

33 be nonunique and highly unstable with respect to the perturbation in the noisy data y^δ ,
 34 regularization is often needed for their stable and accurate numerical solutions, and many ef-
 35 fective techniques have been proposed over the past few decades (see, e.g., [5, 15, 23, 12, 24]).
 36 Among existing techniques, iterative regularization represents a very powerful class of solvers
 37 for problem (1.1), including Landweber method, (regularized) Gauss-Newton method, con-
 38 jugate gradient methods, and Leverberg-Marquardt method etc; see the monographs [15]
 39 and [24] for overviews on iterative regularization methods in Hilbert spaces and Banach
 40 spaces, respectively. In this work, we are interested in the convergence analysis of stochastic
 41 gradient descent (SGD) for problem (1.1) with noisy data y^δ . The basic version of SGD
 42 reads: given the initial guess $x_1^\delta = x_1$, update the iterate x_k^δ by

$$43 \quad (1.3) \quad x_{k+1}^\delta = x_k^\delta - \eta_k F'_{i_k}(x_k^\delta)^* (F_{i_k}(x_k^\delta) - y_{i_k}^\delta); \quad k = 1, 2, \dots,$$

44 where the index i_k is drawn uniformly from the index set $\{1, \dots, n\}$, and $\eta_k > 0$ is the
 45 corresponding step size. SGD was pioneered by Robbins and Monro in statistical inference
 46 [22] (see the monograph [17] for asymptotic convergence results). It has demonstrated
 47 encouraging numerical results on diffuse optical tomography [2]. Further, a variant of SGD,
 48 i.e., randomized Kaczmarz method (RKM), has been successful in the computed tomography
 49 community [9, 10] with revived interest in linear regression and phase retrieval [25, 27].
 50 Algorithmically, SGD is a randomized version of the classical Landweber method [18]

$$51 \quad (1.4) \quad x_{k+1}^\delta = x_k^\delta - \eta_k F'(x_k^\delta)^* (F(x_k^\delta) - y^\delta),$$

52 which may be obtained from gradient descent applied to the functional

$$53 \quad (1.5) \quad J(x) = \frac{1}{2} \|F(x) - y^\delta\|^2 = \frac{1}{n} \sum_{i=1}^n \frac{1}{2} \|F_i(x) - y_i^\delta\|^2.$$

54 Compared with the Landweber method, SGD requires only evaluating one randomly se-
 55 lected (nonlinear) equation at each iteration, instead of the whole nonlinear system, which
 56 substantially reduces the computational cost per iteration and enables excellent scalability
 57 to truly massive data sets (i.e., large n), which are increasingly common in practice due to
 58 advances in data acquisition technologies. This highly desirable property has attracted much
 59 recent interest in machine learning, where currently SGD and its variants are the workhorse
 60 for many challenging training tasks involving deep neural networks [32, 26, 16, 1].

61 Note that due to the ill-posed nature of problem (1.1) (in the sense that the minimizer
 62 depends sensitively on the data perturbation), the minimization problem (1.5) is also *ill-*
 63 *posed*, and due to the inevitable presence of noise in the observational data y^δ , the global
 64 minimizer (if it exists at all!) often represents a poor approximation to the exact solution x^\dagger
 65 and thus is not of interest. The goal of iterative regularization is to iteratively construct an
 66 approximate minimizer that converges to the exact solution x^\dagger as the noise level $\delta \rightarrow 0^+$, and
 67 further, to derive convergence rates in terms of δ . This is achieved by equipping an iterative
 68 algorithm, e.g., Landweber method or SGD, with an early stopping strategy. Early stopping
 69 allows properly balancing the deleterious effect of the perturbation δ and the approximation
 70 error of the iterates for the perturbed data y^δ , which respectively grows and decreases as the
 71 iteration proceeds. Thus the setting differs greatly from *well-posed* optimization problems
 72 that are extensively studied in the optimization and machine learning literature.

73 For a class of nonlinear inverse problems, the Landweber method is relatively well un-
 74 derstood in terms of the regularizing property, since the influential work [8] (see also [20, 30]
 75 for linear inverse problems), and the results were refined and extended in different aspects
 76 [15]. In contrast, the stochastic counterparts, e.g., SGD, remains largely under-explored

77 for inverse problems, despite their computational appeals. The theoretical analysis of sto-
78 chastic iterative methods for inverse problems has just started, and some first theoretical
79 results were obtained in [13, 14] for linear inverse problems. The regularizing property of
80 SGD for linear inverse problems was proved in [14], by drawing on relevant developments in
81 statistical learning theory [31, 4, 19], whereas in [13], the preasymptotic convergence behav-
82 ior of RKM was analyzed. In this work, we study in depth the regularizing property and
83 convergence rates of SGD for a class of nonlinear inverse problems, under an *a priori* choice
84 of the stopping index and standard assumptions on the nonlinear operator F ; see section 2
85 for further details and discussions. The analysis borrows techniques from the works [14, 8],
86 i.e., handling iteration noise [14] and coping with the nonlinearity of forward map [8]. To
87 the best of our knowledge, this work gives a first thorough analysis of SGD for nonlinear
88 ill-posed inverse problems in the lens of iterative regularization.

89 There is a vast literature on the convergence of SGD and its variants in optimization and
90 machine learning; see [1, Section 4] for a comprehensive overview; see also [7] and references
91 therein for recent results and [6] for recent results in a Hilbert space setting. For general
92 nonconvex optimization problems, most of the results are concerned with the convergence
93 in terms of either expected optimality gap or expected norm of its gradient, with respect to
94 the iteration index k . However, these works focus on *well-posed* optimization problems, and
95 the ultimate goal is to find a global minimizer. This differs substantially from the setting of
96 *ill-posed* problems, e.g., (1.5). In particular, the existing convergence results of SGD cannot
97 be applied directly to deduce convergence (and rate) for problem (1.5), due to its least-
98 squares structure and different assumptions (on the forward map, instead of the objective
99 functional J ; see Remark 2.1 below for further discussions. More closely related to this
100 work are the works [31, 28, 4, 19] on generalization error in statistical learning. Ying and
101 Pontil [31] studied an online least-squares gradient descent algorithm in a reproducing kernel
102 Hilbert space (RKHS), and derived bounds on the generalization error. Lin and Rosasco
103 [19] analyzed the influence of batch size on the convergence of mini-batch SGD. See also
104 the recent work [4] on averaged SGD for nonparametric regression in RKHS. There are also
105 major differences between these interesting works and this study. First, in these prior works,
106 the noise arises mainly due to finite sampling, whereas for inverse problems, it arises from
107 imperfect data acquisition process and enters into the data y^δ directly. Second, the main
108 focus of these works is to bound the generalization error, instead of error estimates on the
109 iterate. Third, these prior works analyzed only linear problems (similar to [14]), instead of
110 nonlinear problems of this work. Nonetheless, our proof strategy of decomposing the mean
111 squared error into the bias and variance components shares similarity with these works.

112 Throughout, we denote the iterate for the exact data y^\dagger by x_k . The notation \mathcal{F}_k denotes
113 the filtration generated by the random indices $\{i_1, \dots, i_{k-1}\}$ up to the $(k-1)$ th iteration.
114 The notation c , with or without a subscript, denotes a generic constant, which may differ at
115 each occurrence, but it is always independent of the noise level δ and the iteration number
116 k . We shall abuse $\|\cdot\|$ for the operator norm on Y^n and from X to Y (or Y^n). The rest
117 of the paper is organized as follows. In section 2, we state the main results and provide
118 relevant discussions. Then in section 3 and section 4, we give the proofs on the regularizing
119 property and convergence rate, respectively. The paper concludes with further discussions
120 in section 5. In the appendix, we collect some useful inequalities.

121 **2. Main results and discussions.** To analyze SGD for nonlinear inverse problems,
122 suitable conditions are needed. For example, for Tikhonov regularization, both nonlinearity
123 and source conditions are often employed to derive convergence rates [5, 11, 24, 12]. Below
124 we make a number of assumptions on the nonlinear operators F_i and the reference solution
125 x^\dagger . Since the solution to problem (1.1) may be nonunique, the reference solution x^\dagger is taken

126 to be the minimum norm solution (with respect to the initial guess x_1), which is known to
 127 be unique under [Assumption 2.1\(ii\)](#) below [8].

128 **ASSUMPTION 2.1.** *The following conditions hold:*

- 129 (i) *The operator $F : X \rightarrow Y^n$ is continuous, with a continuous and uniformly bounded*
 130 *Fréchet derivative on X .*
 131 (ii) *There exists an $\eta \in (0, \frac{1}{2})$ such that for any $x, \tilde{x} \in X$,*

132 (2.1)
$$\|F(x) - F(\tilde{x}) - F'(\tilde{x})(x - \tilde{x})\| \leq \eta \|F(x) - F(\tilde{x})\|.$$

- 133 (iii) *There are a family of uniformly bounded operators R_x^i such that for any $x \in X$,*
 134 *$F'_i(x) = R_x^i F'_i(x^\dagger)$ and $R_x = \text{diag}(R_x^i) : Y^n \rightarrow Y^n$, with*

135
$$\|R_x - I\| \leq c_R \|x - x^\dagger\|.$$

- 136 (iv) *The source condition holds: there exist some $\nu \in (0, \frac{1}{2})$ and $w \in X$ such that*

137
$$x^\dagger - x_1 = (F'(x^\dagger)^* F'(x^\dagger))^\nu w.$$

138 The conditions in [Assumption 2.1](#) are standard for analyzing iterative regularization
 139 methods for nonlinear inverse problems [8, 15]. (i) is similar to the λ -smoothness commonly
 140 used in optimization. (ii)–(iii) have been verified for a class of nonlinear inverse problems
 141 [8], e.g., parameter identification for PDEs and nonlinear integral equations. The inequality
 142 (2.1) is often known as tangential cone condition, and it controls the degree of nonlinearity
 143 of the operator F . Roughly speaking, it requires the map F be not far from a linear map; see
 144 Lemma 3.1 for the consequences. The fractional power $(F'(x^\dagger)^* F'(x^\dagger))^\nu$ in (iv) is defined
 145 by spectral decomposition (e.g., via Dunford-Taylor integral). Customarily, it represents a
 146 certain smoothness condition on the exact solution x^\dagger (relative to the initial guess x_1). The
 147 restriction $\nu < \frac{1}{2}$ is due to technical reasons. It is worth noting that most results require only
 148 (i)–(ii), especially the convergence of SGD, whereas (iii)–(iv) are only needed for proving
 149 the convergence rate of SGD.

150 **REMARK 2.1.** *It is instructive to compare Assumption 2.1 with the canonical conditions*
 151 *for the usual finite-sum optimization:*

152 (2.2)
$$\mathcal{F}(x) = n^{-1} \sum_{i=1}^n f_i(x).$$

153 *Clearly problem (1.5) is a special case of (2.2), with the choice $f_i(x) = \frac{1}{2} \|F_i(x) - y_i^\delta\|^2$. In*
 154 *the literature on SGD for problem (2.2), the following two conditions are often adopted*

- 155 • *L-smoothness:* $\|\mathcal{F}'(x) - \mathcal{F}'(\tilde{x})\| \leq L \|x - \tilde{x}\|$
 156 • λ -*convexity:* $\mathcal{F}(x) \geq \mathcal{F}(\tilde{x}) + (\mathcal{F}'(\tilde{x}), x - \tilde{x}) + \frac{\lambda}{2} \|x - \tilde{x}\|^2$.

157 *Under these conditions, various convergence results have been established; see [1, Section 4].*

158 *Assumption 2.1(i) imposes boundness and continuity on the derivative $F'(u)$, which*
 159 *does not imply directly the L-smoothness condition. Nonetheless, the Lipschitz continuity of*
 160 *$F'(u)$ can be verified for a number of inverse problems, which then implies the L-smoothness*
 161 *condition. Assumption 2.1(ii) requires the forward map being not too far from a linear map,*
 162 *and thus one might expect a link with the λ -convexity, which, however, seems not evident.*
 163 *Straightforward computation gives $\nabla^2 J(x) = F'(x)^* F'(x) + \nabla^2 F(x)^* (F(x) - y^\delta)$. First, the*
 164 *map F is not assumed a priori twice differentiable so that $J(x)$ admits a Hessian $\nabla^2 J(x)$.*
 165 *Second, if the Hessian $\nabla^2 F$ does exist, then Taylor expansion gives*

166
$$\|F(x) - F(\tilde{x}) - F'(\tilde{x})(x - \tilde{x})\| = \|\frac{1}{2} \nabla^2 F(\tilde{x})(x - \tilde{x})^2 + \mathcal{O}(\|x - \tilde{x}\|^3)\| \leq \eta \|F(x) - F(\tilde{x})\|.$$

167 Unfortunately it does not imply directly that $\nabla^2 F$ is small. Further, $F'(x)^* F'(x)$ is usually
 168 only positive semidefinite, since the linearized operator $F'(x)$ is degenerate (e.g. compact)
 169 for most ill-posed inverse problems, so even if $\nabla^2 F(\bar{x})$ is small, generally one cannot ensure
 170 $\nabla^2 J(x) \geq 0$, i.e., the convexity. In sum, (2.1) does not imply the λ -convexity condition.
 171 Thus Assumption 2.1 is not directly comparable with standard assumptions for SGD, and
 172 the convergence results in [1] cannot be applied directly.

173 We also need suitable assumptions on the step size schedule $\{\eta_k\}_{k=1}^\infty$. The choice is viable
 174 since $\max_i \sup_{x \in X} \|F'_i(x)\| < \infty$, by Assumption 2.1(i). The choice in Assumption 2.2(i) is
 175 more general than (ii). The latter choice is often known as a polynomially decaying step
 176 size schedule in the literature.

177 ASSUMPTION 2.2. The step sizes $\{\eta_k\}_{k \geq 1}$ satisfy one of the following conditions.

- 178 (i) $\eta_k \max_i \sup_{x \in X} \|F'_i(x)\|^2 < 1$ and $\sum_{k=1}^\infty \eta_k = \infty$.
 179 (ii) $\eta_k = \eta_0 k^{-\alpha}$, with $\alpha \in (0, 1)$ and $\eta_0 \leq (\max_i \sup_{x \in X} \|F'_i(x)\|^2)^{-1}$.

180 Due to the random choice of the index i_k , the SGD iterate x_k^δ is random. There are
 181 several different ways to measure the convergence. We shall employ the mean squared norm
 182 defined by $\mathbb{E}[\|\cdot\|^2]$, where the expectation $\mathbb{E}[\cdot]$ is with respect to the filtration \mathcal{F}_k . Clearly, the
 183 iterate x_k^δ is measurable with respect to \mathcal{F}_k . The first result gives the regularizing property
 184 of SGD for problem (1.1) under *a priori* parameter choice. The notation $\mathcal{N}(\cdot)$ denotes the
 185 kernel of a linear operator.

186 THEOREM 2.1 (convergence for noisy data). Let Assumption 2.1(i)-(ii) and Assump-
 187 tion 2.2(i) be fulfilled. If the stopping index $k(\delta) \in \mathbb{N}$ satisfies $\lim_{\delta \rightarrow 0^+} k(\delta) = \infty$ and
 188 $\lim_{\delta \rightarrow 0^+} \delta^2 \sum_{i=1}^{k(\delta)} \eta_i = 0$, then there exists a solution $x^* \in X$ to problem (1.1) such that

$$189 \quad \lim_{\delta \rightarrow 0^+} \mathbb{E}[\|x_{k(\delta)}^\delta - x^*\|^2] = 0.$$

190 Further, if $\mathcal{N}(F'(x^\dagger)) \subset \mathcal{N}(F'(x))$, then

$$191 \quad \lim_{\delta \rightarrow 0^+} \mathbb{E}[\|x_{k(\delta)}^\delta - x^\dagger\|^2] = 0.$$

192 REMARK 2.2. The conditions on $k(\delta)$ in Theorem 2.1 are identical with that for the
 193 Landweber method [8, Theorem 2.4]. Note that consistency does not require a monotonically
 194 decreasing step size schedule, and holds for a constant step size.

195 Next we make an assumption on the nonlinearity of the operator F in a stochastic sense.

196 ASSUMPTION 2.3. There exist some $\theta \in (0, 1]$ and $c_R > 0$ such that for any function
 197 $G : X \rightarrow Y^n$ and $z_t = tx_k^\delta + (1-t)x^\dagger$, $t \in [0, 1]$, there hold

$$198 \quad \mathbb{E}[\|(I - R_{z_t})G(x_k^\delta)\|^2]^{\frac{1}{2}} \leq c_R \mathbb{E}[\|x_k^\delta - x^\dagger\|^2]^{\frac{\theta}{2}} \mathbb{E}[\|G(x_k^\delta)\|^2]^{\frac{1}{2}},$$

$$199 \quad \mathbb{E}[\|(I - R_{z_t}^*)G(x_k^\delta)\|^2]^{\frac{1}{2}} \leq c_R \mathbb{E}[\|x_k^\delta - x^\dagger\|^2]^{\frac{\theta}{2}} \mathbb{E}[\|G(x_k^\delta)\|^2]^{\frac{1}{2}}.$$

201 Assumption 2.3 is a stochastic version of Assumption 2.1(iii), and strengthens the cor-
 202 responding estimate in the sense of expectation. The case $\theta = 0$ follows trivially from
 203 Assumption 2.1(iii), by the boundedness of the operator R_x , whereas with $\theta = 1$, it recovers
 204 the latter when specialized to a Dirac measure. It will play a role in the convergence rate
 205 analysis, by taking $G(x) = F(x) - y^\delta$ and $G(x) = F'(x^\dagger)(x - x^\dagger)$ (see the proofs in Lemma 4.1
 206 and Lemma 4.6), and it enables bounding the terms involving conditional dependence.

