

CSG: A stochastic gradient method for a wide class of optimization problems appearing in a machine learning or data-driven context

Lukas Pflug*, Max Grieshammer†, Andrian Uihlein†, and Michael Stingl†

Abstract. In a recent article the so called *continuous stochastic gradient method* (CSG) for the efficient solution of a class of stochastic optimization problems was introduced. While the applicability of known stochastic gradient type methods is typically limited to so called expected risk functions, no such limitation exists for CSG. The key to this lies in the computation of design dependent integration weights, which allows for an optimal usage of available information leading to stronger convergence properties. However, due to the nature of the formula for these integration weights, the practical applicability was essentially limited to problems, in which stochasticity enters via a low-dimensional and sufficiently simple probability distribution. In this paper the scope of the CSG method is significantly extended presenting new ways of calculating the integration weights. A full convergence analysis for this new variant of the CSG method is presented and its efficiency is demonstrated in comparison to more classical stochastic gradient methods by means of a number of problem classes, relevant in stochastic optimization and machine learning.

1. Introduction. In the context of optimization problems in which the expected-value of a cost function j is minimized, i.e.,

$$(1.1) \quad \min_{\theta \in \Theta^{\text{ad}}} \mathbb{E}[j(\theta, X)] = \int_{\mathcal{X}} j(\theta, x) \mu(\mathrm{d}x)$$

with probability measure μ and the associated random variables X , a variety of different stochastic optimization schemes has been developed in the past, e.g., [7, 12, 13, 22, 24]. Among the most popular algorithms are the *stochastic gradient method* (SG) [17] and its modification the *stochastic average gradient method* (SAG) [18], both of which shine with their low iteration cost and have been analyzed extensively. A variety of different stochastic optimization schemes has been developed in the past. Among the most popular algorithms are the *stochastic gradient method* (SG) [17] and its modification the *stochastic average gradient method* (SAG) [18]. Both of these methods have been analyzed extensively in literature and are characterized by a low cost per iteration.

Nonetheless, SG and SAG face a number of known disadvantages, like the lack of efficient stopping criteria (cf. [15]) or optimal stepsize rules (cf. [14, 19]). To tackle these issues, a whole variety of modified SG methods can be found in the literature. For example, [7] uses a trust-region-type model to normalize the steplengths, whereas the iSARAH algorithm proposed in [13] combines an inner SG scheme with an outer (inexact) full gradient descent method.

Another disadvantage of SG is the quite restrictive setting of (1.1). [22] and [24] suggest inexact proximal stochastic second-order methods and stochastic primal-dual fixed-point methods to allow for a different type of objective function appearing in (1.1). In the case that the constraints include expected-valued functions, a level set method is analyzed in [12].

*Central Institute for Scientific Computing (ZISC), lukas.pflug@fau.de

†Department of Mathematics, Chair of Applied Mathematics (Continuous Optimization), Friedrich-Alexander University Erlangen-Nürnberg (FAU), max.grieshammer@fau.de, andrian.uihlein@fau.de, michael.stingl@fau.de

39 An even wider class of problems, can be solved by the *continuous stochastic gradient*
40 *method* (CSG) proposed in [16]. The reason is that combining the information collected in
41 previous iterations in an optimal way, CSG gains a significantly improved gradient approxi-
42 mation and is able to estimate the current objective function value during the optimization
43 process. For a characterization of the class of problems CSG can solve, we refer to Remark
44 2.3. Here we just note that among them are objective functions with nested expectation values
45 (Section 5.2) and problems with chance constraints (Section 5.3).

46 While this is already known from the original version of CSG [16], there is also a serious
47 drawback: in order to approximate function values and gradients in the above mentioned
48 way, integration weights have to be computed by an analytical formula, which requires full
49 knowledge about the probability measure μ . Moreover the evaluation is based on a Voronoi
50 diagram, whose computation is not tractable, if the dimension of the parameter set \mathcal{X} is larger
51 than 2. As a consequence, in [16] only examples with a one-dimensional uniform distribution
52 were presented.

53 In this contribution, we expand the setting of CSG even further by introducing new meth-
54 ods of calculating the weights used for the gradient and cost function value approximations.
55 This enables us to apply the CSG algorithm to problems of higher dimension, to arbitrary
56 measure μ and even to problems where the measure μ appearing in (1.1) might be unknown,
57 e.g., in a data-driven context.

58 Depending on the concrete setting, i.e., depending on the dimensions of θ, x and on how
59 time-consuming the evaluation of a gradient sample is, the different methods allow us to
60 continuously trade weight-computation time and speed of convergence (w.r.t. number of
61 gradient sample evaluations).

62 In this article we present of a full convergence analysis for the CSG method extended in
63 this way. In particular, we show that the error in the gradient approximation as well as in
64 the objective function value approximation vanish as the number of steps increases. As a
65 consequence these values can be utilized, for instance, to apply stopping criteria based on first
66 order optimality conditions. Moreover, this potentially allows, to combine the CSG method
67 with slightly adapted step length strategies as they are known from the world of deterministic
68 optimization methods, a topic we leave open for future research.

69 The remaining structure of the paper is as follows. In Section 2, the mathematical structure
70 of the problems, we would like to solve by the CSG method is outlined in details. In Section 3,
71 the CSG method with generalized weight computation is presented. Section 4 is devoted to
72 the convergence analysis and in Section 5 we compare the generalized CSG method to more
73 traditional SG-type algorithm using three different classes of test problems.

74 **2. Problem setting and definitions.** Following the classic setup for expected-value ob-
75 jective functions, we introduce the set of admissible designs $\mathcal{P} \subset \mathbb{R}^{d_{\text{des}}}$ and the parameter
76 set $\mathcal{X} \subset \mathbb{R}^{d_{\text{par}}}$, where $d_{\text{des}}, d_{\text{par}} \in \mathbb{N}$. In the optimization process, the drawn random samples
77 x_1, x_2, \dots from the parameter set \mathcal{X} are assumed to be realizations of independent uniformly
78 random variables $X_i \sim \mu$ for all $i \in \mathbb{N}$, i.e., X_1, X_2, \dots are independent and follow an under-
79 lying probability distribution μ , which may be unknown.

