

RIEMANNIAN NATURAL GRADIENT METHODS

JIANG HU^{*}, RUICHENG AO[†], ANTHONY MAN-CHO SO[‡], MINGHAN YANG[§], AND
ZAIWEN WEN[¶]

Abstract. This paper studies large-scale optimization problems on Riemannian manifolds whose objective function is a finite sum of negative log-probability losses. Such problems arise in various machine learning and signal processing applications. By introducing the notion of Fisher information matrix in the manifold setting, we propose a novel Riemannian natural gradient method, which can be viewed as a natural extension of the natural gradient method from the Euclidean setting to the manifold setting. We establish the almost-sure global convergence of our proposed method under standard assumptions. Moreover, we show that if the loss function satisfies certain convexity and smoothness conditions and the input-output map satisfies a Riemannian Jacobian stability condition, then our proposed method enjoys a local linear—or, under the Lipschitz continuity of the Riemannian Jacobian of the input-output map, even quadratic—rate of convergence. We then prove that the Riemannian Jacobian stability condition will be satisfied by a two-layer fully connected neural network with batch normalization with high probability, provided that the width of the network is sufficiently large. This demonstrates the practical relevance of our convergence rate result. Numerical experiments on applications arising from machine learning demonstrate the advantages of the proposed method over state-of-the-art ones.

Key words. Manifold optimization, Riemannian Fisher information matrix, Kronecker-factored approximation, Natural gradient method

AMS subject classifications. 90C06, 90C22, 90C26, 90C56

1 Introduction Manifold constrained learning problems are ubiquitous in machine learning, signal processing, and deep learning ; see, e.g., [6, 14, 32, 40, 17]. In this paper, we focus on manifold optimization problems of the form

$$(1.1) \quad \min_{\Theta \in \mathcal{M}} \Psi(\Theta) := -\frac{1}{|\mathcal{S}|} \sum_{(x,y) \in \mathcal{S}} \log p(y|f(x, \Theta)),$$

where \mathcal{M} is either an embedded submanifold of $\mathbb{R}^{m \times n}$ or a quotient manifold whose total space is an embedded submanifold of $\mathbb{R}^{m \times n}$, $\Theta \in \mathcal{M}$ is the parameter to be estimated, \mathcal{S} is a collection of $|\mathcal{S}|$ data pairs (x, y) with $x \in \mathcal{X}, y \in \mathcal{Y}$, \mathcal{X} and \mathcal{Y} are the input and output spaces, respectively, $f(\cdot, \Theta) : \mathcal{X} \rightarrow \mathcal{Y}$ is a mapping from the input space to the output space, and $p(y|f(x, \Theta))$ is the conditional probability of taking y conditioning on $f(x, \Theta)$. If the conditional distribution is assumed to be Gaussian, the objective function in (1.1) reduces to the square loss. When the conditional distribution $p(y|f(x, \Theta))$ obeys the multinomial distribution, the corresponding objective function is the cross-entropy loss. As an aside, it is worth noting the equivalence between the negative log probability loss and Kullback-Leibler (KL) divergence shown in [38].

Let us take the low-rank matrix completion (LRMC) problem [14, 32] as an example and explain how it can be fitted into the form (1.1). The goal of LRMC

^{*}Department of Systems Engineering and Engineering Management, The Chinese University of Hong Kong, Shatin, NT, Hong Kong (hjiangopt@gmail.com).

[†]School of Mathematical Sciences, Peking University, China (archer_arc@pku.edu.cn).

[‡]Department of Systems Engineering and Engineering Management, The Chinese University of Hong Kong, Shatin, NT, Hong Kong (manchose@se.cuhk.edu.hk).

[§]Beijing International Center for Mathematical Research, Peking University, China (yangminghan@pku.edu.cn).

[¶]Beijing International Center for Mathematical Research, Center for Data Science and College of Engineering, Peking University, Beijing, China (wenzw@pku.edu.cn).

39 is to recover a low-rank matrix from an observed matrix X of size $n \times N$. Denote by
 40 Ω the set of indices of known entries in X , the rank- p LRMC problem amounts to
 41 solving

$$42 \quad (1.2) \quad \min_{U \in \text{Gr}(n,p), A \in \mathbb{R}^{p \times N}} \frac{1}{2} \|\mathcal{P}_\Omega(UA - X)\|^2,$$

43 where $\text{Gr}(n,p)$ is the Grassmann manifold consists of all p -dimensional subspaces in
 44 \mathbb{R}^n . The operator $\mathcal{P}_\Omega(X)$ is defined in an element-wise manner with $\mathcal{P}_\Omega(X_{ij}) = X_{ij}$
 45 if $(i,j) \in \Omega$ and 0 otherwise. Partitioning $X = [x_1, \dots, x_N]$ leads to the following
 46 equivalent formulation

$$47 \quad \min_{U \in \text{Gr}(n,p), a_i \in \mathbb{R}^p} \frac{1}{2N} \sum_{i=1}^N \|\mathcal{P}_{\Omega_{x_i}}(Ua_i - x_i)\|^2,$$

48 where $x_i \in \mathbb{R}^n$ and the j -th element of $\mathcal{P}_{\Omega_{x_i}}(v)$ is v_j if $(i,j) \in \Omega$ and 0 otherwise.
 49 Given U , we can obtain a_i by solving a least squares problem, i.e.,

$$50 \quad a_i = a(U; x_i) := \arg \min_a \|\mathcal{P}_{\Omega_{x_i}}(Ua - x_i)\|^2.$$

51 Then, the LRMC problem can be written as

$$52 \quad (1.3) \quad \min_{U \in \text{Gr}(n,p)} \Psi(U) := \frac{1}{2N} \sum_{i=1}^N \|\mathcal{P}_{\Omega_{x_i}}(Ua(U; x_i) - x_i)\|^2.$$

53 For the Gaussian distribution $p(y|z) = \frac{1}{\sqrt{(2\pi)^n}} \exp(-\frac{1}{2}(y-z)^\top(y-z))$, it holds that
 54 $-\log p(y|z) = \frac{1}{2}\|y-z\|^2 + \frac{n \log(2\pi)}{2}$. Hence, problem (1.3) is a special case of problem
 55 (1.1), in which $\mathcal{S} = \{(x_i, 0)\}_{i=1}^N$, $\mathcal{X} = \mathbb{R}^n$, $\mathcal{Y} = \mathbb{R}^n$, $f(x, U) = \mathcal{P}_{\Omega_x}(Ua(U; x) - x)$,
 56 $\mathcal{M} = \text{Gr}(n,p)$, and $p(y|z) = \frac{1}{\sqrt{(2\pi)^n}} \exp(-\frac{1}{2}(y-z)^\top(y-z))$. Other applications that
 57 can be fitted into the form (1.1) will be introduced in Section 4.

58 **1.1 Motivation of this work** Since the calculation of the gradient of Ψ in
 59 (1.1) can be expensive when the dataset \mathcal{S} is large, various approximate or stochastic
 60 methods for solving (1.1) have been proposed. On the side of first-order methods, we
 61 have the stochastic gradient method [47], stochastic variance-reduced gradient method
 62 [31], and adaptive gradient methods [19, 35] for solving (1.1) in the Euclidean setting
 63 (i.e., $\mathcal{M} = \mathbb{R}^{m \times n}$). We refer the reader to the book [37] for variants of these algorithms
 64 and a comparison of their performance. For the general manifold setting, by utilizing
 65 manifold optimization techniques [1, 26, 13], Riemannian versions of the stochastic
 66 gradient method [11], stochastic variance-reduced gradient method [52, 67, 29], and
 67 adaptive gradient methods [10] have been developed.

68 On the side of second-order methods, existing algorithms for solving (1.1) in
 69 the Euclidean setting (i.e., $\mathcal{M} = \mathbb{R}^{m \times n}$) can be divided into two classes. The first
 70 is based on approximate Newton or quasi-Newton techniques; see, e.g., [48, 44, 15,
 71 60, 61, 21, 45]. The second is the natural gradient-type methods, which are based
 72 on the Fisher information matrix (FIM) [4]. When the FIM can be approximated
 73 by a Kronecker-product form, the natural gradient direction can be computed us-
 74 ing relatively low computational cost. It is well known that second-order methods
 75 can accelerate convergence by utilizing curvature information. In particular, natural
 76 gradient-type methods can perform much better than the stochastic gradient method

77 [39, 63, 7, 62, 9, 42] in the Euclidean setting. The connections between natural gra-
 78 dient methods and second-order methods have been established in [38]. Compared
 79 with the approximate Newton/quasi-Newton-type methods, methods based on FIM
 80 are shown to be more efficient when tackling large-scale learning problems. For the
 81 general manifold setting, Riemannian stochastic quasi-Newton-type and Newton-type
 82 methods [34, 33, 65] have been proposed by utilizing the second-order manifold ge-
 83 ometry and variance reduction techniques. However, to the best of our knowledge,
 84 there is currently no Riemannian natural gradient-type method for solving (1.1). In
 85 view of the efficiency of Euclidean natural gradient-type methods, we are motivated
 86 to develop their Riemannian analogs for solving (1.1).

87 **1.2 Our contributions** In this paper, we develop a new Riemannian natural
 88 gradient method for solving (1.1). Our main contributions are summarized as follows.

- 89 • We introduce the Riemannian FIM (RFIM) and Riemannian empirical FIM
 90 (REFIM) to approximate the Riemannian Hessian. These notions extend the
 91 corresponding ones for the Euclidean setting [4, 38] to the manifold setting.
 92 Then, we propose an adaptive regularized Riemannian natural gradient de-
 93 scent (RNGD) method. We show that for some representative applications,
 94 Kronecker-factorized approximations of RFIM and REFIM can be construc-
 95 ted, which reduce the computational cost of the Riemannian natural gradient
 96 direction. Our experiment results demonstrate that although RNGD is a
 97 second-order-type method, it has low per-iteration cost and enjoys favorable
 98 numerical performances.
- 99 • Under some mild conditions, we prove that RNGD globally converges to a
 100 stationary point of (1.1) almost surely. Moreover, if the loss function satisfies
 101 certain convexity and smoothness conditions and the input-output map f
 102 satisfies a Riemannian Jacobian stability condition, then we can establish the
 103 local linear—or, under the Lipschitz continuity of the Riemannian Jacobian of
 104 f , even quadratic—rate of convergence of the method by utilizing the notion
 105 of second-order retraction. We then show that for a two-layer neural network
 106 with batch normalization, the Riemannian Jacobian stability condition will
 107 be satisfied with high probability when the width of the network is sufficiently
 108 large.

109 **1.3 Notation** For an $m \times n$ matrix Θ , we denote its Frobenius norm by $\|\Theta\|$
 110 and its vectorization by $\theta = \text{vec}(\Theta) \in \mathbb{R}^{mn}$. For a smooth function $h : \mathbb{R}^{m \times n} \rightarrow \mathbb{R}$, we
 111 use $\nabla h(\Theta) \in \mathbb{R}^{m \times n}$ to denote its Euclidean gradient at $\Theta \in \mathbb{R}^{m \times n}$. For simplicity, we
 112 set $r = mn$. When no confusion can arise, we use $\nabla h(\theta)$ to denote the vectorization
 113 of $\nabla h(\Theta)$. We use $\nabla^2 h(\theta) \in \mathbb{R}^{r \times r}$ to denote the Euclidean Hessian of h at $\theta \in \mathbb{R}^r$.
 114 We denote the tangent space to \mathcal{M} at Θ by $T_\Theta \mathcal{M}$. We write $d \in T_\theta \mathcal{M}$ to mean
 115 $\text{mat}(d) \in T_\Theta \mathcal{M}$, where $d \in \mathbb{R}^r$ and $\text{mat}(d)$ converts d into a m -by- n matrix. For a
 116 retraction R defined on \mathcal{M} , we write $R_\theta(d) := \text{vec}(R_\Theta(D))$ for $D \in T_\Theta \mathcal{M}$, $\theta = \text{vec}(\Theta)$,
 117 and $d = \text{vec}(D)$. We shall use θ and Θ interchangeably when no confusion can arise.
 118 Basically, Θ is used when we want to utilize the manifold structure, while θ is used
 119 when we want to utilize the vector space structure of the ambient space.

120 **1.4 Organization** We begin with the preliminaries on manifold optimization
 121 and natural gradient methods in Section 2. In Section 3, we introduce the RFIM and
 122 its empirical version REFIM and derive some of their properties. Then, we present our
 123 proposed RNGD method by utilizing the RFIM and REFIM. In Section 4, we discuss
 124 practical implementations of the RNGD method when problem (1.1) enjoys certain

125 Kronecker-product structure. In Section 5, we study the convergence behavior of the
 126 RNGD method under various assumptions. Finally, we present numerical results in
 127 Section 6.

128 2 Preliminaries

129 **2.1 Manifold optimization** Consider the optimization problem

$$130 \quad (2.1) \quad \min_{\Theta \in \mathcal{M}} h(\Theta),$$

131 where \mathcal{M} is either an embedded submanifold of $\mathbb{R}^{m \times n}$ or a quotient manifold whose
 132 total space is an embedded submanifold of $\mathbb{R}^{m \times n}$ and $h : \mathbb{R}^{m \times n} \rightarrow \mathbb{R}$ is a smooth
 133 function. For every $\Theta \in \mathcal{M}$, we endow the tangent space $T_{\Theta}\mathcal{M}$ with a general
 134 Riemannian metric $\langle U, V \rangle_{\Theta} := \text{vec}(U)^{\top} D(\theta) \text{vec}(V)$, where $D(\theta) \in \mathbb{R}^{r \times r}$ is symmetric
 135 and positive definite on $T_{\theta}\mathcal{M}$. The design and analysis of numerical algorithms for
 136 tackling (2.1) have been extensively studied over the years; see, e.g., [1, 26, 13] and the
 137 references therein. One of the key constructs in the design of manifold optimization
 138 algorithms is the retraction operator. A smooth mapping $R : T\mathcal{M} := \cup_{\Theta \in \mathcal{M}} T_{\Theta}\mathcal{M} \rightarrow$
 139 \mathcal{M} is called a retraction operator if

- 140 • $R_{\Theta}(0) = \Theta$,
- 141 • $\text{DR}_{\Theta}(0)[\xi] := \frac{d}{dt} R_{\Theta}(t\xi) |_{t=0} = \xi$, for all $\xi \in T_{\Theta}\mathcal{M}$.

142 We call R a second-order retraction [1, Proposition 5.5.5] if $\mathcal{P}_{T_{\Theta}\mathcal{M}} \left(\frac{d^2}{dt^2} R_{\Theta}(t\xi) |_{t=0} \right)$
 143 $= 0$ for all $\Theta \in \mathcal{M}$ and $\xi \in T_{\Theta}\mathcal{M}$. Some examples of second-order retraction can be
 144 found in [3, Theorem 22]. Another key concept is the Riemannian gradient. Given
 145 $\Theta \in \mathcal{M}$, the vectorization of the Riemannian gradient $\widetilde{\text{grad}} h(\Theta) \in \mathbb{R}^{m \times n}$ of h at Θ is
 146 given by

$$147 \quad \widetilde{\text{grad}} h(\theta) = D(\theta)^{-1} \mathcal{P}_{T_{\theta}\mathcal{M}}(\nabla h(\theta)) \in \mathbb{R}^r,$$

148 where $\mathcal{P}_{T_{\theta}\mathcal{M}}(\cdot)$ is the orthogonal projection operator onto $T_{\theta}\mathcal{M}$. The retraction-based
 149 methods for solving (2.1) perform updates of the form

$$150 \quad (2.2) \quad \Theta^{k+1} = R_{\Theta^k}(td^k),$$

151 where d^k is a descent direction in the tangent space $T_{\Theta^k}\mathcal{M}$ and $t > 0$ is the step size.
 152 The retraction operator R constrains the iterates on \mathcal{M} . For the case where \mathcal{M} is
 153 an embedded submanifold, we always take the Euclidean metric as the Riemannian
 154 metric (i.e., $\langle U, V \rangle_{\Theta} = \text{vec}(U)^{\top} \text{vec}(V)$ for any $\Theta \in \mathcal{M}$) and use $\text{grad} h(\theta) \in \mathbb{R}^r$
 155 and $\text{Hess} h(\theta) \in \mathbb{R}^{r \times r}$ to denote the Riemannian gradient and Riemannian Hessian
 156 of h under the Euclidean metric, respectively. For the case where \mathcal{M} is a quotient
 157 manifold, we use a Riemannian metric that satisfies the horizontally invariant property
 158 in [1, Equation (3.38)], so that the Riemannian norm of a vector on $T_{\theta}\mathcal{M}$ does not
 159 depend on the representative element of θ in \mathcal{M} . We also assume that the total space
 160 has a retraction satisfying the projection property in [1, Equation (4.9)], so that the
 161 retraction R on \mathcal{M} can be defined according to [1, Equation (4.10)].

