that we need novel convergence analyses to study 157 the behavior of BFGS in the low SNR setting. 158

## 3. Generalized Linear Model with Polynomial Link Function

In this section, we formally present the generalized linear model (GLM) setting that we consider in our paper, and discuss the low and high SNR settings and optimization challenges corresponding to these cases. Consider the case that the feature vectors are denoted by  $X \in \mathbb{R}^d$  and their corresponding labels are denoted by  $Y \in \mathbb{R}$ . Suppose that we have access to *n* sample points  $(Y_1, X_1), (Y_2, X_2), \ldots, (Y_n, X_n)$  that are i.i.d. samples from the following generalized linear model with polynomial link function of power *p* (Carroll et al., 1997), i.e.,

$$Y_i = (X_i^\top \theta^*)^p + \zeta_i, \tag{6}$$

where  $\theta^*$  is a true but unknown parameter,  $p \in \mathbb{N}$  is a given power, and  $\zeta_1, \ldots, \zeta_n$  are independent Gaussian noises with zero mean and variance  $\sigma^2$ . The Gaussian assumption on the noise is for the simplicity of the discussion and similar results hold for the sub-Gaussian i.i.d. noise case. Furthermore, we assume the feature vectors are generated as  $X \sim \mathcal{N}(0, \Sigma)$  where  $\Sigma \in \mathbb{R}^{d \times d}$  is a symmetric positive definite matrix. Here we focus on the settings that  $p \in \mathbb{N}^+$ and  $p \geq 2$ .

The above class of GLMs with polynomial link functions arise in several settings. When p = 1, the model in (6) is the standard linear regression model, and for the case that p = 2, the above setup corresponds to the phase retrieval model (Fienup, 1982; Shechtman et al., 2015; Candes et al., 2011; Netrapalli et al., 2015), which has found applications in optical imaging, x-ray tomography, and audio signal processing. Moreover, the analysis of GLMs with  $p \ge 2$  also serves as the basis of the analysis on other popular statistical models. For example, as shown by Yi & Caramanis (2015); Wang et al. (2015); Xu et al. (2016); Balakrishnan et al. (2017); Daskalakis et al. (2017); Dwivedi et al. (2020a); Kwon et al. (2019); Dwivedi et al. (2020b); Kwon et al. (2021), the local landscape of log-likelihood for Gaussian mixture models and mixture linear regression models are identical to GLMs for p = 2.

In the case that the polynomial link function parameter in the GLM model in (6) is p = 2, by adapting similar arguments from (Kwon et al., 2021) under the symmetric twocomponent location Gaussian mixture, there are essentially three regimes for estimating the true parameter  $\theta^*$ : Low signal-to-noise ratio (SNR) regime:  $\|\theta^*\|/\sigma \leq C_1(d/n)^{1/4}$ where d is the dimension, n is the sample size, and  $C_1$  is a universal constant; Middle SNR regime:  $C_1(d/n)^{1/4} \leq$  $\|\theta^*\|/\sigma \leq C$  where C is a universal constant; High SNR regime:  $\|\theta^*\|/\sigma \geq C$ . The main idea is that with different  $\theta^*$ , the optimization landscape of the parameter estimation problem changes. By generalizing the insights from the case p = 2, we define the following regimes for any  $p \geq 2$ :

- (i) Low SNR regime: ||θ\*||/σ ≤ C<sub>1</sub>(d/n)<sup>1/(2p)</sup> where d is the dimension, n is the sample size, and C<sub>1</sub> is a universal constant;
- (ii) Middle SNR regime:  $C_1(d/n)^{1/(2p)} \leq \|\theta^*\|/\sigma \leq$

C where C is a universal constant;

• (iii) High SNR regime:  $\|\theta^*\|/\sigma \ge C$ .

Note that, the rate  $(d/n)^{1/2p}$  that we use to define the SNR regimes is from the statistical rate of estimating the true parameter  $\theta^*$  when  $\theta^*$  approaches 0. Next, we provide insight into the landscape of the least-square loss function for each regime. In particular, given the GLM in (6), the sample least-square takes the following form:

$$\min_{\theta \in \mathbb{R}^d} \mathcal{L}_n(\theta) := \frac{1}{n} \sum_{i=1}^n \left( Y_i - \left( X_i^\top \theta \right)^p \right)^2.$$
(7)

To obtain insight about the landscape of the loss function  $\mathcal{L}_n$ , a useful approach is to consider an approximation of that function by its population version, which we refer to as population least-square loss function and is given by:

$$\min_{\theta \in \mathbb{R}^d} \mathcal{L}(\theta) := \mathbb{E}[\mathcal{L}_n(\theta)],\tag{8}$$

where the outer expectation is taken with respect to the data.

**High SNR regime.** In the setting that the ground truth parameter has a relatively large norm, i.e.,  $\|\theta^*\| \ge C$  for some constant C > 0 that only depends on  $\sigma$ , the population loss in (8) is locally strongly convex and smooth around  $\theta^*$ . More precisely, when  $\|\theta - \theta^*\|$  is small, we have

$$(X^{\top}\theta^{*})^{p} - (X^{\top}\theta)^{p}$$
  
=  $p(X^{\top}\theta^{*})^{p-1}X^{\top}(\theta - \theta^{*}) + o(\|\theta - \theta^{*}\|)$ 

Hence, in a neighborhood of the optimal solution, the objective in (8) can be approximated as

$$\mathcal{L}(\theta) = p^2 (\theta - \theta^*)^\top \mathbb{E}_X \left[ X (X^\top \theta^*)^{2p-2} X^\top \right] (\theta - \theta^*) + \sigma^2 + o(\|\theta - \theta^*\|^2).$$

Indeed, if  $||\theta^*|| \ge C\sigma$  the above function behaves as a quadratic function that is smooth and strongly convex, assuming that  $o(||\theta - \theta^*||^2)$  is negligible. As a result, the iterates of gradient descent (GD) converge to the solution at a linear rate and it requires  $\kappa \log(1/\epsilon)$  to reach an  $\epsilon$ -accurate solution, where  $\kappa$  depends on the conditioning of the covariance matrix  $\Sigma$  and the norm of  $\theta^*$ . In this case, BFGS converges superlinearly to  $\theta^*$  and the rate would be independent of  $\kappa$ , while the cost per iteration would be  $\mathcal{O}(d^2)$ .

Low SNR regime. As mentioned above, in the high SNR case, GD has a fast linear rate. However, in the low SNR case where  $\|\theta^*\|$  is small and  $\|\theta^*\| \leq C_1(d/n)^{1/(2p)}$ , the strong convexity parameter approaches zero when the sample size *n* goes to infinity and the problem becomes illconditioned. In this case, we deal with a *function that is only convex and its gradient is not Lipschitz continuous*. To better elaborate on this point, let us focus on the case that  $\theta^* = 0$  as a special case of the low SNR setting. Considering the underlying distribution of X, which is  $X \sim \mathcal{N}(0, \Sigma)$ , for such a low SNR case, the population loss can be written as

$$\mathcal{L}(\theta) = \mathbb{E}_X \left[ (X^\top \theta)^{2p} \right] + \sigma^2 = (2p-1)!! \|\Sigma^{1/2}\theta\|^{2p} + \sigma^2.$$
(9)

Since we focus on  $p \ge 2$  it can be verified that  $\mathcal{L}(\theta)$  is not strongly convex in a neighborhood of the solution  $\theta^* = 0$ . For this class of functions, it is well-known that GD with constant step size would converge at a sublinear rate, and hence GD iterates require polynomial number of iterations to reach the final statistical radius. In the next section, we study the behavior of BFGS for solving the low SNR setting and showcase its advantage compared to GD.

Middle SNR regime. Different from the low and high SNR regimes, the middle SNR regime is generally harder to analyze as the landscapes of both the population and sample least-square loss functions are complex. Adapting the insight from middle SNR regime of the symmetric twocomponent location Gaussian mixtures from (Kwon et al., 2021), for the middle SNR setting of the generalized linear model, the eigenvalues of the Hessian matrix of the population least-square loss function approach 0 and their vanishing rates depend on some increasing function of  $\|\theta^*\|$ . The optimal statistical rate of the optimization algorithms, such as gradient descent algorithm, for solving  $\theta^*$  depends strictly on the tightness of these vanishing rates in terms of  $\|\theta^*\|$ , which are non-trivial to obtain. In fact, to the best of our knowledge, there is no result on the convergence of iterative methods (such as GD or its variants) for GLMs with a polynomial link function in the middle SNR. Hence, we leave the study of BFGS for this setting as a future work.

## 4. Convergence Analysis in the Low SNR Regime: Population Loss

In this section, we focus on the convergence properties of BFGS for the population loss in the low SNR case introduced in (9). This analysis provides an intuition for the analysis of the finite sample case that we discuss in Section 5, as we expect these two loss functions to be close to each other when the number of samples n is sufficiently large. Note that the loss function in (9) can be considered as a special case of the following convex optimization problem:

$$\min_{\theta \in \mathbb{R}^d} f(\theta) = \|A\theta - b\|^q, \tag{10}$$

where  $A \in \mathbb{R}^{m \times d}$  is a matrix,  $b \in \mathbb{R}^m$  is a given vector, and q satisfies  $q \ge 4$ . We should note that for  $q \ge 4$ , the considered objective is not strictly convex because the Hessian is singular when  $A\theta = b$ . Indeed, if we set m = dand further let A be  $\Sigma^{1/2}$  and choose  $b = A\theta^* = 0$ , then we recover the problem in (9) for q = 2p. 220 Notice that the problem in (10), which serves as a surrogate 221 for the finite sample problem that we plan to study in the 222 next section, has the same solution set as the quadratic prob-223 lem of minimizing  $||A\theta - b||^2$  with solution  $(A^{\top}A)^{-1}A^{\top}b$ 224 when  $A^{\top}A$  is invertible. Given this point, one might sug-225 gest that instead of minimizing (10) we could directly solve the quadratic problem which is indeed much easier to solve. 227 This point is valid, but the goal of this section is not to effi-228 ciently solve problem (10) itself. Our goal is to understand 229 the convergence properties of the BFGS method for solv-230 ing the problem in (10) with the hope that it will provide 231 some intuitions for the convergence analysis of BFGS for 232 the empirical loss (7) in the low SNR regime. As we will 233 discuss in Section 5, the convergence analysis of BFGS on 234 the population loss which is a special case of (10) is closely 235 related to the one for the empirical loss in (7). 236

One remark is that population loss in (9) holds for the re-237 strictive assumption of  $\theta^* = 0$ , which is only a special case 238 of the general low SNR regime of  $\|\theta^*\| \leq C_1 (d/n)^{1/(2p)}$ . 239 Our ultimate goal is to analyze the convergence behavior of 240 the BFGS method applied to the empirical loss (7) of GLM 241 problems. The errors between gradients and Hessians of 242 the population loss with  $\theta^* = 0$  and  $\|\theta^*\| \leq C_1 (d/n)^{1/(2p)}$ 243 are upper bounded by the corresponding statistical errors 244 between the population loss (8) and the empirical loss (7)245 in the low SNR regime, respectively. Therefore, the errors 246 between iterations of applying BFGS to the population loss 247 with  $\theta^* = 0$  and  $\|\theta^*\| \le C_1 (d/n)^{1/(2p)}$  are negligible com-248 pared to the statistical errors. Instead of directly analyzing 249 BFGS for the population loss (8) in the general low SNR 250 regime, studying the convergence properties of BFGS for 251 the population loss (9) with  $\theta^* = 0$  can equivalently lay 252 foundations for the convergence analysis of BFGS for the 253 empirical loss (7) in the low SNR regime, which is the main 254 target of this paper. The details can be found in the proof of 255 Theorem 5.1. 256

Moreover, our results would be of general interest from an
optimization point of view, since there is no global convergence theory (without line-search) for the BFGS method in
the literature for the case that the objective function is not
strictly convex, and our analysis provides the first result for
such general settings. Before, stating our results, we first
state the assumptions.

Assumption 4.1. There exists  $\hat{\theta} \in \mathbb{R}^d$ , such that  $b = A\hat{\theta}$ . In other words, b is in the range of matrix A.

265

266

267

272

273

274

This assumption implies that the problem in (10) is realizable,  $\hat{\theta}$  is an optimal solution, and the optimal function value is zero. This assumption is satisfied in our considered low SNR setting in (9) as  $\theta^* = 0$  which implies b = 0.

Assumption 4.2. The matrix  $A^{\top}A \in \mathbb{R}^{d \times d}$  is invertible. This is equivalent to  $A^{\top}A \succ 0$ . The above assumption is also satisfied for our considered setting as we assume that the covariance matrix for our input features is positive definite. Combining Assumptions 4.1 and 4.2, we conclude that  $\hat{\theta}$  is the unique optimal solution of problem (10). Next, we state the convergence rate of BFGS for solving problem 10 under the disclosed assumptions.

**Theorem 4.3.** Consider the BFGS method in (3)-(5). Suppose Assumptions 4.1 and 4.2 hold, and the initial Hessian inverse approximation matrix is selected as  $H_0 = \nabla^2 f(\theta_0)^{-1}$ , where  $\theta_0 \in \mathbb{R}^d$  is an arbitrary initial point. If the step size is  $\eta_k = 1$  for all  $k \ge 0$ , then the iterates of BFGS converge to the optimal solution  $\hat{\theta}$  at a linear rate of

$$\|\theta_k - \hat{\theta}\| \le r_{k-1} \|\theta_{k-1} - \hat{\theta}\|, \quad \forall k \ge 1, \qquad (11)$$

where the contraction factors  $r_k \in [0, 1)$  satisfy

$$r_0 = \frac{q-2}{q-1}, \qquad r_k = \frac{1 - r_{k-1}^{q-2}}{1 - r_{k-1}^{q-1}}, \quad \forall k \ge 1.$$
 (12)

The proof of this theorem is available in Appendix A.1. Theorem 4.3 shows that the iterates of BFGS converge globally at a linear rate to the optimal solution of (10). This result is of interest as it illustrates the iterates generated by BFGS converge globally without any line search scheme and the stepsize is fixed as  $\eta_k = 1$  for any  $k \ge 0$ . Moreover, the initial point  $\theta_0$  is arbitrary and there is no restriction on the distance between  $\theta_0$  and the optimal solution  $\hat{\theta}$ . This result is in contrast to most analyses of quasi-Newton methods which require the initial point  $\theta_0$  to be in a local neighborhood of  $\hat{\theta}$  to guarantee the linear or superlinear convergence rate, without line-search.