80 To be precise, we define the following probability space setup:

81 **Definition 2.1 (Probability space setup).** *The probability space $(\Omega, \mathcal{A}, \mathbb{P})$ is given by*

$$82 \quad \begin{aligned} \Omega &:= \mathcal{X}^{\mathbb{N}}, \mathbb{P} := \mu^{\otimes \mathbb{N}}, \\ \mathcal{A} &:= \sigma(\{A_1 \times \dots \times A_n : A_i \in \mathcal{B}(\mathcal{X}), \forall i, n \in \mathbb{N}\}), \end{aligned}$$

where $\mu^{\otimes \mathbb{N}}(A_1 \times \dots \times A_n) = \prod_{i=1}^n \mu(A_i)$ is the product measure, μ is a probability measure on \mathcal{X} and $\sigma(\cdot)$ is the smallest σ -field that contains \cdot . We denote by $\text{supp}(\mu)$ the support of the measure μ , i.e.,

$$\text{supp}(\mu) := \{x \in \mathcal{X} : \mu(B_\varepsilon(x)) > 0 \forall \varepsilon > 0\},$$

83 where $B_\varepsilon(x)$ denotes an open ball of radius $\varepsilon > 0$ around $x \in \mathcal{X}$. We write $X_n : \Omega \rightarrow \mathcal{X}$,
84 $(\omega_k)_{k \in \mathbb{N}} \mapsto \omega_n$ for the projection to the $n = 1, 2, \dots$ coordinate and define $X := X_1$.

85 With this setup, the objective function takes the following form:

86 **Definition 2.2 (Objective function).** *The objective function $J : \mathcal{P} \rightarrow \mathbb{R}$ is given by*

$$87 \quad J(\theta) := \mathbb{E}[j(\theta, X)] = \int_{\mathcal{X}} j(\theta, x) \mu(dx)$$

88 with a measurable function $j \in C^1(\mathcal{P} \times \mathcal{X}; \mathbb{R})$ and random variable X .

89 **Remark 2.3 (Generalization of the setting).** During the optimization process, we may
90 also generate an approximation \hat{J}_n to the exact objective function value $J(\theta_n)$ with almost no
91 additional computational cost. We will show later that $\|\hat{J}_n - \nabla J(\theta_n)\|_{\mathcal{P}} \rightarrow 0$ (see Remark 4.8).

92 This enables us to solve a much broader class of optimization problems, where the objective
93 function may depend non-linearly on the expression above, i.e.,

$$94 \quad \tilde{J}(\theta) := f(\theta, \mathbb{E}[j(\theta, X)]),$$

95 with a Lipschitz continuously differentiable function $f : \mathcal{P} \times \mathbb{R} \rightarrow \mathbb{R}$. Included in the set of
96 possible objective functions are for example tracking functionals

$$97 \quad \tilde{J}(\theta) := \frac{1}{2} \|h(\theta, \cdot) - f(\theta, \mathbb{E}[j(\theta, \cdot, X)])\|_{L^2}^2$$

98 and nested expected values

$$99 \quad \tilde{J}(\theta) := \mathbb{E}_{\mathcal{Y}}[f(Y, \mathbb{E}_{\mathcal{X}}[j(\theta, X)])].$$

100 Notice that such settings can not be solved by SG algorithms.

101 As we are aiming for a gradient based optimization scheme, we further state the derivative of
102 the objective functional:

103 **Lemma 2.4 (Derivative of objective function).** *The gradient of the objective functional J is*
104 *given by $\nabla J(\theta) = \mathbb{E}[\delta(\theta, X)]$, where $X \sim \mu$ and $\delta : \mathcal{P} \times \mathcal{X} \rightarrow \mathbb{R}^{d_{\text{des}}}$ denotes $\nabla_1 j(\theta, x)$.*

105 **Proof.** This is a direct consequence of the linearity of the expectation value and the finite-
106 dimensional derivative of j . Integration and differentiation can be exchanged due to the
107 Lipschitz continuity of the integrand w.r.t. the integration variable. ■

108 In order to state and prove convergence results for the algorithm presented in this work,
 109 we define the norms on the used spaces as follows:

110 **Definition 2.5 (Norms on \mathcal{X} , \mathcal{P} and $\mathcal{P} \times \mathcal{X}$).** *In this contribution, we will use for the norm*
 111 *on the underlying spaces of the parameter space \mathcal{X} and the design space \mathcal{P} the notation $\|\cdot\|_{\mathcal{X}}$*
 112 *and $\|\cdot\|_{\mathcal{P}}$ respectively. Due to norm-equivalence in finite dimensional spaces, the norm used in*
 113 *the spaces \mathcal{P} , \mathcal{X} does not have to be specified and can be chosen problem specific. In addition,*
 114 *we define on $\mathcal{P} \times \mathcal{X}$ the following metric:*

$$115 \quad d((\theta, x), (\hat{\theta}, \hat{x})) := \|(\|\theta - \hat{\theta}\|_{\mathcal{P}}, \|x - \hat{x}\|_{\mathcal{X}})\|_1 \quad \forall (\theta, \hat{\theta}, x, \hat{x}) \in \mathcal{P}^2 \times \mathcal{X}^2.$$

116 *Choosing the 1-norm in the three-dimensional space as “outer”-norm is arbitrary and could*
 117 *for instance - in the other extreme case - be the ∞ -norm and of course could include positive*
 118 *weights for each individual component.*

119 **Assumption 2.6 (Regularity of the δ).** *We assume $\delta : \mathcal{P} \times \mathcal{X} \rightarrow \mathbb{R}^{d_{\text{des}}}$ to be bounded and*
 120 *Lipschitz continuous, i.e., there exist constants $C_{\delta}, L_{\delta} \in \mathbb{R}_{>0}$ s.t.*

$$121 \quad \begin{aligned} & \|\delta(\theta, x)\|_{\mathcal{P}} \leq C_{\delta} \\ & \|\delta(\theta, x) - \delta(\hat{\theta}, \hat{x})\|_{\mathcal{P}} \leq L_{\delta} (\|\theta - \hat{\theta}\|_{\mathcal{P}} + \|x - \hat{x}\|_{\mathcal{X}}) \end{aligned}$$

122 *for all $\theta, \hat{\theta} \in \mathcal{P}$ and $x, \hat{x} \in \mathcal{X}$. A sufficient condition therefore is to assume ∇j to be Lipschitz*
 123 *continuous in both arguments.*

124 For the convergence analysis of Algorithm 3.1, the following assumptions on the sets \mathcal{P} , \mathcal{X}
 125 and the measure μ are an important ingredient.