162 **2.2 Natural gradient descent method** The natural gradient descent (NGD)
 163 method was originally proposed in [4] to solve (1.1) in the Euclidean setting (i.e.,
 164 $\mathcal{M} = \mathbb{R}^{m \times n}$). Suppose that y follows the conditional distribution $P_{y|f(x, \Theta)}$. Consider
 165 the population loss under $P_{y|x}(\Theta) := P_{y|f(x, \Theta)}$, i.e.,

$$166 \quad (2.3) \quad \Phi(\Theta) := -\mathbb{E}_{P_x} \left[\mathbb{E}_{P_{y|x}(\Theta)} \log p(y|f(x, \Theta)) \right].$$

When $P_{y|x}(\Theta)$ and P_x are replaced by their empirical counterparts defined using \mathcal{S} , the population loss $\Phi(\Theta)$ reduces to the empirical loss $\Psi(\Theta)$. Now, the FIM associated with Φ is defined as

$$F(\theta) := \mathbb{E}_{P_x} [\mathbb{E}_{P_{y|x}(\theta)} [\nabla \log p(y|f(x, \theta)) \nabla \log p(y|f(x, \theta))^\top]] \in \mathbb{R}^{r \times r}.$$

Under certain regularity condition [20], we can interchange the order of expectation and derivative to obtain $F(\theta) = \nabla^2 \Phi(\theta)$. In what follows, we assume that such a regularity condition holds. Since the distribution of x is unknown, we set P_x to be the empirical distribution defined by \mathcal{S} . In practice, we may only be able to get hold of an empirical counterpart of $P_{y|x}(\Theta)$. The empirical FIM (EFIM) associated with Ψ is then defined by replacing $P_{y|x}(\Theta)$ with its empirical counterpart [53], i.e.,

$$\bar{F}(\theta) := \frac{1}{|\mathcal{S}|} \sum_{(x,y) \in \mathcal{S}} \nabla \log p(y|f(x, \theta)) \nabla \log p(y|f(x, \theta))^\top.$$

With the FIM, the natural gradient direction is given by

$$\tilde{\nabla} \Phi(\theta) := (F(\theta))^{-1} \nabla \Phi(\theta) \in \mathbb{R}^r.$$

It is shown in [5, Theorem 1] and [43, Proposition 1] that $\tilde{\nabla} \Phi(\theta)$ is the steepest descent direction in the sense that

$$-\frac{\tilde{\nabla} \Phi(\theta)}{\|\nabla \Phi(\theta)\|_{(F(\theta))^{-1}}} = \lim_{\epsilon \rightarrow 0} \frac{1}{\epsilon} \arg \min_{d \in \mathbb{R}^r: \text{KL}(P_{x,y}(\theta+d) \| P_{x,y}(\theta)) \leq \epsilon^2/2} \Phi(\theta + d),$$

where $\|\nabla \Phi(\theta)\|_{(F(\theta))^{-1}} := \sqrt{\nabla \Phi(\theta) (F(\theta))^{-1} \nabla \Phi(\theta)}$.

In the k -th iteration, the iterative scheme of NGD for minimizing (2.3) is

$$\theta^{k+1} = \theta^k - t_k \tilde{\nabla} \Phi(\theta^k),$$

where $t_k > 0$ is a step size. In the case where $F(\theta)$ is computationally expensive or inaccessible, we use the EFIM instead of the FIM. The connections between NGD and second-order methods are presented in [38].

3 Riemannian natural gradient method

3.1 Fisher information matrix on manifold When the parameter to be estimated Θ lies on \mathcal{M} , the Euclidean natural gradient direction need not lie on the tangent space to \mathcal{M} at Θ and thus cannot be used as a search direction in retraction-based methods. To overcome this difficulty, we first introduce the RFIM, which is defined as

$$(3.1) \quad F^R(\theta) := \mathbb{E}_{P_x} \left[\mathbb{E}_{P_{y|x}(\theta)} \left[\text{grad} \log p(y|f(x, \theta)) \text{grad} \log p(y|f(x, \theta))^\top \right] \right] \in \mathbb{R}^{r \times r},$$

where $\text{grad} \log p(y|f(x, \theta))$ is the Riemannian gradient of $\log p(y|f(x, \theta))$ with respect to θ under the Euclidean metric.¹ Note that the generalization of FIM in the manifold setting has been developed in [55, 12]. The RFIM defined in (3.1) can be regarded as an extrinsic representation (i.e., an r -by- r matrix) of the said generalization. Such

¹The RFIM should not be confused with the Riemannian Fisher information metric. For any two tangent vectors $u, v \in T_\theta \mathcal{M}$, the Riemannian Fisher information metric associated with the RFIM (3.1) is given by $u^\top F^R(\theta) v$.

196 an extrinsic representation relies on the Euclidean representation of the Riemannian
 197 gradient in the total space and presents a straightforward way to compute RFIM. It
 198 is easy to see that the range of $F^R(\theta)$ is included in $T_\theta\mathcal{M}$. Assuming that $F^R(\theta)$ is
 199 positive definite on $T_\theta\mathcal{M}$, we define the Riemannian natural gradient direction $d^R(\theta)$
 200 as

$$201 \quad (3.2) \quad d^R(\theta) := (F^R(\theta))^{-1} \text{grad} \Phi(\theta) \in \mathbb{R}^r,$$

202 which is a vector on $T_\theta\mathcal{M}$. The following theorem justifies our definition of RFIM. It
 203 extends the corresponding results on FIM given in [5, Theorem 1] and [43, Proposition
 204 1].

205 **THEOREM 3.1.** *Let \mathcal{M} be either an embedded submanifold of $\mathbb{R}^{m \times n}$ or a quotient*
 206 *manifold whose total space is an embedded submanifold of $\mathbb{R}^{m \times n}$, and $\Phi : \mathcal{M} \rightarrow \mathbb{R}$ be*
 207 *the function given in (2.3). Given $\Theta \in \mathcal{M}$, suppose that $F^R(\theta)$ is positive definite on*
 208 *$T_\theta\mathcal{M}$. Then, for any second-order retraction R on \mathcal{M} , the steepest descent direction*
 209 *in the tangent space to \mathcal{M} at Θ is given by $-d^R(\theta)$ in (3.2), i.e.,*

$$210 \quad (3.3) \quad \frac{-d^R(\theta)}{\|\text{grad} \Phi(\theta)\|_{(F^R(\theta))^{-1}}} = \lim_{\epsilon \rightarrow 0} \frac{1}{\epsilon} \arg \min_{d \in T_\theta\mathcal{M} : \mathbb{E}_{P_x} [\text{KL}(P_{y|x}(R_\theta(d)) \| P_{y|x}(\theta))] \leq \epsilon^2/2} \Phi(R_\theta(d)),$$

211 where $\|\text{grad} \Phi(\theta)\|_{(F^R(\theta))^{-1}} = \sqrt{\text{grad} \Phi(\theta)^\top (F^R(\theta))^{-1} \text{grad} \Phi(\theta)}$.

Proof. For $\Theta \in \mathcal{M}$, from the definition

$$\text{KL}(P_{y|x}(\theta) \| P_{y|x}(R_\theta(td))) = \mathbb{E}_{P_{y|x}(\theta)} \log p(y|f(x, \theta)) - \mathbb{E}_{P_{y|x}(\theta)} \log p(y|f(x, R_\theta(td))),$$

212 we have

$$213 \quad \begin{aligned} \frac{d}{dt} \text{KL}(P_{y|x}(\theta) \| P_{y|x}(R_\theta(td))) \Big|_{t=0} &= -\frac{d}{dt} \mathbb{E}_{P_{y|x}(\theta)} \log p(y|f(x, R_\theta(td))) \Big|_{t=0} \\ &= -d^\top \nabla \mathbb{E}_{P_{y|x}(\theta)} \log p(y|f(x, \theta)). \end{aligned}$$

214 By definition of the Riemannian gradient, we obtain

$$215 \quad d^\top \text{grad} \text{KL}(P_{y|x}(\theta) \| P_{y|x}(R_\theta(td))) \Big|_{t=0} = -d^\top \nabla \mathbb{E}_{P_{y|x}(\theta)} \log p(y|f(x, \theta)), \quad \forall d \in T_\theta\mathcal{M},$$

216 where $\text{grad} \text{KL}(P_{y|x}(\theta) \| P_{y|x}(R_\theta(td))) \Big|_{t=0} \in T_\theta\mathcal{M}$. Then, we have

$$217 \quad \text{grad} \text{KL}(P_{y|x}(\theta) \| P_{y|x}(R_\theta(td))) \Big|_{t=0} = -\text{grad} \mathbb{E}_{P_{y|x}(\theta)} \log p(y|f(x, \theta)).$$

218 Accordingly, using the Leibniz integral rule and the property of second-order retrac-
 219 tions [1, Proposition 5.5.5], we have the second-order derivative

$$220 \quad \begin{aligned} &\frac{d^2}{dt^2} \text{KL}(P_{y|x}(\theta) \| P_{y|x}(R_\theta(td))) \Big|_{t=0} \\ &= \mathbb{E}_{P_{y|x}(\theta)} [d^\top \text{grad} \log p(y|f(x, \theta)) (\text{grad} \log p(y|f(x, \theta)))^\top d]. \end{aligned}$$

It follows that $\text{grad} \mathbb{E}_{P_{y|x}(\theta)} \log p(y|f(x, \theta)) = 0$. By the definition of F^R , we conclude that

$$\mathbb{E}_{P_x} \text{KL}(P_{y|x}(\theta) \| P_{y|x}(R_\theta(d))) = \frac{1}{2} d^\top F^R(\theta) d + O(d^3), \quad \forall d \in T_\theta\mathcal{M}.$$

221 From the fact [43, Proposition 1] that

$$222 \quad \frac{-A^{-1} \nabla h(\theta)}{\|\nabla h(\theta)\|_{A^{-1}}} = \lim_{\epsilon \rightarrow 0} \frac{1}{\epsilon} \arg \min_{d: \|d\|_A \leq \epsilon} h(\theta + d),$$

223 where A is a positive definite matrix and $\|d\|_{A^{-1}} = \sqrt{d^\top A^{-1}d}$, we have

$$224 \quad (3.4) \quad \frac{-B^{-1}\nabla(\Phi \circ R_\theta)(0)}{\|\nabla(\Phi \circ R_\theta)(0)\|_{B^{-1}}} = \lim_{\epsilon \rightarrow 0} \frac{1}{\epsilon} \arg \min_{d \in T_\theta \mathcal{M}: \|d\|_B \leq \epsilon} \Phi(R_\theta(d)),$$

225 where $B : T_\theta \mathcal{M} \rightarrow T_\theta \mathcal{M}$ is a positive definite linear operator. Note that for all
226 $u \in T_\theta \mathcal{M}$, it holds that

$$227 \quad \nabla(\Phi \circ R_\theta)(0)[u] = \nabla\Phi(R_\theta(0))[DR_\theta(0)[u]] = u^\top \text{grad}\Phi(\theta).$$

228 This gives

$$229 \quad \nabla(\Phi \circ R_\theta)(0) = \text{grad}\Phi(\theta).$$

230 Substituting the above into (3.4) and letting $B = F^R(\theta)$, we have

$$231 \quad (3.5) \quad \frac{-(F^R(\theta))^{-1}\text{grad}\Phi(\theta)}{\|\text{grad}\Phi(\theta)\|_{(F^R(\theta))^{-1}}} = \lim_{\epsilon \rightarrow 0} \frac{1}{\epsilon} \arg \min_{d \in T_\theta \mathcal{M}: \|d\|_{F^R(\theta)} \leq \epsilon} \Phi(R_\theta(d)).$$

232 Therefore, (3.3) holds for any second-order retraction R . This completes the proof. \square

233 Note that for an embedded submanifold \mathcal{M} endowed with the Euclidean metric,
234 the Riemannian Hessian [2, Equation 7] of Φ at θ along $u \in T_\theta \mathcal{M}$ is given by

$$235 \quad \text{Hess}\Phi(\theta)[u] = \mathcal{P}_{T_\theta \mathcal{M}}(\nabla^2\Phi(\theta)[u]) - \mathcal{P}_{T_\theta \mathcal{M}}D_u(\text{grad}\Phi(\theta)).$$

236 Since $\mathbb{E}_{P_{y|x}(\theta)} \nabla \log p(y|f(x, \theta)) = \int_y \nabla p(y|f(x, \theta))dy = \nabla \int_y p(y|f(x, \theta))dy = 0$, we
237 have $\text{grad}\Phi(\theta) = 0$ and $\text{Hess}\Phi(\theta) = F^R(\theta)$. Due to the uniqueness of the second-
238 order Taylor expansion, the Riemannian Newton's direction at θ does not depend on
239 the Riemannian metric and is equal to $d^R(\theta)$ in (3.2). Hence, it is reasonable to use
240 the Euclidean metric to define the Riemannian natural gradient direction (3.2). For
241 a quotient manifold \mathcal{M} whose total space is an embedded submanifold and whose en-
242 dowed Riemannian metric is horizontally invariant, it follows from [1, Equation (3.39)]
243 that the Riemannian gradient of Φ in the total space is the horizontal lift of the corre-
244 sponding Riemannian gradient in \mathcal{M} . Since the total space is an embedded manifold,
245 we see from [2, Equation 7] and our earlier argument that $F^R(\theta)$ is the Riemannian
246 Hessian of Φ in the total space at the representative element θ . Furthermore, by [1,
247 Proposition 5.3.3], the horizontal lift of the corresponding Riemannian Hessian in \mathcal{M}
248 at the representative element θ equals the horizontal projection of $F^R(\theta)$. Since the
249 Riemannian gradient of Φ in the total space at a representative element θ belongs to
250 the horizontal space at θ , we conclude that $d^R(\theta)$ in (3.2), which lies in the horizontal
251 space at θ , is the Riemannian Newton's direction at θ . As the Riemannian natural
252 gradient direction is independent of the choice of the Riemannian metric, we can use
253 the Euclidean metric to define (3.2), but a horizontally invariant Riemannian metric
254 should be introduced to compare the norms of Riemannian gradients. In summary,
255 the Riemannian natural descent direction (3.2) behaves as the Riemannian Newton's
256 direction whenever \mathcal{M} is an embedded submanifold or a quotient manifold whose total
257 space is an embedded submanifold.