Remark 4.4. The cost of computing the initial Hessian inverse approximation  $H_0 = \nabla^2 f(\theta_0)^{-1}$  is  $\mathcal{O}(d^3)$ , but this cost is only required for the first iteration, and it is not required for  $k \ge 1$  as for those iterates we update the Hessian inverse approximation matrix  $H_k$  based on the update in (4) at a cost of  $\mathcal{O}(d^2)$ .

Note that the result in Theorem 4.3 does not specify the exact complexity of BFGS for solving problem(10), as the contraction factors  $r_k$  are not explicitly given. In the following theorem, we show that for  $q \ge 4$ , the linear rate contraction factors  $\{r_k\}_{k=0}^{\infty}$  also converge linearly to a fixed point contraction factor  $r_*$  determined by the parameter q.

**Theorem 4.5.** Consider the linear convergence factors  $\{r_k\}_{k=0}^{\infty}$  defined in (12) from Theorem 4.3. If  $q \ge 4$ , then the sequence  $\{r_k\}_{k=0}^{\infty}$  converges to  $r_* \in (0, 1)$  that is determined by the equation

$$r_*^{q-1} + r_*^{q-2} = 1, (13)$$

and the rate of convergence is linear with a contraction factor that is at most 1/2, i.e.,

$$|r_k - r_*| \le (1/2)^k |r_0 - r_*|, \quad \forall k \ge 0.$$
 (14)

275 The proof is in Appendix A.2. Based on Theorem 4.5, the 276 iterates of BFGS eventually converge to the optimal solution 277 at the linear rate of  $r_*$  determined by (13). Specifically, the 278 factors  $\{r_k\}_{k=0}^{\infty}$  converge to the fixed point  $r_*$  at a linear 279 rate with the contraction factor of 1/2. Further, the linear 280 convergence factors  $\{r_k\}_{k=0}^{\infty}$  and their limit  $r_*$  are all only 281 determined by q, and they are independent of the dimension 282 d and the condition number  $\kappa_A$  of the matrix A. Hence, the performance of BFGS is not influenced by high-dimensional 284 or ill-conditioned problems. This result is independently im-285 portant from an optimization point of view, as it provides the first global linear convergence of BFGS without line-search 287 for a setting that is not strictly or strongly convex, and interestingly the constant of linear convergence is independent 289 of dimension or condition number.

We illustrate the convergence of factors  $\{r_k\}_{k=0}^{\infty}$  to the fixed point  $r_*$  in Figure 4 of the appendix B for q = 4and q = 100. In plots (a) and (b), we observe that  $r_k$ becomes close to  $r_*$  after only 5 iterations. Hence, the linear convergence rate of BFGS is approximately  $r_*$  after only a few iterations. We further observe in plots (c) and (d) that the factors  $\{r_k\}_{k=0}^{\infty}$  converge to the fixed point  $r_*$  at a linear rate upper bounded by 1/2. Note that  $r_*$  is the solution of (13). These plots verify our results in Theorem 4.5.

#### 4.1. Comparison with Newton's Method

Next, we compare the convergence results of BFGS for solving problem (10) with the one for Newton's method. The following theorem characterizes the global linear convergence of Newton's method with unit step size applied to the objective function in (10).

**Theorem 4.6.** Consider applying Newton's method to optimization problem (10) and suppose Assumptions 4.1 and 4.2 hold. Moreover, suppose the step size is  $\eta_k = 1$  for any  $k \ge 0$ . Then, the iterates of Newton's method converge to the optimal solution  $\hat{\theta}$  at a linear rate of

$$\|\theta_k - \hat{\theta}\| = \frac{q-2}{q-1} \|\theta_{k-1} - \hat{\theta}\|, \quad \forall k \ge 1.$$
 (15)

Moreover, this linear convergence rate  $\frac{q-2}{q-1}$  is smaller than the fixed point  $r_*$  defined in (13) of the BFGS quasi-Newton method, i.e.,  $\frac{q-2}{q-1} < r_* < \frac{2q-3}{2q-2}$  for all  $q \ge 4$ .

The proof is available in Appendix A.3. The convergence results of Newton's method are also global without any line search method, and the linear rate  $\frac{q-2}{q-1}$  is independent of dimension d and condition number  $\kappa_A$ . Furthermore, the condition  $\frac{q-2}{q-1} < r_*$  implies that iterates of Newton's method converge faster than BFGS, but the gap is not substantial as  $r_* < \frac{2q-3}{2q-2}$ . On the other hand, the computational cost per iteration of Newton's method is  $\mathcal{O}(d^3)$  which is worse than the  $\mathcal{O}(d^2)$  of BFGS. Moving back to our main problem, one important implication of the above convergence results is that in the low SNr setting the iterates of BFGS converge linearly to the optimal solution of the population loss function, while the contraction coefficient of BFGS is comparable to that of Newton's method which is (2p-2)/(2p-1). For instance, for p = 2, 3, 5, 10, the linear rate contraction factor of Newton's method are 0.667, 0.8, 0.889, 0.947 and the approximate linear rate contraction factor of BFGS denoted by  $r_*$ are 0.755, 0.857, 0.922, 0.963, respectively.

## 5. Convergence Analysis in the Low SNR Regime: Finite Sample Setting

Thus far, we have demonstrated that the BFGS iterates converge linearly to the true parameter  $\theta^*$  when minimizing the population loss function  $\mathcal{L}$  of the GLM in (9) in the low SNR regime. In this section, we study the statistical behavior of the BFGS iterates for the finite sample case by leveraging the insights developed in the previous section about the convergence rate of BFGS in the infinite sample case, i.e., population loss. More precisely, we focus on the application of BFGS for solving the least-square loss function  $\mathcal{L}_n$  defined in (7) for the low SNR setting. The iterates of BFGS in this case follow the update rule

$$\theta_{k+1}^n = \theta_k^n - \eta_k H_k^n \nabla \mathcal{L}_n(\theta_k^n), \tag{16}$$

where  $H_k^n$  is updated using the gradient information of the loss  $\mathcal{L}_n$  by the BFGS rule.

We next show that the BFGS iterates (16)  $\{\theta_k^n\}_{k\geq 0}$  converge to the final statistical radius within a logarithmic number of iterations under the low SNR regime of the GLMs. To prove this claim, we track the difference between the iterates  $\{\theta_k^n\}_{k\geq 0}$  generated based on the empirical loss and the iterates  $\{\theta_k\}_{k\geq 0}$  generated according to the population loss. Assuming that they both start from the same initialization  $\theta_0$ , with the concentration of the gradient  $\|\nabla \mathcal{L}_n(\theta) - \nabla \mathcal{L}(\theta)\|$ and the Hessian  $\|\nabla^2 \mathcal{L}_n(\theta) - \nabla^2 \mathcal{L}(\theta)\|_{op}$  from Mou et al. (2019); Ren et al. (2022b) we control the deviation between these two sequences. Using this bound and the convergence results of the iterates generated based on the population loss discussed in the previous section, we prove the following result for the finite sample setting.

**Theorem 5.1.** Consider the low SNR regime of the GLM in (6) namely,  $\|\theta^*\| \leq C_1(d/n)^{1/(2p)}$ . Apply the BFGS method to the empirical loss (7) with the initial Hessian inverse approximation matrix as

$$H_0 = \nabla^2 f(\theta_0^n)^{-1},$$

where  $\theta_0^n \in \mathbb{B}(\theta^*, r)$  for some r > 0 and step size  $\eta_k = 1$ . For any failure probability  $\delta \in (0, 1)$ , if the number of samples is  $n \ge C_2(d \log(d/\delta))^{2p}$ , and the number of iterations

satisfies  $T \ge \log(n/d(\log(1/\delta)))$ , then with probability  $1 - \delta$ , we have

$$\min_{t \in [T]} \|\theta_t^n - \theta^*\| \le C_3 \left(\frac{d \log(1/\delta)}{n}\right)^{\frac{1}{2p+2}}, \quad (17)$$

where  $C_1$ ,  $C_2$ , and  $C_3$  are universal constants independent of n and d, and  $C_3$  has a polynomial dependency on r.

The proof is available in Appendix A.4. Theorem 5.1 shows that BFGS achieves the statistical accuracy in  $O(\log n)$  iterations, which is faster than the sublinear convergence  $O(n^{\frac{p-1}{p}})$  of GD shown in (Ren et al., 2022a). A few comments about Theorem 5.1 are in order.

344 Comparing to GD, GD with Polyak step size, and 345 Newton's method: Theorem 5.1 indicates that under the 346 low SNR regime, the BFGS iterates reach the final statis-347 tical radius  $\mathcal{O}(n^{-1/(2p+2)})$  within the true parameter  $\theta^*$ 348 after  $\mathcal{O}(\log(n))$  number of iterations. The statistical ra-349 dius is slightly worse than the optimal statistical radius 350  $O(n^{-1/(2p)})$ . However, we conjecture that this is due to 351 the proof technique and BFGS can still reach the optimal 352  $O(n^{-1/(2p)})$  in practice. In our experiments, in the next 353 section, we observe that when d = 4 and p = 2, the sta-354 tistical radius of BFGS is closer to the optimal radius of 355  $O(n^{-1/4})$  instead of  $O(n^{-1/6})$  suggested by our analysis. 356 We leave an improvement of the statistical analysis as the 357 future work. On the other hand, the overall computational 358 complexity of BFGS, which is  $\mathcal{O}(\log(n))$ , is indeed better 359 than the polynomial number of iterations of GD, which is at 360 the order of  $\mathcal{O}(n^{(p-1)/p})$  (Corollary 3 in (Ho et al., 2020)). 361

Moreover, the complexity of BFGS is better than the one for 362 GD with Polyak step size which is  $\mathcal{O}(\kappa \log(n))$  iterations 363 (Corollary 1 in (Ren et al., 2022a)), where  $\kappa$  is the condition 364 number of the covariance matrix  $\Sigma$ . Note that while the 365 iteration complexity of BFGS is comparable to that of GD with Polyak step size in terms of the sample size, the BFGS 367 overcomes the need to approximate the optimal value of the sample least-square loss  $\mathcal{L}_n$ , which can be unstable in 369 practice, and also removes the dependency on the condition 370 number that appears in the complexity bound of GD with 371 Polyak step size. Finally, the iteration complexity of BFGS is comparable to the  $\mathcal{O}(\log(n))$  of Newton's method (Corol-373 374 lary 3 in (Ho et al., 2020)), while per iteration cost of BFGS is lower than Newton's method. 375

376 On the minimum number of iterations: The results in 377 Theorem 5.1 involve the minimum number of iterations, 378 namely, this result holds for some  $1 \le t \le T$ . It suggests 379 that the BFGS iterates may diverge after they reach the 380 final statistical radius. As highlighted in (Ho et al., 2020), 381 such instability behavior of BFGS is inherent to fast and 382 unstable methods. While it may sound limited, this issue 383 can be handled via an early stopping scheme using the cross-384

validation approaches. We illustrate such early stopping of the BFGS iterates for the low SNR regime in Figure 2.

#### 6. Numerical Experiments

We divide our experiments into two sections where the first one focuses on the behavior of iterative methods on the population loss of GLM with polynomial link functions and the second one focuses on the finite sample setting.

Experiments for the population loss function. In this section, we compare the performance of Newton's method, BFGS, GD with constant step size, and GD with Polyak step size applied to (10) which corresponds to the population loss. We choose different values of parameter m, dimension d and the exponential parameter q in (10). We generate a random matrix  $A \in \mathbb{R}^{m \times d}$  and a random vector  $\hat{\theta} \in \mathbb{R}^d$ , and compute the vector  $b = A\hat{\theta} \in \mathbb{R}^d$ . The initial point  $\theta_0 \in \mathbb{R}^d$  is also generated randomly. The GD constant step size  $\eta$  is tuned by hand to achieve the best performance of GD on each problem. We present the logarithmic scale of  $||\theta_k - \hat{\theta}||$  versus the number of iteration k for different algorithms. All the values of different parameters m, d, q and  $\eta$  as well as the numerical results of our experiments are shown in Figure 1.

We observe that GD with constant step converges slowly due to its sub-linear convergence rate. The performance of GD with Polyak step size is also poor when dimension is large or the parameter q is huge. This is due to the fact that as dimension increases the problem becomes more illconditioned and hence the linear convergence contraction factor approaches 1. We observe that both Newton's method and BFGS generate iterations with linear convergence rates, and their linear convergence rates are only affected by the parameter q, i.e., the dimension d has no impact over the performance of BFGS and Newton's method. Although the convergence speed of Newton's method is faster than BFGS, their gap is not significantly large as we expected from our theoretical results in Section 4.

**Experiments for the empirical loss function.** We next study the statistical and computational complexities of BFGS on the empirical loss. In our experiments, we first consider the case that d = 4 and the power of the link function is p = 2, namely, we consider the multivariate setting of the phase retrieval problem. The data is generated by first sampling the inputs according to  $\{X_i\}_{i=1}^n \sim \mathcal{N}(0, \operatorname{diag}(\sigma_1^2, \dots, \sigma_4^2))$  where  $\sigma_k = (0.5)^{k-1}$ , and then generating their labels based on  $Y_i = (X_i^\top \theta^*)^2 + \zeta_i$  where  $\{\zeta_i\}_{i=1}^n$  are i.i.d. samples from  $\mathcal{N}(0, 0.01)$ . In the low SNR regime, we set  $\theta^* = 0$ , and in the high SNR regime we select  $\theta^*$  uniformly at random from the unit sphere. Further, for GD, we set the step size as  $\eta = 0.1$ , while for Newton's method and BFGS, we use the unit stepsize  $\eta = 1$ .