126 **Assumption 2.7 (Regularity of \mathcal{P} , \mathcal{X} and the measure μ).** *The set $\mathcal{P} \subset \mathbb{R}^{d_{\text{des}}}$ is compact*
 127 *and convex. $\text{supp}(\mu) \subset \mathcal{X}$ with $\mathcal{X} \subset \mathbb{R}^{d_{\text{par}}}$ is open and bounded. In addition, there exists*
 128 *$M_1, M_2, M_3 > 0$ s.t. $\forall \varepsilon \in (0, M_3)$ there exists $\mathcal{X}_{\varepsilon} \subset \mathcal{X}$ satisfying $\mu(\mathcal{X}_{\varepsilon}) \geq 1 - M_1 \varepsilon$ and*

$$129 \quad \inf_{x \in \mathcal{X}_{\varepsilon}} \mu(B_{\varepsilon}(x)) \geq M_2 \varepsilon^{d_{\text{par}}},$$

130 *where $B_{\varepsilon}(x) \subset \mathcal{X}$ is an open ball with radius ε centered in $x \in \mathcal{X}$.*

131 **Remark 2.8 (Examples for Assumption 2.7).** *In most cases, the choice $\mathcal{X}_{\varepsilon} = \mathcal{X}$ is suitable,*
 132 *for example when \mathcal{X} satisfies the uniform cone condition (cf. [1, Definition 4.8]). However,*
 133 *there exist cases where the possibility of choosing $\mathcal{X}_{\varepsilon} \subset \mathcal{X}$ in the condition of Assumption 2.7*
 134 *allows to consider even more general measures and sets.*

135 As example therefore, let $\mathcal{X} := \{1/n : n \in \mathbb{N}\}$ and $\mu := \sum_{k=1}^{\infty} 2^{-k} \delta_{k-1}(s)$. Then, for
 136 $M_n := [n^{-1}, 1] \cap \mathcal{X}$, it holds

$$137 \quad \mu(M_n) = \sum_{k=1}^n 2^{-k} = 1 - 2^{-n}.$$

138 Thus for $\varepsilon = 2^{-n}$ and $c' = 1$ we obtain $\mu(M_n) \geq 1 - c' \varepsilon$ and $\inf_{x \in M_n} \mu(B_{\varepsilon}(x)) \geq 2^{-n} = \varepsilon$. Since
 139 $\inf_{x \in \mathcal{X}} \mu(B_{\varepsilon}(x)) = \mu(B_{\varepsilon}(0)) = 2^{1-2^n}$, there exists no $c > 0$ such that $2^{1-2^n} \geq c 2^{-n} \forall n \in \mathbb{N}$.

140 For a uniform distribution and for all $0 < p \leq \infty$, the open p -Balls

$$141 \quad \mathcal{X}^p := \{x \in \mathbb{R}^{d_{\text{par}}} : \|x\|_p < 1\}$$

142 satisfy Assumption 2.7 as well. While case $1 \leq p \leq \infty$ allows for $\mathcal{X}_\varepsilon^p = \mathcal{X}^p$, for $0 < p < 1$ we
 143 first have to obtain $\mathcal{X}_\varepsilon^p$ by trimming of the spikes of \mathcal{X} .

144 **3. The algorithm.** To state the algorithm, we first define the projection operator which
 145 ensures the sequence of generated designs $(\theta_n)_{n \in \mathbb{N}}$ to be in the set \mathcal{P} .

146 **Definition 3.1 (Orthogonal projection).** We define the – in the sense of $\|\cdot\|_{\mathcal{P}}$ – orthogonal
 147 projection onto the set \mathcal{P} as follows:

$$148 \quad \text{Proj}_{\mathcal{P}}(\theta) := \arg \min_{\hat{\theta} \in \mathcal{P}} \|\theta - \hat{\theta}\|_{\mathcal{P}}.$$

149 Note that the existence and uniqueness of $\text{Proj}_{\mathcal{P}}$ is guaranteed by the projection theorem (see
 150 e.g. [3]) building on the convexity of \mathcal{P} as assumed in Assumption 2.7.

151 **Lemma 3.2 (Properties of $\text{Proj}_{\mathcal{P}}$).** Let $\mathcal{P} \subset \mathbb{R}^{d_{\text{des}}}$ satisfy Assumption 2.7. Then the
 152 following holds for all $x, y \in \mathbb{R}^{d_{\text{des}}}$ and $z \in \mathcal{P}$:

- 153 (a) $(\text{Proj}_{\mathcal{P}}(x) - x)^T (\text{Proj}_{\mathcal{P}}(x) - z) \leq 0$,
- 154 (b) $(\text{Proj}_{\mathcal{P}}(y) - \text{Proj}_{\mathcal{P}}(x))^T (y - x) \geq \|\text{Proj}_{\mathcal{P}}(y) - \text{Proj}_{\mathcal{P}}(x)\|_{\mathcal{P}}^2 \geq 0$,
- 155 (c) $\|\text{Proj}_{\mathcal{P}}(y) - \text{Proj}_{\mathcal{P}}(x)\|_{\mathcal{P}} \leq \|y - x\|_{\mathcal{P}}$.

156 *Proof.* A proof of (a) can be found in [3, Thm. 1.4.1 (ii)], (b) and (c) correspond to (iii)
 157 and (ii) in [3, Prop. 1.4.1] respectively. ■

158 Given θ_1 , $n = 1$ and a sequence x_1, x_2, \dots of inputs, where we assume that they are
 159 realizations of the independent random variables X_1, X_2, \dots introduced in Section 2, the CSG
 method for the (possibly unknown) measure μ is given in Algorithm 3.1.

Algorithm 3.1 CSG method

- 1: **while** Termination condition not met **do**
 - 2: Sample objective function (optional):
 $j_n := j(\theta_n, x_n)$
 - 3: Sample gradient:
 $g_n := \nabla_{\theta} j(\theta_n, x_n)$
 - 4: Calculate weights α_k
 - 5: Calculate search direction:
 $\hat{G}_n := \frac{1}{n} \sum_{k=1}^n \alpha_k g_k$
 - 6: Approximation to objective function value (optional):
 $\hat{J}_n := \frac{1}{n} \sum_{k=1}^n \alpha_k j_k$
 - 7: Choose stepsize τ_n
 - 8: Gradient step:
 $\theta_{n+1} := \text{Proj}_{\mathcal{P}}(\theta_n - \tau_n \hat{G}_n)$
 - 9: Update index:
 $n \leftarrow n + 1$
 - 10: **end while**
-

160

161 **3.1. Calculating the weights.** The quality of the weights α_k appearing in Algorithm 3.1
 162 greatly impacts the accuracy of the gradient approximation \hat{G}_n and therefore directly in-
 163 fluences the overall performance of the CSG method. On the other hand, a more optimal
 164 computation of the weights might be time-consuming. Since the trade-off between the time
 165 spent calculating the weights and the time gained by performing fewer gradient evaluations is
 166 heavily problem-specific, we propose four different methods for the weight-calculation in the
 167 n th step:

168 **Exact.** Following an exact nearest neighbor approximation for the integral

$$169 \quad \nabla J(\theta_n) = \int_{\mathcal{X}} \nabla_{\theta} j(\theta_n, x) \mu(dx),$$

170 for each $k = 1, \dots, n$ we define the set

$$171 \quad M_k := \{x \in \mathcal{X} : d((\theta_n, x), (\theta_k, x_k)) < d((\theta_n, x), (\theta_j, x_j)) \text{ for all } j \in \{1, \dots, n\} \setminus \{k\}\},$$

172 i.e., the set of points $x \in \mathcal{X}$ such that (θ_n, x) is closer to (θ_k, x_k) than to any other previous
 173 evaluation point. Assuming that the measure μ is known, we then set $\alpha_k := \mu(M_k)$. This
 174 method has been thoroughly analyzed in [16] and yields the best possible approximation to
 175 the exact gradient, but is computationally infeasible for problems of high dimensions.