207 The next result gives a convergence rate under *a priori* parameter choice, i.e., bound on
 208 the error $e_k^\delta := x_k^\delta - x^\dagger$, in terms of δ and k etc. The notation $[\cdot]$ denotes taking the integral

209 part of a real number, provided that $\|F'(x^\dagger)^*F'(x^\dagger)\| \leq 1$ and $\eta_0 \leq 1$. The assumptions in
 210 [Theorem 2.2](#) are identical with that for the Landweber method [8], except [Assumption 2.3](#).
 211 The strategy of the error analysis is to split the mean squared error $\mathbb{E}[\|e_k^\delta\|^2]$ using bias-
 212 variance decomposition: with bias $\|\mathbb{E}[e_k^\delta]\|^2$ and variance $\mathbb{E}[\|e_k^\delta - \mathbb{E}[e_k^\delta]\|^2]$,

$$213 \quad (2.3) \quad \mathbb{E}[\|e_k^\delta\|^2] = \|\mathbb{E}[e_k^\delta]\|^2 + \mathbb{E}[\|e_k^\delta - \mathbb{E}[e_k^\delta]\|^2].$$

214 The former contains the approximation error and data error, whereas the latter arises from
 215 the random choice of the index i_k . Due to the nonlinearity of the operator F , the two terms
 216 interact with each other (and also $\mathbb{E}[\|F'(x^\dagger)e_k^\delta\|^2]$); see [Theorem 4.4](#) and [Theorem 4.7](#). This
 217 leads to a coupled system of recursive inequalities for $\mathbb{E}[\|e_k^\delta\|^2]$ and $\mathbb{E}[\|F'(x^\dagger)e_k^\delta\|^2]$, and
 218 thus the analysis differs substantially from that for linear inverse problems in [14] and the
 219 Landweber method for nonlinear inverse problems [8].

220 **THEOREM 2.2.** *Let [Assumption 2.1](#), [Assumption 2.2\(ii\)](#) and [Assumption 2.3](#) be fulfilled*
 221 *with $\|w\|$ and η_0 being sufficiently small, and x_k^δ be the SGD iterate defined in (1.3). Then*
 222 *for all $k \leq k^* = \lceil (\frac{\delta}{\|w\|})^{-\frac{2}{(2\nu+1)(1-\alpha)}} \rceil$ and small $\epsilon \in (0, \frac{\alpha}{2})$, there hold*

$$223 \quad \mathbb{E}[\|e_k^\delta\|^2] \leq c^* k^{-\min(2\nu(1-\alpha), \alpha-\epsilon)} \|w\|^2 \quad \text{and} \quad \mathbb{E}[\|F'(x^\dagger)e_k^\delta\|^2] \leq c^* k^{-\min((1+2\nu)(1-\alpha), 1-\epsilon)} \|w\|^2,$$

224 where the constant c^* depends on ν , α , η_0 , n and θ , but is independent of k and δ .

225 **REMARK 2.3.** *When $\alpha \in (0, 1)$ is close to 1, setting $k = k^*$ gives*

$$226 \quad \mathbb{E}[\|e_{k^*}^\delta\|^2] \leq c^* \|w\|^{\frac{2}{2\nu+1}} \delta^{\frac{4\nu}{2\nu+1}} \quad \text{and} \quad \mathbb{E}[\|F'(x^\dagger)e_{k^*}^\delta\|^2] \leq c^* \|w\|^{\frac{4\nu}{2\nu+1}} \delta^{\frac{2}{2\nu+1}}.$$

227 These rates are comparable with that for the Landweber method for nonlinear inverse prob-
 228 lems [8, Theorem 3.2] and SGD for linear inverse problems [14, Theorem 2.2]. The restric-
 229 tion $O(k^{-(\alpha-\epsilon)})$ is due to the computational variance arising from the random index i_k , and
 230 for small α , the convergence rate may suffer from a loss. It is noteworthy that for $\nu > 1/2$,
 231 the convergence rate is suboptimal, just as the classical Landweber method, and thus SGD
 232 may suffer from a saturation phenomenon. It is an interesting open question to remove the
 233 saturation phenomenon.
 234

235 **REMARK 2.4.** *In practice, the domain $\mathcal{D}(F) \subset X$ is often not the whole space X , es-*
 236 *pecially for parameter identifications for PDEs, where box constraints arise naturally due*
 237 *to physical constraints. When the domain $\mathcal{D}(F) \subset X$ is a closed convex set, e.g., box con-*
 238 *straints, it can be incorporated into the algorithm by a projection operator P [29], i.e.,*

$$239 \quad x_{k+1}^\delta = P(x_k^\delta - \eta_k F'_{i_k}(x_k^\delta)^*(F_{i_k}(x_k^\delta) - y_{i_k}^\delta)).$$

240 However, the presence of the projection P significantly complicates the analysis. The exten-
 241 sion to the constrained case is an interesting open question.

242 **3. Convergence of SGD.** Now we analyze the convergence of SGD, and give the
 243 proof of [Theorem 2.1](#). We first recall a useful characterization of an exact solution x^* [8,
 244 Proposition 2.1].

245 **LEMMA 3.1.** *The following statements hold under [Assumption 2.1\(i\)–\(ii\)](#).*

246 (i) *The following upper and lower bounds hold:*

$$247 \quad \frac{1}{1+\eta} \|F'(x)(x - \tilde{x})\| \leq \|F(x) - F(\tilde{x})\| \leq \frac{1}{1-\eta} \|F'(x)(x - \tilde{x})\|.$$

248 (ii) *If x^* is a solution of problem (1.1), then any other solution \tilde{x}^* satisfies $x^* - \tilde{x}^* \in$
 249 $\mathcal{N}(F'(x^*))$, and vice versa.*

288 The next result shows that the sequence $\{x_k\}_{k \geq 1}$ is a Cauchy sequence.

289 LEMMA 3.3. Under *Assumption 2.1(i)-(ii)* and *Assumption 2.2(i)*, for the exact data
290 y^\dagger , the sequence $\{x_k\}_{k \geq 1}$ generated by SGD (1.3) is a Cauchy sequence.

291 *Proof.* The argument below follows closely [8, Theorem 2.3], which can be traced back
292 to [21]. Let x^* be any solution to problem (1.1), and let $e_k := x_k - x^*$. By Corollary 3.2,
293 $\mathbb{E}[\|e_k\|^2]$ is monotonically decreasing to some $\epsilon \geq 0$. Next we show that the sequence $\{x_k\}_{k \geq 1}$
294 is actually a Cauchy sequence. First we note that $\mathbb{E}[\langle \cdot, \cdot \rangle]$ defines an inner product. For any
295 $j \geq k$, choose an index ℓ with $j \geq \ell \geq k$ such that

$$296 \quad (3.1) \quad \mathbb{E}[\|y^\dagger - F(x_\ell)\|^2] \leq \mathbb{E}[\|y^\dagger - F(x_i)\|^2], \quad \forall k \leq i \leq j.$$

297 By the inequality $\mathbb{E}[\|e_j - e_k\|^2]^{\frac{1}{2}} \leq \mathbb{E}[\|e_j - e_\ell\|^2]^{\frac{1}{2}} + \mathbb{E}[\|e_\ell - e_k\|^2]^{\frac{1}{2}}$ and the identities

$$298 \quad (3.2) \quad \begin{aligned} \mathbb{E}[\|e_j - e_\ell\|^2] &= 2\mathbb{E}[\langle e_\ell - e_j, e_\ell \rangle] + \mathbb{E}[\|e_j\|^2] - \mathbb{E}[\|e_\ell\|^2], \\ \mathbb{E}[\|e_\ell - e_k\|^2] &= 2\mathbb{E}[\langle e_\ell - e_k, e_\ell \rangle] + \mathbb{E}[\|e_k\|^2] - \mathbb{E}[\|e_\ell\|^2], \end{aligned}$$

299 it suffices to prove that both $\mathbb{E}[\|e_j - e_\ell\|^2]$ and $\mathbb{E}[\|e_\ell - e_k\|^2]$ tend to zero as $k \rightarrow \infty$. For
300 $k \rightarrow \infty$, the last two terms on each of the right-hand side of (3.2) tend to $\epsilon - \epsilon = 0$, by
301 the monotone convergence of $\mathbb{E}[\|e_k\|^2]$ to ϵ , cf. Corollary 3.2. Next we show that the term
302 $\mathbb{E}[\langle e_\ell - e_k, e_\ell \rangle]$ also tends to zero as $k \rightarrow \infty$. Actually, by the definition of x_k , we have

$$303 \quad e_\ell - e_k = \sum_{i=k}^{\ell-1} (e_{i+1} - e_i) = \sum_{i=k}^{\ell-1} \eta_i F'_{i_i}(x_i)^* (y_{i_i}^\dagger - F_{i_i}(x_i)).$$

304 By triangle inequality and Cauchy-Schwarz inequality, we have

$$\begin{aligned} 305 \quad |\mathbb{E}[\langle e_\ell - e_k, e_\ell \rangle]| &\leq \sum_{i=k}^{\ell-1} \eta_i |\mathbb{E}[\langle F'_{i_i}(x_i)^* (y_{i_i}^\dagger - F_{i_i}(x_i)), e_\ell \rangle]| \\ 306 \quad &= \sum_{i=k}^{\ell-1} \eta_i |\mathbb{E}[\langle y_{i_i}^\dagger - F_{i_i}(x_i), F'_{i_i}(x_i)(x^* - x_i + x_i - x_\ell) \rangle]| \\ 307 \quad &= \sum_{i=k}^{\ell-1} \eta_i |\mathbb{E}[\langle y^\dagger - F(x_i), F'(x_i)(x^* - x_i + x_i - x_\ell) \rangle]| \\ 308 \quad &\leq \sum_{i=k}^{\ell-1} \eta_i \mathbb{E}[\|y^\dagger - F(x_i)\|^2]^{\frac{1}{2}} \mathbb{E}[\|F'(x_i)(x^* - x_i)\|^2]^{\frac{1}{2}} \\ 309 \quad &\quad + \sum_{i=k}^{\ell-1} \eta_i \mathbb{E}[\|y^\dagger - F(x_i)\|^2]^{\frac{1}{2}} \mathbb{E}[\|F'(x_i)(x_i - x_\ell)\|^2]^{\frac{1}{2}} := \text{I} + \text{II}. \end{aligned}$$

311 By Assumption 2.1(ii) and Lemma 3.1(i), we bound the first term I by

$$\begin{aligned} 312 \quad \text{I} &\leq (1 + \eta) \sum_{i=k}^{\ell-1} \eta_i \mathbb{E}[\|y^\dagger - F(x_i)\|^2]^{\frac{1}{2}} \mathbb{E}[\|F(x^*) - F(x_i)\|^2]^{\frac{1}{2}} \\ 313 \quad &= (1 + \eta) \sum_{i=k}^{\ell-1} \eta_i \mathbb{E}[\|y^\dagger - F(x_i)\|^2]. \\ 314 \end{aligned}$$

315 Likewise, we bound the term II by triangle inequality and the choice of ℓ in (3.1) as:

$$\begin{aligned}
316 \quad \text{II} &\leq (1 + \eta) \sum_{i=k}^{\ell-1} \eta_i \mathbb{E}[\|y^\dagger - F(x_i)\|^2]^{\frac{1}{2}} \mathbb{E}[\|(F(x_\ell) - y^\dagger) + (y^\dagger - F(x_i))\|^2]^{\frac{1}{2}} \\
317 \quad &\leq 2(1 + \eta) \sum_{i=k}^{\ell-1} \eta_i \mathbb{E}[\|y^\dagger - F(x_i)\|^2]. \\
318
\end{aligned}$$

319 The last two estimates together imply $|\mathbb{E}\langle e_\ell - e_k, e_\ell \rangle| \leq 3(1 + \eta) \sum_{i=k}^{\ell-1} \eta_i \mathbb{E}[\|y^\dagger - F(x_i)\|^2]$.
320 Similarly, one can deduce $\mathbb{E}\langle e_j - e_\ell, e_\ell \rangle| \leq 3(1 + \eta) \sum_{i=\ell}^{j-1} \eta_i \mathbb{E}[\|y^\dagger - F(x_i)\|^2]$. These two
321 estimates and Corollary 3.2 imply that the right-hand sides of (3.2) tend to zero as $k \rightarrow \infty$.
322 Hence both $\{e_k\}_{k \geq 1}$ and $\{x_k\}_{k \geq 1}$ are Cauchy sequences. \square

323 LEMMA 3.4. Under Assumption 2.1(i)-(ii) and Assumption 2.2(i), there holds

$$324 \quad \lim_{k \rightarrow \infty} \mathbb{E}[\|F(x_k) - y^\dagger\|^2] = 0.$$

325 *Proof.* Lemma 3.3 implies that $\{x_k\}_{k \geq 1}$ is a Cauchy sequence. By Assumption 2.2(i),
326 $\sup_{x \in X} \|F'(x)\| \leq c_F$ for some $c_F > 0$. Further, for any $x, \tilde{x} \in X$, there holds

$$327 \quad \|F(x) - F(\tilde{x})\| \leq (1 - \eta)^{-1} \|F'(x)(x - \tilde{x})\| \leq c_F(1 - \eta)^{-1} \|x - \tilde{x}\|.$$

328 Thus, $\{F(x_k) - y^\dagger\}_{k \geq 1}$ is a Cauchy sequence, and $\mathbb{E}[\|F(x_k) - y^\dagger\|^2]$ converges. Now we
329 proceed by contradiction, and assume that $\lim_{k \rightarrow \infty} \mathbb{E}[\|F(x_k) - y^\dagger\|^2] > 0$. Then there
330 exist some $\epsilon > 0$ and $k^* \in \mathbb{N}$, such that $\mathbb{E}[\|F(x_k) - y^\dagger\|^2] \geq \epsilon$ for all $k \geq k^*$. Hence, by
331 Assumption 2.2(i),

$$332 \quad \sum_{k=1}^{\infty} \eta_k \mathbb{E}[\|F(x_k) - y^\dagger\|^2] \geq \sum_{k=k^*}^{\infty} \eta_k \mathbb{E}[\|F(x_k) - y^\dagger\|^2] \geq \epsilon \sum_{k=k^*}^{\infty} \eta_k = \infty,$$

333 which contradicts the inequality $\sum_{k=1}^{\infty} \eta_k \mathbb{E}[\|F(x_k) - y^\dagger\|^2] < \infty$ from Corollary 3.2. \square

334 Now we can state the convergence of SGD for the exact data y^\dagger . Below x^\dagger denotes the
335 unique solution to problem (1.1) of minimal distance to x_1 .

336 THEOREM 3.5 (Convergence for exact data). Let Assumption 2.1(i)-(ii) and Assump-
337 tion 2.2(i) be fulfilled. Then for the exact data y^\dagger , the sequence $\{x_k\}_{k \geq 1}$ generated by SGD
338 converges to a solution x^* of problem (1.1):

$$339 \quad \lim_{k \rightarrow \infty} \mathbb{E}[\|x_k - x^*\|^2] = 0.$$

340 Further, if $\mathcal{N}(F'(x^\dagger)) \subset \mathcal{N}(F'(x))$, then

$$341 \quad \lim_{k \rightarrow \infty} \mathbb{E}[\|x_k - x^\dagger\|^2] = 0.$$

342 *Proof.* Since $\{x_k\}_{k \geq 1}$ is a Cauchy sequence, it has a limit, denoted by x^* . Further, x^*
343 is a solution, since by Lemma 3.4, the mean squared residual $\mathbb{E}[\|y^\dagger - F(x_k)\|^2]$ converges
344 to zero as $k \rightarrow \infty$. Note that problem (1.1) has a unique solution of minimal distance to
345 the initial guess x_1 that satisfies $x^\dagger - x_1 \in \mathcal{N}(F'(x^\dagger))^\perp$; see Lemma 3.1. If $\mathcal{N}(F'(x^\dagger)) \subset$
346 $\mathcal{N}(F'(x_k))$ for all $k = 1, 2, \dots$, then clearly, $x_k - x_1 \in \mathcal{N}(F'(x^\dagger))^\perp$, $k = 1, 2, \dots$. Hence,
347 $x^\dagger - x^* = x^\dagger - x_1 + x_1 - x^* \in \mathcal{N}(F'(x^\dagger))^\perp$. This and Lemma 3.1 imply $x^* = x^\dagger$. \square

348 **REMARK 3.2.** *Theorem 3.5 does not impose any constraint on the step size schedule*
349 *$\{\eta_k\}_{k=1}^\infty$ directly, apart from the fact that it should not decay too fast to zero. In particular,*
350 *it can be taken to be a constant step size. This result slightly improves that in [14, Theorem*
351 *2.1], where a decreasing step size is required (for linear inverse problems). The improvement*
352 *is achieved by exploiting the quadratic structure of the functional $J(x)$ in (1.5) (and the*
353 *tangential cone condition in Assumption 2.1(i)), whereas in [14] the consistency is derived*
354 *by means of bias-variance decomposition.*

355 **3.2. Convergence for noisy data.** The next result gives the stability of the SGD
356 iterate x_k^δ with respect to the noise level δ (at $\delta = 0$).