258 Similar to EFIM, we can define REFIM as

$$259 \quad (3.6) \quad \bar{F}^R(\theta) := \frac{1}{|\mathcal{S}|} \sum_{(x,y) \in \mathcal{S}} \text{grad} \log p(y|f(x, \theta)) \text{grad} \log p(y|f(x, \theta))^\top.$$

260 **3.2 Algorithmic framework** To fix ideas, let us first consider the case where
 261 \mathcal{M} is an embedded submanifold. In the k -th iteration, once we obtain an estimate F_k
 262 of the RFIM (3.1) associated with Φ or the REFIM (3.6) associated with Ψ at θ^k , the
 263 Riemannian natural gradient direction in the tangent space to \mathcal{M} at θ^k is computed
 264 by solving the following optimization problem:

$$265 \quad (3.7) \quad d^k = \arg \min_{d \in T_{\theta^k} \mathcal{M}} m_k(d) := \Psi_k + \langle g^k, d \rangle + \frac{1}{2} \langle (F_k + \lambda_k I)d, d \rangle,$$

266 where $\langle u, v \rangle := u^\top v$ for two vectors $u, v \in \mathbb{R}^r$, $F_k d$ is the usual matrix-vector multi-
 267 plication, Ψ_k and g^k are stochastic estimates of $\Psi(\theta^k)$ and $\text{grad } \Psi(\theta^k)$, respectively,
 268 and $\lambda_k > 0$ is usually updated adaptively by a trust region-like strategy. In view
 269 of the finite-sum structure of Ψ (see (1.1)), the stochastic estimates Ψ_k and g_k can
 270 be obtained using, e.g., a mini-batch strategy (i.e., randomly sample a subset of \mathcal{S}
 271 and sum the corresponding terms in Ψ and $\text{grad } \Psi$ to get Ψ_k and g_k , respectively).
 272 Since $F_k + \lambda_k I : T_{\theta^k} \mathcal{M} \rightarrow T_{\theta^k} \mathcal{M}$ is positive definite and $g^k \in T_{\theta^k} \mathcal{M}$, the solution of
 273 (3.7) is $d^k = -(F_k + \lambda_k I)^{-1} g^k$. If the inverse of $F_k + \lambda_k I$ is costly to compute, then
 274 the truncated conjugate gradient method can be utilized [41]. We will introduce the
 275 constructions of a few computationally efficient approximation F_k in Section 4.

276 Once d^k is obtained, we construct a trial point

$$277 \quad (3.8) \quad z^k = R_{\theta^k}(d^k).$$

278 To measure whether z^k leads to a sufficient decrease in the objective value, we first
 279 calculate the ratio ρ_k between the reduction of Ψ and the reduction of m_k . Since the
 280 exact evaluation of Ψ is costly, one popular way [16] is to construct estimates Ψ_k^0 and
 281 $\Psi_k^{z^k}$ of $\Psi(\theta^k)$ and $\Psi(z^k)$, respectively. Then, we compute the ratio as

$$282 \quad (3.9) \quad \rho_k = \frac{\Psi_k^{z^k} - \Psi_k^0}{m_k(d^k) - \Psi_k^0}.$$

283 Here, we take $\Psi_k = \Psi_k^0$ in the calculation of $m_k(d^k)$. Lastly, we perform the update

$$284 \quad (3.10) \quad \theta^{k+1} = \begin{cases} z^k, & \text{if } \rho_k \geq \eta_1 \text{ and } \|g^k\| \geq \frac{\eta_2}{\sigma_k}, \\ \theta^k, & \text{otherwise,} \end{cases}$$

285 where $\eta_1 \in (0, 1)$ and $\eta_2 > 0$ are constants and $\sigma_k > 0$ is used to control the regular-
 286 ization parameter λ_k . Indeed, to ensure the descent property of the original function
 287 Ψ , some assumptions on the accuracy of the estimates of $\Psi(\theta^k)$, $\Psi(z^k)$ and the model
 288 m_k are needed, and they will be introduced later in the convergence analysis. Due to
 289 the error in the estimates, the regularization parameter λ_{k+1} should not only depend
 290 on the ratio ρ_k but also on the norm of the estimated Riemannian gradient g^k . In
 291 particular, we set $\lambda_{k+1} := \sigma_{k+1} \|g^{k+1}\|$ and update σ_{k+1} as

$$292 \quad (3.11) \quad \sigma_{k+1} = \begin{cases} \max \left\{ \sigma_{\min}, \frac{1}{\gamma} \sigma_k \right\}, & \text{if } \rho_k \geq \eta_1 \text{ and } \|g^k\| > \frac{\eta_2}{\sigma_k}, \\ \gamma \sigma_k, & \text{otherwise,} \end{cases}$$

293 where $\eta_1 \in (0, 1)$, $\eta_2 > 0$ are as before and $\sigma_{\min} > 0$, $\gamma > 1$ are parameters. Our
 294 proposed RNGD method is summarized in Algorithm 1. It is worth mentioning that a
 295 trust-region method is developed in [16] to solve stochastic optimization problems. Al-
 296 gorithm 1 can be seen as a combination of the stochastic update rule of the trust-region

Algorithm 1: Riemannian natural gradient descent (RNGD) for solving (1.1).

- 1 Choose an initial point θ^0 and parameters $\sigma_0 > 0$, $\sigma_{\min} > 0$, $\lambda_0 = \sigma_0 \|g^0\|$, $\eta_1 \in (0, 1)$, $\eta_2 > 0$, and $\gamma > 1$. Set $k = 0$.
 - 2 **while** *stopping conditions not met* **do**
 - 3 Compute the estimated Riemannian gradient g^k and the estimated Riemannian Fisher information matrix F_k .
 - 4 Compute the negative natural gradient direction d^k by solving (3.7) and compute the trial point z^k by (3.8).
 - 5 Update θ^{k+1} based on (3.10).
 - 6 Update λ_{k+1} based on (3.11).
 - 7 $k \leftarrow k + 1$.
-

297 radius in [16] and the adaptive regularization technique for manifold optimization in
 298 [27]. Compared with the trust-region subproblem in [16, Equation (2)], the subprob-
 299 lem (3.7) can be efficiently solved if the cost of computing the inverse of $F_k + \lambda_k I$
 300 is low. We remark that regularized subproblems similar to (3.7) have appeared in
 301 [39, 63, 62].

302 Now, for the case where \mathcal{M} is a quotient manifold, we have a horizontally invariant
 303 Riemannian metric $\langle U, V \rangle_{\Theta} := \text{vec}(U)^\top D(\theta) \text{vec}(V)$. The Riemannian gradient in the
 304 k -th iteration is $\tilde{g}^k = D(\theta^k)^{-1} g^k$. Thus, in Algorithm 1, we can still use g_k and F_k in
 305 (3.7) but should replace $\|g^k\|$ in λ_k , (3.10), and (3.11) with $\|\tilde{g}^k\|_{\theta^k} := \sqrt{(\tilde{g}^k)^\top D(\theta^k) \tilde{g}^k}$.
 306

307 **4 Practical Riemannian natural gradient descent methods** From the
 308 definitions of RFIM and REFIM in Section 3, the computational cost of solving sub-
 309 problem (3.7) may be high because of the vectorization of Θ . Fortunately, analogous
 310 to [39], the Riemannian natural gradient direction can be computed with a relatively
 311 low cost if the gradient of a single sample is of low rank, i.e., for a pair of observations
 312 $(x, y) \in \mathcal{S}$ and $\psi(\Theta; x, y) := -\log p(y|f(x, \Theta))$, $\nabla \psi$ takes the form

$$313 \quad (4.1) \quad \nabla \psi(\Theta; x, y) = G(x, y) A(x, y)^\top,$$

314 where $G(x, y) \in \mathbb{R}^{m \times q}$ and $A(x, y) \in \mathbb{R}^{n \times q}$ with $q \ll \min(m, n)$. Let us now elaborate
 315 on this observation.

316 Recall that the Riemannian gradient of ψ is given by

$$317 \quad \text{grad} \psi(\Theta; x, y) = \mathcal{P}_{T_{\Theta} \mathcal{M}}(\nabla \psi(\Theta; x, y)).$$

318 When $\nabla \psi$ has the form (4.1), the linearity of the projection operator implies that

$$319 \quad (4.2) \quad \begin{aligned} F^R(\theta) &= \mathbb{E}_{P_{x,y}(\theta)} [\text{grad} \psi(\theta; x, y) \text{grad} \psi(\theta; x, y)^\top] \\ &\approx \mathcal{P} (\mathbb{E}_{P_{x,y}(\theta)} [A(x, y) A(x, y)^\top] \otimes \mathbb{E}_{P_{x,y}(\theta)} [G(x, y) G(x, y)^\top]) \mathcal{P}, \end{aligned}$$

320 where $P_{x,y}(\theta)$ is the joint distribution of (x, y) given θ , $\mathcal{P} \in \mathbb{R}^{r \times r}$ is the matrix
 321 representation of $\mathcal{P}_{T_{\Theta} \mathcal{M}}$ (note that $\mathcal{P}^\top = \mathcal{P}$ due to the symmetry of orthogonal
 322 projection operators), and the approximation is due to the assumption that $A(x, y)$
 323 and $G(x, y)$ are approximately independent; see also [23, Theorem 1] for a use of such
 324 an assumption to derive a simplified form of the FIM. By replacing $P_{x,y}(\theta)$ with its

325 empirical distribution observed from \mathcal{S} , an approximate REFIM is given by
 (4.3)

$$326 \quad \bar{F}^R(\theta) \approx \mathcal{P} \left(\left[\frac{1}{|\mathcal{S}|} \sum_{(x,y) \in \mathcal{S}} A(x,y)A(x,y)^\top \right] \otimes \left[\frac{1}{|\mathcal{S}|} \sum_{(x,y) \in \mathcal{S}} G(x,y)G(x,y)^\top \right] \right) \mathcal{P}.$$

327 When a direct inverse of $\bar{F}^R(\theta)$ is expensive to compute, the truncated conjugate
 328 gradient method can be used. In preparation for the applications, we now show how
 329 to construct computationally efficient approximations of the RFIM and REFIM on
 330 the Grassmann manifold.

331 **4.1 RFIM and REFIM on Grassmann manifold** If the matrix representa-
 332 tion \mathcal{P} of the projection operator $\mathcal{P}_{T_\Theta \mathcal{M}}$ has dimensions m -by- m or n -by- n , i.e.,

$$333 \quad \text{grad} \psi(\Theta; x, y) = B_1 G(x, y) A(x, y)^\top \quad \text{or} \quad \text{grad} \psi(\Theta; x, y) = G(x, y) A(x, y)^\top B_2$$

334 with $B_1 \in \mathbb{R}^{m \times m}$ and $B_2 \in \mathbb{R}^{n \times n}$, then we can approximate the RFIM in (4.2) by

$$335 \quad F^R(\theta) \approx \mathbb{E}_{P_{x,y}(\theta)} [A(x, y)A(x, y)^\top] \otimes \mathbb{E}_{P_{x,y}(\theta)} [B_1 G(x, y)G(x, y)^\top B_1]$$

336 or

$$337 \quad F^R(\theta) \approx \mathbb{E}_{P_{x,y}(\theta)} [B_2 A(x, y)A(x, y)^\top B_2] \otimes \mathbb{E}_{P_{x,y}(\theta)} [G(x, y)G(x, y)^\top].$$

338 Moreover, if we replace $P_{x,y}(\theta)$ by its empirical distribution observed from \mathcal{S} , then
 339 we can approximate the REFIM in (4.3) by

$$340 \quad \bar{F}^R(\theta) \approx \left(\frac{1}{|\mathcal{S}|} \sum_{(x,y) \in \mathcal{S}} A(x, y)A(x, y)^\top \right) \otimes \left(\frac{1}{|\mathcal{S}|} \sum_{(x,y) \in \mathcal{S}} B_1 G(x, y)G(x, y)^\top B_1 \right)$$

341 or

$$342 \quad \bar{F}^R(\theta) \approx \left(\frac{1}{|\mathcal{S}|} \sum_{(x,y) \in \mathcal{S}} B_2 A(x, y)A(x, y)^\top B_2 \right) \otimes \left(\frac{1}{|\mathcal{S}|} \sum_{(x,y) \in \mathcal{S}} G(x, y)G(x, y)^\top \right).$$

343 Note that the Kronecker product form allows the inverse of $\bar{F}^R(\theta)$ to be calculated
 344 efficiently by inverting two smaller matrices [39]. A typical manifold that yields the
 345 above Kronecker product representations is the Grassmann manifold $\text{Gr}(m, n)$, which
 346 consists of all n (resp., m) dimensional subspaces in \mathbb{R}^m (resp., \mathbb{R}^n) if $m \geq n$ (resp.,
 347 $m < n$). The matrix representation of the projection operator at a point Θ with
 348 $\Theta^\top \Theta = I$ is $B_1 = I_m - \Theta \Theta^\top$ ($m \geq n$) or $B_2 = I_n - \Theta^\top \Theta$ ($m < n$). In what
 349 follows, we derive the RFIMs associated with three concrete applications involving
 350 the Grassmann manifold and explain how they can be computed efficiently.

351 4.2 Applications

352 **4.2.1 Low-rank matrix completion** For simplicity, we derive the RFIM
 353 associated with problem (1.3) for the fully observed case, i.e., $\Omega = \{1, \dots, n\} \times$
 354 $\{1, \dots, N\}$. One can derive the RFIM for the partly observed case in a similar fashion.
 355 By definition, we have $f(x, U) = Ua(U; x) - x$ and $\psi(U; x, y) = -\log p(y|f(x, U)) =$
 356 $\frac{1}{2} \|f(x, U) - y\|^2 + \frac{n \log(2\pi)}{2}$. It follows from [14, Subsection 3.4] that the Jacobian
 357 of a along a tangent vector $H \in T_U \text{Gr}(n, p)$ is given by $J_a(U; x)[H] = H^\top x$ and its

358 adjoint $J_a^\top(U; x)$ satisfies $J_a^\top(U; x)[v] = x^\top v$ for $v \in \mathbb{R}^p$. The Riemannian gradient of
 359 $\psi(\cdot; x, y)$ is

$$360 \quad \begin{aligned} \text{grad} \psi(U; x, y) &= (I - UU^\top)((Ua(U; x) - x - y)a(U; x)^\top) \\ &\quad + (I - UU^\top)x(Ua(U; x) - x - y)^\top U. \end{aligned}$$

361 By assuming that the residual $Ua(U; x) - x$ is close to zero, we have $(I - UU^\top)x \approx$
 362 $(I - UU^\top)Ua(U; x) = 0$. This leads to the following approximate Riemannian gradient
 363 of $\psi(\cdot; x, y)$:

$$364 \quad (4.4) \quad \text{grad} \psi(U; x, y) \approx (I - UU^\top)((Ua(U; x) - x - y)a(U; x)^\top).$$

365 Plugging the above approximation into (4.2) leads to

$$366 \quad \begin{aligned} F^R(u) &= \mathbb{E}_{P_x} \left[\mathbb{E}_{P_{y|x}(u)} \left[\text{grad} \psi(u; x, y) \text{grad} \psi(u; x, y)^\top \right] \right] \\ &\approx \mathbb{E}_{P_x} \left[\mathbb{E}_{P_{y|x}(U)} \left[[a(U; x)a(U; x)^\top] \otimes [(I - UU^\top)(Ua(U; x) - x - y) \right. \right. \\ &\quad \left. \left. (Ua(U; x) - x - y)^\top (I - UU^\top)] \right] \right] \\ &\approx \left[\frac{1}{N} \sum_{i=1}^N a(U; x_i)a(U; x_i)^\top \right] \otimes (I - UU^\top), \end{aligned}$$

367 where $u = \text{vec}(U)$ is the vectorization of U , the second line is due to (4.4), $\text{vec}(uv^\top) =$
 368 $v \otimes u$, $(A \otimes B)^\top = A^\top \otimes B^\top$, and $(A \otimes B)(A^\top \otimes B^\top) = (AA^\top) \otimes (BB^\top)$, and the
 369 last line follows from $\mathbb{E}_{P_{y|x}(U)} [(Ua(U; x) - x - y)(Ua(U; x) - x - y)^\top] = I$ and by
 370 substituting P_x with its empirical distribution. For $H \in T_U \text{Gr}(n, p)$, we have

$$371 \quad (4.5) \quad \begin{aligned} \text{mat}(F^R(u)[\text{vec}(H)]) &\approx \left[\frac{1}{N} \sum_{i=1}^N a(U; x_i)a(U; x_i)^\top \right] \otimes (I - UU^\top) \text{vec}(H) \\ &= H \left[\frac{1}{N} \sum_{i=1}^N a(U; x_i)a(U; x_i)^\top \right], \end{aligned}$$

372 where $\text{mat}(b)$ converts the vector $b \in \mathbb{R}^{np}$ into an n -by- p matrix and the equality
 373 follows from $(I - UU^\top)H = H$. For the partly observed case, the matrix $F^R(u)$
 374 defined in the above equation can serve as a good approximation of the exact RFIM.
 375 Note that $\frac{1}{N} \sum_{i=1}^N a(U; x_i)a(U; x_i)^\top \in \mathbb{R}^{p \times p}$ is of low dimension since the rank p is
 376 usually small. Thus, the Riemannian natural gradient direction can be calculated
 377 with a relatively low cost.