176 **Empirical.** Utilizing the properties of the empirical measure μ_n (see Remark 4.6), we may
 177 replace the exact weights mentioned above by the empirical weights

$$178 \quad \alpha_k := \frac{1}{n} \sum_{i=1}^n 1_{M_k}(x_i) = \mu_n(M_k) \approx \mu(M_k),$$

179 where 1_{M_k} denotes the indicator function of the set M_k . Note that the computation of the
 180 empirical weights requires no knowledge of μ and is also feasible for high-dimensional problems,
 181 but needs many samples x_i to approximate the exact gradient with a high accuracy.

182 **Exact hybrid.** Assuming that the dimension of \mathcal{X} is much smaller than the dimension of
 183 \mathcal{P} , we might treat the designs and parameters separately. Instead of M_k , we now consider the
 184 sets

$$185 \quad \widetilde{M}_i = \{x \in \mathcal{X} : \|x - x_i\|_{\mathcal{X}} \leq \|x - x_j\|_{\mathcal{X}} \text{ for all } j = 1, \dots, n\}, \quad i = 1, \dots, n.$$

186 The α_k are now calculated as a combination of the empirical and exact method

$$187 \quad (3.1) \quad \alpha_k := \sum_{i=1}^n 1_{M_k}(x_i) \mu(\widetilde{M}_i).$$

188 **Inexact hybrid.** As for the exact weights, the calculation of the exact hybrid weights
 189 requires knowledge of μ . If μ is unknown, we may replace the factor $\mu(\widetilde{M}_i)$ in (3.1) by an
 190 empirical approximation. Since this only requires samples of X , which we assume to have a
 191 plenitude of, we can control the quality of this approximation through the number of samples
 192 we draw. The inexact hybrid weights are therefore calculated as follows:

$$193 \quad \alpha_k := \frac{1}{\lfloor n^\beta \rfloor} \sum_{i=1}^n 1_{M_k}(x_{j_i}) \sum_{m=1}^{\lfloor n^\beta \rfloor} 1_{\widetilde{M}_{j_i}}(x_m),$$

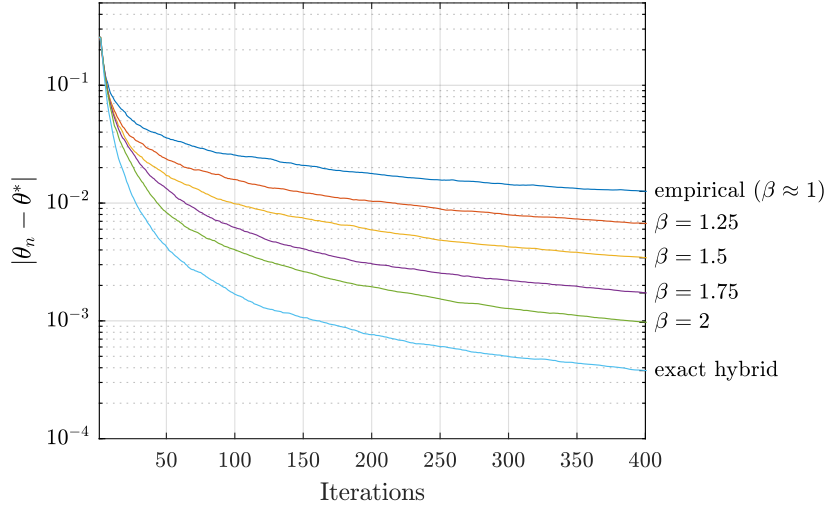


Figure 1. Absolute error $|\theta_n - \theta^*|$ in iteration n for the setting $\mathcal{P} = [-\frac{1}{2}, \frac{1}{2}] = \mathcal{X}$, $j(\theta, x) = \frac{1}{2}(\theta - x)^2$ and $X \sim \mathcal{U}_{\mathcal{X}}$. The curves correspond to the median of 1000 runs with constant stepsizes $\tau_n = 1$ and randomized starting points in \mathcal{P} .

194 where $\beta \geq 1$, $\lfloor n^\beta \rfloor$ is the total number of samples we have drawn until step n and x_{j_i} denote
 195 the samples where $\nabla_{\theta} j(\theta, x)$ has been evaluated at.

196 Figure 1 shows that the inexact hybrid method allows us to interpolate between the purely
 197 empirical method and the exact hybrid variant by choosing β appropriately.

198 *Remark 3.3.* In general, the nearest neighbor approximation, which is used in all methods
 199 mentioned above, worsens as the dimension of $\mathcal{P} \times \mathcal{X}$ increases (cf. [4]). Especially for
 200 problems where $\dim(\mathcal{P}) \ll \dim(\mathcal{X})$, results from Monte Carlo integration ([23]) suggest that
 201 the performance boost gained by better weight calculation starts to become negligible. The
 202 proposed CSG methods are therefore best suited for optimization problems where \mathcal{X} is of
 203 small dimension when compared to \mathcal{P} and the evaluation of $j(\theta, x)$ is time-consuming.

204 Furthermore, the metric d should be chosen problem-specific to ensure the best possible
 205 performance.

206 *Remark 3.4 (SAG and SG as two extreme cases of the algorithm).*

207 As stated in Definition 2.5, our metric d can be chosen as

$$208 \quad d((\theta, x), (\hat{\theta}, \hat{x})) = a_1 \|\theta - \hat{\theta}\|_{\mathcal{P}} + a_2 \|x - \hat{x}\|_{\mathcal{X}},$$

209 where $a_1, a_2 > 0$ are arbitrary. By choosing $a_1 \gg a_2$, the nearest neighbor to (θ_n, x) is almost
 210 exclusively determined by the distance in the design variable. Hence, for the weights α_k we
 211 get $\alpha_n \approx 1$ and $\alpha_1, \dots, \alpha_{n-1} \approx 0$, i.e., the CSG algorithm will behave very similar to the usual
 212 SG algorithm.