357 **LEMMA 3.6.** *Let Assumption 2.1(i) be fulfilled. For any fixed $k \in \mathbb{N}$ and any path*
358 *$(i_1, \dots, i_{k-1}) \in \mathcal{F}_k$, let x_k and x_k^δ be the SGD iterates along the path for exact data y^\dagger and*
359 *noisy data y^δ , respectively. Then*

$$360 \quad \lim_{\delta \rightarrow 0^+} \mathbb{E}[\|x_k^\delta - x_k\|^2] = 0.$$

361 *Proof.* We prove the assertion by mathematical induction. It holds trivially for $k = 1$.
362 Now suppose that it holds for all indices up to k and any path in \mathcal{F}_k . By the definition, for
363 any fixed path (i_1, \dots, i_k) , we have

$$364 \quad x_{k+1}^\delta - x_{k+1} = (x_k^\delta - x_k) - \eta_k \left((F'_{i_k}(x_k^\delta))^* - F'_{i_k}(x_k)^* \right) (F_{i_k}(x_k^\delta) - y_{i_k}^\delta) \\
365 \quad + F'_{i_k}(x_k)^* \left((F_{i_k}(x_k^\delta) - y_{i_k}^\delta) - (F_{i_k}(x_k) - y_{i_k}^\dagger) \right).$$

367 Thus, by triangle inequality,

$$368 \quad (3.3) \quad \|x_{k+1}^\delta - x_{k+1}\| \leq \|x_k^\delta - x_k\| + \eta_k \|F'_{i_k}(x_k^\delta)^* - F'_{i_k}(x_k)^*\| \|F_{i_k}(x_k^\delta) - y_{i_k}^\delta\| \\
369 \quad + \eta_k \|F'_{i_k}(x_k)^*\| \|(F_{i_k}(x_k^\delta) - y_{i_k}^\delta) - (F_{i_k}(x_k) - y_{i_k}^\dagger)\|.$$

371 Next we show that for any fixed k , $\sup_{(i_1, \dots, i_{k-1}) \in \mathcal{F}_k} \|x_k\|$ is bounded. Indeed, by Assump-
372 tion 2.1(i), $\max_i \sup_{x \in X} \|F'_i(x)\| \leq c_F$ for some $c_F > 0$. Then, by Lemma 3.1(i)

$$373 \quad \|x_{k+1} - x^*\| \leq \|x_k - x^*\| + \eta_k \|F'_{i_k}(x_k)^*\| \|F_{i_k}(x_k) - y_{i_k}^\dagger\| \leq (1 + \eta_k \frac{c_F^2}{1-\eta}) \|x_k - x^*\|.$$

375 This and an induction argument show that the claim. Similarly,

$$376 \quad \|F_{i_k}(x_k^\delta) - y_{i_k}^\delta\| \leq \|F_{i_k}(x_k^\delta) - F_{i_k}(x_k)\| + \|F_{i_k}(x_k) - y_{i_k}^\dagger\| + \|y_{i_k}^\dagger - y_{i_k}^\delta\| \\
377 \quad \leq \frac{c_F}{1-\eta} (\|x_k^\delta - x_k\| + \|x_k - x^*\|) + \delta,$$

379 and consequently,

$$380 \quad \|x_{k+1}^\delta - x_{k+1}\| \leq \|x_k^\delta - x_k\| + \eta_k \left(\frac{c_F}{1-\eta} (\|x_k^\delta - x_k\| + \|x_k - x^*\|) + \delta \right) \|F'_{i_k}(x_k^\delta)^* - F'_{i_k}(x_k)^*\| \\
381 \quad + c_F \|((F_{i_k}(x_k^\delta) - y_{i_k}^\delta) - (F_{i_k}(x_k) - y_{i_k}^\dagger))\| \\
382 \quad \leq \|x_k^\delta - x_k\| + 2\eta_k c_F \left(\frac{c_F}{1-\eta} (\|x_k^\delta - x_k\| + \|x_k - x^*\|) + \delta \right) \\
383 \quad + c_F \|((F_{i_k}(x_k^\delta) - y_{i_k}^\delta) - (F_{i_k}(x_k) - y_{i_k}^\dagger))\|,$$

385 This and mathematical induction shows that for any fixed k , $\sup_{(i_1, \dots, i_{k-1}) \in \mathcal{F}_k} \|x_k^\delta - x_k\|$
386 is uniformly bounded. Let $c = \frac{c_F}{1-\eta} \sup_{(i_1, \dots, i_{k-1}) \in \mathcal{F}_k} (\|x_k^\delta - x_k\| + \|x_k - x^*\|) + \delta$. Then it
387 follows from (3.3) that

$$388 \quad \lim_{\delta \rightarrow 0^+} \|x_{k+1}^\delta - x_{k+1}\| \leq \lim_{\delta \rightarrow 0^+} \|x_k^\delta - x_k\|^2 + c\eta_k \lim_{\delta \rightarrow 0^+} \|F'_{i_k}(x_k^\delta)^* - F'_{i_k}(x_k)^*\|$$

389
390

$$+ c_F \lim_{\delta \rightarrow 0^+} \|(F_{i_k}(x_k^\delta) - y_{i_k}^\delta) - (F_{i_k}(x_k) - y_{i_k}^\dagger)\|.$$

391 Then the desired assertion follows from the continuity of the operators F_i and F'_i in **As-**
392 **sumption 2.1(i)**, the induction hypothesis, and taking full expectation. \square

393 Now we can prove **Theorem 2.1** on the regularizing property of SGD.

394 *Proof of Theorem 2.1.* Let $\{\delta_n\}_{n \geq 1} \subset \mathbb{R}$ be a sequence converging to zero, and $y_n := y^{\delta_n}$
395 a corresponding sequence of noisy data. For each pair (δ_n, y_n) , we denote by $k_n = k(\delta_n)$
396 the stopping index. Further, we may assume that k_n increases strictly monotonically with
397 n . By **Proposition 3.1** and Young's inequality $2ab \leq \epsilon a^2 + \epsilon^{-1}b^2$, with the choice $a =$
398 $\mathbb{E}[\|F(x_k^\delta) - y^\delta\|^2]^{\frac{1}{2}}$, $b = (1 + \eta)\delta$ and $\epsilon = 1 - 2\eta > 0$:

$$\begin{aligned} 399 \quad \mathbb{E}[\|x^* - x_{k_{n+1}}^\delta\|^2] - \mathbb{E}[\|x^* - x_{k_n}^\delta\|^2] &\leq - (1 - 2\eta)\eta_k \mathbb{E}[\|F(x_{k_n}^\delta) - y^\delta\|^2] \\ 400 \quad &+ 2\eta_k(1 + \eta)\delta \mathbb{E}[\|F(x_{k_n}^\delta) - y^\delta\|^2]^{\frac{1}{2}} \leq \frac{(1 + \eta)^2}{1 - 2\eta} \eta_k \delta^2. \end{aligned}$$

402 Then for any $m < n$, summing the inequality with $\delta = \delta_n$ from k_m to $k_n - 1$ and applying
403 triangle inequality lead to

$$\begin{aligned} 404 \quad \mathbb{E}[\|x_{k_n}^{\delta_n} - x^*\|^2] &\leq \mathbb{E}[\|x_{k_m}^{\delta_n} - x^*\|^2] + \frac{(1 + \eta)^2}{1 - 2\eta} \delta_n^2 \sum_{j=k_m}^{k_n-1} \eta_j \\ 405 \quad &\leq 2\mathbb{E}[\|x_{k_m}^{\delta_n} - x_{k_m}\|^2] + 2\mathbb{E}[\|x_{k_m} - x^*\|^2] + \frac{(1 + \eta)^2}{1 - 2\eta} \delta_n^2 \sum_{j=1}^{k_n-1} \eta_j. \end{aligned}$$

407 By **Theorem 3.5**, we can fix a large m so that the term $\mathbb{E}[\|x_{k_m} - x^*\|^2]$ is sufficiently
408 small. Since the index k_m is fixed, we may apply **Lemma 3.6** to conclude that the term
409 $\mathbb{E}[\|x_{k_m}^{\delta_n} - x_{k_m}\|^2]$ tends to zero as $n \rightarrow \infty$. The last term also tends to zero under the
410 condition $\lim_{n \rightarrow \infty} \delta_n^2 \sum_{i=1}^{k_n} \eta_i = 0$. This completes the proof of the first assertion. The case
411 $\mathcal{N}(F'(x^\dagger)) \subset \mathcal{N}(F'(x))$ follows similarly as **Theorem 3.5**. \square

412 **4. Convergence rates.** Now we prove convergence rates for SGD under **Assump-**
413 **tion 2.1**, **Assumption 2.2(ii)** and **Assumption 2.3**; see **Theorem 4.8** and **Theorem 2.2** for the
414 results for exact and noisy data, respectively. We employ some shorthand notation. Let

$$415 \quad K_i = F'_i(x^\dagger), \quad K = \frac{1}{\sqrt{n}} \begin{pmatrix} K_1 \\ \vdots \\ K_n \end{pmatrix} \quad \text{and} \quad B = K^*K = \frac{1}{n} \sum_{i=1}^n K_i^* K_i.$$

416 Further, we frequently adopt the shorthand notation

$$417 \quad (4.1) \quad \Pi_j^k(B) = \prod_{i=j}^k (I - \eta_i B),$$

418 with the convention $\Pi_j^k(B) = I$ for $j > k$, and for $s \geq 0$ and $j \in \mathbb{N}$, we define,

$$419 \quad \tilde{s} = s + \frac{1}{2} \quad \text{and} \quad \phi_j^s = \|B^s \Pi_{j+1}^k(B)\|.$$

420 The rest of this section is organized as follows. By bias variance decomposition, we first
421 derive two important recursions for the mean $\|B^s \mathbb{E}[e_k^\delta]\|$ and variance $\mathbb{E}[\|B^s(e_k^\delta - \mathbb{E}[e_k^\delta])\|^2]$,
422 for any $s \geq 0$, in **subsection 4.1** and **subsection 4.2**, respectively, and then use the recursions
423 to derive convergence rates under *a priori* parameter choice in **subsection 4.3**.

424 **4.1. Recursion on the bias.** First, we derive a recursion on the bias of the SGD
 425 iterate x_k^δ . The following bound on the linearization error is useful.

426 LEMMA 4.1. *Under Assumption 2.1(iii), there holds*

$$427 \quad \|F(x) - F(x^\dagger) - K(x - x^\dagger)\| \leq \frac{c_B}{2} \|K(x - x^\dagger)\| \|x - x^\dagger\|.$$

428 Further, under Assumption 2.3, there holds

$$429 \quad \mathbb{E}[\|F(x_k^\delta) - F(x^\dagger) - K(x_k^\delta - x^\dagger)\|^2]^{\frac{1}{2}} \leq \frac{c_R}{1+\theta} \mathbb{E}[\|K(x_k^\delta - x^\dagger)\|^2]^{\frac{1}{2}} \mathbb{E}[\|x_k^\delta - x^\dagger\|^2]^{\frac{\theta}{2}}.$$

430 *Proof.* Let $z_t = tx + (1-t)x^\dagger$. By the mean value theorem and Assumption 2.1(iii),

$$431 \quad \|F(x) - F(x^\dagger) - K(x - x^\dagger)\| \leq \left\| \int_0^1 (F'(z_t) - K)(x - x^\dagger) dt \right\|$$

$$432 \quad \leq \int_0^1 \|(R_{z_t} - I)K(x - x^\dagger)\| dt \leq \frac{c_R}{2} \|K(x - x^\dagger)\| \|x - x^\dagger\|.$$

434 This shows the first estimate. Similarly, using Assumption 2.1(iii) and Assumption 2.3 with
 435 the choice $G(x) = K(x - x^\dagger)$, we obtain

$$436 \quad \mathbb{E}[\|F(x_k^\delta) - F(x^\dagger) - K(x_k^\delta - x^\dagger)\|^2]^{\frac{1}{2}} \leq \int_0^1 \mathbb{E}[\|(R_{z_t} - I)K(x_k^\delta - x^\dagger)\|^2]^{\frac{1}{2}} dt$$

$$437 \quad \leq c_R \mathbb{E}[\|K(x_k^\delta - x^\dagger)\|^2]^{\frac{1}{2}} \int_0^1 \mathbb{E}[\|z_t - x^\dagger\|^2]^{\frac{\theta}{2}} dt \leq \frac{c_R}{1+\theta} \mathbb{E}[\|K(x_k^\delta - x^\dagger)\|^2]^{\frac{1}{2}} \mathbb{E}[\|x_k^\delta - x^\dagger\|^2]^{\frac{\theta}{2}}.$$

439 This completes the proof of the lemma. \square

440 The next result gives a useful representation of the mean $\mathbb{E}[e_k^\delta]$ of the error $e_k^\delta \equiv x_k^\delta - x^\dagger$.

441 LEMMA 4.2. *Under Assumption 2.1(iii), the error e_k^δ satisfies*

$$442 \quad \mathbb{E}[e_{k+1}^\delta] = \Pi_1^k(B)e_1 + \sum_{j=1}^k \eta_j \Pi_{j+1}^k(B) K^*(-(y^\dagger - y^\delta) + \mathbb{E}[v_j]),$$

443 with the vector $v_k \in Y^n$ given by

$$444 \quad (4.2) \quad v_k = -(F(x_k^\delta) - F(x^\dagger) - K(x_k^\delta - x^\dagger)) + (I - R_{x_k^\delta}^*)(F(x_k^\delta) - y^\delta).$$

445 *Proof.* The definition of the SGD iterate x_k^δ in (1.3) and the relation $F'_{i_k}(x_k^\delta)^* =$
 446 $(R_{x_k^\delta}^{i_k} F'_{i_k}(x^\dagger))^* = K_{i_k}^* R_{x_k^\delta}^{i_k}$ from Assumption 2.1(iii) directly imply

$$447 \quad e_{k+1}^\delta = e_k^\delta - \eta_k K_{i_k}^* K_{i_k}(x_k^\delta - x^\dagger) - \eta_k K_{i_k}^*(y_{i_k}^\dagger - y_{i_k}^\delta) + \eta_k K_{i_k}^* v_{k,i_k},$$

449 with the random variable $v_{k,i}$ defined by

$$450 \quad (4.3) \quad v_{k,i} = -(F_i(x_k^\delta) - F_i(x^\dagger) - K_i(x_k^\delta - x^\dagger)) + (I - R_{x_k^\delta}^{i_k})(F_i(x_k^\delta) - y_i^\delta).$$

451 Thus, by the measurability of x_k^δ (and thus e_k^δ) with respect to \mathcal{F}_k , $\mathbb{E}[e_{k+1}^\delta | \mathcal{F}_k]$ is given by

$$452 \quad \mathbb{E}[e_{k+1}^\delta | \mathcal{F}_k] = (I - \eta_k B)e_k^\delta - \eta_k K^*(y^\dagger - y^\delta) + \eta_k K^* v_k.$$

454 Then taking full conditional and applying the recursion repeatedly complete the proof. \square

455 REMARK 4.1. The term v_k in (4.2) includes both the linearization error ($F(x_k^\delta) - F(x^\dagger) -$
456 $K(x_k^\delta - x^\dagger)$) of the nonlinear operator F and the range invariance of the derivative $F'(x)$ in
457 Assumption 2.1(ii)–(iii).

458 The next result gives a useful bound on $\mathbb{E}[v_j]$.

459 LEMMA 4.3. Under Assumption 2.1(i)–(iii), for v_j defined in (4.2), there holds

$$460 \quad \|\mathbb{E}[v_j]\| \leq \frac{(3-\eta)c_R}{2(1-\eta)} \mathbb{E}[\|e_j^\delta\|^2]^{\frac{1}{2}} \mathbb{E}[\|B^{\frac{1}{2}}e_j^\delta\|^2]^{\frac{1}{2}} + c_R \mathbb{E}[\|e_j^\delta\|^2]^{\frac{1}{2}} \delta.$$

461 *Proof.* By the triangle inequality, there holds

$$462 \quad \|\mathbb{E}[v_j]\| \leq \|\mathbb{E}[F(x_j^\delta) - F(x^\dagger) - K(x_j^\delta - x^\dagger)]\| + \|\mathbb{E}[(I - R_{x_j^\delta}^*)(F(x_j^\delta) - y^\delta)]\| := \text{I} + \text{II}.$$

464 The bound on I follows from Lemma 4.1 and Cauchy-Schwarz inequality as

$$465 \quad \text{I} \leq \frac{c_R}{2} \mathbb{E}[\|e_j^\delta\| \|Ke_j^\delta\|] \leq \frac{c_R}{2} \mathbb{E}[\|e_j^\delta\|^2]^{\frac{1}{2}} \mathbb{E}[\|Ke_j^\delta\|^2]^{\frac{1}{2}}.$$

467 For the term II, by triangle inequality, Cauchy-Schwarz inequality and Lemma 3.1,

$$468 \quad \begin{aligned} \text{II} &:= \|\mathbb{E}[(I - R_{x_j^\delta})(y^\delta - F(x_j^\delta))]\| \leq \mathbb{E}[\|(I - R_{x_j^\delta})(y^\delta - F(x_j^\delta))\|] \\ 469 \quad &\leq \frac{c_R}{1-\eta} \mathbb{E}[\|e_j^\delta\| \|Ke_j^\delta\|] + c_R \mathbb{E}[\|e_j^\delta\|] \delta \leq \mathbb{E}[\|e_j^\delta\|^2]^{\frac{1}{2}} \left(\frac{c_R}{1-\eta} \mathbb{E}[\|Ke_j^\delta\|^2]^{\frac{1}{2}} + c_R \delta \right). \end{aligned}$$

471 Combining these estimates with the identity $\|Ke_j^\delta\| = \|B^{\frac{1}{2}}e_j^\delta\|$ gives the assertion. \square

472 Last, we bound the error $\mathbb{E}[e_k^\delta]$ in a weighted norm. The cases $s = 0$ and $s = \frac{1}{2}$ will be
473 employed in the convergence analysis.