Figure 1. Convergence of Newton's method, BFGS, GD with constant step size and GD with Polyak step size for different values of d and q. In plot (a), m = 100 and  $\eta = 10^{-4}$ . In plot (b), m = 100 and  $\eta = 10^{-8}$ . In plot (c), m = 2000 and  $\eta = 10^{-12}$ . In plot (d), m = 2000 and  $\eta = 10^{-15}$ .



*Figure 2.* Illustration of different methods with high SNR regime in (a) and low SNR regime in (b). Illustration of the statistical radius of BFGS with high SNR regime in (c) and low SNR regime in (d).

In plots (a) and (b) of Figure 2, we consider the setting that 414 the sample size is  $n = 10^4$ , and we run GD, GD with Polyak 415 step size, BFGS, and Newton's method to find the optimal 416 solution of the sample least-square loss  $\mathcal{L}_n$ . Furthermore, 417 for both Newton's method and the BFGS algorithm, due to 418 their instability, we also perform cross-validation to choose 419 420 their early stopping. In particular, we split the data into training and the test sets. The training set consists of 90% of 421 the data while the test set has 10% of the data. The yellow 422 points in plots (a) and (b) of Figure 2 show the iterates of 423 BFGS and Newton, respectively, with the minimum valida-424 tion loss. As we observe, under the low SNR regime, the 425 iterates of GD with Polyak step size, BFGS and Newton's 426 method converge geometrically fast to the final statistical 427 radius while those of the GD converge slowly to that radius. 428 Under the high SNR regime, the iterates of all of these meth-429 ods converge geometrically fast to the final statistical radius. 430 The faster convergence of GD with Polyak step size over 431 GD is due to the optimality of step size of, while the faster 432 convergence of BFGS and Newton's method over GD is 433 due to their independence on the problem condition number. 434 Finally, in plots (c) and (d) of Figure 2, we run BFGS when 435 the sample size is ranging from  $10^2$  to  $10^4$  to empirically 436 verify the statistical radius of these methods. As indicated 437 in the plots of that figure, under the high SNR regime, the 438 439

BFGS has statistical radius is  $\mathcal{O}(n^{-1/2})$ , while under the low SNR regime, its statistical radius becomes  $\mathcal{O}(n^{-1/4})$ .

We present additional numerical experiments regarding the linear contraction factors, experiments on the empirical loss with high dimensions and empirical results for the middle SNR regimes in the appendix B.

#### 7. Conclusions

In this paper, we analyzed the convergence rates of BFGS on both population and empirical loss functions of the generalized linear model in the low SNR regime. We showed that in this case, BFGS outperforms GD and performs similar to Newton's method in terms of iteration complexity, while it requires a lower per iteration computational complexity compared to Newton's method. We also provided experiments for both infinite and finite sample loss functions and showed that our empirical results are consistent with our theoretical findings. Perhaps one limitation of the BFGS method is that its computational cost is still not linear in the dimension and scales as  $\mathcal{O}(d^2)$ . One future research direction is to analyze some other iterative methods such as limited memory-BFGS (L-BFGS) which may be able to achieve a fast linear convergence rate in the low SRN setting, while its computational cost per iteration is  $\mathcal{O}(d)$ .

## 440 **References**

473

474

475

476

- Agarwal, A., Negahban, S., and Wainwright, M. J.
  Fast global convergence of gradient methods for highdimensional statistical recovery. *Annals of Statistics*, 40 (5):2452–2482, 2012.
- Balakrishnan, S., Wainwright, M. J., and Yu, B. Statistical
  guarantees for the EM algorithm: From population to
  sample-based analysis. *Annals of Statistics*, 45:77–120,
  2017.
- Broyden, C. G. A class of methods for solving nonlinear
  simultaneous equations. *Mathematics of computation*, 19
  (92):577–593, 1965.
- 454
  455
  456
  456
  457
  457
  Broyden, C. G. The convergence of single-rank quasi-Newton methods. *Mathematics of Computation*, 24(110): 365–382, 1970.
- Broyden, C. G., Jr., J. E. D., Broyden, and More, J. J. On the local and superlinear convergence of quasi-Newton methods. *IMA J. Appl. Math*, 12(3):223–245, June 1973.
- 462 Candes, E. J., Eldar, Y., Strohmer, T., and Voroninski, V.
  463 Phase retrieval via matrix completion, 2011.
- 465 Carroll, R. J., Fan, J., Gijbels, I., and Wand, M. P. General466 ized partially linear single-index models. *Journal of the*467 *American Statistical Association*, 92:477–489, 1997.
- 468
  469
  469
  470
  470
  471
  472
  472
  472
  474
  475
  474
  475
  475
  475
  476
  477
  477
  477
  478
  479
  479
  479
  479
  479
  479
  479
  470
  470
  471
  472
  472
  472
  472
  473
  474
  475
  474
  475
  475
  475
  475
  475
  475
  475
  475
  475
  475
  475
  475
  475
  475
  475
  475
  475
  475
  475
  475
  475
  475
  475
  475
  475
  475
  475
  475
  475
  475
  475
  475
  475
  475
  475
  475
  475
  475
  475
  475
  475
  475
  475
  475
  475
  475
  475
  475
  475
  475
  475
  475
  475
  475
  475
  475
  475
  475
  475
  475
  475
  475
  475
  475
  475
  475
  475
  475
  475
  475
  475
  475
  475
  475
  475
  475
  475
  475
  475
  475
  475
  475
  475
  475
  475
  475
  475
  475
  475
  475
  475
  475
  475
  475
  475
  475
  475
  475
  475
  475
  475
  475
  475
  475
  475
  475
  475
  475
  475
  475
  475
  475
  475
  475
  475
  475
  475
  475
  475
  - Daskalakis, C., Tzamos, C., and Zampetakis, M. Ten steps of EM suffice for mixtures of two Gaussians. In *Proceedings of the 2017 Conference on Learning Theory*, 2017.
- 478
  479
  479
  480
  481
  481
  481
  482
  481
  483
  484
  484
  484
  484
  484
  484
  484
  484
  484
  484
  484
  484
  484
  484
  484
  484
  484
  484
  484
  484
  484
  484
  484
  484
  484
  484
  484
  484
  484
  484
  484
  484
  484
  484
  484
  484
  484
  484
  484
  484
  484
  484
  484
  484
  484
  484
  484
  484
  484
  484
  484
  484
  484
  484
  484
  484
  484
  484
  484
  484
  484
  484
  484
  484
  484
  484
  484
  484
  484
  484
  484
  484
  484
  484
  484
  484
  484
  484
  484
  484
  484
  484
  484
  484
  484
  484
  484
  484
  484
  484
  484
  484
  484
  484
  484
  484
  484
  484
  484
  484
  484
  484
  484
  484
  484
  484
  484
  484
  484
  484
  484
  484
  484
  484
  484
  484
  484
  484
  484
  484
  484
  484
  484
  484
  484
  484
  484
  484
  484
  484
  484
  484
  484
  484
  484
  484
  484
  484
  484
  484
  484
  484
  484
  484
  484
  484
  484
- 482 Dwivedi, R., Ho, N., Khamaru, K., Wainwright, M. J., Jor483 dan, M. I., and Yu, B. Sharp analysis of expectation484 maximization for weakly identifiable models. *AISTATS*,
  485 2020a.
- Dwivedi, R., Ho, N., Khamaru, K., Wainwright, M. J., Jordan, M. I., and Yu, B. Singularity, misspecification, and the convergence rate of EM. *Annals of Statistics*, 44: 2726–2755, 2020b.
- Feiyan Tian, Lei Liu, X. C. Generalized memory approximate message passing. *https://arxiv.org/abs/2110.06069*, 2021.

- Fienup, J. R. Phase retrieval algorithms: a comparison. Appl. Opt., 21(15):2758-2769, Aug 1982. doi: 10.1364/AO.21.002758. URL http://www.osapublishing.org/ao/ abstract.cfm?URI=ao-21-15-2758.
- Fletcher, R. A new approach to variable metric algorithms. *The computer journal*, 13(3):317–322, 1970.
- Fletcher, R. and Powell, M. J. A rapidly convergent descent method for minimization. *The computer journal*, 6(2): 163–168, 1963.
- Gay, D. M. Some convergence properties of Broyden's method. *SIAM Journal on Numerical Analysis*, 16(4): 623–630, 1979.
- Goebel, K. and Kirk, W. A. *Topics in Metric Fixed Point Theory*. Cambridge University Press, 1990.
- Goldfarb, D. A family of variable-metric methods derived by variational means. *Mathematics of computation*, 24 (109):23–26, 1970.
- Hardt, M., Recht, B., and Singer, Y. Train faster, generalize better: Stability of stochastic gradient descent. In Balcan, M. F. and Weinberger, K. Q. (eds.), *Proceedings of The 33rd International Conference on Machine Learning*, volume 48 of *Proceedings of Machine Learning Research*, pp. 1225–1234, New York, New York, USA, 20– 22 Jun 2016. PMLR. URL http://proceedings. mlr.press/v48/hardt16.html.
- Ho, N., Khamaru, K., Dwivedi, R., Wainwright, M. J., Jordan, M. I., and Yu, B. Instability, computational efficiency and statistical accuracy. *Arxiv Preprint Arxiv:* 2005.11411, 2020.
- Jin, Q. and Mokhtari, A. Non-asymptotic superlinear convergence of standard quasi-newton methods. *arXiv preprint arXiv:2003.13607*, 2020.
- Jin, Q., Koppel, A., Rajawat, K., and Mokhtari, A. Sharpened quasi-newton methods: Faster superlinear rate and larger local convergence neighborhood. *The 39th International Conference on Machine Learning (ICML 2022)*, 2022.
- Kwon, J., Qian, W., Caramanis, C., Chen, Y., and Damek, D. Global convergence of the EM algorithm for mixtures of two component linear regression. In *Conference on Learning Theory (COLT)*, 2019.
- Kwon, J. Y., Ho, N., and Caramanis, C. On the minimax optimality of the EM algorithm for learning two-component mixed linear regression. In *AISTATS*, 2021.

- Lin, D., Ye, H., and Zhang, Z. Explicit superlinear convergence of broyden's method in nonlinear equations. *arXiv* preprint arXiv:2109.01974, 2021a.
  - Lin, D., Ye, H., and Zhang, Z. Greedy and random quasinewton methods with faster explicit superlinear convergence. Advances in Neural Information Processing Systems 34, 2021b.
  - Liu, D. C. and Nocedal, J. On the limited memory BFGS method for large scale optimization. *Mathematical programming*, 45(1-3):503–528, 1989.
  - Loh, P.-L. and Wainwright, M. J. Regularized M-estimators
     with nonconvexity: Statistical and algorithmic theory for
     local optima. *Journal of Machine Learning Research*, 16:
     559–616, 2015.
  - Mou, W., Ho, N., Wainwright, M. J., Bartlett, P., and
     Jordan, M. I. A diffusion process perspective on posterior contraction rates for parameters. *arXiv preprint* arXiv:1909.00966, 2019.
  - Netrapalli, P., Jain, P., and Sanghavi, S. Phase retrieval using alternating minimization. *IEEE Transactions on Signal Processing*, 63(18):4814–4826, 2015. doi: 10.1109/TSP. 2015.2448516.
  - Nocedal, J. Updating quasi-Newton matrices with limited storage. *Mathematics of computation*, 35(151):773–782, 1980.
  - Ren, T., Cui, F., Atsidakou, A., Sanghavi, S., and Ho, N.
     Towards statistical and computational complexities of Polyak step size gradient descent. *Artificial Intelligence* and Statistics Conference, 2022a.
  - Ren, T., Zhuo, J., Sanghavi, S., and Ho, N. Improving computational complexity in statistical models with secondorder information. *arXiv preprint arXiv:2202.04219*, 2022b.
  - Rodomanov, A. and Nesterov, Y. Greedy quasi-newton methods with explicit superlinear convergence. *SIAM Journal on Optimization*, 31(1):785–811, 2021a.
  - Rodomanov, A. and Nesterov, Y. Rates of superlinear convergence for classical quasi-newton methods. *Mathematical Programming*, pp. 1–32, 2021b.
  - Rodomanov, A. and Nesterov, Y. New results on superlinear convergence of classical quasi-newton methods. *Journal of Optimization Theory and Applications*, 188(3):744–769, 2021c.
  - Shanno, D. F. Conditioning of quasi-Newton methods for function minimization. *Mathematics of computation*, 24 (111):647–656, 1970.

- Shechtman, Y., Eldar, Y. C., Cohen, O., Chapman, H. N., Miao, J., and Segev, M. Phase retrieval with application to optical imaging: A contemporary overview. *IEEE Signal Processing Magazine*, 32(3):87–109, 2015. doi: 10.1109/MSP.2014.2352673.
- Wang, Z., Gu, Q., Ning, Y., and Liu, H. High-dimensional expectation-maximization algorithm: Statistical optimization and asymptotic normality. In Advances in Neural Information Processing Systems 28, 2015.
- Xu, J., Hsu, D., and Maleki, A. Global analysis of expectation maximization for mixtures of two Gaussians. In *Advances in Neural Information Processing Systems 29*, 2016.
- Ye, H., Lin, D., Zhang, Z., and Chang, X. Explicit superlinear convergence rates of the sr1 algorithm. *arXiv* preprintarXiv:2105.07162, 2021.
- Yi, X. and Caramanis, C. Regularized EM algorithms: A unified framework and statistical guarantees. In *Advances in Neural Information Processing Systems* 28, 2015.
- Yuan, X.-T. and Zhang, T. Truncated power method for sparse eigenvalue problems. *Journal of Machine Learning Research*, 14(Apr):899–925, 2013.