213 Analogously, choosing $a_1 \ll a_2$ will lead to a performance similar to SAG.

214 **4. Convergence analysis.** In this section we will study the convergence of the proposed
 215 algorithm. By the matter of the randomly chosen evaluation point within the algorithm, we

216 will have to study probabilistic convergence behaviour in terms of “almost sure convergence”.
 217 Therefore, we first state first order optimality conditions, assumptions on the regularity of the
 218 involved functions as well as the steplength τ and a suitable probability space setting.

219 **4.1. Optimality conditions and assumptions.** For $h \in C^1(\mathcal{P})$ and \mathcal{P} convex we have the
 220 following equivalent sufficient conditions for first order optimality:

221 **Corollary 4.1 (Optimality conditions).** *For all $\theta^* \in \mathcal{P}$ the following items are equivalent:*

222 (a) $-\nabla h(\theta^*)^T(\theta - \theta^*) \leq 0 \quad \forall \theta \in \mathcal{P}$

223 (b) $\mathcal{P}(\theta^* - t\nabla h(\theta^*)) = \theta^* \quad \forall t \geq 0.$

224 A point $\theta^* \in \mathcal{P}$ satisfying these conditions is called a stationary point.

225 *Proof.* The proof can be found in e.g. [16]. ■

226 In order to guarantee that Algorithm 3.1 generates a convergent subsequence, the stepsizes
 227 have to be damped, i.e., $(\tau_n)_{n \in \mathbb{N}}$ has to be a null series with upper and lower bound as stated
 228 in the following Assumption. However, in contrast to the ordinary stochastic gradient decent
 229 method, if Algorithm 3.1 generates – with stepsizes satisfying $\tau_n \geq \tau > 0 \quad \forall n \in \mathbb{N}$ – a
 230 convergent sequence, the limit point is a stationary point of the objective function too. This
 231 is shown in Theorem 4.11.

232 **Assumption 4.2 (Steplength).** *The steplength $(\tau_n)_{n \in \mathbb{N}}$ in Algorithm 3.1 satisfies the fol-*
 233 *lowing: $\exists N \in \mathbb{N}, \underline{S}, \bar{S} \in \mathbb{R}_{>0}$ and $D \in (0, \frac{1}{\max\{d_{\text{par}}, 2\}})$ s.t.*

$$234 \quad \underline{S}n^{-1} \leq \tau_n \leq \bar{S}n^{-1 + \frac{1}{\max\{d_{\text{par}}, 2\}} - D} \quad \forall n \in \mathbb{N}_{>N}.$$

235 These bounds on the steplength satisfy the conditions stated in [17, Eqns. (6) and (26)] as
 236 well as equivalently in [6, Eqn. (4.19)] in the one-dimensional case and can be seen as a higher
 237 dimensional equivalent.

238 In the following we assume that these assumptions are always satisfied without mentioning
 239 it explicitly.

240 **4.2. Error in the search direction.** In this subsection we analyse the error in the n -th
 241 iteration of the search direction \hat{G}_n and the gradient of the objective functional ∇J_n . For
 242 this, we define the following random variables:

243 **Definition 4.3 (Random variables).** *For $x \in \mathcal{X}$ and $\omega \in \Omega$ the sequence of random variables*
 244 *$(Z_n)_{n \in \mathbb{N}}$ with $Z_n : \Omega \times \mathcal{X} \rightarrow \mathbb{R}_{\geq 0}$ is defined by*

$$245 \quad Z_n(\omega, x) := \min_{k=1, \dots, n} d((\Theta_k(\omega), X_k(\omega)), (\Theta_n(\omega), x)),$$

246 where the designs $\Theta_k \in \mathcal{P}$ for $k > 1$ depend by their construction on the initial design Θ_1 and
 247 all “previous” random variables $X_1(\omega), \dots, X_{k-1}(\omega)$, i.e.,

$$248 \quad \Theta_k(\Theta_1, X_1(\omega), \dots, X_{k-1}(\omega))$$

249 and thus is also a random variable. We shorten this dependency by the notation $\Theta_k(\omega)$.

250 This random variable fulfills the following property:

251 **Lemma 4.4.** For μ almost all $x \in \text{supp}(\mu)$

$$252 \quad \sum_{n=1}^{\infty} \mathbb{P}(Z_n(\cdot, x) > \varepsilon_n) < \infty \quad \text{and} \quad \sum_{n=1}^{\infty} \sup_{x \in \mathcal{X}_{\varepsilon_n}} \mathbb{P}(Z_n(\cdot, x) > \varepsilon_n) < \infty,$$

253 with

$$254 \quad (4.1) \quad \varepsilon_n := \frac{C_\delta \bar{S}}{1 - 2^{-\frac{1}{2 \max\{d_{\text{par}}, 2\}}}} \cdot n^{-\frac{D}{2}} + \tilde{\varepsilon}_n \quad \text{and} \quad \tilde{\varepsilon}_n := n^{\frac{D}{2} - \frac{1}{2 \max\{d_{\text{par}}, 2\}}}.$$

255 Therein, C_δ is defined in Assumption 2.6, $\mathcal{X}_{\varepsilon_n}$ in Assumption 2.7 and \bar{S}, D in Assumption 4.2.

256 *Proof.* We first define $i_0 \in \mathbb{N}$ as an auxiliary index as follows:

$$257 \quad i_0 := \lceil n - a_n + 1 \rceil \quad \text{with} \quad a_n := n^{1 + \frac{D}{2} - \frac{1}{\max\{d_{\text{par}}, 2\}}}.$$

258 By construction, we have

$$\begin{aligned} \mathbb{P}(Z_n(\cdot, x) \geq \varepsilon_n) &\leq \mathbb{P}\left(\min_{k=i_0, \dots, n} d((\Theta_k, X_k), (\Theta_n, x)) \geq \varepsilon_n\right) \\ 259 \quad &\leq \mathbb{P}\left(\sum_{i=i_0}^{n-1} \|\tau_i \hat{G}_i\|_{\mathcal{P}} + \min_{k=i_0, \dots, n} \|X_k - x\|_{\mathcal{X}} \geq \varepsilon_n\right) \\ &\leq \mathbb{P}\left(C_\delta \sum_{i=i_0}^{n-1} \tau_i + \min_{k=i_0, \dots, n} \|X_k - x\|_{\mathcal{X}} \geq \varepsilon_n\right). \end{aligned}$$