474 THEOREM 4.4. Under Assumption 2.1, for any $s \geq 0$, there holds

$$475 \quad \|B^s \mathbb{E}[e_{k+1}^\delta]\| \leq \phi_0^{s+\nu} \|w\| + \sum_{j=1}^k \eta_j \phi_j^s \left(\frac{(3-\eta)c_R}{2(1-\eta)} \mathbb{E}[\|e_j^\delta\|^2]^{\frac{1}{2}} \mathbb{E}[\|B^{\frac{1}{2}}e_j^\delta\|^2]^{\frac{1}{2}} + c_R \mathbb{E}[\|e_j^\delta\|^2]^{\frac{1}{2}} \delta + \delta \right).$$

477 *Proof.* By Lemma 4.2 and triangle inequality,

$$478 \quad \|B^s \mathbb{E}[e_{k+1}^\delta]\| \leq \text{I} + \sum_{j=1}^k \eta_j \text{II}_j.$$

479 with $\text{I} = \|B^s \Pi_1^k(B)(x_1 - x^\dagger)\|$ and $\text{II}_j = \|B^s \Pi_{j+1}^k(B) K^*(\mathbb{E}[v_j] - (y^\dagger - y^\delta))\|$. It suffices to
480 bound the terms I and II_j . By Assumption 2.1(iv),

$$481 \quad \text{I} = \|B^s \Pi_1^k(B) B^\nu w\| \leq \|\Pi_1^k(B) B^{s+\nu}\| \|w\|.$$

483 To bound the terms II_j , we have

$$484 \quad \text{II}_j \leq \|B^s \Pi_{j+1}^k(B) K^*(\mathbb{E}[v_j] - (y^\dagger - y^\delta))\| \leq \|B^{s+\frac{1}{2}} \Pi_{j+1}^k(B)\| (\|\mathbb{E}[v_j]\| + \delta).$$

486 This, Lemma 4.3 and the notation ϕ_j^s complete the proof. \square

487 REMARK 4.2. The bound on $\mathbb{E}[e_k^\delta]$ depends on the variance of the iterate x_k^δ (via the
488 terms like $\mathbb{E}[\|e_k^\delta\|^2]$ etc.), which differs from the linear case [14]. This is one of the com-
489 plications for nonlinear inverse problems. The weighted norm $\|B^s \mathbb{E}[e_k^\delta]\|$ is useful since the
490 upper bound in Theorem 4.4 involves $\mathbb{E}[\|B^{\frac{1}{2}}e_k^\delta\|^2]$, i.e., $s = \frac{1}{2}$. For linear inverse problems,
491 $R_x = I$ and $c_R = 0$, and the recursion simplifies to $\|B^s \mathbb{E}[e_{k+1}^\delta]\| \leq \phi_0^{s+\nu} \|w\| + \sum_{j=1}^k \eta_j \phi_j^s \delta$,
492 i.e., the approximation error and data error, respectively.

493 **4.2. Recursion on variance.** Now we turn to the computational variance $\mathbb{E}[\|B^s(x_k^\delta -$
494 $\mathbb{E}[x_k^\delta])\|^2]$, which arises from the random index i_k . First, we bound on the variance in terms
495 of iteration noises $N_{j,1}$ and $N_{j,2}$ (defined in (4.4) below).

496 LEMMA 4.5. Under *Assumption 2.1(iii)*, for the SGD iterate x_k^δ , there holds

$$497 \quad \mathbb{E}[\|B^s(x_{k+1}^\delta - \mathbb{E}[x_{k+1}^\delta])\|^2] \leq \sum_{j=1}^k \eta_j^2 (\phi_j^\delta)^2 \mathbb{E}[\|N_{j,1}\|^2] + 2 \sum_{i=1}^k \sum_{j=i}^k \eta_i \eta_j \phi_i^\delta \phi_j^\delta \mathbb{E}[\|N_{i,1}\| \|N_{j,2}\|]$$

$$498 \quad + \sum_{i=1}^k \sum_{j=1}^k \eta_i \eta_j \phi_i^\delta \phi_j^\delta \mathbb{E}[\|N_{i,2}\| \|N_{j,2}\|],$$

499

500 with the random variables $N_{j,1}$ and $N_{j,2}$ respectively given by

$$501 \quad (4.4) \quad \begin{aligned} N_{j,1} &= (K(x_j^\delta - x^\dagger) - K_{i_j}(x_j^\delta - x^\dagger)\varphi_{i_j}) + ((y^\dagger - y^\delta) - (y_{i_j}^\dagger - y_{i_j}^\delta)\varphi_{i_j}), \\ N_{j,2} &= -\mathbb{E}[v_j] + v_{j,i_j}\varphi_{i_j}, \end{aligned}$$

502 where v_k and $v_{k,i}$ are given in (4.2) and (4.3), and $\varphi_i = (0, \dots, 0, n^{\frac{1}{2}}, 0, \dots, 0)$ denotes the
503 canonical i th Cartesian basis vector in \mathbb{R}^n scaled by $n^{\frac{1}{2}}$.

504 *Proof.* Similar to the proof of Lemma 4.2, we rewrite the SGD iteration (1.3) as

$$505 \quad (4.5) \quad x_{k+1}^\delta = x_k^\delta - \eta_k K_{i_k}^* K_{i_k}(x_k^\delta - x^\dagger) - \eta_k K_{i_k}^*(y_{i_k}^\dagger - y_{i_k}^\delta) + \eta_k K_{i_k}^* v_{k,i_k},$$

506 with $v_{k,i}$ defined in (4.3). By the definition of v_k in (4.2) and the measurability of x_k^δ with
507 respect to \mathcal{F}_k , we obtain

$$508 \quad \mathbb{E}[x_{k+1}^\delta | \mathcal{F}_k] = x_k^\delta - \eta_k B(x_k^\delta - x^\dagger) - \eta_k K^*(y^\dagger - y^\delta) + \eta_k K^* v_k.$$

510 Taking full conditional yields

$$511 \quad (4.6) \quad \mathbb{E}[x_{k+1}^\delta] = \mathbb{E}[x_k^\delta] - \eta_k B \mathbb{E}[x_k^\delta - x^\dagger] - \eta_k K^*(y^\dagger - y^\delta) + \eta_k K^* \mathbb{E}[v_k].$$

513 Thus, subtracting (4.6) from (4.5) shows that $z_k := x_k^\delta - \mathbb{E}[x_k^\delta]$ satisfies

$$514 \quad (4.7) \quad z_{k+1} = (I - \eta_k B)z_k + \eta_k M_k,$$

516 with $z_1 = 0$ and the iteration noise M_j given by $M_j = M_{j,1} + M_{j,2}$, where

$$517 \quad \begin{aligned} M_{j,1} &= (B(x_j^\delta - x^\dagger) - K_{i_j}^* K_{i_j}(x_j^\delta - x^\dagger)) + (K^*(y^\dagger - y^\delta) - K_{i_j}^*(y_{i_j}^\dagger - y_{i_j}^\delta)), \\ 518 \quad M_{j,2} &= -(K^* \mathbb{E}[v_j] - K_{i_j}^* v_{j,i_j}). \end{aligned}$$

520 Repeatedly applying the recursion (4.7) with $z_1 = 0$ leads to

$$521 \quad z_{k+1} = \sum_{j=1}^k \eta_j \Pi_{j+1}^k(B) M_j.$$

522 With the decomposition of $M_j = M_{j,1} + M_{j,2}$, we directly obtain

$$523 \quad \mathbb{E}[\|B^s z_{k+1}\|^2] = \sum_{i=1}^k \sum_{j=1}^k \eta_i \eta_j \mathbb{E}[\langle B^s \Pi_{i+1}^k(B) M_{i,1}, B^s \Pi_{j+1}^k(B) M_{j,1} \rangle]$$

$$\begin{aligned}
524 \quad & + 2 \sum_{i=1}^k \sum_{j=1}^k \eta_i \eta_j \mathbb{E}[\langle B^s \Pi_{i+1}^k(B) M_{i,1}, B^s \Pi_{j+1}^k(B) M_{j,2} \rangle] \\
525 \quad & + \sum_{i=1}^k \sum_{j=1}^k \eta_i \eta_j \mathbb{E}[\langle B^s \Pi_{i+1}^k(B) M_{i,2}, B^s \Pi_{j+1}^k(B) M_{j,2} \rangle] := \text{I} + \text{II} + \text{III}. \\
526 \quad &
\end{aligned}$$

527 Below we simplify the three terms. Since x_j^δ is measurable with respect to \mathcal{F}_j , we have
528 $\mathbb{E}[M_{j,1} | \mathcal{F}_j] = 0$, which directly implies the independence $\mathbb{E}[\langle B^s M_{i,1}, B^s M_{j,1} \rangle] = 0$, $i \neq j$.
529 Indeed, for $i > j$, $\mathbb{E}[\langle B^s M_{i,1}, B^s M_{j,1} \rangle | \mathcal{F}_i] = \langle B^s \mathbb{E}[M_{i,1} | \mathcal{F}_i], B^s M_{j,1} \rangle = 0$, and taking full
530 conditional yields the claim. Thus, the term I simplifies to

$$531 \quad \text{I} = \sum_{j=1}^k \eta_j^2 \mathbb{E}[\|B^s \Pi_{j+1}^k(B) M_{j,1}\|^2].$$

532 Further, for $i > j$, a similar argument yields $\mathbb{E}[\langle B^s M_{i,1}, B^s M_{j,2} \rangle] = 0$ and thus

$$533 \quad \text{II} = 2 \sum_{i=1}^k \sum_{j=i}^k \eta_i \eta_j \mathbb{E}[\langle B^s \Pi_{i+1}^k M_{i,1}, B^s \Pi_{j+1}^k M_{j,2} \rangle].$$

534 Now we further simplify $M_{j,1}$ and $M_{j,2}$. By the definitions of $N_{j,1}$ and $N_{j,2}$, with $(K^*)^\dagger$
535 being the pseudoinverse of K^* , we have $(K^*)^\dagger M_j = N_{j,1} + N_{j,2}$. Thus, by triangle inequality,

$$\begin{aligned}
536 \quad \mathbb{E}[\|B^s z_{k+1}\|^2] & \leq \sum_{j=1}^k \eta_j^2 \mathbb{E}[\|B^{s+\frac{1}{2}} \Pi_{j+1}^k(B)\|^2 \|N_{j,1}\|^2] \\
537 \quad & + 2 \sum_{i=1}^k \sum_{j=i}^k \eta_i \eta_j \|B^{s+\frac{1}{2}} \Pi_{i+1}^k(B)\| \|B^{s+\frac{1}{2}} \Pi_{j+1}^k(B)\| \mathbb{E}[\|N_{i,1}\| \|N_{j,2}\|] \\
538 \quad & + \sum_{i=1}^k \sum_{j=1}^k \eta_i \eta_j \|B^{s+\frac{1}{2}} \Pi_{i+1}^k(B)\| \|B^{s+\frac{1}{2}} \Pi_{j+1}^k(B)\| \mathbb{E}[\|N_{i,2}\| \|N_{j,2}\|]. \\
539 \quad &
\end{aligned}$$

540 This completes the proof of the lemma. \square

541 The next result bounds the iteration noises $N_{j,1}$ and $N_{j,2}$.

542 **LEMMA 4.6.** *Under [Assumption 2.1](#) (i)–(iii) and [Assumption 2.3](#), for $N_{j,1}$ and $N_{j,2}$ de-*
543 *defined in [\(4.4\)](#), there hold*

$$544 \quad (4.8) \quad \mathbb{E}[\|N_{j,1}\|^2]^{\frac{1}{2}} \leq n^{\frac{1}{2}} (\mathbb{E}[\|B^{\frac{1}{2}} e_j^\delta\|^2]^{\frac{1}{2}} + \delta),$$

$$545 \quad (4.9) \quad \mathbb{E}[\|N_{j,2}\|^2]^{\frac{1}{2}} \leq n^{\frac{1}{2}} \left(\frac{c_R(2+\theta-\eta)}{(1+\theta)(1-\eta)} \mathbb{E}[\|B^{\frac{1}{2}} e_j^\delta\|^2]^{\frac{1}{2}} + c_R \delta \right) \mathbb{E}[\|e_j^\delta\|^2]^{\frac{\theta}{2}}.$$

547 *Proof.* By the measurability of x_j^δ with respect to \mathcal{F}_j , we have $\mathbb{E}[K_{i_j}(x_j^\delta - x^\dagger) \varphi_{i_j} | \mathcal{F}_j] =$
548 $K(x_j^\delta - x^\dagger)$. Then by bias-variance decomposition, we have

$$\begin{aligned}
549 \quad & \mathbb{E}[\|(K(x_j^\delta - x^\dagger) - K_{i_j}(x_j^\delta - x^\dagger)) \varphi_{i_j}\|^2 | \mathcal{F}_j] \leq \mathbb{E}[\|K_{i_j}(x_j^\delta - x^\dagger) \varphi_{i_j}\|^2 | \mathcal{F}_j] \\
550 \quad & = n^{-1} \sum_{i=1}^n \|K_i(x_j^\delta - x^\dagger)\|^2 n = n \|K(x_j^\delta - x^\dagger)\|^2, \\
551 \quad &
\end{aligned}$$

552 and then by taking full expectation, we obtain

$$553 \quad \mathbb{E}[\|(K(x_j^\delta - x^\dagger) - K_{i_j}(x_j^\delta - x^\dagger)\varphi_{i_j})\|^2]^{\frac{1}{2}} \leq n^{\frac{1}{2}}\mathbb{E}[\|K(x_j^\delta - x^\dagger)\|^2]^{\frac{1}{2}}.$$

554 Similarly, $\mathbb{E}[\|(y^\dagger - y^\delta) - (y_{i_j}^\dagger - y_{i_j}^\delta)\varphi_{i_j}\|^2]^{\frac{1}{2}} \leq n^{\frac{1}{2}}\delta$. This and triangle inequality show the
 555 estimate (4.8). Similarly, by the measurability of x_j^δ with respect to \mathcal{F}_j and bias variance
 556 decomposition, we deduce (with $\mathbb{E}_{\mathcal{F}_j}$ denoting taking expectation in \mathcal{F}_j)

$$557 \quad \mathbb{E}[\|(\mathbb{E}[v_j] - v_{j,i_j}\varphi_{i_j})\|^2] \leq \mathbb{E}_{\mathcal{F}_j}[\mathbb{E}[\|v_{j,i_j}\varphi_{i_j}\|^2|\mathcal{F}_j]] = n\mathbb{E}[\|v_j\|^2],$$

559 i.e., $\mathbb{E}[\|(\mathbb{E}[v_j] - v_{j,i_j}\varphi_{i_j})\|^2]^{\frac{1}{2}} \leq n^{\frac{1}{2}}\mathbb{E}[\|v_j\|^2]^{\frac{1}{2}}$. Then by triangle inequality, [Assumption 2.3](#)
 560 and [Lemma 4.1](#),

$$561 \quad \mathbb{E}[\|v_j\|^2]^{\frac{1}{2}} \leq \mathbb{E}[\|(F(x_j^\delta) - F(x^\dagger) - K(x_j^\delta - x^\dagger))\|^2]^{\frac{1}{2}} + \mathbb{E}[\|(I - R_{x_j^\delta}^*)(F(x_j^\delta) - y^\delta)\|^2]^{\frac{1}{2}} \\
 562 \quad \leq \frac{c_R}{1+\theta}\mathbb{E}[\|Ke_j^\delta\|^2]^{\frac{1}{2}}\mathbb{E}[\|e_j^\delta\|^2]^{\frac{\theta}{2}} + c_R(\frac{1}{1-\eta}\mathbb{E}[\|Ke_j^\delta\|^2]^{\frac{1}{2}} + \delta)\mathbb{E}[\|e_j^\delta\|^2]^{\frac{\theta}{2}} \\
 563 \quad = (\frac{(2+\theta-\eta)c_R}{(1+\theta)(1-\eta)}\mathbb{E}[\|Ke_j^\delta\|^2]^{\frac{1}{2}} + c_R\delta)\mathbb{E}[\|e_j^\delta\|^2]^{\frac{\theta}{2}}.$$

565 This completes the proof of the lemma. \square

566 **REMARK 4.3.** Note that the convergence analysis in [14] relies on the independence
 567 $\mathbb{E}[(B^s M_j, B^s M_\ell)] = 0$ for $j \neq \ell$. This identity is no longer valid for nonlinear inverse
 568 problems, although it still holds for the linear part $M_{j,1}$: $\mathbb{E}[(B^s M_{j,1}, B^s M_{\ell,1})] = 0$ for $j \neq \ell$.
 569 The conditional dependence among the iteration noises $M_{j,2}$ poses one big challenge to the
 570 convergence analysis, and the splitting of the conditionally dependent and independent com-
 571 ponents will plays a role in the analysis below. [Assumption 2.3](#) is to compensate the condi-
 572 tional dependence.

573 **REMARK 4.4.** The constants in [Lemma 4.6](#) involve an unpleasant dependence on n as
 574 $n^{\frac{1}{2}}$, due to the variance inflation of the estimated gradient. It can be reduced by various
 575 strategies, e.g., mini-batch or variance reduction.

576 Last, we give a bound on the variance $\mathbb{E}[\|B^s(x_k^\delta - \mathbb{E}[x_k^\delta])\|^2]$. This result will play an
 577 important role in the error analysis in [subsection 4.3](#).