378 **4.2.2 Low-dimension subspace learning** In multi-task learning [6, 40], dif-
 379 ferent tasks are assumed to share the same latent low-dimensional feature represen-
 380 tation. Specifically, suppose that the i -th task has the training set $X_i \in \mathbb{R}^{d_i \times n}$ and
 381 the corresponding label set $y_i \in \mathbb{R}^{d_i}$ for $i = 1, \dots, N$. The multi-task feature learning
 382 problem can then be formulated as

$$383 \quad (4.6) \quad \min_{U \in \text{Gr}(n, p)} \Psi(U) = \frac{1}{2N} \sum_{i=1}^N \|X_i U w(U; X_i, y_i) - y_i\|^2,$$

384 where $w(U; X_i, y_i) = \arg \min_w \frac{1}{2} \|X_i U w - y_i\|^2 + \lambda \|w\|^2$ and $\lambda > 0$ is a regularization
 385 parameter. Suppose that $d_1 = \dots = d_N = d$. Then, problem (4.6) has the form (1.1),

386 where $\mathcal{S} = \{(X_i, y_i), 0\}_{i=1}^N$, $\mathcal{X} = \mathbb{R}^{d \times (n+1)}$, $\mathcal{Y} = \mathbb{R}^d$, $f(X, y, U) = XUw(U; X, y) - y$,
387 and $p(z|f(X, y, U)) = \frac{1}{\sqrt{(2\pi)^d}} \exp(-\frac{1}{2}(z - f(X, y, U))^\top(z - f(X, y, U)))$. By ignoring
388 the constant $\frac{1}{\sqrt{(2\pi)^d}}$ when computing ψ , we denote $\psi(U; X, y, z) = \frac{1}{2}\|XUw(U; X, y) -$
389 $y - z\|^2$. Using the optimality of $w(U; X, y)$, we have $U^\top X^\top (XU$
390 $w(U; X, y) - y) + \lambda w(U; X, y) = 0$. Then, we can compute the Euclidean gradient of
391 $\psi(\cdot; X, y, z)$ as

$$\begin{aligned} & \nabla\psi(U; X, y, z) \\ 392 &= X^\top (XUw(U; X, y) - y - z)w(U; X, y)^\top + J_w^\top(U) [U^\top X^\top (XUw(U; X, y) - y - z)] \\ & \approx X^\top (XUw(U; X, y) - y)w(U; X, y)^\top, \end{aligned}$$

393 where $J_w(U)$ is the Jacobian of $w(U; X, y)$, $J_w^\top(U)$ denotes the adjoint of $J_w(U)$,
394 and the approximation holds for small λ and $\|z\|$. Note that z will lie in a small
395 neighborhood of zero with high probability if $f(X, y, U)$ is close to 0. Besides, z is
396 always zero in the dataset \mathcal{S} . With the above, an approximate Riemannian gradient
397 of $\psi(\cdot; X, y, z)$ is given by

$$398 \quad (4.7) \quad \text{grad}\psi(U; X, y, z) \approx (I - UU^\top)X^\top (XUw(U; X, y) - y - z)w(U; X, y)^\top.$$

399 Consequently, we have

$$\begin{aligned} F^R(u) &= \mathbb{E}_{P_{(X,y)}} \left[\mathbb{E}_{P_{z|(X,y)}(u)} [\text{grad}\psi(u; X, y, z)\text{grad}\psi(u; X, y, z)^\top] \right] \\ & \approx \frac{1}{N} \sum_{i=1}^N (w_i \otimes ((I - UU^\top)X_i^\top)) (w_i \otimes ((I - UU^\top)X_i^\top))^\top \\ 400 \quad (4.8) &= \frac{1}{N} \sum_{i=1}^N [(w_i w_i^\top) \otimes ((I - UU^\top)X_i^\top X_i (I - UU^\top))] \\ & \approx \frac{1}{N} \left[\sum_{i=1}^N w_i w_i^\top \right] \otimes \left[\frac{1}{N} \sum_{i=1}^N (I - UU^\top)X_i^\top X_i (I - UU^\top) \right], \end{aligned}$$

401 where $u = \text{vec}(U)$ is the vectorization of U , $w_i := w(U; X_i, y_i)$, the second line follows
402 from (4.7), $\mathbb{E}_{P_{z|(X,y)}(u)} [(XUw(U; X, y) - y - z)(XUw(U; X, y) - y - z)^\top] = I$, and the
403 empirical approximation of $P_{(X,y)}$, and the last line holds under the same condition
404 as in (4.2). Though the construction of $F^R(u)$ is for the case $d_1 = \dots = d_N$, it can
405 be easily extended to the case where the d_i 's are not equal.

406 **4.2.3 Fully connected network with batch normalization** Consider an
407 L -layer neural network with input $a_0 = x$. In the l -th layer, we have

$$408 \quad (4.9) \quad s_l = W_l a_{l-1} + b_l, \quad t_{l,i} = \frac{s_{l,i} - \mathbb{E}(s_{l,i})}{\text{Var}(s_{l,i})} \times \gamma_{l,i} + \beta_{l,i}, \quad i = 1, \dots, n_l, \quad a_l = \varphi_l(t_l),$$

409 where φ_l is an element-wise activation function, $W_l \in \mathbb{R}^{n_l \times n_{l-1}}$ is the weight, $b_l \in \mathbb{R}^{n_l}$
410 is the bias, $s_{l,i}$ is the i -th component of $s_l \in \mathbb{R}^{n_l}$, $\gamma_{l,i}, \beta_{l,i} \in \mathbb{R}$ are two learnable
411 parameters, $\text{Var}(s_{l,i})$ is the variance of $s_{l,i}$, and $f(x, \Theta) = a_L \in \mathbb{R}^m$ is the output of
412 the network with Θ being the collection of parameters $\{W_l, b_l, \gamma_l, \beta_l\}$. By default, the
413 elements of $\gamma_{l,i}$ are set to 1 and the elements of $\beta_{l,i}$ are set to 0. In [28], $t_{l,i}$ is called
414 the batch normalization of $s_{l,i}$.

415 Given a dataset \mathcal{S} , our goal is to minimize the discrepancy between the network
 416 output $f(x, \Theta)$ and the observed output y , namely,

$$417 \quad (4.10) \quad \min_{\Theta} \Psi(\Theta) = -\frac{1}{|\mathcal{S}|} \sum_{(x,y) \in \mathcal{S}} \log p(y|f(x, \Theta)).$$

418 By [17], each row of W_l lies on the Grassmann manifold $\text{Gr}(1, n_{l-1})$. It follows that W_l
 419 lies on the product of Grassmann manifolds, i.e., $W_l \in \text{Gr}(1, n_{l-1}) \times \cdots \times \text{Gr}(1, n_{l-1}) \in$
 420 $\mathbb{R}^{n_l \times n_{l-1}}$. The remaining parameters lie in the Euclidean space. Rather than batch
 421 normalization, layer normalization [8] and weight normalization [49] have also been
 422 widely investigated in the study of deep neural networks, where $\text{vec}(W_l) \in \text{Gr}(n_l \times$
 423 $n_{l-1}, 1)$ and $W_l \in \text{Sp}(n_{l-1} - 1) \times \cdots \times \text{Sp}(n_{l-1} - 1) \in \mathbb{R}^{n_l \times n_{l-1}}$ with $\text{Sp}(n_{l-1} - 1) :=$
 424 $\{u \in \mathbb{R}^{n_{l-1}} : \|u\| = 1\}$, respectively.

By back-propagation, the Euclidean gradient of Ψ with respect to W_l is given by

$$g_l \leftarrow Da_l \odot \varphi'_l(t_l) \odot Dt_l, \quad \nabla \Psi(W_l) \leftarrow g_l a_{l-1}^\top, \quad Da_{l-1} \leftarrow W_l^\top g_l.$$

425 In particular, we see that $\nabla \Psi(W_l)$ has the Kronecker product form (4.1). Moreover,
 426 note that $\Psi(w_{l,i}) = \Psi(cw_{l,i})$, $\forall c \neq 0$. Now, we compute

$$427 \quad \nabla \Psi(w_{l,i}) w_{l,i}^\top = \lim_{t \rightarrow 0} \frac{\Psi(w_{l,i} + tw_{l,i}) - \Psi(w_{l,i})}{t} = 0.$$

428 By definition of the projection operator defined on the product of Grassmann man-
 429 ifolds, the Riemannian gradient $\text{grad} \Psi(W_l)$ is actually the same as the Euclidean
 430 gradient $\nabla \Psi(W_l)$. Specifically, for the i -th row of $\text{grad} \Psi(W_l)$, we have

$$431 \quad [\text{grad} \Psi(W_l)]_i = \text{grad} \Psi(w_{l,i}) = \nabla \Psi(w_{l,i}) - \nabla \Psi(w_{l,i}) w_{l,i}^\top w_{l,i} = \nabla \Psi(w_{l,i}).$$

432 Therefore, the RFIM coincides with the FIM. The inverse of $F^R(\theta)$ can be computed
 433 easily when the FIM has a Kronecker product form.

434 **5 Convergence Analysis** In this section, we study the convergence behavior
 435 of the RNGD method (Algorithm 1).

436 **5.1 Global convergence to a stationary point** To begin, let us consider
 437 the case where \mathcal{M} is an embedded submanifold and extend some of the definitions
 438 used in the study of Euclidean stochastic trust-region methods (see, e.g., [16]) to this
 439 setting.

440 **DEFINITION 5.1.** Let $\kappa_{\text{ef}}, \kappa_{\text{eg}} > 0$ be given constants. A function m_k is called a
 441 $(\kappa_{\text{ef}}, \kappa_{\text{eg}})$ -fully linear model of Ψ on $B_{\theta^k}(0, 1/\sigma_k)$ if for any $y \in B_{\theta^k}(0, 1/\sigma_k)$,

$$442 \quad (5.1) \quad \|\nabla(\Psi \circ R_{\theta^k})(y) - \nabla m_k(y)\| \leq \frac{\kappa_{\text{eg}}}{\sigma_k} \quad \text{and} \quad |\Psi \circ R_{\theta^k}(y) - m_k(y)| \leq \frac{\kappa_{\text{ef}}}{\sigma_k^2},$$

443 where $B_\theta(0, \rho) := \{d \in T_\theta \mathcal{M} : \|d\| \leq \rho\}$.

444 **DEFINITION 5.2.** Let $\epsilon_F, \sigma_k > 0$ be given constants. The quantities Ψ_k^0 and $\Psi_k^{z^k}$
 445 are called ϵ_F -accurate estimates of $\Psi(\theta^k)$ and $\Psi_k(z^k)$, respectively if

$$446 \quad (5.2) \quad \left| \Psi_k^0 - \Psi(\theta^k) \right| \leq \frac{\epsilon_F}{\sigma_k^2} \quad \text{and} \quad \left| \Psi_k^{z^k} - \Psi_k(z^k) \right| \leq \frac{\epsilon_F}{\sigma_k^2},$$

447 where z^k is defined in (3.8).

448 Analogous to [16, 58], the inequalities (5.1) and (5.2) can be guaranteed when
 449 \mathcal{M} is compact, the number of samples is large enough, and $\nabla(\Psi \circ R)$ is Lipschitz
 450 continuous.

451 Next, we introduce the assumptions needed for our convergence analysis. Their
 452 Euclidean counterparts can be found in, e.g., [16, Assumptions 4.1 and 4.3].

453 **ASSUMPTION 5.3.** *Let $\theta^0 \in \mathbb{R}^r, \sigma_{\min} > 0$ be given. Let $\mathcal{L}(\theta^0)$ denote the set
 454 of iterates generated by Algorithm 1. Then, the function Ψ is bounded from below
 455 on $\mathcal{L}(\theta^0)$. Moreover, the function $\Psi \circ R$ and its gradient $\nabla(\Psi \circ R)$ are L -Lipschitz
 456 continuous on the set*

$$457 \quad \mathcal{L}_{\text{enl}}(\theta^0) = \bigcup_{\theta \in \mathcal{L}(\theta^0)} B_{\theta} \left(0, \frac{1}{\sigma_{\min}} \right).$$

458 **ASSUMPTION 5.4.** *The RFIM or REFIM F_k satisfies $\|F_k\|_{\text{op}} \leq \kappa_{\text{fim}}$ for all $k \geq 0$,
 459 where $\|\cdot\|_{\text{op}}$ is the operator norm.*

460 We remark that Assumptions 5.3 and 5.4 hold for any compact \mathcal{M} and smooth
 461 Ψ . With the above assumptions, we can prove the convergence of Algorithm 1 by
 462 adapting the arguments in [16]. The main difference is that our analysis makes use of
 463 the pull-back function $\Psi \circ R$ and its Euclidean gradient; see Definitions 5.1 and 5.2.

THEOREM 5.5. *Suppose that Assumptions 5.3 and 5.4 hold, m_k is a $(\kappa_{\text{ef}}, \kappa_{\text{eg}})$ -
 fully linear model for some $\kappa_{\text{ef}}, \kappa_{\text{eg}} > 0$, and the estimates Ψ_k^0 and $\Psi_k^{z^k}$ are ϵ_F -
 accurate for some $\epsilon_F > 0$. Furthermore, suppose that $\eta_2 \geq \max \left\{ \kappa_{\text{fim}}, \frac{16\kappa_{\text{ef}}}{1-\eta_1} \right\}$ and
 $\epsilon_F \leq \min \left\{ \kappa_{\text{ef}}, \frac{1}{32}\eta_1\eta_2 \right\}$. Then, the sequence of iterates $\{\theta^k\}$ generated by Algorithm
 1 will almost surely satisfy*

$$\liminf_{k \rightarrow \infty} \|\text{grad} \Psi(\theta^k)\| = 0.$$

464 *Proof.* One can prove the conclusion by following the arguments in [16, Theorem
 465 4.16]. We here present a sketch of the proof. Define \mathcal{F}_k as the σ -algebra generated by
 466 $\Psi_1^0, \Psi_1^{z^1}, \dots, \Psi_k^0, \Psi_k^{z^k}$ and m_1, \dots, m_k . Consider the random function $\Phi_k = v\Psi(\theta^k) +$
 467 $(1-v)/\sigma_k^2$, where $v \in (0, 1)$ is fixed. The idea is to prove that there exists a constant
 468 $\tau > 0$ such that for all k ,

$$469 \quad (5.3) \quad \mathbb{E}[\Phi_{k+1} - \Phi_k \mid \mathcal{F}_{k-1}] \leq -\frac{\tau}{\sigma_k^2} < 0.$$

470 Summing (5.3) over $k \geq 1$ and taking expectations on both sides lead to $\sum_{k=1}^{\infty} 1/\sigma_k^2 <$
 471 ∞ . The inequality (5.3) can be proved in the following steps. Firstly, a decrease on Ψ
 472 of order $-\mathcal{O}(1/\sigma_k^2)$ can be proved using the fully linear model approximation and the
 473 positive definiteness of $F_k + \sigma_k \|g^k\|I$ with a sufficiently large σ_k . Secondly, the trial
 474 point z^k is accepted provided that the estimates Ψ_k^0 and $\Psi_k^{z^k}$ are ϵ_F -accurate with
 475 sufficiently small ϵ_F and large σ_k . In addition, with $\eta_2 \geq \max \left\{ \kappa_{\text{fim}}, \frac{16\kappa_{\text{ef}}}{1-\eta_1} \right\}$, if z^k is
 476 accepted (i.e., $\theta^{k+1} = z^k$), then a decrease of $-\mathcal{O}(1/\sigma_k^2)$ on Ψ can always be guaranteed
 477 when $\epsilon_F \leq \min \left\{ \kappa_{\text{ef}}, \frac{1}{32}\eta_1\eta_2 \right\}$ based on the update scheme (3.11). On the other hand,
 478 if z^k is rejected (i.e., $\theta^{k+1} = \theta^k$), then $\mathbb{E}[\Phi_{k+1} - \Phi_k \mid \mathcal{F}_{k-1}] = (1-v)(1/\gamma^2 - 1)/\sigma_k^2$.
 479 By choosing v to be sufficiently close to 1, the inequality (5.3) holds for any k .