## A. Proofs

**Lemma A.1.** Consider the objective function in (10) satisfying Assumption 4.1 and 4.2. Then, the inverse matrix of its Hessian  $\nabla^2 f(\theta)$  can be expressed as

$$\nabla^2 f(\theta)^{-1} = \frac{(A^\top A)^{-1}}{q \|A\theta - b\|^{q-2}} - \frac{(q-2)(\theta - \hat{\theta})(\theta - \hat{\theta})^\top}{q(q-1)\|A\theta - b\|^q}.$$
(18)

*Proof.* Notice that the Hessian of objective function (10) can be expressed as

$$\nabla^2 f(\theta) = q \|A\theta - b\|^{q-2} A^\top A + q(q-2) \|A\theta - b\|^{q-4} A^\top (A\theta - b) (A\theta - b)^\top A.$$
(19)

We use the Sherman–Morrison formula. Suppose that  $X \in \mathbb{R}^{d \times d}$  is an invertible matrix and  $a, b \in \mathbb{R}^d$  are two vectors satisfying that  $1 + b^\top X^{-1}a \neq 0$ . Then, the matrix  $X + ab^\top$  is invertible and

$$(X+ab^{\top})^{-1} = X^{-1} - \frac{X^{-1}ab^{\top}X^{-1}}{1+b^{\top}X^{-1}a}.$$
(20)

Applying the Sherman–Morrison formula with  $X = q ||A\theta - b||^{q-2}A^{\top}A$ ,  $a = q(q-2)||A\theta - b||^{q-4}A^{\top}(A\theta - b)$  and  $b = A^{\top}(A\theta - b)$ . Notice that  $A^{\top}A$  is invertible, hence X is invertible and

$$1 + b^{\top} X^{-1} a$$

$$= 1 + (A\theta - b)^{\top} A \frac{(A^{\top} A)^{-1}}{q \| A\theta - b \|^{q-2}} q(q-2) \| A\theta - b \|^{q-4} A^{\top} (A\theta - b)$$

$$= 1 + (q-2) (A\theta - b)^{\top} A \frac{(A^{\top} A)^{-1} A^{\top} A (\theta - \hat{\theta})}{\| A\theta - b \|^{2}}$$

$$= 1 + (q-2) \frac{(A\theta - b)^{\top} (A\theta - b)}{\| A\theta - b \|^{2}}$$

$$= q - 1 \neq 0. \qquad (\text{Since } q \ge 4.)$$
(21)

Therefore, we obtain that

 $\nabla^2 f(\theta)^{-1}$ 

$$= \frac{(A^{\top}A)^{-1}}{q\|A\theta - b\|^{q-2}} - \frac{\frac{(A^{\top}A)^{-1}}{q\|A\theta - b\|^{q-2}}q(q-2)\|A\theta - b\|^{q-4}A^{\top}(A\theta - b)(A^{\top}(A\theta - b))^{\top}\frac{(A^{\top}A)^{-1}}{q\|A\theta - b\|^{q-2}}}{q-1}$$

$$= \frac{(A^{\top}A)^{-1}}{q\|A\theta - b\|^{q-2}} - \frac{(q-2)}{q(q-1)\|A\theta - b\|^{q}}(A^{\top}A)^{-1}AA^{\top}(\theta - \hat{\theta})(\theta - \hat{\theta})^{\top}AA^{\top}(A^{\top}A)^{-1}$$

$$= \frac{(A^{\top}A)^{-1}}{q\|A\theta - b\|^{q-2}} - \frac{(q-2)(\theta - \hat{\theta})(\theta - \hat{\theta})^{\top}}{q(q-1)\|A\theta - b\|^{q}}.$$
(22)

591 As a consequence, we obtain the conclusion of the lemma.

Lemma A.2. Banach's Fixed-Point Theorem. Consider the differentiable function  $f: D \subset \mathbb{R} \to D \subset \mathbb{R}$ . Suppose that there exists  $C \in (0,1)$  such that  $|f'(x)| \leq C$  for any  $x \in D$ . Now let  $x_0 \in D$  be arbitrary and define the sequence  $\{x_k\}_{k=0}^{\infty}$  as

$$x_{k+1} = f(x_k), \qquad \forall k \ge 0. \tag{23}$$

597 Then, the sequence  $\{x_k\}_{k=0}^{\infty}$  converges to the unique fixed point  $x_*$  defined as

$$x_* = f(x_*), \tag{24}$$

with linear convergence rate of

$$|x_k - x_*| \le C^k |x_0 - x_*|, \qquad \forall k \ge 0.$$
(25)

Proof. Check (Goebel & Kirk, 1990).

## 605 A.1. Proof of Theorem 4.3

 We use induction to prove the convergence results in (11) and (12). Note that  $b = A\hat{\theta}$  by Assumption 4.1 and the gradient and Hessian of the objective function in (10) are explicitly given by

$$\nabla f(\theta) = q \|A\theta - b\|^{q-2} A^{\top} (A\theta - b) = q \|A\theta - b\|^{q-2} A^{\top} A(\theta - \hat{\theta}),$$
(26)

$$\nabla^2 f(\theta) = q \|A\theta - b\|^{q-2} A^\top A + q(q-2) \|A\theta - b\|^{q-4} A^\top (A\theta - b) (A\theta - b)^\top A.$$
(27)

Applying Lemma A.1, we can obtain that

$$\nabla^2 f(\theta)^{-1} = \frac{(A^\top A)^{-1}}{q \|A\theta - b\|^{q-2}} - \frac{(q-2)(\theta - \hat{\theta})(\theta - \hat{\theta})^\top}{q(q-1)\|A\theta - b\|^q}.$$
(28)

First, we consider the initial iteration

$$\theta_1 = \theta_0 - H_0 \nabla f(\theta_0) = \theta_0 - \nabla f(\theta_0)^{-1} \nabla f(\theta_0), \tag{29}$$

$$\theta_1 - \hat{\theta} = \theta_0 - \hat{\theta} - \nabla f(\theta_0)^{-1} \nabla f(\theta_0).$$
(30)

623 Notice that  $b = A\hat{\theta}$  by Assumption 4.1 and

$$\nabla^{2} f(\theta_{0})^{-1} \nabla f(\theta_{0}) = \left[ \frac{(A^{\top} A)^{-1}}{q \| A \theta_{0} - b \|^{q-2}} - \frac{(q-2)(\theta_{0} - \hat{\theta})(\theta_{0} - \hat{\theta})^{\top}}{q(q-1) \| A \theta_{0} - b \|^{q}} \right] q \| A \theta - b \|^{q-2} A^{\top} A(\theta_{0} - \hat{\theta}) = \theta_{0} - \hat{\theta} - \frac{q-2}{q-1} \frac{(\theta_{0} - \hat{\theta})^{\top} A^{\top} A(\theta_{0} - \hat{\theta})}{\| A \theta_{0} - b \|^{2}} (\theta_{0} - \hat{\theta}) = \theta_{0} - \hat{\theta} - \frac{q-2}{q-1} \frac{(A \theta_{0} - b)^{\top} (A \theta_{0} - b)}{\| A \theta_{0} - b \|^{2}} (\theta_{0} - \hat{\theta}) = \theta_{0} - \hat{\theta} - \frac{q-2}{q-1} \frac{(\theta_{0} - \theta)^{\top} (A \theta_{0} - b)}{\| A \theta_{0} - b \|^{2}} (\theta_{0} - \hat{\theta}) = \theta_{0} - \hat{\theta} - \frac{q-2}{q-1} (\theta_{0} - \hat{\theta}).$$
(31)

Therefore, we obtain that

$$\theta_1 - \hat{\theta} = \theta_0 - \hat{\theta} - \nabla f(\theta_0)^{-1} \nabla f(\theta_0) = \frac{q-2}{q-1} (\theta_0 - \hat{\theta}).$$
(32)

Condition (11) holds for k = 1 with  $r_0 = \frac{q-2}{q-1}$ . Now we assume that condition (11) holds for k = t where  $t \ge 1$ , i.e.,

$$\theta_t - \hat{\theta} = r_{t-1}(\theta_{t-1} - \hat{\theta}). \tag{33}$$

Considering the condition  $b = A\hat{\theta}$  in Assumption 4.1 and the condition in (33), we further have

$$A\theta_t - b = A(\theta_t - \hat{\theta}) = r_{t-1}A(\theta_{t-1} - \hat{\theta}) = r_{t-1}(A\theta_{t-1} - b),$$
(34)

which implies that

$$\nabla f(\theta_t) = q r_{t-1}^{q-1} \| A(\theta_{t-1} - \hat{\theta}) \|^{q-2} A^\top A(\theta_{t-1} - \hat{\theta}).$$
(35)

We further show that the variable displacement and gradient difference can be written as

$$s_{t-1} = \theta_t - \theta_{t-1} = \theta_t - \hat{\theta} - \theta_{t-1} + \hat{\theta} = (r_{t-1} - 1)(\theta_{t-1} - \hat{\theta}),$$
(36)

654 and

$$u_{t-1} = \nabla f(\theta_t) - \nabla f(\theta_{t-1}) = q(r_{t-1}^{q-1} - 1) \|A(\theta_{t-1} - \hat{\theta})\|^{q-2} A^{\top} A(\theta_{t-1} - \hat{\theta}).$$
(37)

Considering these expressions, we can show that the rank-1 matrix in the update of BFGS  $u_{t-1}s_{t-1}^{\top}$  is given by

$$u_{t-1}s_{t-1}^{\top} = q(r_{t-1}^{q-1} - 1)(r_{t-1} - 1) \|A(\theta_{t-1} - \hat{\theta})\|^{q-2} A^{\top} A(\theta_{t-1} - \hat{\theta})(\theta_{t-1} - \hat{\theta})^{\top},$$
(38)

and the inner product  $s_{t-1}^{\top} u_{t-1}$  can be written as

$$s_{t-1}^{\top} u_{t-1} = q(r_{t-1}^{q-1} - 1)(r_{t-1} - 1) \|A(\theta_{t-1} - \hat{\theta})\|^{q-2} (\theta_{t-1} - \hat{\theta})^{\top} A^{\top} A(\theta_{t-1} - \hat{\theta}) = q(r_{t-1}^{q-1} - 1)(r_{t-1} - 1) \|A(\theta_{t-1} - \hat{\theta})\|^{q}.$$
(39)

These two expressions allow us to simplify the matrix  $I - \frac{u_{t-1}s_{t-1}^{\top}}{s_{t-1}^{\top}u_{t-1}}$  in the update of BFGS as

$$I - \frac{u_{t-1}s_{t-1}^{\top}}{s_{t-1}^{\top}u_{t-1}} = I - \frac{A^{\top}A(\theta_{t-1} - \hat{\theta})(\theta_{t-1} - \hat{\theta})^{\top}}{\|A(\theta_{t-1} - \hat{\theta})\|^2}.$$
(40)

An important property of the above matrix is that its null space is the set of the vectors that are parallel to  $u_{t-1}$ . Considering the expression for  $u_{t-1}$ , any vector that is parallel to  $A^{\top}A(\theta_{t-1} - \hat{\theta})$  is in the null space of the above matrix. We can easily observe that the gradient defined in (35) satisfies this condition and therefore

- $\begin{pmatrix} I \frac{u_{t-1}s_{t-1}^{\top}}{s_{t-1}^{\top}u_{t-1}} \end{pmatrix} \nabla f(\theta_{t})$  $= qr_{t-1}^{q-1} \|A(\theta_{t-1} - \hat{\theta})\|^{q-2} \left( I - \frac{A^{\top}A(\theta_{t-1} - \hat{\theta})(\theta_{t-1} - \hat{\theta})^{\top}}{\|A(\theta_{t-1} - \hat{\theta})\|^{2}} \right) A^{\top}A(\theta_{t-1} - \hat{\theta})$  $= qr_{t-1}^{q-1} \|A(\theta_{t-1} - \hat{\theta})\|^{q-2} \left( A^{\top}A(\theta_{t-1} - \hat{\theta}) - \frac{A^{\top}A(\theta_{t-1} - \hat{\theta})\|A(\theta_{t-1} - \hat{\theta})\|^{2}}{\|A(\theta_{t-1} - \hat{\theta})\|^{2}} \right)$ = 0.(41)
- This important observation shows that if the condition in (33) holds, then the BFGS descent direction  $H_t \nabla f(\theta_t)$  can be simplified as
- $H_t \nabla f(\theta_t)$  $= \left(I - \frac{s_{t-1}u_{t-1}^{\top}}{s_{t-1}^{\top}u_{t-1}}\right)H_{t-1}\left(I - \frac{u_{t-1}s_{t-1}^{\top}}{s_{t-1}^{\top}u_{t-1}}\right)\nabla f(\theta_{t}) + \frac{s_{t-1}s_{t-1}^{\top}}{s_{t-1}^{\top}u_{t-1}}\nabla f(\theta_{t})$  $= \frac{s_{t-1}s_{t-1}}{s_{t-1}^{\top}u_{t-1}}\nabla f(\theta_t)$  $= \frac{(r_{t-1}-1)^2(\theta_{t-1}-\hat{\theta})(\theta_{t-1}-\hat{\theta})^\top}{q(r_{t-1}^{q-1}-1)(r_{t-1}-1)\|A(\theta_{t-1}-\hat{\theta})\|^q} qr_{t-1}^{q-1} \|A(\theta_{t-1}-\hat{\theta})\|^{q-2} A^\top A(\theta_{t-1}-\hat{\theta})$ (42) $= \frac{1 - r_{t-1}}{1 - r_{t-1}^{q-1}} r_{t-1}^{q-1}(\theta_{t-1} - \hat{\theta}) \frac{\|A(\theta_{t-1} - \hat{\theta})\|^{q-2}(\theta_{t-1} - \hat{\theta})^{\top} A^{\top} A(\theta_{t-1} - \hat{\theta})}{\|A(\theta_{t-1} - \hat{\theta})\|^{q}}$  $= \frac{1 - r_{t-1}}{1 - r_{t-1}^{q-1}} r_{t-1}^{q-1} (\theta_{t-1} - \hat{\theta}).$

This simplification implies that for the new iterate  $\theta_{t+1}$ , we have

$$\theta_{t+1} - \hat{\theta} = \theta_t - H_t \nabla f(\theta_t) - \hat{\theta} = \theta_t - \hat{\theta} - \frac{1 - r_{t-1}}{1 - r_{t-1}^{q-1}} r_{t-1}^{q-1} \frac{(\theta_t - \theta)}{r_{t-1}} = \frac{1 - r_{t-1}^{q-2}}{1 - r_{t-1}^{q-1}} (\theta_t - \hat{\theta}) = r_t (\theta_t - \hat{\theta}),$$
(43)

where

$$r_t = \frac{1 - r_{t-1}^{q-2}}{1 - r_{t-1}^{q-1}}.$$
(44)

Therefore, we prove that condition (11) holds for k = t + 1. By induction, we prove the linear convergence results in (11) and (12).