578 **THEOREM 4.7.** Let [Assumption 2.1\(i\)-\(iii\)](#) and [Assumption 2.3](#) be fulfilled. Then for
 579 any $s \in [0, \frac{1}{2}]$, there holds

$$580 \quad \mathbb{E}[\|B^s(\mathbb{E}[x_{k+1}^\delta] - x_{k+1}^\delta)\|^2] \leq n \sum_{j=1}^k \eta_j^2 (\phi_j^s)^2 (\mathbb{E}[\|B^{\frac{1}{2}}e_j^\delta\|^2]^{\frac{1}{2}} + \delta)^2 \\
 581 \quad + 2n \sum_{i=1}^k \sum_{j=i}^k \eta_i \eta_j \phi_i^s \phi_j^s (\mathbb{E}[\|B^{\frac{1}{2}}e_i^\delta\|^2]^{\frac{1}{2}} + \delta) (\frac{(2+\theta-\eta)c_R}{(1+\theta)(1-\eta)}\mathbb{E}[\|B^{\frac{1}{2}}e_j^\delta\|^2]^{\frac{1}{2}} + c_R\delta) \mathbb{E}[\|e_j^\delta\|^2]^{\frac{\theta}{2}} \\
 582 \quad + n \left(\sum_{j=1}^k \eta_j \phi_j^s (\frac{(2+\theta-\eta)c_R}{(1+\theta)(1-\eta)}\mathbb{E}[\|B^{\frac{1}{2}}e_j^\delta\|^2]^{\frac{1}{2}} + c_R\delta) \mathbb{E}[\|e_j^\delta\|^2]^{\frac{\theta}{2}} \right)^2.$$

584 *Proof.* The assertion follows directly from [Lemma 4.5](#) and [Lemma 4.6](#). \square

585 **4.3. Convergence rates.** This part is devoted to convergence rates analysis of SGD
 586 under [Assumption 2.1\(ii\)](#). We analyze the cases of exact and noisy data separately. For exact
 587 data, the bounds involve constants that are more transparent in terms of their dependence
 588 on various algorithmic parameters. First we analyze the case of exact data y^\dagger , and the bound

589 boils down to the approximation error and computational variance. Further, we assume that
 590 $\|B\| \leq 1$ and $\eta_0 \leq 1$ below, which can be easily achieved by rescaling the operator F and
 591 the data y^\dagger/y^δ . The analysis relies heavily on various technical estimates in [Appendix A](#),
 592 especially [Proposition A.1](#).

593 **THEOREM 4.8.** *Let [Assumption 2.1](#), [Assumption 2.2\(ii\)](#) and [Assumption 2.3](#) be fulfilled*
 594 *with $\|w\|$, θ and η_0 being sufficiently small. Then the error $e_k = x_k - x^\dagger$ satisfies*

$$595 \quad \mathbb{E}[\|e_k\|^2] \leq c^* \|w\|^2 k^{-\min(2\nu(1-\alpha), \alpha-\epsilon)}, \quad \mathbb{E}[\|B^{\frac{1}{2}} e_k\|^2] \leq c^* \|w\|^2 k^{-\min((1+2\nu)(1-\alpha), 1-\epsilon)},$$

597 where $\epsilon \in (0, \frac{\alpha}{2})$ is small, and c^* is independent of k , but depends on α, ν, η_0, n , and θ .

598 *Proof.* For any $s \geq 0$, [Theorem 4.4](#) and [Theorem 4.7](#) give (with $c_0 = \frac{(2+\theta-\eta)c_R}{(1+\theta)(1-\eta)}$)

$$599 \quad \mathbb{E}[\|B^s e_{k+1}\|^2] \leq \left(c_0 \sum_{j=1}^k \eta_j \phi_j^{\bar{s}} \mathbb{E}[\|e_j\|^2]^{\frac{1}{2}} \mathbb{E}[\|B^{\frac{1}{2}} e_j\|^2]^{\frac{1}{2}} + \phi_0^{s+\nu} \|w\| \right)^2$$

$$600 \quad (4.10) \quad + 2nc_0 \left(\sum_{i=1}^k \eta_i \phi_i^{\bar{s}} \mathbb{E}[\|B^{\frac{1}{2}} e_i\|^2]^{\frac{1}{2}} \right) \left(\sum_{j=1}^k \eta_j \phi_j^{\bar{s}} \mathbb{E}[\|B^{\frac{1}{2}} e_j\|^2]^{\frac{1}{2}} \mathbb{E}[\|e_j\|^2]^{\frac{\theta}{2}} \right)$$

$$601 \quad + nc_0^2 \left(\sum_{j=1}^k \eta_j \phi_j^{\bar{s}} \mathbb{E}[\|B^{\frac{1}{2}} e_j\|^2]^{\frac{1}{2}} \mathbb{E}[\|e_j\|^2]^{\frac{\theta}{2}} \right)^2 + n \sum_{j=1}^k \eta_j^2 (\phi_j^{\bar{s}})^2 \mathbb{E}[\|B^{\frac{1}{2}} e_j\|^2].$$

603 Under [Assumption 2.2\(ii\)](#), [Lemma A.1](#) and [Lemma A.2](#) directly give

$$604 \quad \phi_0^{s+\nu} \leq \frac{(s+\nu)^{s+\nu}}{e^{s+\nu} (\sum_{i=1}^k \eta_i)^{s+\nu}} \leq \frac{(s+\nu)^{s+\nu} (1-\alpha)^{\nu+s}}{e^{s+\nu} \eta_0^{\nu+s} (1-2\alpha)^{\nu+s}} (k+1)^{-(1-\alpha)(\nu+s)}.$$

606 Note that the function $\frac{s^s}{e^s}$ is decreasing in s over the interval $[0, 1]$, and the function $\frac{1-\alpha}{1-2\alpha-1}$
 607 is decreasing in α over the interval $[0, 1]$ (and upper bounded by 2). Thus, for $\eta_0 \leq 1$ and
 608 any $0 \leq \nu, s \leq \frac{1}{2}$, there holds (with $c_\nu = \frac{2\nu^\nu}{\eta_0 e^\nu}$)

$$609 \quad (4.11) \quad \phi_0^{s+\nu} \leq c_\nu (k+1)^{-(\nu+s)(1-\alpha)}.$$

611 Let $a_j \equiv \mathbb{E}[\|e_j\|^2]$ and $b_j \equiv \mathbb{E}[\|B^{\frac{1}{2}} e_j\|^2]$. Since $\|B\| \leq 1$, we have $\phi_j^s \leq \phi_j^{\bar{s}}$ for any $0 \leq \bar{s} \leq s$.
 612 Then setting $s = 0$ and $s = 1/2$ in the recursion [\(4.10\)](#) and applying [\(4.11\)](#) lead to

$$613 \quad a_{k+1} \leq \left(c_0 \sum_{j=1}^k \eta_j \phi_j^{\frac{1}{2}} a_j^{\frac{1}{2}} b_j^{\frac{1}{2}} + c_\nu \|w\| (k+1)^{-\nu(1-\alpha)} \right)^2 + n \sum_{j=1}^k \eta_j^2 (\phi_j^{\frac{1}{2}})^2 b_j$$

$$614 \quad (4.12) \quad + 2nc_0 \left(\sum_{i=1}^k \eta_i \phi_i^{\frac{1}{2}} b_i^{\frac{1}{2}} \right) \left(\sum_{j=1}^k \eta_j \phi_j^{\frac{1}{2}} b_j^{\frac{1}{2}} a_j^{\frac{\theta}{2}} \right) + nc_0^2 \left(\sum_{j=1}^k \eta_j \phi_j^{\frac{1}{2}} b_j^{\frac{1}{2}} a_j^{\frac{\theta}{2}} \right)^2,$$

$$615 \quad b_{k+1} \leq \left(c_0 \sum_{j=1}^k \eta_j \phi_j^1 a_j^{\frac{1}{2}} b_j^{\frac{1}{2}} + c_\nu \|w\| (k+1)^{-(\frac{1}{2}+\nu)(1-\alpha)} \right)^2 + n \left(\sum_{j=1}^{\lfloor \frac{k}{2} \rfloor} \eta_j^2 (\phi_j^r)^2 b_j \right.$$

$$616 \quad (4.13) \quad \left. + \sum_{j=\lfloor \frac{k}{2} \rfloor + 1}^k \eta_j^2 (\phi_j^{\frac{1}{2}})^2 b_j \right) + 2nc_0 \left(\sum_{i=1}^k \eta_i \phi_i^1 b_i^{\frac{1}{2}} \right) \left(\sum_{j=1}^k \eta_j \phi_j^1 b_j^{\frac{1}{2}} a_j^{\frac{\theta}{2}} \right) + nc_0^2 \left(\sum_{j=1}^k \eta_j \phi_j^1 b_j^{\frac{1}{2}} a_j^{\frac{\theta}{2}} \right)^2,$$

617 with $r = \min(\frac{1}{2} + \nu, \frac{1-\epsilon}{2(1-\alpha)}) \in (\frac{1}{2}, 1)$. The rest of the proof is to prove

$$619 \quad (4.14) \quad a_k \leq c^* \|w\|^2 k^{-\beta} \quad \text{and} \quad b_k \leq c^* \|w\|^2 k^{-\gamma}.$$

620

621 where $\beta = \min(2\nu(1 - \alpha), \alpha - \epsilon)$ and $\gamma = \min((1 + 2\nu)(1 - \alpha), 1 - \epsilon)$, and $c^* > 0$ is to be
 622 specified. The proof is based on mathematical induction. When $k = 1$, (4.14) holds trivially
 623 for any large c^* . Now we assume that (4.14) holds up to the case k , and prove it for the case
 624 $k + 1$. Actually, it follows from (4.12) and the induction hypothesis that (with $\varrho = c^* \|w\|^2$)

$$\begin{aligned}
 625 \quad a_{k+1} &\leq \left(c_0 \varrho \sum_{j=1}^k \eta_j \phi_j^{\frac{1}{2}} j^{-\frac{\beta+\gamma}{2}} + c_\nu \|w\| (k+1)^{-\nu(1-\alpha)} \right)^2 + n \varrho \sum_{j=1}^k \eta_j^2 (\phi_j^{\frac{1}{2}})^2 j^{-\gamma} \\
 626 \quad (4.15) \quad &+ 2nc_0 \varrho^{1+\frac{\theta}{2}} \left(\sum_{i=1}^k \eta_i \phi_i^{\frac{1}{2}} i^{-\frac{\gamma}{2}} \right) \left(\sum_{j=1}^k \eta_j \phi_j^{\frac{1}{2}} j^{-\frac{\gamma+\theta\beta}{2}} \right) + nc_0^2 \varrho^{1+\theta} \left(\sum_{j=1}^k \eta_j \phi_j^{\frac{1}{2}} j^{-\frac{\gamma+\theta\beta}{2}} \right)^2. \\
 627
 \end{aligned}$$

628 Next we bound the terms on the right-hand side. By Proposition A.1, we have

$$\begin{aligned}
 629 \quad \sum_{j=1}^k \eta_j \phi_j^{\frac{1}{2}} j^{-\frac{\gamma}{2}} &\leq c_1 (k+1)^{-\frac{\beta}{2}} \quad \text{and} \quad \sum_{j=1}^k \eta_j^2 (\phi_j^{\frac{1}{2}})^2 j^{-\gamma} \leq c_2 (k+1)^{-\beta}, \\
 630
 \end{aligned}$$

631 with $c_1 = 2^{\frac{\beta}{2}} \eta_0^{\frac{1}{2}} (2^{-1} B(\frac{1}{2}, \zeta) + 1)$, $\zeta = (\frac{1}{2} - \nu)(1 - \alpha) > 0$, and $c_2 = 2^\beta \eta_0 (\alpha^{-1} + 2)$. Then we
 632 derive from (4.15) that

$$\begin{aligned}
 633 \quad (4.16) \quad a_{k+1} &\leq ((c_0 c_1 \varrho + c_\nu \|w\|)^2 + nc_2 \varrho + 2nc_0 c_1^2 \varrho^{1+\frac{\theta}{2}} + nc_0^2 c_1^2 \varrho^{1+\theta}) (k+1)^{-\beta}.
 \end{aligned}$$

635 Next we bound b_k similarly. It follows from (4.13) (with $r = \min(\frac{1}{2} + \nu, \frac{1-\epsilon}{2(1-\alpha)}) \in (\frac{1}{2}, 1)$)
 636 and the induction hypothesis that

$$\begin{aligned}
 637 \quad b_{k+1} &\leq \left(c_0 \varrho \sum_{j=1}^k \eta_j \phi_j^1 j^{-\frac{\beta+\gamma}{2}} + c_\nu \|w\| (k+1)^{-(\frac{1}{2}+\nu)(1-\alpha)} \right)^2 \\
 638 \quad (4.17) \quad &+ n \varrho \left(\sum_{j=1}^{\lfloor \frac{k}{2} \rfloor} \eta_j^2 (\phi_j^r)^2 j^{-\gamma} + \sum_{j=\lfloor \frac{k}{2} \rfloor + 1}^k \eta_j^2 (\phi_j^{\frac{1}{2}})^2 j^{-\gamma} \right) \\
 639 \quad &+ 2nc_0 \varrho^{1+\frac{\theta}{2}} \left(\sum_{i=1}^k \eta_i \phi_i^1 i^{-\frac{\gamma}{2}} \right) \left(\sum_{j=1}^k \eta_j \phi_j^1 j^{-\frac{\gamma+\theta\beta}{2}} \right) + nc_0^2 \varrho^{1+\theta} \left(\sum_{j=1}^k \eta_j \phi_j^1 j^{-\frac{\gamma+\theta\beta}{2}} \right)^2. \\
 640
 \end{aligned}$$

641 By Proposition A.1, there hold

$$\begin{aligned}
 642 \quad \sum_{j=1}^k \eta_j \phi_j^1 j^{-\frac{\beta+\gamma}{2}} &\leq c'_1 (k+1)^{-\frac{\gamma}{2}}, \quad \sum_{j=1}^{\lfloor \frac{k}{2} \rfloor} \eta_j^2 (\phi_j^r)^2 j^{-\gamma} + \sum_{j=\lfloor \frac{k}{2} \rfloor + 1}^k \eta_j^2 (\phi_j^{\frac{1}{2}})^2 j^{-\gamma} \leq c'_2 (k+1)^{-\gamma}, \\
 643 \quad \left(\sum_{i=1}^k \eta_i \phi_i^1 i^{-\frac{\gamma}{2}} \right) \left(\sum_{j=1}^k \eta_j \phi_j^1 j^{-\frac{\gamma+\theta\beta}{2}} \right) &\leq c'_3 (k+1)^{-\gamma}, \quad \sum_{j=1}^k \eta_j \phi_j^1 j^{-\frac{\gamma+\theta\beta}{2}} \leq c'_4 (k+1)^{-\frac{\gamma}{2}}, \\
 644
 \end{aligned}$$

645 with $c'_1 = 2^{\frac{\gamma}{2}} (\zeta^{-1} + 2\beta^{-1} + 1)$, $c'_2 = 2^\gamma \eta_0^{2-2r} (3\alpha^{-1} + 1)$, $c'_3 = 2^{\frac{\gamma}{2}} (((\frac{1}{2} - \nu - \theta\nu)(1 - \alpha))^{-1} +$
 646 $4(\theta\beta)^{-1} + 1)$ and $c'_4 = 2^{\frac{\gamma}{2}} (\zeta^{-1} + 2(\theta\beta)^{-1} + 1)$. These estimates and (4.17) yield

$$\begin{aligned}
 647 \quad (4.18) \quad b_{k+1} &\leq ((c_0 c'_1 \varrho + c_\nu \|w\|)^2 + nc'_2 \varrho + 2nc_0 c'_3 \varrho^{1+\frac{\theta}{2}} + nc_0^2 c'_4 \varrho^{1+\theta}) (k+1)^{-\gamma}. \\
 648
 \end{aligned}$$

649 In view of (4.16) and (4.18), upon dividing by ϱ , assertion (4.14) holds if we can show the
650 existence of a $c^* > 0$ such that

$$651 \quad (c_0 c_1 \varrho^{\frac{1}{2}} + c_\nu c^{*-\frac{1}{2}})^2 + n c_2 + 2n c_0 c_1^2 \varrho^{\frac{\theta}{2}} + n c_0^2 c_1^2 \varrho^\theta \leq 1,$$

$$652 \quad (c_0 c'_1 \varrho^{\frac{1}{2}} + c_\nu c^{*-\frac{1}{2}})^2 + n c'_2 + 2n c_0 c_3'^2 \varrho^{\frac{\theta}{2}} + n c_0^2 c_4'^2 \varrho^\theta \leq 1.$$

654 Since the constants c_2 and c'_2 are proportional to η_0 and η_0^{2-2r} (with the exponent $1 >$
655 $2 - 2r > 0$), respectively, for sufficiently small η_0 , there holds $n \max(c_2, c'_2) < 1$. Now for
656 sufficiently small $\|w\|$ and large c^* such that ρ is small, the above two inequalities hold. This
657 completes the induction step and the proof of the theorem. \square

658 **REMARK 4.5.** $\mathbb{E}[\|B^{\frac{1}{2}} e_k\|^2]$ decays as $\mathbb{E}[\|B^{\frac{1}{2}} e_k\|^2] \leq ck^{-\min((1+2\nu)(1-\alpha), 1-\epsilon)}$, which, for α
659 close to unit, is comparable with that for the Landweber method [8]: $\|B^{\frac{1}{2}} e_k\| \leq ck^{-(\nu+\frac{1}{2})(1-\alpha)}$.
660 The factor $k^{-(1-\epsilon)}$ limits the fastest possible rate. This restriction arises from the compu-
661 tational variance, due to the random selection of the row index i_k , which limits the conver-
662 gence rate $\mathbb{E}[\|e_k\|^2]$ to $O(k^{-\min(2\nu(1-\alpha), \alpha-\epsilon)})$. Thus for order optimality, the largest possible
663 smoothness index is $\nu = \frac{1}{2}$, beyond which SGD suffers from suboptimality, similar to the
664 Landweber method for nonlinear inverse problems [8]. Further, it shows the impact of the
665 exponent α : a smaller α may restrict the error $\mathbb{E}[\|e_k\|^2]$ to $O(k^{-(\alpha-\epsilon)})$.

666 **REMARK 4.6.** The exponent α in the step size schedule in Assumption 2.2(ii) enters
667 into the constant c^* via the constants c_1, \dots, c'_4 etc, and the constant c_0 is independent of α .
668 The constants c_1, \dots, c'_4 blow up either like $(1-\alpha)^{-1}$ as $\alpha \rightarrow 1^-$, according to the well-known
669 asymptotic behavior of the Beta function, or like α^{-1} as $\alpha \rightarrow 0^+$. These dependencies partly
670 exhibit the delicacy of choosing a proper step size schedule for SGD.