480 Now, we will have $\sigma_k \rightarrow \infty$ as $k \rightarrow \infty$ with probability 1. If there exist $\epsilon > 0$
 481 and $k_0 \geq 1$ such that $\|\text{grad} \Psi(\theta^k)\| \geq \epsilon$ for all $k \geq k_0$, then the trial point will be
 482 accepted eventually because the estimates Ψ_k^0 and $\Psi_k^{z^k}$ are ϵ_F -accurate. Recall that

483 σ_k is decreasing in the case of accepting z^k . This means that σ_k is bounded above,
 484 which leads to a contradiction. Hence, we conclude that $\liminf_{k \rightarrow \infty} \|\text{grad } \Psi(\theta^k)\| = 0$
 485 will hold almost surely. \square

486 **REMARK 5.6.** *Analogous to [16, Theorem 4.18], one can show that $\lim_{k \rightarrow \infty} \|\text{grad } \Psi(\theta^k)\|$
 487 $= 0$ will hold almost surely by assuming the Lipschitz continuity of $\text{grad } \Psi$.*

488 **REMARK 5.7.** *For the case where \mathcal{M} is a quotient manifold, we modify Algorithm
 489 1 according to the approach mentioned in the last paragraph of Section 3.2. The iterate
 490 θ^k and the tangent space at θ^k should be understood as a representative element and the
 491 horizontal space at θ^k , respectively. Due to the horizontal invariance of the Riemannian
 492 metric, the almost sure convergence result of $\liminf_{k \rightarrow \infty} \|D(\theta^k)^{-1} \text{grad } \Psi(\theta^k)\|_{\theta^k}$
 493 $\rightarrow 0$ also holds.*

494 **5.2 Convergence rate analysis of RNGD** In this subsection, we study the
 495 local convergence rate of a deterministic version of the RNGD method. Let us start
 496 with some definitions. Let

$$497 \quad L(z, y) := -\log p(y|z)$$

498 and suppose that P_x is the empirical distribution defined by \mathcal{S} . We define $\mathcal{S}_x := \{x :$
 499 $(x, y) \in \mathcal{S}\}$, $\mathcal{S}_y := \{y : (x, y) \in \mathcal{S}\}$, $F_L(x, \theta) := \mathbb{E}_{P_{y|x}(\theta)}[\nabla_z \log p(y|z) \nabla_z \log p(y|z)^\top]_{z=f(x, \theta)}$,
 500 and write $J^R(x, \theta) := [\text{grad } f_1(x, \theta), \dots, \text{grad } f_q(x, \theta)]^\top$ for the Riemannian Jaco-
 501 bian of $f(x, \theta) = [f_1(x, \theta), \dots, f_q(x, \theta)]^\top$ with respect to θ . Furthermore, we write
 502 $\mathcal{S} = \{(x_i, y_i)\}_{i=1}^N$ with $N = |\mathcal{S}|$ and $u(\theta) = [f(x_1, \theta), \dots, f(x_N, \theta)]^\top$. Let $J^R(\theta) :=$
 503 $[J^R(x_1, \theta), \dots, J^R(x_N, \theta)]$ and $H_L(u(\theta)) := \text{blkdiag}(H_L(u(\theta)_1), \dots, H_L(u(\theta)_N))$. For
 504 simplicity, we write $u^k := u(\theta^k)$.

505 **5.2.1 Convergence rate** Throughout this subsection, we make the following
 506 assumptions on the loss function L .

507 **ASSUMPTION 5.8.** *For any $y \in \mathcal{S}_y$, the loss function $L(\cdot, y)$ is smooth and μ -
 508 strongly convex and has κ_L -Lipschitz gradient and κ_H -Lipschitz Hessian, namely,*

$$509 \quad \mu I \preceq \nabla_{zz}^2 L(z, y) \preceq \kappa_L I, \quad \|\nabla_{zz}^2 L(z, y) - \nabla_{zz}^2 L(x, y)\| \leq \kappa_H \|z - x\|, \quad \forall z, x \in \mathbb{R}^n.$$

510 *In addition, the following condition holds:*

$$511 \quad (5.4) \quad F_L(x, \theta) = \nabla_{zz}^2 L(z, y)|_{z=f(x, \theta)} := H_L(f(x, \theta)).$$

512 We remark that the equality (5.4) holds if $\nabla_{zz}^2 L(z, y)|_{z=f(x, \theta)}$ does not depend on
 513 y , which is the case for the square loss $L(z, y) = \|z - y\|^2$ and the cross-entropy loss
 514 $L(y, z) = -\sum_j y_j \log z_j + \log(\sum_j \exp(z_j))$. We refer the reader to [38, Section 9.2]
 515 for other loss functions that satisfy (5.4). We remark that the square loss $L(z, y) =$
 516 $\|z - y\|^2$, which appears in both the LRMC and low-dimension subspace learning
 517 problems, satisfies Assumption 5.8.

518 According to the definition of RFIM in (3.1) and the chain rule, we obtain

$$519 \quad F^R(\theta) = \frac{1}{|\mathcal{S}_x|} \sum_{x \in \mathcal{S}_x} J^R(x, \theta)^\top F_L(x, \theta) J^R(x, \theta).$$

520 Based on Assumption 5.8, we have $F^R(\theta) = J^R(\theta)^\top H_L(u(\theta)) J^R(\theta)$. Note that $F^R(\theta)$
 521 may be singular when $J^R(\theta)$ is not of full column rank. In this case, provided that
 522 $(J^R(\theta^k) J^R(\theta^k)^\top)^{-1}$ exists, we can use the pseudo-inverse

$$523 \quad F^R(\theta^k)^\dagger = J^R(\theta^k)^\top (J^R(\theta^k) J^R(\theta^k)^\top)^{-1} H_L(u^k)^{-1} (J^R(\theta^k) J^R(\theta^k)^\top)^{-1} J^R(\theta^k)$$

524 for computation. As mentioned at the beginning of this subsection, we focus on a
 525 deterministic version of the RNGD method, in which we adopt a fixed step size $t > 0$
 526 and perform the update

$$527 \quad (5.5) \quad d^k = (F^R(\theta^k))^\dagger J^R(\theta^k)^\top \nabla L(u^k, y), \quad \theta^{k+1} = R_{\theta^k}(-td^k).$$

528 For concreteness, let us take R to be the exponential map for \mathcal{M} in our subsequent
 529 development. Our convergence rate analysis of this deterministic RNGD method can
 530 be divided into two steps. The first step is to prove that the iterates $\{\theta^k\}$ always stay
 531 in a neighborhood of θ^0 if J^R satisfies certain stability condition. The second step
 532 is to establish the convergence rate of the method by utilizing the strong convexity
 533 of L . We remark that the zero acceleration property of the exponential map [1,
 534 Equation (5.24)] is essential to our analysis. As such, we can only handle the case
 535 where the retraction is the exponential map. The analysis for the case of a more
 536 general retraction is left as an open problem. Motivated by [66], we now formulate
 537 the aforementioned stability condition on J^R .

538 **ASSUMPTION 5.9.** *For any θ satisfying $\|\theta - \theta^0\| \leq 4\kappa_L(\mu\sigma_0)^{-1}\|u^0 - y\|$, where
 539 $\sigma_0 := \sqrt{\lambda_{\min}(J^R(\theta^0)J^R(\theta^0)^\top)} > 0$, it holds that*

$$540 \quad (5.6) \quad \|J^R(\theta) - J^R(\theta^0)\| \leq \min \left\{ \frac{1}{2}, \frac{\mu}{6\kappa_L} \right\} \sigma_0.$$

541 As will be seen in Section 5.2.2, Assumption 5.9 is satisfied by the Riemannian
 542 Jacobian that arises in a two-layer fully connected neural network with batch nor-
 543 malization and sufficiently large width. We are now ready to prove the following
 544 theorem.

545 **THEOREM 5.10.** *Let R be the exponential map for \mathcal{M} . Suppose that Assumptions
 546 5.8 and 5.9 hold. Let $\{\theta^k\}$ be the iterates generated by (5.5).*

547 (a) *There exists a constant $\kappa_R > 0$ such that if $\|u^0 - y\| < \frac{\mu}{3\kappa_H}$ and $t \leq$*

$$548 \quad \min \left\{ 1, \left(\frac{1}{6\|u^0 - y\|} - \frac{\kappa_H}{2\mu} \right) \cdot \frac{3\mu^2\sigma_0}{8\kappa_R\kappa_L^2} \right\}, \text{ then}$$

$$549 \quad (5.7) \quad \|u^{k+1} - y\| \leq \left(1 - \frac{t}{2} \right) \|u^k - y\|.$$

550 (b) *Suppose further that J^R is κ_J -Lipschitz continuous with respect to θ , i.e.,*

$$551 \quad (5.8) \quad \|J^R(\theta) - J^R(\nu)\| \leq \kappa_J \|\theta - \nu\|, \quad \forall \theta, \nu \in \mathbb{R}^r.$$

552 *The rate of convergence is quadratic when $t = 1$, namely, there is a constant
 553 $\kappa_q > 0$ such that*

$$554 \quad (5.9) \quad \|u^{k+1} - y\| \leq \kappa_q \|u^k - y\|^2.$$

555 *Proof.* (a). We proceed by induction. Assume that for $j \leq k$, we have

$$556 \quad \|\theta^j - \theta^0\| \leq 4\kappa_L(\mu\sigma_0)^{-1}\|u^0 - y\|, \quad \|u^j - y\| \leq \left(1 - \frac{\eta}{2} \right) \|u^{j-1} - y\|.$$

557 By the definition of d^k in (5.5),

$$558 \quad (5.10) \quad \begin{aligned} \|d^k\| &\leq \|J^R(\theta^k)^\top (J^R(\theta^k)J^R(\theta^k)^\top)^{-1}\| \|H_L(\theta^k)^{-1}\| \|\nabla_u L(u^k, y) - \nabla_u L(y, y)\| \\ &\leq \mu^{-1} \kappa_L \sigma_{\min}^{-1}(J^R(\theta^k)) \|u^k - y\| \\ &\leq 2\kappa_L(\mu\sigma_0)^{-1} \|u^k - y\|, \end{aligned}$$

559 where the first inequality is due to $\nabla L(y, y) = 0$ and the last inequality is from
 560 Assumption 5.9. Now, define the map $c_k : [0, 1] \rightarrow \mathcal{M}$ as $c_k(s) = R_{\theta^k}(-std^k)$. Note
 561 that for the exponential map R , the geodesic distance between θ and $R_\theta(\xi)$ is equal
 562 to $\|\xi\|$ [1, Equation (7.25)], and inequality (2.2) holds with $\alpha = 1$ when we take the
 563 Euclidean metric as the Riemannian metric on \mathcal{M} . Thus, for any $s \in [0, 1]$,

$$\begin{aligned} \|c_k(s) - \theta^0\| &\leq \|c_k(s) - \theta^k\| + \sum_{j=0}^{k-1} \|\theta^{j+1} - \theta^j\| \leq t \sum_{j=0}^k \|d^j\| \\ &\leq 2\kappa_L(\mu\sigma_0)^{-1} t \sum_{j=0}^k \|u^j - y\|. \end{aligned}$$

564

565 Since $\|u^j - y\| \leq (1 - \frac{\eta}{2})\|u^{j-1} - y\|$ for all $j \leq k$, we have $\|c_k(s) - \theta^0\| \leq 4\kappa_L(\mu\sigma_0)^{-1}\|u^0 - y\|$
 566 $\|y\|$ for all $s \in (0, 1]$. This gives $\|\theta^{k+1} - \theta^0\| \leq 4\mu\kappa_L\sigma_0^{-1}\|u^0 - y\|$. To prove (5.7), we
 567 split $\|u^{k+1} - y\|$ into three terms, namely,

$$\begin{aligned} u^{k+1} - y &= u^{k+1} - u^k + u^k - y = \int_0^1 J^R(c_k(s))c'_k(s)ds + u^k - y \\ &= \underbrace{\int_0^1 J^R(c_k(s))(c'_k(s) - td^k)ds}_{b_1} + \underbrace{t \int_0^1 (J^R(c_k(s)) - J^R(\theta^k))d^k ds}_{b_2} \\ &\quad + \underbrace{t \int_0^1 J^R(\theta^k)d^k ds}_{b_3} + u^k - y. \end{aligned} \tag{5.11}$$

569 For the exponential map R [1, Equation (5.24)], it holds that

$$570 \quad (5.12) \quad c'_k(s) - td^k = c''_k(s)[-std^k] + \tilde{\kappa}_R s^2 t^2 \|d^k\|^2,$$

571 where $c''_k(s)[-std^k]$ belongs to the normal space to \mathcal{M} at $c_k(s)$ and $\tilde{\kappa}_R > 0$ is the
 572 smoothness constant. Plugging (5.12) into (5.11), we have

$$\begin{aligned} \|b_1\| &\leq \int_0^1 (\|J^R(\theta^0)\| + \|J^R(c_k(s)) - J^R(\theta^0)\|) \tilde{\kappa}_R s^2 t^2 \|d^k\|^2 ds \\ &\leq \int_0^1 2\sigma_0 \kappa_R s^2 t^2 \|d^k\|^2 ds = \frac{2}{3} \sigma_0 \kappa_R t^2 \|d^k\|^2, \end{aligned}$$

573

574 where $\kappa_R := \tilde{\kappa}_R \cdot (1/4 + \|J^R(\theta^0)\|/(2\sigma_0))$. By (5.6) and (5.10), we have

$$575 \quad \|b_2\| \leq t \int_0^1 \min \left\{ \frac{1}{2}, \frac{\mu}{6\kappa_L} \right\} \sigma_0 \cdot 2\kappa_L(\mu\sigma_0)^{-1} \|u^k - y\| ds \leq \frac{t}{3} \|u^k - y\|.$$

576 Now, the update (5.5) yields $J^R(u^k)d^k = H_L(u^k)^{-1}\nabla L(u^k, y)$. It follows that

$$\begin{aligned}
\|b_3\| &= \|u^k - y - tH_L(u^k)^{-1}(\nabla L(u^k, y) - \nabla L(y, y))\| \\
&= \|H_L(u^k)^{-1}(H_L(u^k)(u^k - y) - t(\nabla L(u^k, y) - \nabla L(y, y)))\| \\
&= \left\| H_L(u^k)^{-1} \left(H_L(u^k)(u^k - y) - t \int_0^1 H_L(u^k + s(y - u^k))(u^k - y) \, ds \right) \right\| \\
577 &= \left\| H_L(u^k)^{-1} \left[\int_0^1 (H_L(u^k) - tH_L(u^k + s(y - u^k))) \, ds \right] (u^k - y) \right\| \\
&\leq \int_0^1 (1 - t + t\mu^{-1}\kappa_H s \|u^k - y\|) \, ds \cdot \|u^k - y\| \\
&= \left(1 - t + \frac{\kappa_H t}{2\mu} \|u^k - y\| \right) \|u^k - y\|,
\end{aligned}$$

578 where the first inequality is due to Assumption 5.8. Combining the estimates on
579 b_1, b_2, b_3 , we conclude that

(5.13)

$$\begin{aligned}
\|u^{k+1} - y\| &\leq \left(1 - \frac{2t}{3} + \frac{\kappa_H t}{2\mu} \|u^k - y\| \right) \|u^k - y\| + \frac{8}{3}\mu^{-2}\kappa_R\kappa_L^2\sigma_0^{-1}t^2\|u^k - y\|^2 \\
580 &\leq \left(1 - \frac{t}{2} \right) \|u^k - y\|
\end{aligned}$$

581 whenever $\|u^k - y\| < \frac{\mu}{3\kappa_H}$ and $t \leq \left(\frac{1}{6\|u^k - y\|} - \frac{\kappa_H}{2\mu} \right) \cdot \frac{3\mu^2\sigma_0}{8\kappa_R\kappa_L^2}$. Therefore, the inequality
582 (5.7) holds by using the inductive hypothesis $\|u^k - y\| \leq \|u^0 - y\|$.