## A.2. Proof of Theorem 4.5

Recall that we have

$$r_0 = \frac{q-2}{q-1}, \qquad r_k = \frac{1 - r_{k-1}^{q-2}}{1 - r_{k-1}^{q-1}}, \qquad \forall k \ge 1.$$
(45)

Consider that  $q \ge 4$  and define the function g(r) as

$$g(r) := \frac{1 - r^{q-2}}{1 - r^{q-1}}, \qquad r \in [0, 1].$$
(46)

Suppose that  $r_* \in (0, 1)$  satisfying that  $r_* = g(r_*)$ , which is equivalent to

$$r_*^{q-1} + r_*^{q-2} = 1. (47)$$

Notice that

$$g'(r) = \frac{(q-1)r^{q-2} - r^{2q-4} - (q-2)r^{q-3}}{(1-r^{q-1})^2},$$
(48)

$$(q-1)r^{q-2} - r^{2q-4} - (q-2)r^{q-3}$$

$$= r^{q-3}[(q-1)(r-1) - (r^{q-1} - 1)]$$

$$= r^{q-3}(r-1)(q-1 - \frac{r^{q-1} - 1}{r-1})$$

$$= r^{q-3}(r-1)(q-1 - \sum_{i=0}^{q-2} r^{i}).$$
(49)

Since  $r \in [0, 1]$ , we have that

$$r^{q-3} \ge 0, \quad r-1 \le 0, \quad \sum_{i=0}^{q-2} r^i \le \sum_{i=0}^{q-2} 1 = q-1.$$
 (50)

Therefore, we obtain that

$$(q-1)r^{q-2} - r^{2q-4} - (q-2)r^{q-3} \le 0, (51)$$

53 and

$$|g'(r)| = \frac{r^{2q-4} + (q-2)r^{q-3} - (q-1)r^{q-2}}{(1-r^{q-1})^2}.$$
(52)

Our target is to prove that for any  $r \in [0, 1]$ ,

$$|g'(r)| \le \frac{1}{2}.\tag{53}$$

First, we present the plots of |g'(r)| for  $r \in [0,1]$  with  $4 \le q \le 11$  in Figure 3. We observe that for  $4 \le q \le 11$ ,  $|g'(r)| \le 1/2$  always holds.

Next, we prove that for  $q \ge 12$  and any  $r \in [0, 1]$ , we have

$$|g'(r)| = \frac{(q-1)r^{q-2} - r^{2q-4} - (q-2)r^{q-3}}{(1-r^{q-1})^2} \le \frac{1}{2},$$
(54)

which is equivalent to

$$r^{2q-2} - 2r^{2q-4} - 2r^{q-1} + 2(q-1)r^{q-2} - 2(q-2)r^{q-3} + 1 \ge 0, \qquad \forall r \in [0,1].$$
(55)



We observe that for  $q \ge 14$ , we have l(q) > 0 and we calculate that l(12) = 84 > 0 and l(13) = 236 > 0. Hence, we obtain that l(q) > 0 for all  $q \ge 12$ , which indicates that for all  $r \in [0, 1]$ ,

$$r^{2} \leq 1 < \frac{2(q-2)^{2}(q-3)}{q(q+1)(q-1)},$$

(65)

 $q(q+1)(q-1)r^2 - 2(q-2)^2(q-3) < 0.$ (66)

Therefore, for all  $r \in [0, 1]$ ,  $h^{(4)}(r)$  defined in (63) satisfies that  $h^{(4)}(r) < 0$  and from (62) we know that  $\frac{dh^{(3)}(r)}{dr} < 0$ . Hence,  $h^{(3)}(r)$  defined in (61) is decreasing function and  $h^{(3)}(r) <= h^{(3)}(0) = -2 < 0$ . We know that  $\frac{dh^{(2)}(r)}{dr} = h^{(3)}(r) < 0$ , which implies that  $h^{(2)}(r)$  defined in (60) is decreasing function. So we have that  $h^{(2)}(r) \ge h^{(2)}(1) = 1 > 0$ . From (59) we know that  $\frac{dh^{(1)}(r)}{dr} > 0$  and  $h^{(1)}(r)$  defined in (58) is increasing function for  $r \in [0, 1]$ . Hence, we get that  $h^{(1)}(r) \le h^{(1)}(1) = 0$  and from (57) we obtain that h(r) defined in(56) is decreasing function for  $r \in [0, 1]$ . Therefore, we have that  $h(r) \ge h(1) = 0$  and condition (55) holds for all  $r \in [0, 1]$ , which is equivalent to  $|g'(r)| \le 1/2$ .

In summary, we proved that for any  $q \ge 12$ , we have  $|g'(r)| \le 1/2$ . Combining this with the results from Figure 3, we obtain that  $|g'(r)| \le 1/2$  holds for all  $q \ge 4$ . Applying Banach's Fixed-Point Theorem from Lemma A.2, we prove the final conclusion (14).

#### A.3. Proof of Theorem 4.6

Notice that the gradient and the Hessian of the objective function (10) can be expressed as

$$\nabla f(\theta) = q \|A\theta - b\|^{q-2} A^{\top} (A\theta - b) = q \|A\theta - b\|^{q-2} A^{\top} A(\theta - \hat{\theta}), \tag{67}$$

$$\nabla^2 f(\theta) = q \|A\theta - b\|^{q-2} A^\top A + q(q-2) \|A\theta - b\|^{q-4} A^\top (A\theta - b) (A\theta - b)^\top A.$$
(68)

Applying Lemma A.1, we can obtain that

$$\nabla^2 f(\theta)^{-1} = \frac{(A^\top A)^{-1}}{q \| A\theta - b \|^{q-2}} - \frac{(q-2)(\theta - \hat{\theta})(\theta - \hat{\theta})^\top}{q(q-1) \| A\theta - b \|^q}.$$
(69)

Hence, we have that for any  $k \ge 1$ ,

$$\theta_k = \theta_{k-1} - \nabla f(\theta_{k-1})^{-1} \nabla f(\theta_{k-1}), \tag{70}$$

$$\theta_k - \hat{\theta} = \theta_{k-1} - \hat{\theta} - \nabla f(\theta_{k-1})^{-1} \nabla f(\theta_{k-1}).$$
(71)

Notice that  $b = A\hat{\theta}$  by Assumption 4.1 and

$$\nabla f(\theta_{k-1})^{-1} \nabla f(\theta_{k-1}) = \left[\frac{(A^{\top}A)^{-1}}{q\|A\theta_{k-1} - b\|^{q-2}} - \frac{(q-2)(\theta_{k-1} - \hat{\theta})(\theta_{k-1} - \hat{\theta})^{\top}}{q(q-1)\|A\theta_{k-1} - b\|^{q}}\right] q\|A\theta - b\|^{q-2}A^{\top}A(\theta_{k-1} - \hat{\theta}) = \theta_{k-1} - \hat{\theta} - \frac{q-2}{q-1} \frac{(\theta_{0} - \hat{\theta})^{\top}A^{\top}A(\theta_{k-1} - \hat{\theta})}{\|A\theta_{k-1} - b\|^{2}}(\theta_{k-1} - \hat{\theta}) = \theta_{k-1} - \hat{\theta} - \frac{q-2}{q-1} \frac{(A\theta_{k-1} - b)^{\top}(A\theta_{k-1} - b)}{\|A\theta_{k-1} - b\|^{2}}(\theta_{k-1} - \hat{\theta}) = \theta_{k-1} - \hat{\theta} - \frac{q-2}{q-1} \frac{(A\theta_{k-1} - b)^{\top}(A\theta_{k-1} - b)}{\|A\theta_{k-1} - b\|^{2}}(\theta_{k-1} - \hat{\theta}) = \theta_{k-1} - \hat{\theta} - \frac{q-2}{q-1}(\theta_{k-1} - \hat{\theta}).$$
(72)

870 Therefore, we prove the conclusion that for any  $k \ge 1$ ,

$$\theta_k - \hat{\theta} = \theta_{k-1} - \hat{\theta} - \nabla f(\theta_{k-1})^{-1} \nabla f(\theta_{k-1}) = \frac{q-2}{q-1} (\theta_{k-1} - \hat{\theta}).$$
(73)

We observe that the iterations generated by Newton's method also satisfy the parallel property, i.e., all vectors  $\{\theta_k - \hat{\theta}\}_{k=0}^{\infty}$ are parallel to each other with the same direction.

Notice that the function  $h(r) = r^{q-1} + r^{q-2}$  is strictly increasing and  $h(\frac{q-2}{q-1}) < 1$ ,  $h(r_*) = 1$  as well as  $h(\frac{2q-3}{2q-2}) > 1$ . Hence, we know that  $\frac{q-2}{q-1} < r_* < \frac{2q-3}{2q-2}$ .

#### A.4. Proof of Theorem 5.1

In the following proof, we use  $\mathcal{L}$  and  $\mathcal{L}_n$  to denote the population objective and empirical objective, and  $\theta_t$  and  $\theta_t^n$ ,  $H_t^n$ and  $H_t$  to denote the corresponding iterations and Hessian approximating matrices in the population update and empirical update accordingly. For the ease of presentation, we use  $C_p$  to denote any constant that is independent of d, n, and  $C_p$ can be varied case by case to simply the proof. From Mou et al. (2019) and Ren et al. (2022b), we have that, as long as  $n = \Omega((d \log d/\delta)^{2p})$ , we have the following two uniform concentration results:

$$\sup_{\theta \in \mathbb{B}(\theta^*, r)} \|\nabla \mathcal{L}_n(\theta)) - \nabla \mathcal{L}(\theta))\| \leq C_p (\|\theta^*\| + r)^{p-1} \sqrt{d \log(1/\delta)/n},$$

$$\sup_{\theta \in \mathbb{B}(\theta^*, r)} \|\nabla^2 \mathcal{L}_n(\theta)) - \nabla^2 \mathcal{L}(\theta))\| \leq C_p (\|\theta^*\| + r)^{p-2} \sqrt{d \log(1/\delta)/n}.$$
(74)

Notice that we have

$$\mathcal{L}(\theta) = \mathbb{E}[(Y - (X^{\top}\theta)^p)^2] = \mathbb{E}[((X^{\top}\theta^*)^p + \zeta - (X^{\top}\theta)^p)^2]$$
(75)

$$= \mathbb{E}[((X^{\top}\theta^{*})^{p} - (X^{\top}\theta)^{p})^{2}] + \mathbb{E}[2\zeta((X^{\top}\theta^{*})^{p} - (X^{\top}\theta)^{p})^{2}] + \mathbb{E}[\zeta^{2}]$$
(76)

$$= \mathbb{E}[((X^{\top}\theta^{*})^{p} - (X^{\top}\theta)^{p})^{2}] + \sigma^{2},$$
(77)

where we use the fact that  $\zeta$  is independent of X and  $\mathbb{E}[\zeta] = 0$ ,  $\mathbb{E}[\zeta^2] = \sigma^2$ . Hence, we have that 

$$\nabla \mathcal{L}(\theta) = 2p\mathbb{E}[((X^{\top}\theta^{*})^{p} - (X^{\top}\theta)^{p})(X^{\top}\theta)^{p-1}X].$$
(78)

$$\nabla^{2} \mathcal{L}(\theta) = -2p^{2} \mathbb{E}[(X^{\top} \theta)^{2p-2} X X^{\top}] + 2p(p-1) \mathbb{E}[((X^{\top} \theta^{*})^{p} - (X^{\top} \theta)^{p})(X^{\top} \theta)^{p-2}) X X^{\top}].$$
(79)

We denote  $\mathcal{L}^0$  as the population loss function with respect to the assumption of  $\theta^* = 0$ . Therefore, we have that

$$\mathcal{L}^{0}(\theta) = \mathbb{E}[(X^{\top}\theta)^{p})^{2}] + \sigma^{2}.$$
(80)

$$\nabla \mathcal{L}^{0}(\theta) = -2p\mathbb{E}[(X^{\top}\theta)^{p})(X^{\top}\theta)^{p-1}X].$$
(81)

$$\nabla^2 \mathcal{L}^0(\theta) = -2p^2 \mathbb{E}[(X^\top \theta)^{2p-2} X X^\top] - 2p(p-1)\mathbb{E}[(X^\top \theta)^p (X^\top \theta)^{p-2}) X X^\top].$$
(82)

Hence, we have

$$\|\nabla \mathcal{L}(\theta) - \nabla \mathcal{L}^{0}(\theta)\| = C_{p}\mathbb{E}[(X^{\top}\theta^{*})^{p}(X^{\top}\theta)^{p-1}X] \le C_{p}\mathbb{E}[\|X\|^{2p}]\|\theta^{*}\|^{p}\|\theta\|^{p-1}.$$
(83)

Recall that X is a Gaussian or sub-Gaussian random variable with  $\mathbb{E}[||X||^{2p}] < +\infty$ . For any  $\theta \in \mathbb{B}(\theta^*, r)$ ,  $||\theta|| \le ||\theta^*|| + r$ and in the low SNR regime, we have  $\|\theta^*\| \leq C_1(\frac{d}{n})^{\frac{1}{2p}}$ . Hence, we have 

$$\sup_{\theta \in \mathbb{B}(\theta^*, r)} \|\nabla \mathcal{L}(\theta) - \nabla \mathcal{L}^0(\theta)\| \le C_p (\|\theta^*\| + r)^{p-1} \sqrt{d \log(1/\delta)/n}.$$
(84)