671 **REMARK 4.7.** We briefly comment on the “smallness” conditions on w , η_0 and θ in
672 the analysis. The smallness assumption on w in the source condition in Assumption 2.1(iv)
673 appears also for the classical Landweber method [8] and the standard Tikhonov regularization
674 [5, 11], and thus it is not surprising. The smallness condition on η_0 is to control the influence
675 of the computational variance, and in a slightly different context of statistical learning theory,
676 similar conditions also appear in the convergence analysis of variants of SGD. The smallness
677 condition on θ is only to facilitate the analysis, i.e., a concise form of the constant c'_3 , and the
678 assumption can be removed at the expense of a less transparent (but more benign) expression
679 for c'_3 ; see the proof in Proposition A.1 and also Remark A.1.

680 Last, we prove the main result in this work, i.e., Theorem 2.2, which gives the conver-
681 gence rate of SGD (1.3) for noisy data y^δ .

682 *Proof of Theorem 2.2.* The main proof strategy is similar to that of Theorem 4.8. Let
683 $a_j \equiv \mathbb{E}[\|e_j^\delta\|^2]$ and $b_j \equiv \mathbb{E}[\|B^{\frac{1}{2}} e_j^\delta\|^2]$. Then with $c_0 = \frac{(2+\theta-\eta)c_R}{(1+\theta)(1-\eta)}$, repeating the argument of
684 Theorem 4.8 leads to

$$685 \quad a_{k+1} \leq \left(\sum_{j=1}^k \eta_j \phi_j^{\frac{1}{2}} (c_0 a_j^{\frac{1}{2}} b_j^{\frac{1}{2}} + c_R a_j^{\frac{1}{2}} \delta + \delta) + c_\nu \|w\| (k+1)^{-\nu(1-\alpha)} \right)^2$$

$$686 \quad + n \sum_{j=1}^k \eta_j^2 (\phi_j^{\frac{1}{2}})^2 (b_j^{\frac{1}{2}} + \delta)^2 + n \left(\sum_{j=1}^k \eta_j \phi_j^{\frac{1}{2}} (c_0 b_j^{\frac{1}{2}} + c_R \delta) a_j^{\frac{\theta}{2}} \right)^2$$

$$687 \quad + 2n \left(\sum_{i=1}^k \eta_i \phi_i^{\frac{1}{2}} (b_i^{\frac{1}{2}} + \delta) \right) \left(\sum_{j=1}^k \eta_j \phi_j^{\frac{1}{2}} (c_0 b_j^{\frac{1}{2}} + c_R \delta) a_j^{\frac{\theta}{2}} \right),$$

$$\begin{aligned}
688 \quad b_{k+1} &\leq \left(\sum_{j=1}^k \eta_j \phi_j^1 (c_0 a_j^{\frac{1}{2}} b_j^{\frac{1}{2}} + c_R a_j^{\frac{1}{2}} \delta + \delta) + c_\nu \|w\| (k+1)^{-(\nu+\frac{1}{2})(1-\alpha)} \right)^2 \\
689 \quad &+ n \sum_{j=1}^k \eta_j^2 (\phi_j^1)^2 (b_j^{\frac{1}{2}} + \delta)^2 + n \left(\sum_{j=1}^k \eta_j \phi_j^1 (c_0 b_j^{\frac{1}{2}} + c_R \delta) a_j^{\frac{\theta}{2}} \right)^2 \\
690 \quad &+ 2n \left(\sum_{i=1}^k \eta_i \phi_i^1 (b_i^{\frac{1}{2}} + \delta) \right) \left(\sum_{j=1}^k \eta_j \phi_j^1 (c_0 b_j^{\frac{1}{2}} + c_R \delta) a_j^{\frac{\theta}{2}} \right). \\
691
\end{aligned}$$

692 Like in the proof of [Theorem 4.8](#), the goal is to show

$$693 \quad (4.19) \quad a_k \leq c^* \|w\|^2 k^{-\beta} \quad \text{and} \quad b_k \leq c^* \|w\|^2 k^{-\gamma},$$

694 for all $k \leq k^* = \lfloor (\frac{\delta}{\|w\|})^{-\frac{2}{(2\nu+1)(1-\alpha)}} \rfloor$, with $\beta = \min(2\nu(1-\alpha), \alpha - \epsilon)$ and $\gamma = \min((1+2\nu)(1 -$
695 $\alpha), 1 - \epsilon)$, and the constant $c^* > 0$ to be specified. By the choice of k^* , for any $k \leq k^*$,

$$696 \quad (4.20) \quad k^{\frac{1-\alpha}{2}} \delta \leq k^{-\nu(1-\alpha)} \|w\|.$$

697 Now the proof proceeds by mathematical induction. When $k = 1$, [\(4.19\)](#) holds trivially for
698 any sufficiently large c^* . Now we assume [\(4.19\)](#) holds up to some $k < k^*$, and prove it for
699 $k+1 \leq k^*$. Upon substituting the induction hypothesis, with $\varrho = c^* \|w\|^2$, we obtain

$$\begin{aligned}
700 \quad a_{k+1} &\leq \left(\sum_{j=1}^k \eta_j \phi_j^{\frac{1}{2}} (c_0 \varrho j^{-\frac{\beta+\gamma}{2}} + c_R \varrho^{\frac{1}{2}} j^{-\frac{\beta}{2}} \delta + \delta) + c_\nu \|w\| (k+1)^{-\nu(1-\alpha)} \right)^2 \\
701 \quad (4.21) \quad &+ n \sum_{j=1}^k \eta_j^2 (\phi_j^{\frac{1}{2}})^2 (\varrho^{\frac{1}{2}} j^{-\frac{\gamma}{2}} + \delta)^2 + 2n \left(\sum_{i=1}^k \eta_i \phi_i^{\frac{1}{2}} (\varrho^{\frac{1}{2}} i^{-\frac{\gamma}{2}} + \delta) \right) \\
702 \quad &\times \left(\sum_{j=1}^k \eta_j \phi_j^{\frac{1}{2}} (c_0 \varrho^{\frac{1}{2}} j^{-\frac{\gamma}{2}} + c_R \delta) \varrho^{\frac{\theta}{2}} j^{-\frac{\theta\beta}{2}} \right) + n \left(\sum_{j=1}^k \eta_j \phi_j^{\frac{1}{2}} (c_0 \varrho^{\frac{1}{2}} j^{-\frac{\gamma}{2}} + c_R \delta) \varrho^{\frac{\theta}{2}} j^{-\frac{\theta\beta}{2}} \right)^2. \\
703
\end{aligned}$$

704 Next, using [Proposition A.2](#), we obtain

$$\begin{aligned}
705 \quad (4.22) \quad a_{k+1} &\leq \left((c_1 (c_0 \varrho + (c_R \varrho^{\frac{1}{2}} + 1) \|w\|) + c_\nu \|w\|)^2 + 2n (c_2 \varrho + c_3 \|w\|^2) \right. \\
706 \quad &\left. + 2n c_1^2 (\varrho^{\frac{1}{2}} + \|w\|) (c_0 \varrho^{\frac{1}{2}} + c_R \|w\|) \varrho^{\frac{\theta}{2}} + n c_1^2 (c_0 \varrho^{\frac{1}{2}} + c_R \|w\|)^2 \varrho^\theta \right) (k+1)^{-\beta}, \\
707
\end{aligned}$$

708 with the constants c_1, \dots, c_3 given in [Proposition A.2](#). Similarly, it follows from the induction
709 hypothesis that

$$\begin{aligned}
710 \quad b_{k+1} &\leq \left(\sum_{j=1}^k \eta_j \phi_j^1 (c_0 \varrho j^{-\frac{\beta+\gamma}{2}} + c_R \varrho^{\frac{1}{2}} j^{-\frac{\beta}{2}} \delta + \delta) + c_\nu \|w\| (k+1)^{-(1-\alpha)(\nu+\frac{1}{2})} \right)^2 \\
711 \quad (4.23) \quad &+ n \sum_{j=1}^k \eta_j^2 (\phi_j^1)^2 (\varrho^{\frac{1}{2}} j^{-\frac{\gamma}{2}} + \delta)^2 + 2n \left(\sum_{i=1}^k \eta_i \phi_i^1 (\varrho^{\frac{1}{2}} i^{-\frac{\gamma}{2}} + \delta) \right) \\
712 \quad &\times \left(\sum_{j=1}^k \eta_j \phi_j^1 (c_0 \varrho^{\frac{1}{2}} j^{-\frac{\gamma}{2}} + c_R \delta) \varrho^{\frac{\theta}{2}} j^{-\frac{\theta\beta}{2}} \right) + n \left(\sum_{j=1}^k \eta_j \phi_j^1 (c_0 \varrho^{\frac{1}{2}} j^{-\frac{\gamma}{2}} + c_R \delta) \varrho^{\frac{\theta}{2}} j^{-\frac{\theta\beta}{2}} \right)^2, \\
713
\end{aligned}$$

714 from which and [Proposition A.2](#), it follows that

$$715 \quad b_{k+1} \leq \left((c_0 c'_1 \varrho + c'_5 (c_R \varrho^{\frac{1}{2}} + 1) \|w\| + c_\nu \|w\|)^2 + 2n(c'_2 \varrho + c_3 \|w\|^2) \right. \\ 716 \quad (4.24) \quad \left. + 2n(c'_3 \varrho^{\frac{1}{2}} + c'_5 \|w\|)(c_0 c'_3 \varrho^{\frac{1}{2}} + c_R c'_5 \|w\|) \varrho^{\frac{\theta}{2}} + n(c_0 c'_4 \varrho^{\frac{1}{2}} + c_R c'_5 \|w\|)^2 \varrho^\theta \right) (k+1)^{-\gamma}, \\ 717$$

718 with the constants c'_1, \dots, c'_5 given in [Proposition A.2](#). In view of [\(4.22\)](#) and [\(4.24\)](#), for small
719 $\|w\|$ and η_0 , repeating the argument for [Theorem 4.8](#) (and noting that c_1, c_2, c_3, c'_2 tend to
720 zero as $\eta_0 \rightarrow 0^+$) concludes the existence of a $c^* > 0$ such that [\(4.19\)](#) hold. This completes
721 the induction step and the proof of [Theorem 2.2](#). \square

722 **5. Concluding remarks.** In this work, we have provided a convergence analysis of
723 stochastic gradient descent for a class of nonlinear ill-posed inverse problems. The method
724 employs an unbiased estimate of the gradient, computed from one randomly selected equa-
725 tion of the nonlinear system, and admits excellent scalability to the problem size. We
726 proved that it is regularizing under the traditional tangential cone condition with *a priori*
727 parameter choice, and also showed a convergence rate under canonical source condition and
728 range invariance condition (and its stochastic variant), for a polynomially decaying step size
729 schedule. The analysis combines techniques from both nonlinear regularization theory and
730 stochastic calculus, and the results extend the existing works [\[8\]](#) and [\[14\]](#).

731 There are several avenues in both theoretical and practical aspects for further research.
732 First, it is important to verify the assumptions for concrete nonlinear inverse problems,
733 especially nonlinearity conditions in [Assumption 2.1\(ii\)–\(iii\)](#) and [Assumption 2.3](#), for e.g.,
734 parameter identifications for PDEs, which would justify the usage of SGD. Several important
735 inverse problems in medical imaging are of the form [\(1.1\)](#), e.g., electrical impedance tomog-
736 raphy and diffuse optical spectroscopy. These applications often involve natural physical
737 constraints, e.g., positivity, which the algorithm should be adapted to preserve. Second, the
738 source condition employed in the work is canonical, and alternative approaches, e.g., varia-
739 tional inequalities and conditional stability, should also be studied for convergence rates [\[24\]](#),
740 or the Frechét differentiability of the forward operator in [Assumption 2.1](#) may be relaxed
741 [\[3\]](#). Third, the influence of various algorithmic parameters, e.g., mini-batch, random sam-
742 pling, step size schedules (including adaptive rules) and *a posteriori* stopping rule, should
743 be analyzed to provide useful practical guidelines.

744 **Acknowledgements.** The authors are grateful to the associate editor, Professor Frank
745 E. Curtis, and two anonymous referees for helpful comments.

746 **Appendix A. Auxiliary estimates.** In this appendix, we collect several auxiliary
747 inequalities that have been used in the convergence rates analysis. Most estimates follow
748 from routine but rather tedious computations. We begin with a well known estimate on
749 operator norms (see, e.g., [\[19\]](#) [\[14, Lemma A.1\]](#)).

750 **LEMMA A.1.** *For any $j < k$, and any symmetric and positive semidefinite operator S*
751 *and step sizes $\eta_j \in (0, \|S\|^{-1}]$ and $p \geq 0$, there holds*

$$752 \quad \left\| \prod_{i=j}^k (I - \eta_i S) S^p \right\| \leq \frac{p^p}{e^p (\sum_{i=j}^k \eta_i)^p}.$$

753 Below we need the Beta function $B(a, b) = \int_0^1 s^{a-1} (1-s)^{b-1} ds$ for any $a, b > 0$. Note
754 that for fixed a , the function $B(a, \cdot)$ is monotonically decreasing.

755 LEMMA A.2. For $\eta_j = \eta_0 j^{-\alpha}$ with $\alpha \in (0, 1)$, $r \in [0, 1)$, $\beta \in [0, 1]$, and $\gamma = \alpha + \beta$, the
 756 following estimates hold

$$757 \quad \sum_{i=1}^k \eta_i \geq (1 - 2^{\alpha-1})(1 - \alpha)^{-1} \eta_0 (k+1)^{1-\alpha},$$

$$758 \quad \sum_{j=1}^{k-1} \frac{\eta_j}{(\sum_{\ell=j+1}^k \eta_\ell)^r} j^{-\beta} \leq \eta_0^{1-r} B(1-r, 1-\gamma) k^{r\alpha+1-r-\gamma}, \quad r \in [0, 1), \gamma < 1,$$

$$759 \quad \sum_{j=1}^{k-1} \frac{\eta_j}{\sum_{\ell=j+1}^k \eta_\ell} j^{-\beta} \leq \begin{cases} 2^\gamma (1-\gamma)^{-1} k^{-\beta}, & \gamma < 1, \\ 4k^{\alpha-1} \ln k, & \gamma = 1, \\ 2^\gamma (\gamma-1)^{-1} k^{\alpha-1}, & \gamma > 1, \end{cases} + 2^{1+\gamma} k^{-\beta} \ln k.$$

761 *Proof.* The first estimate follows from the fact $1 - (k+1)^{\alpha-1} \geq 1 - 2^{\alpha-1}$ for $k \geq 1$ that

$$762 \quad \sum_{i=1}^k \eta_i \geq \eta_0 \int_1^{k+1} s^{-\alpha} ds = \eta_0 (1-\alpha)^{-1} ((k+1)^{1-\alpha} - 1) \geq \eta_0 (1-\alpha)^{-1} (1 - 2^{\alpha-1}) (k+1)^{1-\alpha}. \blacksquare$$

764 To prove the second estimate, we note $\eta_i \geq \eta_0 k^{-\alpha}$ for any $i = j+1, \dots, k$ and thus

$$765 \quad (\text{A.1}) \quad \eta_0^{-1} \sum_{i=j+1}^k \eta_i \geq k^{-\alpha} (k-j).$$

766 Thus, if $\gamma = \alpha + \beta < 1$ and $r < 1$,

$$767 \quad \sum_{j=1}^{k-1} \frac{\eta_j}{(\sum_{\ell=j+1}^k \eta_\ell)^r} j^{-\beta} \leq \eta_0^{1-r} k^{r\alpha} \sum_{j=1}^{k-1} (k-j)^{-r} j^{-\gamma} \leq \eta_0^{1-r} k^{r\alpha} \int_0^k (k-s)^{-r} s^{-\gamma} ds$$

$$768 \quad = \eta_0^{1-r} B(1-r, 1-\gamma) k^{r\alpha+1-r-\gamma}.$$

770 Similarly, if $r = 1$, it follows from (A.1) that

$$771 \quad \sum_{j=1}^{k-1} \frac{\eta_j}{\sum_{\ell=j+1}^k \eta_\ell} j^{-\beta} \leq k^\alpha \sum_{j=1}^{k-1} (k-j)^{-1} j^{-\gamma}$$

$$772 \quad = k^\alpha \sum_{j=1}^{\lfloor \frac{k}{2} \rfloor} j^{-\gamma} (k-j)^{-1} + k^\alpha \sum_{j=\lfloor \frac{k}{2} \rfloor + 1}^{k-1} j^{-\gamma} (k-j)^{-1}$$

$$773 \quad \leq 2k^{\alpha-1} \sum_{j=1}^{\lfloor \frac{k}{2} \rfloor} j^{-\gamma} + 2^\gamma k^{-\beta} \sum_{j=\lfloor \frac{k}{2} \rfloor + 1}^{k-1} (k-j)^{-1}.$$

775 Simple computation gives

$$776 \quad (\text{A.2}) \quad \sum_{j=\lfloor \frac{k}{2} \rfloor + 1}^{k-1} (k-j)^{-1} \leq 2 \ln k \quad \text{and} \quad \sum_{j=1}^{\lfloor \frac{k}{2} \rfloor} j^{-\gamma} \leq \begin{cases} (1-\gamma)^{-1} (\frac{k}{2})^{1-\gamma}, & \gamma \in [0, 1), \\ 2 \ln k, & \gamma = 1, \\ \gamma(\gamma-1)^{-1}, & \gamma > 1. \end{cases}$$