583 (b). The proof is similar to that for (a). Substituting $t = 1$ into (5.11), we obtain

$$\begin{aligned}
\|u^{k+1} - y\| &\leq \frac{\kappa_H}{2\mu}\|u^k - y\|^2 + \frac{1}{2}\kappa_J\|d^k\|^2 + \frac{8}{3}\mu^{-2}\kappa_R\kappa_L^2\sigma_0^{-1}\|u^k - y\|^2 \\
584 &\leq \left[\frac{\kappa_H}{2\mu} + 2\kappa_L^2(\mu\sigma_0)^{-2} \left(\kappa_J + \frac{4}{3}\sigma_0\kappa_R \right) \right] \|u^k - y\|^2,
\end{aligned}$$

585 where we use (5.8) to get

$$\begin{aligned}
&\left\| \int_0^1 (J^R(c_k(s)) - J^R(\theta^k))d^k \, ds \right\| \\
586 &\leq \kappa_J \int_0^1 \|c_k(s) - \theta^k\| \|d^k\| \, ds \leq \frac{1}{2}\kappa_J\|d^k\|^2 \leq 2\kappa_J\kappa_L^2(\mu\sigma_0)^{-2}\|u^k - y\|^2.
\end{aligned}$$

587 The verification of the neighborhood condition for θ^k is similar to that in (a). This
588 completes the proof. \square

589 5.2.2 Jacobian stability of two-layer neural network with batch nor- 590 malization

591 **Problem setting** From the previous subsection, we see that the Jacobian stability
592 condition in Assumption 5.9 plays an important role in the convergence rate analysis
593 of the RNGD method. Let us now show that such a condition is satisfied by a two-layer
594 neural network with batch normalization, thereby demonstrating its relevance. The
595 difference between our setting and that of [66] lies in the use of batch normalization.

596 To begin, consider the input-output map f given by

$$597 \quad (5.14) \quad f(x, \theta, a) = \frac{1}{\sqrt{m}} \sum_{j=1}^m a_j \phi \left(\frac{\theta_j^\top (x - \mathbb{E}[x])}{\sqrt{\theta_j^\top V \theta_j}} \right),$$

598 where $x \in \mathbb{R}^n$ is the (random) input vector, $V = \mathbb{E}[(x - \mathbb{E}[x])(x - \mathbb{E}[x])^\top]$ is the
599 covariance matrix, $\theta = [\theta_1^\top, \theta_2^\top, \dots, \theta_m^\top]^\top \in \mathbb{R}^{mn}$ is the weight vector of the first layer,
600 $a_j \in \mathbb{R}$ is the output weight of hidden unit j , and ϕ is the ReLU activation function.
601 This represents a single-output two-layer neural network with batch normalization.
602 We fix the a_j 's throughout as in [66] and apply the RNGD method with a fixed
603 step size on θ , in which each weight vector θ_j is assumed to be normalized. For the
604 Grassmann manifold $\text{Gr}(1, n)$, we choose d with $\|d\| = 1$ as the representative element
605 of the one-dimensional subspace $\{cd : c \neq 0\}$. With a slight abuse of notation, we
606 write $\text{Gr}(1, n) := \{d \in \mathbb{R}^n : \|d\| = 1\}$. Then, we can regard the vector θ as lying on a
607 Cartesian product of $\text{Gr}(1, n)$'s.

608 **Jacobian stability** It is well known that if θ_j is a standard Gaussian random
609 vector, then the random vector $\theta_j/\|\theta_j\|$ is uniformly distributed on $\text{Gr}(1, n)$. We draw
610 each θ_j uniformly from $\text{Gr}(1, n)$ and each a_j uniformly from $\{-1, +1\}$. As mentioned
611 in Section 4.2.3, we have $J^R(\theta) = J(\theta)$. Thus, our goal now is to establish the stability
612 of J . To begin, let $\mathcal{S} = \{(x_i, y_i)\}_{i=1}^N$ denote the dataset and $u(\theta) = [f(x_1, \theta, a),$
613 $f(x_2, \theta, a), \dots, f(x_N, \theta, a)]^\top$ denote the output vector. Following [18, 59, 66], we make
614 the following assumption on \mathcal{S} .

615 ASSUMPTION 5.11. For any $(x, y) \in \mathcal{S}$, it holds that $\|x\| = 1$ and $|y| = \mathcal{O}(1)$.
616 For any $x_i, x_j \in \mathcal{S}_x$ with $i \neq j$, it holds that $x_i \neq \pm x_j$. In addition, the input vector
617 x satisfies $\mathbb{E}[x] = \mathbf{0}$ and the covariance matrix $V = \mathbb{E}[xx^\top]$ is positive definite with
618 minimum eigenvalue $\sigma_V > 0$.

619 The above assumptions on the dataset \mathcal{S} are mild as explained in [66, Assump-
620 tion 1]. The positive-definite property of the variance V is used to ensure the
621 well-posedness of the input-output map (5.14). If V is just positive semidefinite,
622 one can replace it by the shift matrix $V + \sigma_V I$ in (5.14) and remove the assump-
623 tion on V . Motivated by the shift matrix [66], we use $[x_i^\top \theta_j^0]_{k-}$ to represent the k -th smallest en-
624 try of $[x_i^\top \theta_1^0, x_i^\top \theta_2^0, \dots, x_i^\top \theta_m^0]$ in absolute value. Since V is positive definite and
625 $\text{Gr}(1, n) = \{d \in \mathbb{R}^n : \|d\| = 1\}$ is compact, for $i = 1, \dots, N$, the function $u \mapsto$
626 $\varphi_i(u) = \frac{x_i}{\sqrt{u^\top V u}} - \frac{V u u^\top x_i}{(u^\top V u)^{3/2}}$ is L -Lipschitz on $\text{Gr}(1, n)$ for some constant $L > 0$, i.e.,
627 $\|\varphi_i(u) - \varphi_i(v)\| \leq L\|u - v\|$ for any $u, v \in \text{Gr}(1, n)$. To prove the desired Jacobian
628 stability result, we need the following lemmas. They extend those in [66], which are
629 developed for the Euclidean setting, to the Grassmann manifold setting. In what
630 follows, we use δ_A to denote the indicator function of an event A , i.e., δ_A takes the
631 value 1 if the event A happens and 0 otherwise.

632 LEMMA 5.12. Let $\theta_j, \theta_j^0 \in \text{Gr}(1, n)$, where $j = 1, \dots, m$, be given. Suppose that
633 for some $k \in \{1, \dots, m\}$, we have $\|\theta - \theta^0\| \leq \sqrt{k} [x_i^\top \theta_j^0]_{k-}$ for $i = 1, 2, \dots, N$ and
634 $j = 1, 2, \dots, m$. Then, we have

$$635 \quad (5.15) \quad \|J(\theta) - J(\theta^0)\|^2 \leq \frac{2NkM + NkL}{m},$$

636 where $M = \max_{i \in \{1, \dots, N\}} \left(\max_{u \in \text{Gr}(1, n)} \left\| \frac{x_i}{\sqrt{u^\top V u}} - \frac{V u u^\top x_i}{(u^\top V u)^{3/2}} \right\|^2 \right)$.

637 *Proof.* Let $A_{i,j}$ denote the event that the signs of $x_i^\top \theta_j$ and $x_i^\top \theta_j^0$ are different.
638 We claim that, for $i = 1, 2, \dots, N$, there are at most $2k$ non-zero entries of $\{\delta_{A_{i,j}}\}_{j=1}^m$.
639 Otherwise, there exists an $i \in \{1, \dots, N\}$ such that

$$\begin{aligned}
640 \quad & \|\theta - \theta^0\|^2 \geq \sum_{j=1}^m |x_i^\top \theta_j - x_i^\top \theta_j^0|^2 \\
641 \quad & \geq \sum_{j \in \{j: \delta_{A_{i,j}}=1\}} |x_i^\top \theta_j - x_i^\top \theta_j^0|^2 \geq \sum_{j \in \{j: \delta_{A_{i,j}}=1\}} |x_i^\top \theta_j^0|^2 > k[x_i^\top \theta_j^0]_{k-}^2,
\end{aligned}$$

which contradicts our assumption. Now, the generalized Jacobian of f with respect to θ is given by

$$J(\theta) = \frac{1}{\sqrt{m}} \sum_{j=1}^m \sum_{i=1}^N a_j \left[\delta_{x_i^\top \theta_1 \geq 0} \cdot \varphi_i(\theta_1)^\top, \dots, \delta_{x_i^\top \theta_m \geq 0} \cdot \varphi_i(\theta_m)^\top \right].$$

642 When $x_i^\top \theta_j$ and $x_i^\top \theta_j^0$ have the same sign, the difference $\delta_{x_i^\top \theta_j \geq 0} \cdot \frac{a_j}{\sqrt{m}} \varphi_i(\theta_j) - \delta_{x_i^\top \theta_j^0 \geq 0} \cdot$
643 $\frac{a_j}{\sqrt{m}} \varphi_i(\theta_j^0)$ is either $\mathbf{0}$ or $\frac{a_j}{\sqrt{m}} (\varphi_i(\theta_j) - \varphi_i(\theta_j^0))$. Splitting $\|J(\theta) - J(\theta^0)\|^2$ into two parts
644 according to the event $A_{i,j}$ yields

$$\begin{aligned}
645 \quad & \|J(\theta) - J(\theta^0)\|^2 \\
646 \quad & \leq \frac{M}{m} \sum_{(x_i, y_i) \in \mathcal{S}} \sum_{j=1}^m \delta_{A_{i,j}} + \frac{L}{m} \sum_{(x_i, y_i) \in \mathcal{S}} \sum_{j=1}^m \|\theta_j - \theta_j^0\|^2 \\
647 \quad & \leq \frac{2NkM}{m} + \frac{L}{m} \sum_{(x_i, y_i) \in \mathcal{S}} \|\theta - \theta^0\|^2 \\
648 \quad & \leq \frac{2NkM + NkL}{m},
\end{aligned}$$

649 where the last inequality follows from the assumption on $\|\theta - \theta^0\|$ and the fact that
650 $|[x_i^\top \theta_j^0]_{k-}| \leq 1$ for $i = 1, \dots, N$ and $j = 1, \dots, m$. \square

651 The next lemma gives an upper bound on the probability of the event $\{|x_i^\top \theta_j| \leq \gamma\}$
652 for all $\gamma > 0$, which will be used to estimate $[x_i^\top \theta_j^0]_{k-}$ in Lemma 5.14.

653 **LEMMA 5.13.** *Let v be uniformly distributed on $\text{Gr}(1, n)$, $x \in \text{Gr}(1, n)$ be a given*
654 *unit-norm vector, and $\gamma > 0$ be a given positive number, where $n \geq 2$. Then, we have*
655 $\mathbb{P}(|x^\top v| \leq \gamma) \leq \sqrt{\pi n} \gamma$. *Moreover, the dependence on n in the bound is optimal up to*
656 *constant factors.*

657 *Proof.* Without loss of generality, we may assume that $x = (1, 0, \dots, 0)$ since
658 the Euclidean inner product and the distribution of v are invariant under orthogo-
659 nal transformation. Then, we have $x^\top v = v_1$. Let Z_1, \dots, Z_n be standard Gauss-
660 ian random variables. Then, the random variable $x^\top v$ has the same distribution as
661 $B := \frac{Z_1}{\sqrt{Z_1^2 + \dots + Z_n^2}}$. It is well known that B^2 follows the distribution $\text{Beta}(\frac{1}{2}, \frac{n-1}{2})$ [30,
662 Section 25.2]. As a result, the density function h of B can be explicitly written as

$$663 \quad (5.16) \quad h(r) = \frac{\Gamma(\frac{n}{2})}{\sqrt{\pi} \Gamma(\frac{n-1}{2})} (1 - r^2)^{\frac{n-3}{2}}, \quad |r| < 1.$$

664 It follows directly that

$$665 \quad (5.17) \quad \mathbb{P}(|x^\top v| \leq \gamma) = \mathbb{P}(|B| \leq \gamma) = \int_{-\gamma}^{\gamma} h(r) dr \leq \frac{\gamma \Gamma(\frac{n}{2})}{\sqrt{\pi} \Gamma(\frac{n-1}{2})} \leq \sqrt{\pi n} \gamma,$$

666 where the last step uses the classic result $\Gamma(\frac{n}{2}) \leq \pi\sqrt{n}\Gamma(\frac{n-1}{2})$ in calculus.

667 To see the optimality of the dependence on n in the bound, note that for $\gamma \leq \frac{1}{\sqrt{n}}$,
 668 we have

$$669 \quad \mathbb{P}(|x^\top v| \leq \gamma) = \mathbb{P}(|B| \leq \gamma) = \int_{-\gamma}^{\gamma} h(r) dr \geq \frac{\gamma\Gamma(\frac{n}{2})}{2\sqrt{\pi}\Gamma(\frac{n-1}{2})} \geq \frac{5}{12\sqrt{2e}}\sqrt{n},$$

670 where the third step uses $(1 - r^2)^{\frac{n-3}{2}} \geq 1 - \frac{n-3}{2}r^2$ and the fact that $\gamma \leq \frac{1}{\sqrt{n}}$, and
 671 the last step follows from an application of Stirling's formula; see, e.g., [56, Eq. (33)].
 672 Hence, the dependence on n in the bound is optimal up to constant factors. \square

673 Using the above lemmas, we show that Assumption 5.9 will hold with high prob-
 674 ability.

675 **LEMMA 5.14.** *Let $\theta_j, \theta_j^0 \in \text{Gr}(1, n)$, where $j = 1, \dots, m$, be given. For any given*
 676 *$Q, \epsilon > 0$, if $\|\theta - \theta^0\| \leq Q$, then with probability at least $1 - \epsilon$, we will have*

$$677 \quad (5.18) \quad \|J(\theta) - J(\theta^0)\|^2 \leq \frac{2(\pi n)^{\frac{1}{3}} N^{\frac{5}{3}} M Q^{\frac{2}{3}}}{\epsilon^{\frac{2}{3}} m^{\frac{1}{3}}} + \frac{(\pi n)^{\frac{1}{3}} N^{\frac{5}{3}} L Q^{\frac{2}{3}}}{\epsilon^{\frac{2}{3}} m^{\frac{1}{3}}}.$$

678 *Proof.* For given integers $k \in \{1, \dots, m\}$ and $i \in \{1, 2, \dots, N\}$, we prove that with
 679 probability at least $1 - \epsilon/N$, there will be at most $k - 1$ hidden units θ_j^0 such that
 680 $|x_i^\top \theta_j^0| \leq \frac{k\epsilon}{Nm\sqrt{\pi n}}$. For $\tau > 0$, let γ_τ be the positive number such that $\mathbb{P}(|g| \leq \gamma_\tau) = \tau$,
 681 where g follows the same distribution as $x_i^\top \theta_j^0$. It follows from Lemma 5.13 that
 682 $\gamma_\tau \geq \frac{1}{\sqrt{\pi n}}\tau$. Let $\tau = \frac{k\epsilon}{Nm}$. Then, we have

$$683 \quad (5.19) \quad \mathbb{E} \left[\sum_{j=1}^m \delta_{|x_i^\top \theta_j^0| \leq \gamma_\tau} \right] = \sum_{j=1}^m \mathbb{P} [|x_i^\top \theta_j^0| \leq \gamma_\tau] \leq \frac{k\epsilon}{N}.$$

684 Applying the Markov inequality yields

$$685 \quad (5.20) \quad \mathbb{P} \left[\sum_{j=1}^m \delta_{|x_i^\top \theta_j^0| \leq \gamma_\tau} \geq k \right] \leq \frac{\epsilon}{N}.$$

686 Therefore, by taking $k = \frac{Q^{\frac{2}{3}} m^{\frac{2}{3}} (\pi n)^{\frac{1}{3}} N^{\frac{2}{3}}}{\epsilon^{\frac{2}{3}}}$, the inequalities $\sqrt{k}[x_i^\top \theta^0]_{k-} \geq \frac{k^{\frac{3}{2}} \epsilon}{Nm\sqrt{\pi n}} = Q$
 687 will hold simultaneously for $i = 1, \dots, N$ with probability at least $1 - \epsilon$. The desired
 688 conclusion then follows from Lemma 5.12. \square

689 **Linear convergence of RNGD** With the help of Lemma 5.14, we are now ready
 690 to establish the convergence rate of the RNGD method when applied to the two-layer
 691 neural network with batch normalization.