Similarly, we have that

$$\|\nabla^{2}\mathcal{L}(\theta) - \nabla^{2}\mathcal{L}^{0}(\theta)\| = C_{p}\mathbb{E}[(X^{\top}\theta^{*})^{p}(X^{\top}\theta)^{p-2}XX^{\top}] \le C_{p}\mathbb{E}[\|X\|^{2p}]\|\theta^{*}\|^{p}\|\theta\|^{p-2}.$$
(85)

$$\sup_{\theta \in \mathbb{B}(\theta^*, r)} \|\nabla^2 \mathcal{L}(\theta) - \nabla^2 \mathcal{L}^0(\theta)\| \le C_p (\|\theta^*\| + r)^{p-2} \sqrt{d \log(1/\delta)/n}.$$
(86)

Leveraging (74), (84) and (86) we obtain that 

ŧ

$$\sup_{\theta \in \mathbb{B}(\theta^*, r)} \|\nabla \mathcal{L}_n(\theta)) - \nabla \mathcal{L}^0(\theta))\| \le C_p (\|\theta^*\| + r)^{p-1} \sqrt{d \log(1/\delta)/n},$$
(87)

932  
933  

$$\sup_{\theta \in \mathbb{B}(\theta^*, r)} \|\nabla^2 \mathcal{L}_n(\theta)) - \nabla^2 \mathcal{L}^0(\theta))\| \le C_p (\|\theta^*\| + r)^{p-2} \sqrt{d \log(1/\delta)/n}.$$
(88)

This explains the remark of assumption  $\theta^* = 0$  in section 4. The errors between gradients and Hessians of the population loss with  $\theta^* = 0$  and  $\|\theta^*\| \le C_1 (d/n)^{1/(2p)}$  are upper bounded by the corresponding statistical errors between the population loss (8) and the empirical loss (7) in the low SNR regime, respectively. Hence,  $\mathcal{L}^0$  and  $\mathcal{L}$  can be treated as equivalent. In the following proof, we replace all  $\mathcal{L}^0$  with  $\mathcal{L}$  in the general low SNR regime. We can assume that the results of Theorem 4.3 and 4.5 also hold for the  $\mathcal{L}$ , i.e., there exists linear convergence rates  $\{r_t\}_{t=0}^{\infty}$  with 

$$\|\theta_{t+1} - \theta^*\| \le r_t \|\theta_t - \theta^*\| \qquad \forall t \ge 0$$

Recall that  $\mathcal{L}^0(\theta) = \mathbb{E}_X \left[ (X^\top \theta)^{2p} \right] = C_p \|\Sigma^{1/2} \theta\|^{2p}$ . Since  $\mathcal{L}^0$  is equivalent to  $\mathcal{L}$ , we can also assume that

$$\lambda_{\max}(\nabla^2 \mathcal{L}(\theta))) \le C_p \|\theta - \theta^*\|^{2p-2}, \quad \lambda_{\min}(\nabla^2 \mathcal{L}(\theta))) \ge C_p \|\theta - \theta^*\|^{2p-2}, \quad \|\nabla \mathcal{L}(\theta)\| \le C_p \|\theta - \theta^*\|^{2p-1}.$$
(89)

We assume that  $\|\theta^*\| \leq \|\theta - \theta^*\|$  for any  $\theta$ . Otherwise, we have that  $\|\theta - \theta^*\| < \|\theta^*\| \leq C_1(\frac{d}{n})^{\frac{1}{2p}}$  by the definition of the low SNR regime, which indicates that  $\theta$  has already achieved the optimal statistical radius. Hence, for any  $\theta \in \mathbb{B}(\theta^*, r)$ , we have that  $\|\theta^*\| \leq \|\theta - \theta^*\| \leq r$  and thus, we have the following conditions 

$$\sup_{\theta \in \mathbb{B}(\theta^*, r)} \|\nabla \mathcal{L}_n(\theta)) - \nabla \mathcal{L}(\theta))\| \le C_p r^{p-1} \sqrt{d \log(1/\delta)/n},\tag{90}$$

$$\sup_{\theta \in \mathbb{B}(\theta^*, r)} \|\nabla^2 \mathcal{L}_n(\theta)) - \nabla^2 \mathcal{L}(\theta))\| \le C_p r^{p-2} \sqrt{d \log(1/\delta)/n}.$$
(91)

In this proof, we use induction and apply conditions (89), (90) and (91) to prove the final results and start with the base case. Assume  $\theta_0^n = \theta^0$ , and  $\|\theta_0 - \theta^*\|^p \ge C_p \varepsilon(n, \delta)$ . For the first step, we have that 

$$\|\theta_{1}^{n}-\theta_{1}\| = \|(\nabla^{2}\mathcal{L}_{n}(\theta_{0}))^{-1}\nabla\mathcal{L}_{n}(\theta_{0}) - (\nabla^{2}\mathcal{L}(\theta_{0}))^{-1}\nabla\mathcal{L}(\theta_{0})\| \\ \leq \left\|\left(\nabla^{2}\mathcal{L}_{n}(\theta_{0})\right)^{-1} - (\nabla^{2}\mathcal{L}(\theta_{0}))^{-1}\right\| \|\nabla\mathcal{L}_{n}(\theta_{0})\| + \left\|\left(\nabla^{2}\mathcal{L}(\theta_{0})\right)^{-1}\right\| \|\nabla\mathcal{L}_{n}(\theta_{0}) - \nabla\mathcal{L}(\theta_{0})\|$$

$$(92)$$

We observe that for invertible matrices A and B,

$$(A^{-1} - B^{-1}) = A^{-1}(B - A)B^{-1}.$$
(93)

We can bound

$$\begin{aligned} \|\nabla^{2}\mathcal{L}_{n}(\theta_{0})^{-1} - \nabla^{2}\mathcal{L}(\theta_{0})^{-1}\| \\ &\leq \|(\nabla^{2}\mathcal{L}_{n}(\theta_{0}))^{-1}\| \|\nabla^{2}\mathcal{L}_{n}(\theta_{0}) - \nabla^{2}\mathcal{L}(\theta_{0})\| \|(\nabla^{2}\mathcal{L}(\theta_{0}))^{-1}\| \\ &\leq C_{p}\|\theta_{0} - \theta^{*}\|^{2-3p}\sqrt{d\log(1/\delta)/n}, \end{aligned}$$
(94)

which leads to the bound

$$\|\theta_1^n - \theta_1\| \le C_p \|\theta_0 - \theta^*\|^{1-p} \sqrt{d\log(1/\delta)/n}.$$
(95)

We now show that,  $\forall t < T = O(\log(n/d))$ , we have

$$\|\theta_t^n - \theta_t\| \le c_t \|\theta_{t-1} - \theta^*\|^{1-p} \sqrt{d\log(1/\delta)/n},$$
(96)

where  $c_t = \Theta(\exp(t))$ . Before we start, we assume  $\forall t < T$ , we have 

$$c_t \|\theta_t - \theta^*\|^{-p} \sqrt{d\log(1/\delta)/n} \le \frac{1}{C_p}.$$
(97)

Otherwise, our results simply follow. We now prove (96) by induction. Notice that from (95), we know that (96) holds for t = 1. Assume for  $k \le t$ , we have  $\|\theta_k^n - \theta_k\| \le c_k \|\theta_k - \theta^*\|^{1-p} \sqrt{d \log(1/\delta)/n}$ . Note that 

$$\|\theta_{t+1}^n - \theta_{t+1}\| \le \|\theta_t^n - \theta_t\| + \|H_t^n \nabla \mathcal{L}_n(\theta_t^n) - H_t \nabla \mathcal{L}(\theta_t)\|.$$
(98)

 $H_t^n = \left(I - \frac{s_{t-1}^n (u_{t-1}^n)^\top}{(u_{t-1}^n)^\top s_{t-1}^n}\right) H_{t-1}^n \left(I - \frac{u_{t-1}^n (s_{t-1}^n)^\top}{(s_{t-1}^n)^\top u_{t-1}^n}\right) + \frac{s_{t-1}^n (s_{t-1}^n)^\top}{(s_{t-1}^n)^\top u_{t-1}^n},$ 

 $H_{t} = \left(I - \frac{s_{t-1}(u_{t-1})^{\top}}{(u_{t-1})^{\top}s_{t-1}}\right) H_{t-1} \left(I - \frac{u_{t-1}(s_{t-1})^{\top}}{(s_{t-1})^{\top}u_{t-1}}\right) + \frac{s_{t-1}(s_{t-1})^{\top}}{(s_{t-1})^{\top}u_{t-1}},$ 

We apply the update of 990

93

- 994
- 995

- 997
- 998

and

- 999 1000
- 1001

 $\left(I - \frac{u_{t-1}s_{t-1}^{\top}}{s_{t-1}^{\top}u_{t-1}}\right)\nabla\mathcal{L}(\theta_t) = 0,$ (101)

(99)

(100)

(105)

1003 which follows equation (41) from the population analysis. Then, we have the following decomposition: 1004

$$\begin{aligned} & \|H_{t}^{n} \nabla \mathcal{L}_{n}(\theta_{t}^{n}) - H_{t} \nabla f(\theta_{t})\| = \left\| H_{t}^{n} \nabla \mathcal{L}_{n}(\theta_{t}^{n}) - \frac{s_{t-1} s_{t-1}^{\top}}{s_{t-1}^{\top} u_{t-1}} \nabla \mathcal{L}(\theta_{t}) \right\| \\ & \| \left( I - \frac{s_{t-1}^{n} (u_{t-1}^{n})^{\top}}{(u_{t-1}^{n})^{\top} s_{t-1}^{n}} \right) H_{t-1}^{n} \left( I - \frac{u_{t-1}^{n} (s_{t-1}^{n})^{\top}}{(s_{t-1}^{n})^{\top} u_{t-1}^{n}} \right) \nabla \mathcal{L}_{n}(\theta_{t}^{n}) \right\| + \left\| \frac{s_{t-1}^{n} (s_{t-1}^{n})^{\top}}{(s_{t-1}^{n})^{\top} u_{t-1}^{n}} \nabla \mathcal{L}_{n}(\theta_{t}^{n}) - \frac{s_{t-1} s_{t-1}^{\top}}{s_{t-1}^{\top} u_{t-1}} \nabla \mathcal{L}(\theta_{t}) \right\|, \end{aligned}$$

$$(102)$$

Now we bound these two terms separately. The first term can be bounded as 1012

$$\left\| \left( I - \frac{s_{t-1}^{n} (u_{t-1}^{n})^{\top}}{(u_{t-1}^{n})^{\top} s_{t-1}^{n}} \right) H_{t-1}^{n} \left( I - \frac{u_{t-1}^{n} (s_{t-1}^{n})^{\top}}{(s_{t-1}^{n})^{\top} u_{t-1}^{n}} \right) \nabla \mathcal{L}_{n}(\theta_{t}^{n}) \right\| \\
\leq \left\| I - \frac{s_{t-1}^{n} (u_{t-1}^{n})^{\top}}{(u_{t-1}^{n})^{\top} s_{t-1}^{n}} \right\| \| H_{t-1}^{n} \| \left\| \left( I - \frac{u_{t-1}^{n} (s_{t-1}^{n})^{\top}}{(s_{t-1}^{n})^{\top} u_{t-1}^{n}} \right) \nabla \mathcal{L}_{n}(\theta_{t}^{n}) \right\| \\
\leq \left\| H_{t-1}^{n} \right\| \left\| \left( I - \frac{u_{t-1}^{n} (s_{t-1}^{n})^{\top}}{(s_{t-1}^{n})^{\top} u_{t-1}^{n}} \right) \nabla \mathcal{L}_{n}(\theta_{t}^{n}) \right\|,$$
(103)

1021 and 1022

1016

1018 1019

$$\begin{aligned} & \left\| \left( I - \frac{u_{t-1}^{n}(s_{t-1}^{n})^{\top}}{(s_{t-1}^{n})^{\top}u_{t-1}^{n}} \right) \nabla \mathcal{L}_{n}(\theta_{t}^{n}) \right\| = \left\| \left( I - \frac{u_{t-1}^{n}(s_{t-1}^{n})^{\top}}{(s_{t-1}^{n})^{\top}u_{t-1}^{n}} \right) \nabla \mathcal{L}_{n}(\theta_{t}^{n}) - \left( I - \frac{u_{t-1}s_{t-1}^{\top}}{s_{t-1}^{\top}u_{t-1}} \right) \nabla \mathcal{L}(\theta_{t}) \right\| \\ & \leq \left\| \nabla \mathcal{L}_{n}(\theta_{t}^{n}) - \nabla \mathcal{L}(\theta_{t}) \right\| + \left\| \frac{u_{t-1}^{n}s_{t-1}^{n}^{\top} \nabla \mathcal{L}_{n}(\theta_{t}^{n})}{(s_{t-1}^{n})^{\top}u_{t-1}^{n}} - \frac{u_{t-1}^{n}s_{t-1}^{\top} \nabla \mathcal{L}(\theta_{t})}{s_{t-1}^{\top}u_{t-1}} \right\| + \left\| \frac{u_{t-1}^{n}s_{t-1}^{\top} \nabla \mathcal{L}(\theta_{t})}{s_{t-1}^{\top}u_{t-1}} - \frac{u_{t-1}s_{t-1}^{\top} \nabla \mathcal{L}(\theta_{t})}{s_{t-1}^{\top}u_{t-1}} \right\| \\ & \leq \left\| \nabla \mathcal{L}_{n}(\theta_{t}^{n}) - \nabla \mathcal{L}(\theta_{t}) \right\| + \left\| u_{t-1}^{n} \right\| \left\| \frac{(s_{t-1}^{n})^{\top} \nabla \mathcal{L}_{n}(\theta_{t}^{n})}{(s_{t-1}^{n})^{\top}u_{t-1}^{n}} - \frac{s_{t-1}^{\top} \nabla \mathcal{L}(\theta_{t})}{s_{t-1}^{\top}u_{t-1}} \right\| + \left\| u_{t-1}^{n} - u_{t-1} \right\| \frac{s_{t-1}^{\top} \nabla \mathcal{L}(\theta_{t})}{s_{t-1}^{\top}u_{t-1}}, \end{aligned}$$

$$\tag{104}$$

where we use the fact that  $\left\|I - \frac{s_{t-1}^n(u_{t-1}^n)^\top}{(u_{t-1}^n)^\top s_{t-1}^n}\right\| \le 1$  and  $\left(I - \frac{u_{t-1}s_{t-1}^\top}{s_{t-1}^\top u_{t-1}}\right) \nabla \mathcal{L}(\theta_t) = 0$ . The second term can be bounded as