777 Combining the last three estimates gives the assertion for the case $r = 1$.

778 Next we recall two useful estimates.

779 LEMMA A.3. For $\eta_j = \eta_0 j^{-\alpha}$ with $\alpha \in (0, 1)$, $\beta \in [0, 1]$ and $r \geq 0$, there hold

$$780 \quad \sum_{j=1}^{\lfloor \frac{k}{2} \rfloor} \frac{\eta_j^2}{(\sum_{\ell=j+1}^k \eta_\ell)^r} j^{-\beta} \leq c_{\alpha, \beta, r} k^{-r(1-\alpha) + \max(0, 1-2\alpha-\beta)},$$

$$781 \quad \sum_{j=\lfloor \frac{k}{2} \rfloor + 1}^{k-1} \frac{\eta_j^2}{(\sum_{\ell=j+1}^k \eta_\ell)^r} j^{-\beta} \leq c'_{\alpha, \beta, r} k^{-((2-r)\alpha + \beta) + \max(0, 1-r)},$$

782 where we slightly abuse the notation $k^{-\max(0,0)}$ for $\ln k$, and $c_{\alpha, \beta, r}$ and $c'_{\alpha, \beta, r}$ are given by

$$784 \quad c_{\alpha, \beta, r} = 2^r \eta_0^{2-r} \begin{cases} \frac{2\alpha+\beta}{2\alpha+\beta-1}, & 2\alpha+\beta > 1, \\ 2, & 2\alpha+\beta = 1, \\ \frac{2^{2\alpha+\beta-1}}{1-2\alpha-\beta}, & 2\alpha+\beta < 1, \end{cases} \quad \text{and} \quad c'_{\alpha, \beta, r} = 2^{2\alpha+\beta} \eta_0^{2-r} \begin{cases} \frac{r}{r-1}, & r > 1, \\ 2, & r = 1, \\ \frac{2^{r-1}}{1-r}, & r < 1. \end{cases}$$

786 *Proof.* The proof is based on (A.1) and (A.2) and essentially given in [14, Lemma A.3],
787 but the constants are corrected. \square

788 The next result collects some lengthy estimates needed in the proof of Theorem 4.8.

789 PROPOSITION A.1. Let $\beta = \min(2\nu(1-\alpha), \alpha - \epsilon)$, $\gamma = \min((1+2\nu)(1-\alpha), 1-\epsilon)$ and
790 $r = \min(\frac{1}{2} + \nu, \frac{1-\epsilon}{2(1-\alpha)})$. Then under the conditions in Theorem 4.8, i.e., $\|B\| \leq 1$, $\eta_0 \leq 1$
791 and θ being sufficiently small, with $\zeta = (\frac{1}{2} - \nu)(1-\alpha)$, the following estimates hold:

$$792 \quad (\text{A.3}) \quad \sum_{j=1}^k \eta_j \phi_j^{\frac{1}{2}} j^{-\frac{\gamma}{2}} \leq c_1 (k+1)^{-\frac{\beta}{2}}, \quad \sum_{j=1}^k \eta_j^2 (\phi_j^{\frac{1}{2}})^2 j^{-\gamma} \leq c_2 (k+1)^{-\beta},$$

$$793 \quad (\text{A.4}) \quad \sum_{j=1}^{\lfloor \frac{k}{2} \rfloor} \eta_j^2 (\phi_j^r)^2 j^{-\gamma} + \sum_{j=\lfloor \frac{k}{2} \rfloor + 1}^k \eta_j^2 (\phi_j^{\frac{1}{2}})^2 j^{-\gamma} \leq c_3 (k+1)^{-\gamma}, \quad \sum_{j=1}^k \eta_j \phi_j^1 j^{-\frac{\beta+\gamma}{2}} \leq c_4 (k+1)^{-\frac{\gamma}{2}},$$

$$794 \quad (\text{A.5}) \quad \left(\sum_{i=1}^k \eta_i \phi_i^1 i^{-\frac{\gamma}{2}} \right) \left(\sum_{j=1}^k \eta_j \phi_j^1 j^{-\frac{\gamma+\theta\beta}{2}} \right) \leq c_5 (k+1)^{-\gamma}, \quad \sum_{j=1}^k \eta_j \phi_j^1 j^{-\frac{\gamma+\theta\beta}{2}} \leq c_6 (k+1)^{-\frac{\gamma}{2}}.$$

796 with $c_1 = 2^{\frac{\beta}{2}} \eta_0^{\frac{1}{2}} (2^{-1} B(\frac{1}{2}, \zeta) + 1)$, $c_2 = 2^\beta \eta_0 (\alpha^{-1} + 2)$, $c_3 = 2^\gamma \eta_0^{2-2r} (3\alpha^{-1} + 1)$, $c_4 = 2^{\frac{\gamma}{2}} (\zeta^{-1} +$
797 $2\beta^{-1} + 1)$, $c_5 = 2^\gamma (((\frac{1}{2} - \nu - \theta\nu)(1-\alpha))^{-1} + 4(\theta\beta)^{-1} + 1)^2$ and $c_6 = 2^{\frac{\gamma}{2}} (\zeta^{-1} + 2(\theta\beta)^{-1} + 1)$.

798 *Proof.* It follows from Lemma A.1 and the condition $\|B\| \leq 1$ that

$$799 \quad \sum_{j=1}^k \eta_j \phi_j^{\frac{1}{2}} j^{-\frac{\gamma}{2}} \leq (2e)^{-\frac{1}{2}} \sum_{j=1}^{k-1} \frac{\eta_j}{(\sum_{\ell=1}^k \eta_\ell)^{\frac{1}{2}}} j^{-\frac{\gamma}{2}} + \eta_0 k^{-\alpha - \frac{\gamma}{2}}$$

$$800 \quad \leq (\eta_0^{\frac{1}{2}} 2^{-1} B(\frac{1}{2}, 1 - \alpha - \frac{\gamma}{2}) + \eta_0) k^{\frac{1-\alpha}{2} - \frac{\gamma}{2}}.$$

802 By the definitions of β and γ , we have $\frac{1-\alpha}{2} - \frac{\gamma}{2} = -\frac{\beta}{2}$, and $1 - \alpha - \frac{\gamma}{2} \geq (\frac{1}{2} - \nu)(1-\alpha) := \zeta$.
803 Thus, the monotonicity of the Beta function, and the inequality $2k \geq k+1$ for $k \geq 1$ imply
804 the first inequality of (A.3). Now by Lemma A.1 and Lemma A.3,

$$805 \quad (\text{A.6}) \quad \sum_{j=1}^k \eta_j^2 (\phi_j^{\frac{1}{2}})^2 j^{-\gamma} \leq (2e)^{-1} \sum_{j=1}^{k-1} \frac{\eta_j^2}{\sum_{\ell=j+1}^k \eta_\ell} j^{-\gamma} + \eta_0^2 \|B^{\frac{1}{2}}\|^2 k^{-2\alpha-\gamma}$$

$$\leq \eta_0 \left((2e)^{-1} \frac{2(2\alpha + \gamma)}{2\alpha + \gamma - 1} k^{-(1-\alpha)} + (2e)^{-1} 2^{1+2\alpha+\gamma} k^{-\alpha-\gamma} \ln k + \eta_0 \|B^{\frac{1}{2}}\|^2 k^{-2\alpha-\gamma} \right).$$

Now, for any $r > 0$, there holds

$$(A.7) \quad s^{-r} \ln s \leq (er)^{-1}, \quad \forall s \geq 0,$$

and thus $k^{-\alpha-\gamma} \ln k = k^{-\beta}(k^{-1} \ln k) \leq e^{-1} k^{-\beta}$. Further, by the definition of γ , $2\alpha + \gamma \leq \min(2, 1 + 2\alpha) \leq 2$, and since $\epsilon < \frac{\alpha}{2}$, $2\alpha + \gamma - 1 \geq \alpha$,

$$(A.8) \quad \frac{2\alpha+\gamma}{2\alpha+\gamma-1} = 1 + \frac{1}{2\alpha+\gamma-1} \leq 1 + \alpha^{-1}.$$

Then, the last three estimates (with $\|B\| \leq 1$) imply

$$\sum_{j=1}^k \eta_j^2 (\phi_j^{\frac{1}{2}})^2 j^{-\gamma} \leq 2^\beta \eta_0 (\alpha^{-1} + 2) (k+1)^{-\beta}.$$

This proves the second inequality in (A.3).

Next, by letting $r = \min(\frac{1}{2} + \nu, \frac{1-\epsilon}{2(1-\alpha)}) \in (\frac{1}{2}, 1)$, and using (A.7) and (A.8), Lemma A.1 and Lemma A.3 and the monotonicity of $\frac{s^\epsilon}{e^s}$ for $s \in [0, 1]$, the first part of (A.4) follows from

$$\begin{aligned} & \sum_{j=1}^{\lfloor \frac{k}{2} \rfloor} \eta_j^2 (\phi_j^r)^2 j^{-\gamma} + \sum_{j=\lfloor \frac{k}{2} \rfloor + 1}^k \eta_j^2 (\phi_j^{\frac{1}{2}})^2 j^{-\gamma} \\ & \leq (2e)^{-1} \left(\sum_{j=1}^{\lfloor \frac{k}{2} \rfloor} \frac{\eta_j^2}{(\sum_{\ell=1}^j \eta_\ell)^{2r}} j^{-\gamma} + \sum_{j=\lfloor \frac{k}{2} \rfloor + 1}^{k-1} \frac{\eta_j^2}{\sum_{\ell=j+1}^k \eta_\ell} j^{-\gamma} \right) + \eta_0^2 k^{-2\alpha-\gamma} \\ & \leq \eta_0^{2-2r} \frac{2^{2r}(2\alpha + \gamma)}{2e(2\alpha + \gamma - 1)} k^{-\gamma} + \frac{2^{1+2\alpha+\gamma}}{2e} \eta_0 k^{-(\alpha+\gamma)} \ln k + \eta_0^2 k^{-2\alpha-\gamma} \leq c_3 (k+1)^{-\gamma}. \end{aligned}$$

Now, we bound the sum $\sum_{j=1}^k \eta_j \phi_j^1 j^{-\sigma}$ for any $\sigma \in [\frac{\gamma}{2}, \frac{\gamma+\beta}{2}]$, and then set σ to $\frac{\gamma}{2}$, $\frac{\gamma+\theta\beta}{2}$ and $\frac{\gamma+\beta}{2}$ to complete the proof. By Lemma A.1 and Lemma A.2, there hold

$$(A.9) \quad \sum_{j=1}^{\lfloor \frac{k}{2} \rfloor} \eta_j \phi_j^1 j^{-\sigma} \leq e^{-1} \begin{cases} \frac{2^{\alpha+\sigma}}{1-\alpha-\sigma} k^{-\sigma}, & \alpha + \sigma < 1, \\ 4k^{\alpha-1} \ln k, & \alpha + \sigma = 1, \\ \frac{2(\alpha+\sigma)}{\alpha+\sigma-1} k^{\alpha-1}, & \alpha + \sigma > 1, \end{cases}$$

$$(A.10) \quad \sum_{j=\lfloor \frac{k}{2} \rfloor + 1}^k \eta_j \phi_j^1 j^{-\sigma} \leq e^{-1} 2^{1+\alpha+\sigma} k^{-\sigma} \ln k + \eta_0 k^{-\sigma}.$$

First, we choose $\sigma = \frac{\beta+\gamma}{2}$. By (A.7), since $(1 - \alpha - \frac{\gamma}{2})^{-1} \leq \zeta^{-1}$, $\alpha + \frac{\gamma}{2} < 1$, $\|B\| \leq 1$ and $\eta_0 \leq 1$, we obtain

$$\begin{aligned} & \sum_{j=1}^k \eta_j \phi_j^1 j^{-\frac{\beta+\gamma}{2}} \leq \sum_{j=1}^{\lfloor \frac{k}{2} \rfloor} \eta_j \phi_j^1 j^{-\frac{\gamma}{2}} + \sum_{j=\lfloor \frac{k}{2} \rfloor + 1}^k \eta_j \phi_j^1 j^{-\frac{\beta+\gamma}{2}} \\ & \leq 2^{\alpha+\frac{\gamma}{2}} e^{-1} (1 - \alpha - \frac{\gamma}{2})^{-1} k^{-\frac{\gamma}{2}} + 2^{1+\alpha+\frac{\gamma+\beta}{2}} e^{-1} k^{-\frac{\gamma+\beta}{2}} \ln k + \eta_0 k^{-\frac{\gamma}{2}} \leq c_4 (k+1)^{-\frac{\gamma}{2}}, \end{aligned}$$

832 due to the inequality $2^{1+\alpha+\frac{\beta+\gamma}{2}} < e^2$, from the definitions of the exponents β and γ . This
833 shows the second inequality of (A.4). Since θ is small, we may assume $\theta < \frac{1}{2\nu} - 1 \leq$
834 $\frac{1-\alpha}{\beta} - 1$. Then by the relations $\gamma = 1 - \alpha + \beta$ and $\beta \leq 2\nu(1 - \alpha)$, direct computation shows
835 $1 - \alpha - \frac{\gamma+\theta\beta}{2} \geq (\frac{1}{2} - \nu - \theta\nu)(1 - \alpha) > 0$. Further, since $\theta < \frac{1-\alpha}{\beta} - 1$, $\min(\frac{\theta\beta}{2}, 1 - \alpha - \frac{\gamma}{2}) = \frac{\theta\beta}{2}$.
836 Hence, it follows from (A.9) and (A.10), with $\sigma = \frac{\gamma}{2}$ and $\frac{\gamma+\theta\beta}{2}$ that

$$837 \quad \left(\sum_{i=1}^k \eta_i \phi_i^1 i^{-\frac{\gamma}{2}} \right) \left(\sum_{j=1}^k \eta_j \phi_j^1 j^{-\frac{\gamma+\theta\beta}{2}} \right) \leq \left(\frac{2^{\alpha+\frac{\gamma}{2}}}{e(1-\alpha-\frac{\gamma}{2})} + \frac{2^{1+\alpha+\frac{\gamma}{2}}}{e} \ln k + 1 \right)$$

$$838 \quad \times \left(\frac{2^{\alpha+\frac{\gamma+\theta\beta}{2}}}{e(1-\alpha-\frac{\gamma+\theta\beta}{2})} k^{-\min(\frac{\theta\beta}{2}, 1-\alpha-\frac{\gamma}{2})} + \frac{2^{1+\alpha+\frac{\gamma+\theta\beta}{2}}}{e} k^{-\frac{\theta\beta}{2}} \ln k + k^{-\frac{\theta\beta}{2}} \right) k^{-\gamma}.$$

840 Then we move one factor $k^{-\frac{\theta\beta}{4}}$ from the second bracket to the first and bound by (A.7):

$$841 \quad \left(\sum_{i=1}^k \eta_i \phi_i^1 i^{-\frac{\gamma}{2}} \right) \left(\sum_{j=1}^k \eta_j \phi_j^1 j^{-\frac{\gamma+\theta\beta}{2}} \right) \leq \left(\frac{2^{\alpha+\frac{\gamma}{2}}}{e(1-\alpha-\frac{\gamma}{2})} + \frac{2^{1+\alpha+\frac{\gamma}{2}}}{e} k^{-\frac{\theta\beta}{4}} \ln k + 1 \right)$$

$$842 \quad \times \left(\frac{2^{\alpha+\frac{\gamma+\theta\beta}{2}}}{e(1-\alpha-\frac{\gamma+\theta\beta}{2})} + \frac{2^{1+\alpha+\frac{\gamma+\theta\beta}{2}}}{e} k^{-\frac{\theta\beta}{4}} \ln k + 1 \right) k^{-\gamma}$$

$$843 \quad \leq 2^\gamma \left(\left(\frac{1}{2} - \nu - \theta\nu \right) (1 - \alpha) \right)^{-1} + 4(\theta\beta)^{-1} + 1 \Big)^2 (k + 1)^{-\gamma},$$

845 proving the first inequality of (A.5). The other estimate in (A.5) follows similarly by choosing
846 $\sigma = \frac{\gamma+\theta\beta}{2}$, and hence omitted. \square

847 **REMARK A.1.** *The proof of Proposition A.1 implies $\sum_{j=1}^{k-1} \eta_j \phi_j^1 j^{-\frac{\gamma}{2}} \leq (\zeta^{-1} + 2 \ln k) k^{-\frac{\gamma}{2}}$.*
848 *The log factor $\ln k$ seems not removable, and precludes a direct application of mathematical*
849 *induction in the proof of Theorem 4.8. The extra factor $j^{-\frac{\theta\beta}{2}}$ due to Assumption 2.3*
850 *gracefully compensates the log factor $\ln k$ using (A.7). The smallness condition on θ can*
851 *be removed at the expense of less transparent dependence. Specifically, by Lemma A.2, with*
852 $\sigma = \alpha + \frac{\gamma+\theta\beta}{2}$, *there holds*

$$853 \quad \sum_{j=1}^k \eta_j \phi_j^1 j^{-\frac{\gamma+\theta\beta}{2}} \leq \frac{1}{ek^{\frac{\gamma}{2}}} \begin{cases} \frac{2^\sigma}{1-\sigma} k^{-\frac{\theta\beta}{2}}, & \sigma < 1 \\ 4k^{-(1-\alpha-\frac{\gamma}{2})} \ln k, & \sigma = 1 \\ \frac{2^\sigma}{\sigma-1} k^{-(1-\alpha-\frac{\gamma}{2})}, & \sigma > 1 \end{cases} + 2^{1+\sigma} e^{-1} k^{-\frac{\gamma}{2}-\frac{\theta\beta}{2}} \ln k + k^{-(\alpha+\frac{\gamma+\theta\beta}{2})}.$$

855 *Instead of applying (A.7) directly, we rearrange the terms and discuss the cases $\sigma < 1$, $\sigma = 1$*
856 *and $\sigma > 1$ separately with the argument in the proof of Proposition A.1 and obtain*

$$857 \quad \left(\sum_{i=1}^k \eta_i \phi_i^1 i^{-\frac{\gamma}{2}} \right) \left(\sum_{j=1}^k \eta_j \phi_j^1 j^{-\frac{\gamma+\theta\beta}{2}} \right) \leq c_\sigma 2^\gamma (k + 1)^{-\gamma},$$

859 *with the constant c_σ given by*

$$860 \quad c_\sigma = \begin{cases} (1 - \sigma)^{-1} + 4(\theta\beta)^{-1} + 1, & \sigma < 1, \\ \zeta^{-1} + 8(\theta\beta)^{-1} + 1, & \sigma = 1, \\ 2(\sigma - 1)^{-1} + 3\zeta^{-1} + 1, & \sigma > 1. \end{cases}$$

862 The next result gives some basic estimates used in the proof of Theorem 2.2.

863 PROPOSITION A.2. Under the induction hypothesis of [Theorem 2.2](#) and [\(4.20\)](#), there
864 hold

$$865 \quad a_{k+1} \leq \left((c_1(c_0\varrho + (c_R\varrho^{\frac{1}{2}} + 1)\|w\|) + c_\nu\|w\|)^2 + 2n(c_2\varrho + c_3\|w\|^2) \right. \\ 866 \quad \left. + 2nc_1^2(\varrho^{\frac{1}{2}} + \|w\|)(c_0\varrho^{\frac{1}{2}} + c_R\|w\|)\varrho^{\frac{\theta}{2}} + nc_1^2(c_0\varrho^{\frac{1}{2}} + c_R\|w\|)^2\varrho^\theta \right) (k+1)^{-\beta}, \\ 867 \quad b_{k+1} \leq \left((c_0c'_1\varrho + c'_5(c_R\varrho^{\frac{1}{2}} + 1)\|w\| + c_\nu\|w\|)^2 + 2n(c'_2\varrho + c_3\|w\|^2) \right. \\ 868 \quad \left. + 2n(c'_3\varrho^{\frac{1}{2}} + c'_5\|w\|)(c_0c'_3\varrho^{\frac{1}{2}} + c'_5c_R\|w\|)\varrho^{\frac{\theta}{2}} + n(c_0c'_4\varrho^{\frac{1}{2}} + c'_5c_R\|w\|)^2\varrho^\theta \right) (k+1)^{-\gamma},$$

870 where the constants c_1, c_2, c_3 and c'_1, \dots, c'_5 are given in the proof.