260 Observe that for $n > 2$ we obtain for all $\kappa \in (0, 1)$

$$\begin{aligned} \sum_{i=i_0}^{n-1} \frac{1}{i^\kappa} &\leq \int_{i_0-1}^n \frac{1}{s^\kappa} ds = \frac{1}{1-\kappa} \cdot (n^{1-\kappa} - ([n - a_n]^{1-\kappa})) \leq \frac{1}{1-\kappa} \cdot (n^{1-\kappa} - (n - a_n)^{1-\kappa}) \\ 261 \quad &= \frac{1}{1-\kappa} \cdot \left(\frac{n}{n^\kappa} - \frac{n - a_n}{(n - a_n)^\kappa}\right) = \frac{n}{1-\kappa} \cdot \left(\frac{(n - a_n)^\kappa - n^\kappa}{n^\kappa \cdot (n - a_n)^\kappa}\right) + \frac{a_n}{(1-\kappa)(n - a_n)^\kappa} \\ &= \frac{n}{1-\kappa} \cdot \left(\frac{n^\kappa(1 - a_n/n)^\kappa - n^\kappa}{n^\kappa \cdot (n - a_n)^\kappa}\right) + \frac{a_n}{(1-\kappa)(n - a_n)^\kappa}. \end{aligned}$$

262 Applying Bernoulli's inequality in the first term, we conclude

$$\begin{aligned} 263 \quad \sum_{i=i_0}^{n-1} \frac{1}{i^\kappa} &\leq \frac{n}{1-\kappa} \cdot \left(\frac{n^\kappa(1 - \kappa \cdot \frac{a_n}{n}) - n^\kappa}{n^\kappa \cdot (n - a_n)^\kappa}\right) + \frac{a_n}{(1-\kappa)(n - a_n)^\kappa} \\ 264 \quad &= \frac{1}{1-\kappa} \cdot \left(\frac{-\kappa a_n}{(n - a_n)^\kappa}\right) + \frac{a_n}{(1-\kappa)(n - a_n)^\kappa} \\ 265 \quad (4.2) \quad &= \frac{a_n}{(n - a_n)^\kappa} = \frac{a_n}{n^\kappa} \left(1 - \frac{a_n}{n}\right)^{-\kappa}. \\ 266 \end{aligned}$$

267 Combining Assumption 4.2 and (4.2) yields

$$268 \quad \sum_{i=i_0}^{n-1} \tau_i \leq \bar{S} \sum_{i=i_0}^{n-1} \frac{1}{i^{1+D-\frac{1}{\max\{d_{\text{par}},2\}}}} \leq \bar{S} \frac{a_n}{n^\kappa} (1 - \frac{a_n}{n})^{-\kappa},$$

269 with $D \in \left(0, \frac{1}{\max\{d_{\text{par}},2\}}\right)$ and $\kappa := 1 + D - \frac{1}{\max\{d_{\text{par}},2\}} \in (0, 1)$. Hence, for $n \geq 2$ we obtain

$$270 \quad \left(1 - \frac{a_n}{n}\right)^{-\kappa} = \left(1 - n^{\frac{D}{2} - \frac{1}{\max\{d_{\text{par}},2\}}}\right)^{-\kappa} \leq \left(1 - n^{-\frac{1}{2\max\{d_{\text{par}},2\}}}\right)^{-\kappa} \leq \left(1 - 2^{-\frac{1}{2\max\{d_{\text{par}},2\}}}\right)^{-1}.$$

271 Collecting these results, we see

$$272 \quad \sum_{i=i_0}^{n-1} \tau_i \leq \frac{\bar{S}}{1 - 2^{-\frac{1}{2\max\{d_{\text{par}},2\}}}} \frac{a_n}{n^{1+D-\frac{1}{\max\{d_{\text{par}},2\}}}} \leq \frac{\bar{S}}{1 - 2^{-\frac{1}{2\max\{d_{\text{par}},2\}}}} n^{-\frac{D}{2}}.$$

273 Consequently,

$$274 \quad \sum_{i=i_0}^{n-1} \|\tau_i \hat{G}_i\|_{\mathcal{P}} \leq \frac{C_\delta \bar{S}}{1 - 2^{-\frac{1}{2\max\{d_{\text{par}},2\}}}} n^{-\frac{D}{2}} = \varepsilon_n - \tilde{\varepsilon}_n.$$

275 By Assumption 2.7, $\mu(\mathcal{X} \setminus \mathcal{X}_{\varepsilon_n}) \rightarrow 0$. Hence, for μ almost all $x \in \text{supp}(\mu)$, there exists
276 $n \in \mathbb{N}$ large enough, such that $x \in \mathcal{X}_{\varepsilon_n}$. Therefore,

$$\begin{aligned} & \mathbb{P}(Z_n(\cdot, x) \geq \varepsilon_n) \\ & \leq \mathbb{P}\left(\min_{k=i_0, \dots, n-1} \|X_k - x\|_{\mathcal{X}} \geq \tilde{\varepsilon}_n\right) \\ & \leq \mathbb{P}\left(\|X_k - x\|_{\mathcal{X}} \geq \tilde{\varepsilon}_n \quad \forall k \in \{i_0, \dots, n-1\}\right) \\ 277 & = \prod_{k=i_0}^{n-1} \mathbb{P}\left(\|X_k - x\|_{\mathcal{X}} \geq \tilde{\varepsilon}_n\right) = \prod_{k=i_0}^{n-1} (1 - \mu(B_{\tilde{\varepsilon}_n}(x))) \\ & \leq \left(1 - \min\{M_2(\tilde{\varepsilon}_n)^{d_{\text{par}}}, 1\}\right)^{a_n}. \end{aligned}$$

278 As $\tilde{\varepsilon}_n \rightarrow 0$, there exists $N \in \mathbb{N}$ s.t. for $n \geq N$ we obtain

$$279 \quad \mathbb{P}(Z_n(\cdot, x) \geq \varepsilon_n) \leq \left(1 - M_2 n^{\frac{d_{\text{par}} D}{2} - \frac{d_{\text{par}}}{2\max\{2, d_{\text{par}}\}}}\right)^{a_n}.$$

280 For simplicity, we define

$$281 \quad c_1 := \frac{d_{\text{par}} D}{2} - \frac{d_{\text{par}}}{2\max\{2, d_{\text{par}}\}}$$