1034

1037 1

$$\begin{aligned} & \left\| \frac{s_{t-1}^{n}(s_{t-1}^{n})^{\top}}{(s_{t-1}^{n})^{\top}u_{t-1}^{n}} \nabla \mathcal{L}_{n}(\theta_{t}^{n}) - \frac{s_{t-1}^{n}s_{t-1}^{\top}}{s_{t-1}^{\top}u_{t-1}} \nabla \mathcal{L}(\theta_{t}) \right\| + \left\| \frac{s_{t-1}^{n}s_{t-1}^{\top}}{s_{t-1}^{\top}u_{t-1}} \nabla \mathcal{L}(\theta_{t}) - \frac{s_{t-1}s_{t-1}^{\top}}{s_{t-1}^{\top}u_{t-1}} \nabla \mathcal{L}(\theta_{t}) \right\| \\ & \left\| \frac{s_{t-1}^{n}u_{t-1}^{n}}{(s_{t-1}^{n})^{\top}\nabla \mathcal{L}_{n}(\theta_{t}^{n})} - \frac{s_{t-1}^{\top}\nabla \mathcal{L}(\theta_{t})}{s_{t-1}^{\top}u_{t-1}} \right\| + \left\| s_{t-1} - s_{t-1}^{n} \right\| \frac{s_{t-1}^{\top}\nabla \mathcal{L}(\theta_{t})}{s_{t-1}^{\top}u_{t-1}}. \end{aligned}$$

 $\left\|\frac{s_{t-1}^n(s_{t-1}^n)^\top}{(s_{t-1}^n)^\top u_{t-1}^n}\nabla\mathcal{L}_n(\theta_t^n) - \frac{s_{t-1}s_{t-1}^\top}{s_{t-1}^\top u_{t-1}}\nabla\mathcal{L}(\theta_t)\right\|$ 

1042

1043 We start from the bound on several basic terms. 1044

• Bounds related to  $s_{t-1}^n$ :

$$\begin{aligned} \|s_{t-1}^{n} - s_{t-1}\| &\leq \|\theta_{t}^{n} - \theta_{t}\| + \|\theta_{t-1}^{n} - \theta_{t-1}\| \\ &\leq (c_{t}\|\theta_{t-1} - \theta^{*}\|^{1-p} + c_{t-1}\|\theta_{t-2} - \theta^{*}\|^{1-p})\sqrt{d\log(1/\delta)/n} \\ &\leq (c_{t} + c_{t-1}r_{t-1}^{p-1})\|\theta_{t-1} - \theta^{*}\|^{1-p}\sqrt{d\log(1/\delta)/n}, \end{aligned}$$
(106)

which also gives

$$\begin{aligned} \|s_{t-1}^{n}\| \leq \|s_{t-1}\| + \|s_{t-1}^{n} - s_{t-1}\| \\ &= (1 - r_{t-1})\|\theta_{t-1} - \theta^{*}\| + (c_{t}\|\theta_{t-1} - \theta^{*}\|^{1-p} + c_{t-1}\|\theta_{t-2} - \theta^{*}\|^{1-p})\sqrt{d\log(1/\delta)/n} \\ &\leq (1 - r_{t-1})\|\theta_{t-1} - \theta^{*}\| + \frac{1 + 1/r_{t-2}}{C_{p}}\|\theta_{t-1} - \theta^{*}\| \\ &\leq \|\theta_{t-1} - \theta^{*}\|, \end{aligned}$$
(107)

where we use the fact that  $c_t \|\theta_{t-1} - \theta^*\|^{-p} \sqrt{d \log(1/\delta)/n} \le \frac{1}{C_p}$  and we can choose sufficiently large  $C_p$  to make the last inequality holds.

## • Bounds related to $\nabla \mathcal{L}_n(\theta_n^t)$ :

$$\begin{aligned} \|\nabla \mathcal{L}_{n}(\theta_{t}^{n}) - \nabla \mathcal{L}(\theta_{t})\| \\ \leq \|\nabla \mathcal{L}_{n}(\theta_{t}^{n}) - \nabla \mathcal{L}(\theta_{t}^{n})\| + \|\nabla \mathcal{L}(\theta_{t}^{n}) - \nabla \mathcal{L}(\theta_{t})\| \\ \leq C_{p} \|\theta_{t}^{n} - \theta^{*}\|^{p-1} \sqrt{d\log(1/\delta)/n} + C_{p} \|\theta_{t}^{n} - \theta_{t}\| \left( \|\theta_{t}^{n} - \theta_{t}\| + \|\theta_{t} - \theta^{*}\| \right)^{2p-2} \\ \leq C_{p} \left( \|\theta_{t}^{n} - \theta_{t}\| + \|\theta_{t} - \theta^{*}\| \right)^{p-1} \sqrt{d\log(1/\delta)/n} \left( 1 + c_{t} \|\theta_{t-1} - \theta^{*}\|^{1-p} \left( \|\theta_{t}^{n} - \theta_{t}\| + \|\theta_{t} - \theta^{*}\| \right)^{p-1} \right) \\ \leq (C_{p} + C_{p}c_{t}) \|\theta_{t-1} - \theta^{*}\|^{p-1} \sqrt{d\log(1/\delta)/n}, \end{aligned}$$
(108)

which also gives

$$\|\nabla \mathcal{L}_n(\theta_t^n)\| \le \|\nabla \mathcal{L}(\theta_t)\| + \|\nabla \mathcal{L}_n(\theta_t^n) - \nabla \mathcal{L}(\theta_t)\| \le C_p \|\theta_t - \theta^*\|^{2p-1},$$
(109)

where we still use the fact that  $c_t \|\theta_{t-1} - \theta^*\|^{-p} \sqrt{d \log(1/\delta)/n} \le \frac{1}{C_p}$  and we can choose sufficiently large  $C_p$  to make the noise term negligible.

• Bounds related to  $u_{t-1}^n$ :

$$\|u_{t-1}^{n} - u_{t-1}\|$$

$$\leq \|\nabla \mathcal{L}_{n}(\theta_{t}^{n}) - \nabla \mathcal{L}(\theta_{t})\| + \|\nabla \mathcal{L}_{n}(\theta_{t-1}^{n}) - \nabla \mathcal{L}(\theta_{t-1})\|$$

$$\leq (C_{p} + C_{p}c_{t} + C_{p}c_{t-1})\|\theta_{t-1} - \theta^{*}\|^{p-1}\sqrt{d\log(1/\delta)/n}.$$
(110)

With the same technique, we can show that

$$\|u_{t-1}^n\| \le C_p \|\theta_{t-1} - \theta^*\|^{2p-1}.$$
(111)

• Bound of  $|(s_{t-1}^n)^\top \nabla \mathcal{L}_n(\theta_t^n) - s_{t-1}^\top \nabla \mathcal{L}(\theta_t)|$ :

$$|(s_{t-1}^{n})^{\top} \nabla \mathcal{L}_{n}(\theta_{t}^{n}) - s_{t-1}^{\top} \nabla \mathcal{L}(\theta_{t})|$$

$$\leq ||s_{t-1}^{n} - s_{t-1}|| ||\nabla \mathcal{L}_{n}(\theta_{t}^{n})|| + ||s_{t-1}|| ||\nabla \mathcal{L}_{n}(\theta_{t}^{n}) - \nabla \mathcal{L}(\theta_{t})||$$

$$\leq (C_{p} + C_{p}c_{t} + C_{p}c_{t-1}) ||\theta_{t-1} - \theta^{*}||^{p} \sqrt{d\log(1/\delta)/n}.$$
(112)

And with the same technique, we have

$$C_p \|\theta_{t-1} - \theta^*\|^{2p} \le (s_{t-1}^n)^\top \nabla \mathcal{L}_n(\theta_t^n) \le C_p \|\theta_{t-1} - \theta^*\|^{2p}.$$



$$|(s_{t-1}^{n})^{\top}(u_{t-1}^{n}) - s_{t-1}^{\top}u_{t-1}| \leq ||s_{t-1}^{n} - s_{t-1}|| ||u_{t-1}^{n}|| + ||s_{t-1}|| ||u_{t-1}^{n} - u_{t-1}|| \leq (C_{p} + C_{p}c_{t} + C_{p}c_{t-1})||\theta_{t-1} - \theta^{*}||^{p}\sqrt{d\log(1/\delta)/n},$$

$$(113)$$

which also gives

$$C_p \|\theta_{t-1} - \theta^*\|^{2p} \le (s_{t-1}^n)^\top u_{t-1}^n \le C_p \|\theta_{t-1} - \theta^*\|^{2p}.$$
(114)

• Bound for 
$$\left| \frac{(s_{t-1}^{n})^{\top} \nabla \mathcal{L}_{n}(\theta_{t}^{n})}{(s_{t-1}^{n})^{\top} u_{t-1}^{n}} - \frac{s_{t-1}^{\top} \nabla \mathcal{L}(\theta_{t})}{s_{t-1}^{\top} u_{t-1}} \right|$$
:  

$$\left| \frac{\left| (s_{t-1}^{n})^{\top} \nabla \mathcal{L}_{n}(\theta_{t}^{n})}{(s_{t-1}^{n})^{\top} u_{t-1}^{n}} - \frac{s_{t-1}^{\top} \nabla \mathcal{L}(\theta_{t})}{s_{t-1}^{\top} u_{t-1}} \right|$$

$$\leq \frac{\left| (s_{t-1}^{n})^{\top} \nabla \mathcal{L}_{n}(\theta_{t}^{n}) - s_{t-1}^{\top} \nabla \mathcal{L}(\theta_{t}) \right| s_{t-1}^{\top} u_{t-1} + s_{t-1}^{\top} \nabla \mathcal{L}(\theta_{t}) \left| s_{t-1}^{\top} u_{t-1} - (s_{t-1}^{n})^{\top} u_{t-1}^{n} \right|}{(s_{t-1}^{n})^{\top} u_{t-1}^{n} s_{t-1}^{\top} u_{t-1}}$$

$$\leq (C_{p} + C_{p}c_{t} + C_{p}c_{t-1}) \| \theta_{t-1} - \theta^{*} \|^{-p} \sqrt{d \log(1/\delta)/n}.$$
(115)

• Straightforward computation shows that  $s_{t-1}^{\top} \nabla \mathcal{L}(\theta_t) = C_p \|\theta_t - \theta^*\|^2 p$ ,  $s_{t-1}^{\top} u_{t-1} = C_p \|\theta - \theta^*\|^{2p}$ .

1121 To summarize, we have that 1122

$$\left\| \left( I - \frac{u_{t-1}^n (s_{t-1}^n)^\top}{(s_{t-1}^n)^\top u_{t-1}^n} \right) \nabla \mathcal{L}_n(\theta_t^n) \right\| \le (C_p + C_p c_t + C_p c_{t-1}) \|\theta_{t-1} - \theta^*\|^{p-1} \sqrt{d \log(1/\delta)/n},$$
(116)

$$\left\|\frac{s_{t-1}^{n}(s_{t-1}^{n})^{\top}}{(s_{t-1}^{n})^{\top}u_{t-1}^{n}}\nabla\mathcal{L}_{n}(\theta_{t}^{n}) - \frac{s_{t-1}s_{t-1}^{\top}}{s_{t-1}^{\top}u_{t-1}}\nabla\mathcal{L}(\theta_{t})\right\| \leq (C_{p} + C_{p}c_{t} + C_{p}c_{t-1})\|\theta_{t-1} - \theta^{*}\|^{1-p}\sqrt{d\log(1/\delta)/n}.$$
(117)

1129 Now we bound  $||H_t^n||$ . Recall the update of  $H_t^n$ 

$$H_t^n = \left(I - \frac{s_{t-1}^n (u_{t-1}^n)^\top}{(u_{t-1}^n)^\top s_{t-1}^n}\right) H_{t-1}^n \left(I - \frac{u_{t-1}^n (s_{t-1}^n)^\top}{(s_{t-1}^n)^\top u_{t-1}^n}\right) + \frac{s_{t-1}^n (s_{t-1}^n)^\top}{(s_{t-1}^n)^\top u_{t-1}^n}.$$
(118)