871 *Proof.* First, it follows directly from [Lemma A.1](#), [Lemma A.2](#), and [Lemma A.3](#) and the
872 assumptions $\|B\| \leq 1$ and $\eta_0 \leq 1$ that for any $\sigma \in [0, 1 - \alpha)$,

$$873 \quad (\text{A.11}) \quad \sum_{j=1}^k \eta_j \phi_j^{\frac{1}{2}} j^{-\sigma} \leq \eta_0^{\frac{1}{2}} (2^{-1}B(\frac{1}{2}, 1 - \alpha - \sigma) + 1) k^{\frac{1-\alpha}{2} - \sigma},$$

$$874 \quad (\text{A.12}) \quad \sum_{j=1}^k \eta_j^2 (\phi_j^{\frac{1}{2}})^2 \leq \eta_0 (|1 - 2\alpha|^{-1} + \alpha^{-1} + 1) := c_3, \\ 875$$

876 where we have abused the writing 0^{-1} for 1. Meanwhile, by [Proposition A.1](#), we have

$$877 \quad (\text{A.13}) \quad \sum_{j=1}^k \eta_j \phi_j^{\frac{1}{2}} j^{-\frac{\gamma}{2}} \leq c_1 (k+1)^{-\frac{\beta}{2}} \quad \text{and} \quad \sum_{j=1}^k \eta_j^2 (\phi_j^{\frac{1}{2}})^2 j^{-\gamma} \leq c_2 (k+1)^{-\beta}, \\ 878$$

879 with $c_1 = 2^{\frac{\beta}{2}} \eta_0^{\frac{1}{2}} (2^{-1}B(\frac{1}{2}, \zeta) + 1)$, $\zeta = (\frac{1}{2} - \nu)(1 - \alpha)$ and $c_2 = 2^\beta \eta_0 (\alpha^{-1} + 2)$. By [\(A.11\)](#)-[\(A.13\)](#)
880 and the monotonicity of the Beta function, and $k+1 \leq k^*$ (cf. [\(4.20\)](#)), we obtain

$$881 \quad \sum_{j=1}^k \eta_j \phi_j^{\frac{1}{2}} (c_0\varrho j^{-\frac{\beta+\gamma}{2}} + c_R\varrho^{\frac{1}{2}} j^{-\frac{\beta}{2}} \delta + \delta) \leq c_0 c_1 \varrho (k+1)^{-\frac{\beta}{2}} + (c_R\varrho^{\frac{1}{2}} + 1) c_1 (k+1)^{\frac{1-\alpha}{2}} \delta \\ 882 \quad \leq c_1 (c_0\varrho + (c_R\varrho^{\frac{1}{2}} + 1)\|w\|) (k+1)^{-\frac{\beta}{2}},$$

$$883 \quad \sum_{j=1}^k \eta_j^2 (\phi_j^{\frac{1}{2}})^2 (\varrho^{\frac{1}{2}} j^{-\frac{\gamma}{2}} + \delta)^2 \leq 2(c_2\varrho + c_3\|w\|^2) (k+1)^{-\beta}. \\ 884$$

885 Likewise, by the monotonicity of the Beta function, we deduce

$$886 \quad \left(\sum_{i=1}^k \eta_i \phi_i^{\frac{1}{2}} (\varrho^{\frac{1}{2}} i^{-\frac{\gamma}{2}} + \delta) \right) \left(\sum_{j=1}^k \eta_j \phi_j^{\frac{1}{2}} (c_0\varrho^{\frac{1}{2}} j^{-\frac{\gamma}{2}} + c_R\delta) \varrho^{\frac{\theta}{2}} j^{-\frac{\theta\beta}{2}} \right) \\ 887 \quad \leq c_1^2 (\varrho^{\frac{1}{2}} + \|w\|) (c_0\varrho^{\frac{1}{2}} + c_R\|w\|) \varrho^{\frac{\theta}{2}} (k+1)^{-\beta}, \\ 888 \quad \sum_{j=1}^k \eta_j \phi_j^{\frac{1}{2}} (c_0\varrho^{\frac{1}{2}} j^{-\frac{\gamma}{2}} + c_R\delta) \varrho^{\frac{\theta}{2}} j^{-\frac{\theta\beta}{2}} \leq c_1 (c_0\varrho^{\frac{1}{2}} + c_R\|w\|) \varrho^{\frac{\theta}{2}} (k+1)^{-\frac{\beta}{2}}. \\ 889$$

890 The last four estimates give [\(4.21\)](#). Now we prove [\(4.23\)](#). By [Proposition A.1](#), we have

$$891 \quad \sum_{j=1}^k \eta_j \phi_j^1 j^{-\frac{\beta+\gamma}{2}} \leq c'_1 (k+1)^{-\frac{\gamma}{2}}, \quad \sum_{j=1}^{\lfloor \frac{k}{2} \rfloor} \eta_j^2 (\phi_j^r)^2 j^{-\gamma} + \sum_{j=\lfloor \frac{k}{2} \rfloor + 1}^k \eta_j^2 (\phi_j^{\frac{1}{2}})^2 j^{-\gamma} \leq c'_2 (k+1)^{-\gamma},$$

$$\left(\sum_{i=1}^k \eta_i \phi_i^1 i^{-\frac{\gamma}{2}}\right) \left(\sum_{j=1}^k \eta_j \phi_j^1 j^{-\frac{\gamma+\theta\beta}{2}}\right) \leq c_3^2 (k+1)^{-\gamma}, \quad \sum_{j=1}^k \eta_j \phi_j^1 j^{-\frac{\gamma+\theta\beta}{2}} \leq c_4' (k+1)^{-\frac{\gamma}{2}},$$

with $c_1' = 2^{\frac{\gamma}{2}}(\zeta^{-1} + 2\beta^{-1} + 1)$, $c_2' = 2\gamma\eta_0^{2-2r}(3\alpha^{-1} + 1)$, $c_3' = 2^{\frac{\gamma}{2}}((\frac{1}{2} - \nu - \theta\nu)(1 - \alpha))^{-1} + 4(\theta\beta)^{-1} + 1$ and $c_4' = 2^{\frac{\gamma}{2}}(\zeta^{-1} + 2(\theta\beta)^{-1} + 1)$. Further, by (A.9) and (A.10), for any $\sigma \in [0, \frac{\gamma}{2}]$,

$$k^{-\nu(1-\alpha)} \sum_{j=1}^k \eta_j \phi_j^1 j^{-\sigma} \leq \zeta^{-1} + 2(\nu(1-\alpha))^{-1} + 1 := c_5'.$$

With these estimates and (4.20), we deduce

$$\sum_{j=1}^k \eta_j \phi_j^1 (c_0 \varrho j^{-\frac{\beta+\gamma}{2}} + c_R \varrho^{\frac{1}{2}} j^{-\frac{\beta}{2}} \delta + \delta) \leq (c_0 c_1' \varrho + c_5' (c_R \varrho^{\frac{1}{2}} + 1) \|w\|) (k+1)^{-\frac{\gamma}{2}},$$

$$\sum_{j=1}^k \eta_j^2 (\phi_j^1)^2 (\varrho^{\frac{1}{2}} j^{-\frac{\gamma}{2}} + \delta)^2 \leq 2(c_2' \varrho + c_3 \|w\|^2) (k+1)^{-\gamma},$$

$$\sum_{j=1}^k \eta_j \phi_j^1 (c_0 \varrho^{\frac{1}{2}} j^{-\frac{\gamma}{2}} + c_R \delta) \varrho^{\frac{\theta}{2}} j^{-\frac{\theta\beta}{2}} \leq (c_0 c_4' \varrho^{\frac{1}{2}} + c_5' c_R \|w\|) \varrho^{\frac{\theta}{2}} (k+1)^{-\frac{\gamma}{2}},$$

where the second line is due to (A.12) and the inequality $\sum_{j=1}^k \eta_j^2 (\phi_j^1)^2 \leq \sum_{j=1}^k \eta_j^2 (\phi_j^{\frac{1}{2}})^2$ (since $\|B\| \leq 1$). Last, repeating the argument in Proposition A.1 gives

$$\left(\sum_{i=1}^k \eta_i \phi_i^1 (\varrho^{\frac{1}{2}} i^{-\frac{\gamma}{2}} + \delta)\right) \left(\sum_{j=1}^k \eta_j \phi_j^1 (c_0 \varrho^{\frac{1}{2}} j^{-\frac{\gamma}{2}} + c_R \delta) \varrho^{\frac{\theta}{2}} j^{-\frac{\theta\beta}{2}}\right) \leq (c_3' \varrho^{\frac{1}{2}} + c_5' \|w\|) (c_0 c_3' \varrho^{\frac{1}{2}} + c_5' c_R \|w\|) \varrho^{\frac{\theta}{2}} (k+1)^{-\gamma}.$$

Then combining the last four estimates yields the desired bound on b_{k+1} . \square

909

REFERENCES

- [1] L. BOTTOU, F. E. CURTIS, AND J. NOCEDAL, *Optimization methods for large-scale machine learning*, SIAM Rev., 60 (2018), pp. 223–311.
- [2] K. CHEN, Q. LI, AND J.-G. LIU, *Online learning in optical tomography: a stochastic approach*, Inverse Problems, 34 (2018), pp. 075010, 26 pp.
- [3] C. CLASON AND V. H. NHU, *Bouligand–Landweber iteration for a non-smooth ill-posed problem*, Numer. Math., (2019), p. in press.
- [4] A. DIEULEVEUT AND F. BACH, *Nonparametric stochastic approximation with large step-sizes*, Ann. Statist., 44 (2016), pp. 1363–1399.
- [5] H. W. ENGL, M. HANKE, AND A. NEUBAUER, *Regularization of Inverse Problems*, Kluwer, Dordrecht, 1996.
- [6] C. GEIERSBACH AND G. C. PFLUG, *Projected stochastic gradients for convex constrained problems in Hilbert spaces*, SIAM J. Optim., 29 (2019), pp. 2079–2099.
- [7] R. M. GOWER, N. LOIZOU, X. QIAN, A. SAILANBAYEV, E. SHULGIN, AND P. RICHTÁRIK, *SGD: general analysis and improved rates*, in Proceedings of the 36 th International Conference on Machine Learning, PMLR 97, K. Chaudhuri and R. Salakhutdinov, eds., Long Beach, California, 2019, pp. 5200–5209.
- [8] M. HANKE, A. NEUBAUER, AND O. SCHERZER, *A convergence analysis of the Landweber iteration for nonlinear ill-posed problems*, Numer. Math., 72 (1995), pp. 21–37.
- [9] G. T. HERMAN, A. LENT, AND P. H. LUTZ, *Relaxation method for image reconstruction*, Comm. ACM, 21 (1978), pp. 152–158.

929

- 930 [10] G. T. HERMAN AND L. B. MEYER, *Algebraic reconstruction techniques can be made computationally*
931 *efficient*, IEEE Trans. Medical Imag., 12 (1993), pp. 600–609.
- 932 [11] K. ITO AND B. JIN, *A new approach to nonlinear constrained Tikhonov regularization*, Inverse Prob-
933 *lems*, 27 (2011), pp. 105005, 23 pp.
- 934 [12] ———, *Inverse Problems: Tikhonov Theory and Algorithms*, World Scientific, Hackensack, NJ, 2015.
- 935 [13] Y. JIAO, B. JIN, AND X. LU, *Preasymptotic convergence of randomized Kaczmarz method*, Inverse
936 *Problems*, 33 (2017), pp. 125012, 21 pp.
- 937 [14] B. JIN AND X. LU, *On the regularizing property of stochastic gradient descent*, Inverse Problems, 35
938 (2019), pp. 015004, 27 pp.
- 939 [15] B. KALTENBACHER, A. NEUBAUER, AND O. SCHERZER, *Iterative Regularization Methods for Nonlinear*
940 *Ill-posed Problems*, Walter de Gruyter, Berlin, 2008.
- 941 [16] D. P. KINGMA AND J. BA, *Adam: a method for stochastic optimization*, in Proceedings of the 3rd
942 *International Conference on Learning Representations (ICLR)*, 2015.
- 943 [17] H. J. KUSHNER AND G. G. YIN, *Stochastic Approximation and Recursive Algorithms and Applications*,
944 Springer-Verlag, New York, second ed., 2003.
- 945 [18] L. LANDWEBER, *An iteration formula for Fredholm integral equations of the first kind*, Amer. J. Math.,
946 73 (1951), pp. 615–624.
- 947 [19] J. LIN AND L. ROSASCO, *Optimal rates for multi-pass stochastic gradient methods*, J. Mach. Learn.
948 *Res.*, 18 (2017), pp. 1–47.
- 949 [20] A. K. LOUIS, *Inverse und Schlecht Gestellte Probleme*, B. G. Teubner, Stuttgart, 1989.
- 950 [21] S. F. MCCORMICK AND G. H. RODRIGUE, *A uniform approach to gradient methods for linear operator*
951 *equations*, J. Math. Anal. Appl., 49 (1975), pp. 275–285.
- 952 [22] H. ROBBINS AND S. MONRO, *A stochastic approximation method*, Ann. Math. Stat., 22 (1951), pp. 400–
953 407.
- 954 [23] O. SCHERZER, M. GRASMAIR, H. GROSSAUER, M. HALTMEIER, AND F. LENZEN, *Variational Methods*
955 *in Imaging*, Springer, New York, 2009.
- 956 [24] T. SCHUSTER, B. KALTENBACHER, B. HOFMANN, AND K. S. KAZIMIERSKI, *Regularization Methods in*
957 *Banach Spaces*, Walter de Gruyter, Berlin, 2012.
- 958 [25] T. STROHMER AND R. VERSHYNIN, *A randomized Kaczmarz algorithm with exponential convergence*,
959 *J. Fourier Anal. Appl.*, 15 (2009), pp. 262–278.
- 960 [26] I. SUTSKEVER, J. MARTENS, G. DAHL, AND G. E. HINTON, *On the importance of initialization and*
961 *momentum in deep learning*, in Proceedings of the 30th International Conference on Machine
962 *Learning (ICML-13)*, S. Dasgupta and D. Mcallester, eds., Atlanta, GA, 2013, pp. 1139–1147.
- 963 [27] Y. S. TAN AND R. VERSHYNIN, *Phase retrieval via randomized Kaczmarz: theoretical guarantees*, Inf.
964 *Inference*, 8 (2019), pp. 97–123.
- 965 [28] P. TARRÈS AND Y. YAO, *Online learning as stochastic approximation of regularization paths: optimality*
966 *and almost-sure convergence*, IEEE Trans. Inform. Theory, 60 (2014), pp. 5716–5735.
- 967 [29] V. V. VASIN, *Iterative methods for solving ill-posed problems with a priori information in Hilbert*
968 *spaces*, Zh. Vychisl. Mat. i Mat. Fiz., 28 (1988), pp. 971–980, 1117.
- 969 [30] G. M. VAĬNIKKO AND A. Y. VERETENNIKOV, *Iteration Procedures in Ill-posed Problems*, “Nauka”,
970 Moscow, 1986.
- 971 [31] Y. YING AND M. PONTIL, *Online gradient descent learning algorithms*, Found. Comput. Math., 8
972 (2008), pp. 561–596.
- 973 [32] T. ZHANG, *Solving large scale linear prediction problems using stochastic gradient descent algorithms*,
974 in Proceedings of the Twenty First International Conference on Machine Learning, C. Brodley,
975 ed., Banff, Alberta, Canada, 2004, pp. 919–926.