282 and recall that $\log(1-x) \leq -x$ for all $x < 1$. Since $c_1 < 0$, for n large enough it holds

$$\begin{aligned}
& \left(1 - M_2 n^{\frac{d_{\text{par}} D}{2} - \frac{d_{\text{par}}}{2 \max\{2, d_{\text{par}}\}}}\right)^{a_n} = (1 - M_2 n^{c_1})^{a_n} = \exp(a_n \log(1 - M_2 n^{c_1})) \\
& \leq \exp(-a_n M_2 n^{c_1}) = \exp\left(-M_2 n^{1 + \frac{D}{2} - \frac{1}{\max\{2, d_{\text{par}}\}} + \frac{d_{\text{par}} D}{2} - \frac{d_{\text{par}}}{2 \max\{2, d_{\text{par}}\}}}\right) \\
& = \begin{cases} \exp\left(-M_2 n^{1+D-\frac{1}{2}-\frac{1}{4}}\right) & d_{\text{par}} = 1 \\ \exp\left(-M_2 n^{1+\frac{D}{2}-\frac{1}{d_{\text{par}}}+\frac{d_{\text{par}} D}{2}-\frac{1}{2}}\right) & d_{\text{par}} \geq 2 \end{cases} \leq \exp(-M_2 n^D).
\end{aligned}$$

284 Recall that there is $N \in \mathbb{N}$ such that $\exp(-x) \leq x^{-\frac{2}{D}}$ for all $x \geq N$. It follows that for all n
285 large enough: $\exp(-M_2 n^D) \leq M_2^{-\frac{2}{D}} n^{-2}$. Hence,

$$\sum_{n=N}^{\infty} \left(1 - M_2 n^{\frac{d_{\text{par}} D}{2} - \frac{d_{\text{par}}}{2 \max\{2, d_{\text{par}}\}}}\right)^{a_n} \leq \sum_{n=N}^{\infty} \exp(-M_2 n^D) \leq \sum_{n=N}^{\infty} M_2^{-\frac{2}{D}} n^{-2}$$

287 and thus

$$\sum_{n=1}^{\infty} \mathbb{P}(Z_n(\cdot, x) > \varepsilon_n) < \infty.$$

289 Finally, note that Assumption 2.7 gives

$$\sup_{x \in \mathcal{X}^{\varepsilon_n}} \mathbb{P}(Z_n(\cdot, x) \geq \varepsilon_n) \leq \left(1 - c \cdot \frac{n^{\frac{D}{2}}}{a_n}\right)^{a_n} = \left(1 - c \cdot n^{\frac{1}{\max\{d_{\text{par}}, 2\}}}\right)^{a_n}$$

291 with $c > 0$. By the same steps as above, we obtain

$$\sum_{n=1}^{\infty} \sup_{x \in \mathcal{X}^{\varepsilon_n}} \mathbb{P}(Z_n(\cdot, x) > \varepsilon_n) < \infty.$$

293 As a direct consequence of the latter result we get almost sure convergence.

294 **Corollary 4.5.** For μ almost all $x \in \text{supp}(\mu)$

$$Z_n(\cdot, x) \xrightarrow{a.s.} 0 \quad \text{for } n \rightarrow \infty.$$

296 *Proof.* The result follows by Lemma 4.4 and the Borel-Cantelli Lemma (see for example
297 Theorem 2.7 in [11]).

298 **Remark 4.6 (Empirical distribution).** The empirical measure defined as

$$(4.3) \quad \mu_n := \frac{1}{n} \sum_{i=1}^n \delta_{X_i}$$

300 satisfies $\mu_n \Rightarrow \mu$ as $n \rightarrow \infty$ almost surely, see [20, Theorem 3]. Here \Rightarrow denotes the weak
301 convergence of measures which is the weak-* convergence in dual space theory, i.e.,

$$\mu_n \Rightarrow \mu \quad \text{iff} \quad \int_{\mathcal{X}} f(x) \mu_n(dx) \rightarrow \int_{\mathcal{X}} f(x) \mu(dx) \quad \forall f \in C_b(\mathcal{X}, \mathbb{R}).$$

303 See for instance [5] for the empirical distribution and [8, Section 7.3] for a functional
 304 analytical perspective on weak-* convergence in the discussed function space setting.

305 Since this property of μ_n is all we need in the following proofs and since the measures

$$306 \quad \mu_n^{\text{eh}} := \sum_{i=1}^n \delta_{X_i} \mu(\widetilde{M}_i) \quad \text{and} \quad \mu_n^{\text{ih}} := \sum_{i=1}^n \delta_{X_{j_i} \mu_{[n^\beta]}(\widetilde{M}_{j_i})},$$

307 corresponding to exact hybrid weights and inexact hybrid weights respectively, satisfy $\mu_n^{\text{eh}} \Rightarrow \mu$
 308 and $\mu_n^{\text{ih}} \Rightarrow \mu$ as well, we will w.l.o.g. work with empirical weights only.

309 Thus, due to the Lipschitz continuity of δ as defined in Theorem 2.6, the expected value
 310 $\nabla J(\theta) = \mathbb{E}[\delta(\theta, X)]$ is for $n \rightarrow \infty$ better and better approximated by \hat{G}_n :

311 **Theorem 4.7 (Error in gradient approximation).**

312 *The norm of the difference between the search direction \hat{G}_n and the gradient of the objective*
 313 *functional $\mathbb{E}[\delta(\Theta_n, X)]$ vanishes for $n \rightarrow \infty$, i.e.,*

$$314 \quad \|\hat{G}_n - \mathbb{E}[\delta(\Theta_n, X)]\|_{\mathcal{P}} \xrightarrow{a.s.} 0 \quad \text{and} \quad \lim_{n \rightarrow \infty} \mathbb{E} \left[\|\hat{G}_n - \mathbb{E}[\delta(\Theta_n, X)]\|_{\mathcal{P}} \right] = 0.$$

315 *Proof.* For $x \in \text{supp}(\mu)$ define

$$316 \quad k^n(\omega; x) := \arg \min_{k=1, \dots, n} d((\Theta_k(\omega), X_k(\omega)), (\Theta_n(\omega), x)).$$

317 For \hat{G}_n as generated by Algorithm 3.1 with $n \in \mathbb{N}$ arbitrary but fixed the following holds:

$$\begin{aligned} 318 \quad & \|\hat{G}_n - \mathbb{E}[\delta(\Theta_n, X)]\|_{\mathcal{Y}} \\ 319 \quad &= \left\| \sum_{i=1}^n \int_{\mathcal{X}} \delta_{k^n(\omega; x)}(i) \delta(\Theta_i(\omega), x_i) \mu_n(dx) - \int_{\mathcal{X}} \delta(\Theta_n(\omega), x) \mu(dx) \right\|_{\mathcal{Y}} \\ 320 \quad &\leq \left\| \int_{\mathcal{X}} \sum_{i=1}^n \delta_{k^n(\omega; x)}(i) \delta(\Theta_i(\omega), x_i) - \delta(\Theta_n(\omega), x) \mu_n(dx) \right\|_{\mathcal{Y}} \\ 321 \quad &\quad + \left\| \int_{\mathcal{X}} \delta(\Theta_n(\omega), x) \mu_n(dx) - \int_{\mathcal{X}} \delta(\Theta_n(\omega), x) \mu(dx) \right\|_{\mathcal{Y}} \\ 322 \quad (4.4) \quad &\leq \mathsf{L}_\delta \int_{\mathcal{X}} Z_n(\omega, x) \mu_n(dx) + \left\| \int_{\mathcal{X}} \delta(\Theta_n(\omega), x) \mu_n(dx) - \int_{\mathcal{X}} \delta(\Theta_n(\omega), x) \mu(dx) \right\|_{\mathcal{Y}}, \\ 323 \end{aligned}$$

324 where μ_n is the empirical measure given in Remark (4.3) and L_δ the Lipschitz constant defined
 325 in Assumption 2.6. We need to prove that both terms in (4.4) vanish for $n \rightarrow \infty$.