1134 With the property of  $\|I - \frac{s_{t-1}^n (u_{t-1}^n)^\top}{(u_{t-1}^n)^\top s_{t-1}^n}\| \le 1$  and  $\|I - \frac{u_{t-1}^n (s_{t-1}^n)^\top}{(s_{t-1}^n)^\top u_{t-1}^n}\| \le 1$ , we have that 

$$\begin{aligned} \|H_{t}^{n}\| \leq \|I - \frac{s_{t-1}^{n}(u_{t-1}^{n})^{\top}}{(u_{t-1}^{n})^{\top}s_{t-1}^{n}}\|\|H_{t-1}^{n}\|\|I - \frac{u_{t-1}^{n}(s_{t-1}^{n})^{\top}}{(s_{t-1}^{n})^{\top}u_{t-1}^{n}}\| + \|\frac{s_{t-1}^{n}(s_{t-1}^{n})^{\top}u_{t-1}^{n}}{(s_{t-1}^{n})^{\top}u_{t-1}^{n}}\| \\ \leq \|H_{t-1}^{n}\| + \frac{\|s_{t-1}^{n}\|^{2}}{(s_{t-1}^{n})^{\top}u_{t-1}^{n}} \\ < \|(\nabla^{2}\mathcal{L}_{n}(\theta_{0}))^{-1}\| + \sum^{t-1}\frac{\|s_{t}^{n}\|^{2}}{(s_{t-1}^{n})^{\top}(s_{t-1}^{n})^{\top}(s_{t-1}^{n})^{\top}(s_{t-1}^{n})^{\top}(s_{t-1}^{n})^{\top}(s_{t-1}^{n})^{\top}(s_{t-1}^{n})^{\top}(s_{t-1}^{n})^{\top}(s_{t-1}^{n})^{\top}(s_{t-1}^{n})^{\top}(s_{t-1}^{n})^{\top}(s_{t-1}^{n})^{\top}(s_{t-1}^{n})^{\top}(s_{t-1}^{n})^{\top}(s_{t-1}^{n})^{\top}(s_{t-1}^{n})^{\top}(s_{t-1}^{n})^{\top}(s_{t-1}^{n})^{\top}(s_{t-1}^{n})^{\top}(s_{t-1}^{n})^{\top}(s_{t-1}^{n})^{\top}(s_{t-1}^{n})^{\top}(s_{t-1}^{n})^{\top}(s_{t-1}^{n})^{\top}(s_{t-1}^{n})^{\top}(s_{t-1}^{n})^{\top}(s_{t-1}^{n})^{\top}(s_{t-1}^{n})^{\top}(s_{t-1}^{n})^{\top}(s_{t-1}^{n})^{\top}(s_{t-1}^{n})^{\top}(s_{t-1}^{n})^{\top}(s_{t-1}^{n})^{\top}(s_{t-1}^{n})^{\top}(s_{t-1}^{n})^{\top}(s_{t-1}^{n})^{\top}(s_{t-1}^{n})^{\top}(s_{t-1}^{n})^{\top}(s_{t-1}^{n})^{\top}(s_{t-1}^{n})^{\top}(s_{t-1}^{n})^{\top}(s_{t-1}^{n})^{\top}(s_{t-1}^{n})^{\top}(s_{t-1}^{n})^{\top}(s_{t-1}^{n})^{\top}(s_{t-1}^{n})^{\top}(s_{t-1}^{n})^{\top}(s_{t-1}^{n})^{\top}(s_{t-1}^{n})^{\top}(s_{t-1}^{n})^{\top}(s_{t-1}^{n})^{\top}(s_{t-1}^{n})^{\top}(s_{t-1}^{n})^{\top}(s_{t-1}^{n})^{\top}(s_{t-1}^{n})^{\top}(s_{t-1}^{n})^{\top}(s_{t-1}^{n})^{\top}(s_{t-1}^{n})^{\top}(s_{t-1}^{n})^{\top}(s_{t-1}^{n})^{\top}(s_{t-1}^{n})^{\top}(s_{t-1}^{n})^{\top}(s_{t-1}^{n})^{\top}(s_{t-1}^{n})^{\top}(s_{t-1}^{n})^{\top}(s_{t-1}^{n})^{\top}(s_{t-1}^{n})^{\top}(s_{t-1}^{n})^{\top}(s_{t-1}^{n})^{\top}(s_{t-1}^{n})^{\top}(s_{t-1}^{n})^{\top}(s_{t-1}^{n})^{\top}(s_{t-1}^{n})^{\top}(s_{t-1}^{n})^{\top}(s_{t-1}^{n})^{\top}(s_{t-1}^{n})^{\top}(s_{t-1}^{n})^{\top}(s_{t-1}^{n})^{\top}(s_{t-1}^{n})^{\top}(s_{t-1}^{n})^{\top}(s_{t-1}^{n})^{\top}(s_{t-1}^{n})^{\top}(s_{t-1}^{n})^{\top}(s_{t-1}^{n})^{\top}(s_{t-1}^{n})^{\top}(s_{t-1}^{n})^{\top}(s_{t-1}^{n})^{\top}(s_{t-1}^{n})^{\top}(s_{t-1}^{n})^{\top}(s_{t-1}^{n})^{\top}(s_{t-1}^{n})^{\top}(s_{t-1}^{n})^{\top}(s_{t-1}^{n})^{\top}(s_{t-1}^{n})^{\top}(s_{t-1}^{n})^{\top}(s_{t-1}^{n})^{\top}(s_{t-1}^{n})^{\top}(s_{t-1}^{n})^{\top}(s_{$$

$$\leq \left\| \left( \nabla^2 \mathcal{L}_n(\theta_0) \right)^{-1} \right\| + \sum_{i=0}^{i-1} \frac{\|s_i^n\|^2}{(s_i^n)^\top(u_i^n)}$$

1144 From the previous computation, we know

$$\frac{\|s_i^n\|^2}{(s_i^n)^\top (u_i^n)} \le C_p \|\theta_i - \theta^*\|^{2-2p}.$$
(120)

From Theorem 4.5 of the population analysis, we know that there exists  $0 < r_l < r_h < 1$ , such that for all the linear convergence rates  $\{r_t\}_{t=0}^{\infty}$ , we have  $r_l \leq r_t \leq r_h$  for all  $t \geq 0$  and  $\|\theta_{t+1} - \theta^*\| = r_t \|\theta_t - \theta^*\|$  for all  $t \geq 0$ . Hence, we have that for all  $0 \leq i \leq t - 1$ ,

$$\|\theta_{t-1} - \theta^*\| = \prod_{j=i}^{t-2} r_j \|\theta_i - \theta^*\| \le r_h^{t-1-i} \|\theta_i - \theta^*\|,$$
(121)

$$\|\theta_{i} - \theta^{*}\| \ge r_{h}^{i+1-t} \|\theta_{t-1} - \theta^{*}\|,$$
(122)

$$\|\theta_i - \theta^*\|^{2-2p} \le r_h^{(2p-2)(t-1-i)} \|\theta_{t-1} - \theta^*\|^{2-2p},$$
(123)

$$\sum_{i=0}^{t-1} \|\theta_i - \theta^*\|^{2-2p} \le \sum_{i=0}^{t-1} r_h^{(2p-2)(t-1-i)} \|\theta_{t-1} - \theta^*\|^{2-2p} \le \frac{1}{1 - r_h^{2p-2}} \|\theta_{t-1} - \theta^*\|^{2-2p}.$$
(124)

## 1164 Hence, we obtain that

$$\|H_{t}^{n}\| \leq \left\| \left( \nabla^{2} \mathcal{L}_{n}(\theta_{0}) \right)^{-1} \right\| + \sum_{i=0}^{t-1} \frac{\|s_{i}^{n}\|^{2}}{(s_{i}^{n})^{\top}(u_{i}^{n})} \\ \leq \frac{1}{\lambda_{min}(\nabla^{2} \mathcal{L}_{n}(\theta_{0}))} + C_{p} \sum_{i=0}^{t-1} \frac{\|s_{i}^{n}\|^{2}}{(s_{i}^{n})^{\top}(u_{i}^{n})} \\ \leq \frac{1}{2p\|\theta_{0} - \theta^{*}\|^{2p-2}} + C_{p} \frac{1}{1 - r_{h}^{2p-2}} \|\theta_{t-1} - \theta^{*}\|^{2-2p} \\ \leq \left( \frac{r_{h}^{(2p-2)(t-1)}}{2p} + \frac{C_{p}}{1 - r_{h}^{2p-2}} \right) \|\theta_{t-1} - \theta^{*}\|^{2-2p}.$$

$$(125)$$

1177 As long as  $C_p$  is large enough, we obtain that

$$\|H_t^n\| \le C_p \|\theta_{t-1} - \theta^*\|^{2-2p},\tag{126}$$

 $^{1180}_{1181} \ \, \text{and} \ \,$ 

$$\|H_{t-1}^n\| \le C_p \|\theta_{t-2} - \theta^*\|^{2-2p} \le C_p \frac{1}{r_{t-2}^{2-2p}} \|\theta_{t-1} - \theta^*\|^{2-2p} \le C_p r_h^{2p-2} \|\theta_{t-1} - \theta^*\|^{2-2p}.$$
(127)

1185 Combining the previous results of (102), (103), (116), (117) and (127), we have that

$$\|H_t^n \nabla \mathcal{L}(\theta_t^n) - H_t \nabla \mathcal{L}(\theta_t)\| \le (C_p + C_p c_t + C_p c_{t-1}) \|\theta_{t-1} - \theta^*\|^{1-p} \sqrt{d\log(1/\delta)/n}.$$
(128)

1189 Notice that by induction, we observe that

$$\begin{aligned} &\|\theta_{t}^{n} - \theta_{t}\| \leq c_{t} \|\theta_{t-1} - \theta^{*}\|^{1-p} \sqrt{d \log(1/\delta)/n} \\ &= c_{t} \frac{1}{r_{t-1}^{1-p}} \|\theta_{t} - \theta^{*}\|^{1-p} \sqrt{d \log(1/\delta)/n} \\ &\leq c_{t} r_{h}^{p-1} \|\theta_{t} - \theta^{*}\|^{1-p} \sqrt{d \log(1/\delta)/n} \\ &\leq (C_{p} + C_{p}c_{t} + C_{p}c_{t-1}) \|\theta_{t} - \theta^{*}\|^{1-p} \sqrt{d \log(1/\delta)/n}. \end{aligned}$$

$$(129)$$

Therefore for  $C_p$  large enough, we have that

$$\begin{aligned} \|\theta_{t+1}^{n} - \theta_{t+1}\| &\leq \|\theta_{t}^{n} - \theta_{t}\| + \|H_{t}^{n} \nabla \mathcal{L}_{n}(\theta_{t}^{n}) - H_{t} \nabla \mathcal{L}(\theta_{t})\| \\ &\leq (C_{p} + C_{p}c_{t} + C_{p}c_{t-1})\|\theta_{t} - \theta^{*}\|^{1-p} \sqrt{d\log(1/\delta)/n}. \end{aligned}$$
(130)

 $\frac{1202}{1203}$  We define that

$$c_{t+1} = C_p + C_p c_t + C_p c_{t-1}, (131)$$

1206 and we prove that 1207

$$\|\theta_{t+1}^n - \theta_{t+1}\| \le c_{t+1} \|\theta_t - \theta^*\|^{1-p} \sqrt{d\log(1/\delta)/n}.$$
(132)

With the standard recursion, we know that  $c_t \leq (C_p)^t$  for  $C_p$  large enough. Therefore, using induction we proved that (96) holds:

  $\|\theta_t^n - \theta_t\| < c_t \|\theta_{t-1} - \theta^*\|^{1-p} \sqrt{d \log(1/\delta)/n},$ 

where  $c_t = \Theta(\exp(t))$ . Notice that 

$$\|\theta_t^n - \theta^*\| \le \|\theta_t^n - \theta^t\| + \|\theta_t - \theta^*\| \le (C_p)^t \|\theta_t - \theta^*\|^{1-p} \sqrt{d\log(1/\delta)/n} + \|\theta_t - \theta^*\|.$$
(134)

The optimal T with minimum  $\|\theta_T^n - \theta^*\|$  should satisfy that 

$$C_{p}^{T} \|\theta_{T} - \theta^{*}\|^{-p} \sqrt{d \log(1/\delta)/n} = C_{p},$$
(135)

(133)

for which we obtain  $T = \frac{C \log(n/d \log(1/\delta))}{2(p+1)}$  for some C = polylog(p), and  $\|\theta_T^n - \theta^*\| \le C_p (d \log(1/\delta)/n)^{1/(2p+2)}$ , that finishes the proof. 

## **B.** Additional Experiments Results

#### **B.1. Experiments for Linear Factors**



Figure 4. Convergence of factors  $\{r_k\}_{k=0}^{\infty}$  to  $r_*$ .

#### **B.2.** Additional Experiments for the Empirical Loss



Figure 5. Convergence and statistical results in d = 50. Convergence of different methods for high SNR regime are shown in (a) and low SNR regime in (b). Statistical radius of BFGS in high SNR regime and low SNR regime are shown in (c) and (d). 

To show that BFGS can also be applied to high dimension scenarios, we conduct additional experiments on the generalized linear model with input d = 50,100,500 and the power of link function p = 2. The inputs are generated by  $\{X_i\}_{i=1}^n \sim$  $\mathcal{N}(0, \operatorname{diag}(\sigma_1^2, \cdot, \sigma_d^2))$  where  $\sigma_k = (0.96)^{k-1}$ , and the remaining setting and hyper-parameters are set identical to the low dimension scenarios. The results are shown in Figure 5 and 6. As the results show, the performance of BFGS in high dimensional scenarios are nearly identical to the low dimensional scenarios. 

Manuscript under review by ICML 2023



(a) High SNR regime (d = 100). (b) Low SNR regime (d = 100). (c) High SNR regime (d = 500). (d) Low SNR regime (d = 500).

Figure 6. Convergence of different methods with d = 100 for high SNR regime are shown in (a) and low SNR regime in (b). Convergence of different methods with d = 500 for high SNR regime are shown in (c) and low SNR regime in (d).

#### **B.3. Experiments in Middle SNR Regime**

Here we briefly illustrate the behavior of BFGS in Middle SNR regime. We consider the generalized linear model with d = 50, 100, 500 and p = 2. The inputs are still generated by  $\{X_i\}_{i=1}^n$ , but  $\theta^*$  now is uniformly sampled from the sphere with radius  $n^{-1/6}$ .



Figure 7. Convergence results and statistical results for medium SNR regime with d = 50 are shown in (a) and (b). Convergence of different methods with d = 100 and d = 500 for medium SNR regime are shown in (c) and (d).

The results are shown in Figure 7. We can see BFGS still converges fast, and the statistical radius of middle SNR regime lies between the Hign SNR and Low SNR. A rigorous characterization of the statistical radius of middle SNR regime will be left as future work.