THEOREM 5.15. *Suppose that Assumptions 5.8 and 5.11 hold. Let $\epsilon > 0$ be a
 given constant. Suppose that the number m of hidden units satisfies*

$$m = \Omega \left(\frac{128(L + 2M)^3 \pi n N^6 \kappa_L^2}{\mu^2 \sigma_0^8 \sigma_V \epsilon^3 \min \left\{ \frac{1}{2}, \frac{\mu}{6\kappa_L} \right\}^6} \right),$$

692 where the constants $L, M, \kappa_L, \mu, \sigma_0, \sigma_V$ are defined previously. If we draw θ_j^0 uni-
 693 formly from $\text{Gr}(1, n)$ and a_j uniformly from $\{-1, +1\}$ for $j = 1, 2, \dots, m$, then the

694 *Riemannian Jacobian stability condition in Assumption 5.9 will hold with probability*
 695 *at least $1 - \epsilon$. Furthermore, when $m \geq \frac{16(L+2M)^3 \pi n N^5 \kappa_L^2}{9\sigma_0^8 \kappa_H^2 \epsilon^2 \min\{\frac{1}{2}, \frac{\mu}{6\kappa_L}\}^6}$, $\|u^0 - y\| \leq \frac{\mu}{3\kappa_H}$, and*

696 $\eta \leq \min \left\{ 1, \left(\frac{1}{6\|u^0 - y\|} - \frac{\kappa_H}{2\mu} \right) \cdot \frac{3\mu^2 \sigma_0}{8\kappa_R \kappa_L^2} \right\}$, *with probability at least $1 - \epsilon$, we will have*

$$697 \quad (5.21) \quad \|u^{k+1} - y\| \leq \left(1 - \frac{1}{2}\eta \right) \|u^k - y\|.$$

698 *Proof.* By Assumption 5.11 and the fact that a_j is drawn uniformly from $\{-1, +1\}$,
 699 we have $\mathbb{E}[u^0] = \mathbf{0}$ and

$$700 \quad \begin{aligned} \mathbb{E}[(u_j^0)^2] &= \mathbb{E} \left[\frac{1}{m} \left(\sum_{j=1}^m a_j \phi \left(\frac{(\theta_j^0)^\top (x - \mathbb{E}[x])}{\sqrt{(\theta_j^0)^\top V \theta_j^0}} \right) \right)^2 \right] \\ &= \mathbb{E} \left[\frac{1}{m} \sum_{j=1}^m \phi \left(\frac{(\theta_j^0)^\top x}{\sqrt{(\theta_j^0)^\top V \theta_j^0}} \right)^2 \right] = \mathcal{O} \left(\frac{1}{\sigma_V} \right), \quad j = 1, \dots, N. \end{aligned}$$

701 This gives

$$702 \quad (5.22) \quad \mathbb{E}[\|u^0 - y\|^2] = \|y\|^2 + 2y^\top \mathbb{E}[u^0] + \mathbb{E}[\|u^0\|^2] = \mathcal{O} \left(\frac{N}{\sigma_V} \right).$$

703 Applying the Markov inequality, we see that $\|u^0 - y\|^2 = \mathcal{O} \left(\frac{2N}{\epsilon \sigma_V} \right)$ will hold with
 704 probability at least $1 - \frac{1}{2}\epsilon$. This, together with the result of Lemma 5.14 with $Q =$
 705 $4\kappa_L(\mu\sigma_0)^{-1}\|u^0 - y\|$, implies that Assumption 5.9 will hold with probability at least

$$706 \quad 1 - \epsilon \text{ for } m = \Omega \left(\frac{128(L+2M)^3 \pi n N^6 \kappa_L^2}{\mu^2 \sigma_0^8 \sigma_V \epsilon^3 \min\{\frac{1}{2}, \frac{\mu}{6\kappa_L}\}^6} \right).$$

707 To establish the convergence rate result, observe from Theorem 5.10 that $\|\theta^k -$
 708 $\theta^0\| \leq 4\kappa_L(\mu\sigma_0)^{-1}\|u^0 - y\|$ when $\|u^0 - y\| \leq \frac{\mu}{3\kappa_H}$ and $\eta \leq \min \left\{ 1, \left(\frac{1}{6\|u^0 - y\|} - \frac{\kappa_H}{2\mu} \right) \cdot \frac{3\mu^2 \sigma_0}{8\kappa_R \kappa_L^2} \right\}$.

709 By taking $Q = 4\kappa_L \sigma_0^{-1} / (3\kappa_H)$ in Lemma 5.14, we see that Assumption 5.9 will hold
 710 with probability at least $1 - \epsilon$ if $m \geq \frac{16(L+2M)^3 \pi n N^5 \kappa_L^2}{9\sigma_0^8 \kappa_H^2 \epsilon^2 \min\{\frac{1}{2}, \frac{\mu}{6\kappa_L}\}^6}$. Following the proof of

711 Theorem 5.10, we conclude that (5.21) will hold for all $k \geq 0$ with probability at least
 712 $1 - \epsilon$. This completes the proof. \square

713 **6 Numerical results** In this section, we demonstrate the efficacy of our proposed
 714 method via numerical experiments on three problems: Low-rank matrix completion,
 715 low-dimension subspace learning, and deep learning model training. Our code
 716 is available at <https://github.com/hujiangpku/RNGD>.

717 **6.1 Low-rank matrix completion** We compare our proposed RNGD method
 718 with the Riemannian stochastic gradient descent (RSGD) method [11], the Riemannian
 719 adaptive stochastic gradient algorithm (RASA) [32], the Riemannian stochastic
 720 variance-reduced gradient (RSVRG) method [52], and the Riemannian conjugate gra-
 721 dient (RCG) method without preconditioner [14, 46, 51]. All algorithms are initial-
 722 ized by the QR decomposition of a random n -by- p matrix whose entries are generated
 723 from the standard Gaussian distribution. We consider two real datasets. One is taken

724 from the Jester joke recommender system,² which contains ratings (with scores from
725 -10.00 to $+10.00$) of 100 jokes from 24983 users. The other is the movie rating dataset
726 MovieLens-1M,³ which contains ratings (with stars from 1 to 5) of 3952 movies from
727 6040 users. In the experiments, each dataset is randomly divided into 2 sets, one
728 for training and the other for testing. We utilize the implementations of RSGD and
729 RSVRG given in the RSOpt package⁴ and the implementation of RCG given in the
730 Manopt package.⁵ For RASA, the LR-type variant is adopted due to its efficiency.
731 The default parameters therein are used. For RNGD, the same variance reduction
732 technique as that in RSVRG is adopted to update both the estimated gradient and
733 the approximate RFIM (4.5). Specifically, we compute $a_i(U)$ for all i in each outer
734 iteration and update $a_i(U)$ if the i -th sample is used in the estimation of the gradient.
735 We use fixed step sizes for RNGD and RSVRG. For RSGD, the step size η_k is set to
736 $\eta_k = \frac{\eta_0}{1+\eta_0 k/10}$. As suggested in [32], the step size $\eta_k = \eta_0/\sqrt{k}$ is used for RASA. We
737 search in the set $\{2, 1, 0.5, \dots, 2 \times 10^{-8}, 10^{-8}, 5 \times 10^{-9}\}$ to find the best initial step
738 size η_0 for RSGD and RASA and the best step size for RSVRG. The step size for
739 RNGD is set to 0.05 for both datasets.

740 Figure 6.1 reports the mean squared error (MSE) on both the training and testing
741 datasets, which are defined as $\|\mathcal{P}_{\Omega_{\text{train}}}(UA - X)\|^2/|\Omega_{\text{train}}|$ and $\|\mathcal{P}_{\Omega_{\text{test}}}(UA - X)\|^2/|\Omega_{\text{test}}|$,
742 respectively, where Ω_{train} and Ω_{test} are the sets of known indices in the
743 training and testing datasets, respectively. The label “#grad/ N ” means the number
744 of epochs, which is defined as the number of cycles through the full dataset. The
745 label “time” represents the wall-clock time. We run all algorithms with a specified
746 number of epochs for different datasets. On the Jester dataset, we see that RNGD,
747 RSVRG, and RCG achieve lower MSEs than the other two methods. Furthermore,
748 RNGD converges faster than RSVRG and RCG. In the case of the MovieLens-1M
749 dataset, RASA and RSVRG exhibit fast reductions of MSEs in the early iterations.
However, RNGD returns a point with the lowest MSE.

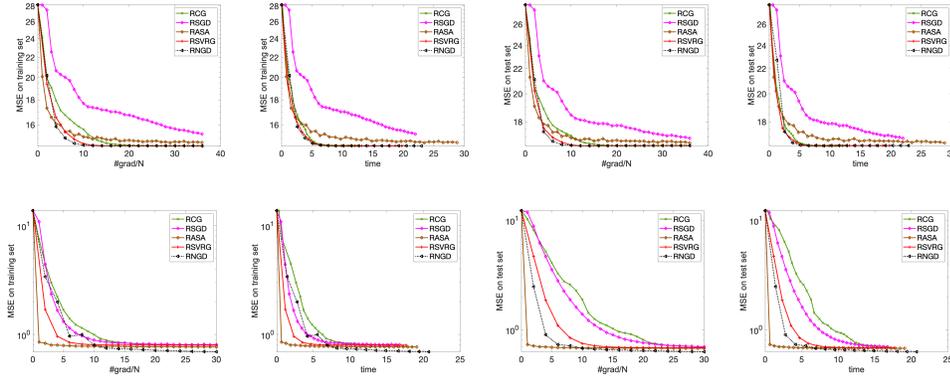


FIG. 6.1. Numerical results for LRMC on the Jester dataset (first row) and the MovieLens-1M dataset (second row).

750

²The dataset Jester can be downloaded from <https://grouplens.org/datasets/jester>

³The dataset MovieLens-1M can be downloaded from <https://grouplens.org/datasets/movielens>

⁴The code of RSOpt can be downloaded from <https://github.com/hiroyuki-kasai/RSOpt>

⁵The code of Manopt can be downloaded from <https://github.com/NicolasBoumal/manopt>

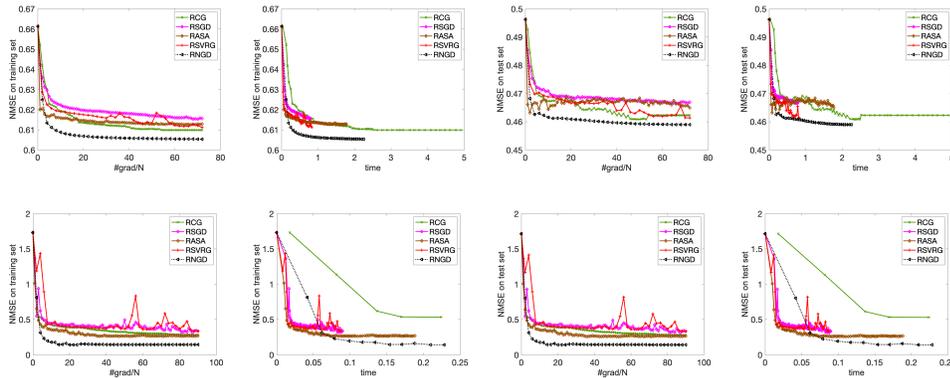


FIG. 6.2. Numerical results for multitask learning on the *School* dataset (first row) and the *Sarcos* dataset (second row).

751 **6.2 Low-dimension subspace learning** We compare our proposed RNGD
 752 with RCG, RSGD, RASA, and RSVRG on two real-world datasets: *School* [22] and
 753 *Sarcos* [57]. The dimension p is set to be 6 for both datasets. We choose the best step
 754 sizes for RSVRG, RASA, and RSGD from the set $\{1, 0.5, 0.2, 0.1, 0.05, 0.02, \dots, 10^{-8}, 5$
 755 $\times 10^{-9}, 2 \times 10^{-9}, 10^{-9}\}$. We use the step size 4 (resp., 1) on the *School* (resp., *Sarcos*)
 756 dataset for RNGD. All the codes are implemented within the RSOpt framework and
 757 the other parameters of the algorithms are set to the default values therein.

758 Figure 6.2 reports the normalized MSE (NMSE) [40] on both datasets, which is the
 759 mean of the normalized squared error of all tasks. For both datasets, RNGD returns a
 760 point with the lowest NMSE. Especially for the *Sarcos* dataset, a significant difference
 761 in the NMSE between RNGD and other methods is observed. Another noteworthy
 762 phenomenon is that RSGD and RSVRG tend to be less efficient than RCG. This
 763 demonstrates the advantage of using the Fisher information.

764 **6.3 Deep learning model training** Batch normalization and momentum-
 765 based optimizer are standard techniques to train state-of-the-art image classification
 766 models [24, 50, 54]. We evaluate the proposed method with Kronecker-factorized
 767 approximate RFIM described in Section 4, denoted by MKFAC, on VGG16BN [54]
 768 and WRN-16-4 [64] while the benchmark datasets CIFAR-10/100 [36] are used. The
 769 detailed network structures are described in [54, 64]. In VGG16BN, batch normaliza-
 770 tion layers are added before every ReLU activation layer. Additionally, we change the
 771 number of neurons in fully connected layers from 4096 to 512 and remove the middle
 772 layer of the last three in VGG due to memory allocation problems (otherwise, one
 773 has to compute the inverse of 4096^2 -by- 4096^2 matrices). This setting is also adopted
 774 in [17, 63].

775 The baseline algorithms are SGD, Adam, KFAC [39], AdamP, and SGDP [25].
 776 The tangential projections are used to control the increase in norms of the weight
 777 parameters in AdamP and SGDP. These methods can be seen as approximate Rie-
 778 mannian first-order methods. We fine tune the initial learning rates of the base-
 779 line algorithms by searching in the set $\{0.5, 0.2, 0.1, 0.05, 0.02, 0.01, \dots, 5 \times 10^{-5}, 2 \times$
 780 $10^{-5}, 10^{-5}\}$. The learning rate decays in epoch 30, 60, and 90 with a decay rate 0.1,
 781 where an epoch is defined as one cycle through the full training dataset. We choose
 782 the parameters β_1, β_2 in Adam and AdamP from the set $\{0.9, 0.99, 0.999\}$. We search
 783 in the set $\{0.05, 0.1, 0.2, 0.5, 1, 2\}$ to determine the damping parameter λ used in cal-

TABLE 6.1
Classification accuracy of various networks on CIFAR-10/100 (median of five runs).

Dataset	CIFAR-10		CIFAR-100	
	WRN-16-4	VGG16BN	WRN-16-4	VGG16BN
SGD	93.84	92.88	74.30	71.79
SGDP	93.42	92.49	73.67	71.54
Adam	92.53	89.88	71.64	62.79
AdamP	92.55	91.43	71.23	58.88
KFAC	93.90	94.36	74.31	76.38
MKFAC	94.06	94.76	74.55	77.28

874 culating the natural direction $(F_k + \lambda I)^{-1}g^k$ and update the KFAC matrix in epoch
875 30, 60, and 90. The initial damping parameter of KFAC is set to 2 in all four tasks.
876 We set the weight decay to 5×10^{-4} for all algorithms. Each mini-batch contains
877 128 samples. The maximum number of epochs is set to 100 for all algorithms. For
878 MKFAC, we use RNGD for parameters constrained on the Grassmann manifold and
879 SGD for the remaining parameters. Let η, η_g denote the learning rates for the Euclid-
890 ean space and Grassmann manifold, respectively. For the dataset CIFAR-10, we set
891 $\eta_g = 0.25$ and $\eta = 0.05$ with decay rates 0.2 and 0.1, respectively. The weight decay
892 is only applied to the unconstrained weights with parameter 5×10^{-4} . The initial
893 MKFAC damping parameters for WRN-16-4 and VGG16BN are set to 1 and 2 with
894 decay rates 0.8 and 0.5, respectively, when the preconditioners update in epoch 30,
895 60, and 90. For the dataset CIFAR-100, we set $\eta_g = 0.3$ for WRN-16-4, $\eta_g = 0.15$
896 for VGG16BN, and $\eta = 0.05$ for both. The learning rate η_g has a decay rate 0.15 for
897 WRN-16-4 and 0.2 for VGG16BN, while η has a decay rate 0.1 for both of them. The
898 initial MKFAC damping parameters for VGG16BN and WRN16-4 are set to 0.5 and
899 1 with decay rates 0.5 and 0.8, respectively. Other settings are the same as KFAC.

900 Table 6.1 presents the comparison of the baseline and the proposed algorithms
901 on CIFAR-10 and CIFAR-100 datasets. We list the best classification accuracy in
902 100 epochs, where the results are obtained from the median of 5 runs. The per-
903 formance of our proposed MKFAC method is the best in all four tasks. Compared
904 with the second-order type method KFAC, our MKFAC method reaches higher accu-
905 racy, though KFAC has a much better behavior than SGD on these tasks. Compared
906 with the manifold geometry-based first-order algorithms SGDP and AdamP, we see
907 that using second-order information can give better accuracy than using first-order
908 information alone.

909 **7 Conclusion** In this paper, we developed a novel efficient RNGD method for
910 tackling the problem of minimizing a sum of negative log-probability losses over a
911 manifold. Key to our development is a new notion of FIM on manifolds, which we
912 introduced in this paper and could be of independent interest. We established the
913 global convergence of RNGD and the local convergence rate of a deterministic ver-
914 sion of RNGD. Our numerical results on representative machine learning applications
915 demonstrate the efficiency and efficacy of the proposed method.

816

REFERENCES

- 817 [1] P.-A. ABSIL, R. MAHONY, AND R. SEPULCHRE, *Optimization Algorithms on Matrix Manifolds*,
818 Princeton University Press, Princeton, NJ, 2008.
819 [2] P.-A. ABSIL, R. MAHONY, AND J. TRUMPF, *An extrinsic look at the Riemannian Hessian*, in
820 Geometric Science of Information, Springer, 2013, pp. 361–368.

- 821 [3] P.-A. ABSIL AND J. MALICK, *Projection-like retractions on matrix manifolds*, SIAM Journal
822 on Optimization, 22 (2012), pp. 135–158.
- 823 [4] S.-I. AMARI, *Neural learning in structured parameter spaces — natural Riemannian gradient*,
824 in Advances in Neural Information Processing Systems, vol. 9, 1996, pp. 127–133.
- 825 [5] S.-I. AMARI, *Natural gradient works efficiently in learning*, Neural Computation, 10 (1998),
826 pp. 251–276.
- 827 [6] R. K. ANDO, T. ZHANG, AND P. BARTLETT, *A framework for learning predictive structures*
828 *from multiple tasks and unlabeled data*, Journal of Machine Learning Research, 6 (2005),
829 pp. 1817–1853.
- 830 [7] R. ANIL, V. GUPTA, T. KOREN, K. REGAN, AND Y. SINGER, *Scalable second order optimization*
831 *for deep learning*, arXiv:2002.09018, (2020).
- 832 [8] J. L. BA, J. R. KIROS, AND G. E. HINTON, *Layer normalization*, in Advances in Neural
833 Information Processing Systems - Deep Learning Symposium, 2016, p. arXiv preprint
834 arXiv:1607.06450.
- 835 [9] A. BAHAMOU, D. GOLDFARB, AND Y. REN, *A mini-block natural gradient method for deep*
836 *neural networks*, arXiv:2202.04124, (2022).
- 837 [10] G. BÉCIGNEUL AND O.-E. GANEA, *Riemannian adaptive optimization methods*, in International
838 Conference on Learning Representations, 2019.
- 839 [11] S. BONNABEL, *Stochastic gradient descent on Riemannian manifolds*, IEEE Transactions on
840 Automatic Control, 58 (2013), pp. 2217–2229.
- 841 [12] F. BOUCHARD, A. BRELOY, A. RENAUX, AND G. GINOLHAC, *Riemannian geometry and Cramér-*
842 *Rao bound for blind separation of Gaussian sources*, in ICASSP 2020-2020 IEEE International
843 Conference on Acoustics, Speech and Signal Processing (ICASSP), IEEE, 2020,
844 pp. 4717–4721.
- 845 [13] N. BOUMAL, *An Introduction to Optimization on Smooth Manifolds*, Available online, May,
846 2020.
- 847 [14] N. BOUMAL AND P.-A. ABSIL, *Low-rank matrix completion via preconditioned optimization on the*
848 *Grassmann manifold*, Linear Algebra and its Applications, 475 (2015), pp. 200–239.
- 849 [15] R. H. BYRD, S. L. HANSEN, J. NOCEDAL, AND Y. SINGER, *A stochastic quasi-Newton method*
850 *for large-scale optimization*, SIAM Journal on Optimization, 26 (2016), pp. 1008–1031.
- 851 [16] R. CHEN, M. MENICKELLY, AND K. SCHEINBERG, *Stochastic optimization using a trust-region*
852 *method and random models*, Mathematical Programming, 169 (2018), pp. 447–487.
- 853 [17] M. CHO AND J. LEE, *Riemannian approach to batch normalization*, in Advances in Neural
854 Information Processing Systems, vol. 30, 2017, pp. 5231–5241.
- 855 [18] S. S. DU, X. ZHAI, B. POCZOS, AND A. SINGH, *Gradient descent provably optimizes over-*
856 *parameterized neural networks*, in International Conference on Learning Representations,
857 2019.
- 858 [19] J. DUCHI, E. HAZAN, AND Y. SINGER, *Adaptive subgradient methods for online learning and*
859 *stochastic optimization.*, Journal of Machine Learning Research, 12 (2011).
- 860 [20] H. FLANDERS, *Differentiation under the integral sign*, The American Mathematical Monthly,
861 80 (1973), pp. 615–627.
- 862 [21] D. GOLDFARB, Y. REN, AND A. BAHAMOU, *Practical quasi-Newton methods for training deep*
863 *neural networks*, in Advances in Neural Information Processing Systems, vol. 33, 2020,
864 pp. 2386–2396.
- 865 [22] H. GOLDSTEIN, *Multilevel modelling of survey data*, Journal of the Royal Statistical Society.
866 Series D (The Statistician), 40 (1991), pp. 235–244.
- 867 [23] R. GROSSE AND J. MARTENS, *A Kronecker-factored approximate Fisher matrix for convolution*
868 *layers*, in International Conference on Machine Learning, 2016, pp. 573–582.
- 869 [24] K. HE, X. ZHANG, S. REN, AND J. SUN, *Deep residual learning for image recognition*, in
870 Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, 2016,
871 pp. 770–778.
- 872 [25] B. HEO, S. CHUN, S. J. OH, D. HAN, S. YUN, G. KIM, Y. UH, AND J.-W. HA, *AdamP: Slowing*
873 *down the slowdown for momentum optimizers on scale-invariant weights*, in International
874 Conference on Learning Representations, 2021.
- 875 [26] J. HU, X. LIU, Z.-W. WEN, AND Y.-X. YUAN, *A brief introduction to manifold optimization*,
876 Journal of the Operations Research Society of China, 8 (2020), pp. 199–248.
- 877 [27] J. HU, A. MILZAREK, Z. WEN, AND Y. YUAN, *Adaptive quadratically regularized Newton method*
878 *for Riemannian optimization*, SIAM Journal on Matrix Analysis and Applications, 39
879 (2018), pp. 1181–1207.
- 880 [28] S. IOFFE AND C. SZEGEDY, *Batch normalization: Accelerating deep network training by reducing*
881 *internal covariate shift*, in International Conference on Machine Learning, 2015, pp. 448–
882 456.

- 883 [29] B. JIANG, S. MA, A. M.-C. SO, AND S. ZHANG, *Vector transport-free SVRG with general re-*
884 *traction for Riemannian optimization: Complexity analysis and practical implementation*,
885 arXiv:1705.09059, (2017).
- 886 [30] N. L. JOHNSON, S. KOTZ, AND N. BALAKRISHNAN, *Continuous Univariate Distributions, volume*
887 *2*, vol. 289, John Wiley & Sons, 1995.
- 888 [31] R. JOHNSON AND T. ZHANG, *Accelerating stochastic gradient descent using predictive variance*
889 *reduction*, in Advances in Neural Information Processing Systems, vol. 26, 2013, pp. 315–
890 323.
- 891 [32] H. KASAI, P. JAWANPURIA, AND B. MISHRA, *Riemannian adaptive stochastic gradient al-*
892 *gorithms on matrix manifolds*, in International Conference on Machine Learning, 2019,
893 pp. 3262–3271.
- 894 [33] H. KASAI AND B. MISHRA, *Inexact trust-region algorithms on Riemannian manifolds.*, in
895 NeurIPS, 2018, pp. 4254–4265.
- 896 [34] H. KASAI, H. SATO, AND B. MISHRA, *Riemannian stochastic quasi-Newton algorithm with*
897 *variance reduction and its convergence analysis*, in International Conference on Artificial
898 Intelligence and Statistics, 2018, pp. 269–278.
- 899 [35] D. P. KINGma AND J. BA, *Adam: A method for stochastic optimization*, International Confer-
900 ence for Learning Representations, (2015).
- 901 [36] A. KRIZHEVSKY, G. HINTON, ET AL., *Learning multiple layers of features from tiny images*,
902 (2009).
- 903 [37] Y. LECUN, Y. BENGIO, AND G. HINTON, *Deep learning*, Nature, 521 (2015), p. 436.
- 904 [38] J. MARTENS, *New insights and perspectives on the natural gradient method*, The Journal of
905 Machine Learning Research, 21 (2020), pp. 5776–5851.
- 906 [39] J. MARTENS AND R. GROSSE, *Optimizing neural networks with Kronecker-factored approximate*
907 *curvature*, in International Conference on Machine Learning, 2015, pp. 2408–2417.
- 908 [40] B. MISHRA, H. KASAI, P. JAWANPURIA, AND A. SAROOP, *A Riemannian gossip approach to*
909 *subspace learning on Grassmann manifold*, Machine Learning, 108 (2019), pp. 1783–1803.
- 910 [41] J. NOCEDAL AND S. J. WRIGHT, *Numerical Optimization*, Springer Series in Operations Re-
911 search and Financial Engineering, Springer, New York, second ed., 2006.
- 912 [42] L. NURBEKYAN, W. LEI, AND Y. YANG, *Efficient natural gradient descent methods for large-*
913 *scale optimization problems*, arXiv:2202.06236, (2022).
- 914 [43] Y. OLLIVIER, L. ARNOLD, A. AUGER, AND N. HANSEN, *Information-geometric optimization*
915 *algorithms: A unifying picture via invariance principles*, Journal of Machine Learning
916 Research, 18 (2017), pp. 1–65.
- 917 [44] M. PILANCI AND M. J. WAINWRIGHT, *Newton sketch: A near linear-time optimization algorithm*
918 *with linear-quadratic convergence*, SIAM Journal on Optimization, 27 (2017), pp. 205–245.
- 919 [45] Y. REN AND D. GOLDFARB, *Kronecker-factored quasi-Newton methods for convolutional neural*
920 *networks*, arXiv:2102.06737, (2021).
- 921 [46] W. RING AND B. WIRTH, *Optimization methods on Riemannian manifolds and their application*
922 *to shape space*, SIAM Journal on Optimization, 22 (2012), pp. 596–627.
- 923 [47] H. ROBBINS AND S. MONRO, *A stochastic approximation method*, The Annals of Mathematical
924 Statistics, (1951), pp. 400–407.
- 925 [48] F. ROOSTA-KHORASANI AND M. W. MAHONEY, *Sub-sampled Newton methods*, Mathematical
926 Programming, 174 (2019), pp. 293–326.
- 927 [49] T. SALIMANS AND D. P. KINGMA, *Weight normalization: A simple reparameterization to ac-*
928 *celerate training of deep neural networks*, in Advances in Neural Information Processing
929 Systems, vol. 29, 2016, pp. 901–909.
- 930 [50] M. SANDLER, A. HOWARD, M. ZHU, A. ZHMOGINOV, AND L.-C. CHEN, *MobileNetV2: Inverted*
931 *residuals and linear bottlenecks*, in Proceedings of the IEEE Conference on Computer
932 Vision and Pattern Recognition, 2018, pp. 4510–4520.
- 933 [51] H. SATO, *Riemannian Optimization and Its Applications*, Springer, 2021.
- 934 [52] H. SATO, H. KASAI, AND B. MISHRA, *Riemannian stochastic variance reduced gradient algo-*
935 *rithm with retraction and vector transport*, SIAM Journal on Optimization, 29 (2019),
936 pp. 1444–1472.
- 937 [53] N. N. SCHRAUDOLPH, *Fast curvature matrix-vector products for second-order gradient descent*,
938 Neural Computation, 14 (2002), pp. 1723–1738.
- 939 [54] K. SIMONYAN AND A. ZISSERMAN, *Very deep convolutional networks for large-scale image recog-*
940 *nition*, arXiv:1409.1556, (2014).
- 941 [55] S. T. SMITH, *Covariance, subspace, and intrinsic Cramér-Rao bounds*, IEEE Transactions on
942 Signal Processing, 53 (2005), pp. 1610–1630.
- 943 [56] A. M.-C. SO, *Non-asymptotic performance analysis of the semidefinite relaxation detector in*
944 *digital communications*. Preprint, 2010.

- 945 [57] S. VIJAYAKUMAR, A. D’SOUZA, T. SHIBATA, J. CONRADT, AND S. SCHAAL, *Statistical learning*
946 *for humanoid robots*, *Autonomous Robots*, 12 (2002), pp. 55–69.
- 947 [58] X. WANG AND Y.-X. YUAN, *Stochastic trust region methods with trust region radius depending*
948 *on probabilistic models*, arXiv:1904.03342, (2019).
- 949 [59] X. WU, S. S. DU, AND R. WARD, *Global convergence of adaptive gradient methods for an*
950 *over-parameterized neural network*, arXiv:1902.07111, (2019).
- 951 [60] M. YANG, A. MILZAREK, Z. WEN, AND T. ZHANG, *A stochastic extra-step quasi-Newton method*
952 *for nonsmooth nonconvex optimization*, *Mathematical Programming*, (2021), pp. 1–47.
- 953 [61] M. YANG, D. XU, H. CHEN, Z. WEN, AND M. CHEN, *Enhance curvature information by struc-*
954 *tured stochastic quasi-Newton methods*, in *Proceedings of the IEEE/CVF Conference on*
955 *Computer Vision and Pattern Recognition*, 2021, pp. 10654–10663.
- 956 [62] M. YANG, D. XU, Q. CUI, Z. WEN, AND P. XU, *An efficient Fisher matrix approximation*
957 *method for large-scale neural network optimization*, *IEEE Transactions on Pattern Analysis*
958 *and Machine Intelligence*, (2022).
- 959 [63] M. YANG, D. XU, Z. WEN, M. CHEN, AND P. XU, *Sketch-based empirical natural gradient*
960 *methods for deep learning*, *Journal of Scientific Computing*, 92 (2022), pp. 1–29.
- 961 [64] S. ZAGORUYKO AND N. KOMODAKIS, *Wide residual networks*, arXiv:1605.07146, (2016).
- 962 [65] D. ZHANG AND S. D. TAJBAKHSH, *Riemannian stochastic variance-reduced cubic regularized*
963 *Newton method*, arXiv:2010.03785, (2020).
- 964 [66] G. ZHANG, J. MARTENS, AND R. GROSSE, *Fast convergence of natural gradient descent for*
965 *overparameterized neural networks*, in *Advances in Neural Information Processing Systems*,
966 2019, pp. 8082–8093.
- 967 [67] H. ZHANG, S. J. REDDI, AND S. SRA, *Riemannian SVRG: Fast stochastic optimization on*
968 *Riemannian manifolds*, in *Advances in Neural Information Processing Systems*, 2016,
969 pp. 4592–4600.