326 For the first term, the uniform (in n) Lipschitz continuity of $Z_n(\omega, \cdot)$ yields

$$\begin{aligned} 327 \quad & \int_{\mathcal{X}} Z_n(\omega, x) \mu_n(dx) = \int_{\mathcal{X}} Z_n(\omega, x) \mu(dx) + \int_{\mathcal{X}} Z_n(\omega, x) \mu_n(dx) - \int_{\mathcal{X}} Z_n(\omega, x) \mu(dx) \\ & \leq \int_{\mathcal{X}} Z_n(\omega, x) \mu(dx) + \mathsf{L}_Z d_W(\mu_n, \mu), \end{aligned}$$

328 where d_W denotes the Wasserstein distance of the measure μ_n and μ (see [10]). Since \mathcal{X} is
329 bounded, [10, Theorem 6] gives that the Wasserstein distance metrizes the weak topology on
330 the set of probability measures on \mathcal{X} . Since $\mu_n \Rightarrow \mu$ almost surely, this gives $d_W(\mu_n, \mu) \rightarrow 0$
331 almost surely. Furthermore, by Assumption 2.6, there exists $C > 0$ s.t. $0 \leq Z_n \leq C$. Using
332 Corollary 4.5, we obtain $Z_n(\omega, x) \rightarrow 0$ for almost all $\omega \in \Omega$. Therefore, Lebesgue's dominated
333 convergence theorem yields

$$334 \quad \int_{\mathcal{X}} Z_n(\omega, x) \mu(dx) \rightarrow 0 \quad \text{for almost all } \omega \in \Omega.$$

335 In order to show that the second part of (4.4) vanishes, observe that

$$\begin{aligned} & \left\| \int_{\mathcal{X}} \delta(\Theta_n(\omega), x) \mu_n(dx) - \int_{\mathcal{X}} \delta(\Theta_n(\omega), x') \mu(dx') \right\|_{\mathcal{P}} \\ 336 \quad &= \left\| \int_{\mathcal{X} \times \mathcal{X}} \delta(\Theta_n(\omega), x) - \delta(\Theta_n(\omega), x') Q_n(d(x, x')) \right\|_{\mathcal{P}} \\ &\leq L_{\delta} \int_{\mathcal{X} \times \mathcal{X}} \|x - x'\|_{\mathcal{X}} Q_n(d(x, x')), \end{aligned}$$

337 where $Q_n(\cdot \times \mathcal{X}) = \mu_n$ and $Q_n(\mathcal{X} \times \cdot) = \mu$ is an arbitrary but fixed coupling of μ_n and μ .
338 By taking the infimum of all such couplings, we again obtain the Wasserstein distance d_W of
339 the measure μ_n and μ , i.e.,

$$340 \quad (4.5) \quad \left\| \int_{\mathcal{X}} \delta(\Theta_n(\omega), x) \mu_n(dx) - \int_{\mathcal{X}} \delta(\Theta_n(\omega), x) \mu(dx) \right\|_{\mathcal{P}} \leq L_{\delta} d_W(\mu_n, \mu).$$

341 By the same arguments as mentioned earlier, $d_W(\mu_n, \mu) \rightarrow 0$ almost surely. Combining all
342 the above facts gives

$$343 \quad \|\hat{G}_n - \mathbb{E}[\delta(\Theta_n, X)]\|_{\mathcal{P}} \rightarrow 0$$

344 almost surely. Since the above quantities are bounded, the almost sure convergence also
345 implies the convergence in expectation via Lebesgue's dominated convergence theorem. ■

346 *Remark 4.8.* Due to the regularity of J , we can show

$$347 \quad \|\hat{J}_n - J(\theta_n)\|_{\mathcal{P}} \rightarrow 0$$

348 analogously to the proof of Theorem 4.7.

349 **Theorem 4.9 (Sum of error in gradient approximation).** *The expectation value of the summed*
350 *norm of the difference between the search direction \hat{G}_n and the gradient of the reduced objective*
351 *functional ∇J weighted by the respected stepsize τ_n vanishes for $n \rightarrow \infty$, i.e.,*

$$352 \quad (4.6) \quad \sum_{n=1}^{\infty} \tau_n \mathbb{E} \left[\|\hat{G}_n - \mathbb{E}[\delta(\Theta, X)]\|_{\mathcal{P}} \right] < \infty.$$

353 *Proof.* Recall from the proof of Theorem 4.7 that

$$\begin{aligned}
354 \quad \mathbb{E} \left[\|\hat{G}_n - \mathbb{E}[\delta(\Theta_n, X)]\|_{\mathcal{P}} \right] &\leq L_\delta \mathbb{E} \left[\int_{\mathcal{X}} Z_n(\omega, x) \mu_n(dx) \right] \\
355 \quad (4.7) \quad &+ \mathbb{E} \left[\left\| \int_{\mathcal{X}} \delta(\Theta_n(\omega), x) \mu_n(dx) - \int_{\mathcal{X}} \delta(\Theta_n(\omega), x) \mu(dx) \right\|_{\mathcal{P}} \right]. \\
356
\end{aligned}$$

357 We start with deriving an upper bound for the first term on the right hand side of the latter
358 inequality. Recall the definition of $\tilde{\varepsilon}_n$ in Lemma 4.4, i.e.,

$$359 \quad \tilde{\varepsilon}_n := n^{-\frac{1}{\max\{2, d_{\text{par}}\} + \frac{D}{2}}}$$

360 with D as defined in Assumption 4.2. Then, analogue to Lemma 4.4 (cf. the proof and the
361 notation there), together with

$$362 \quad D := \sup_{(\tilde{\theta}, \tilde{x}), (\hat{\theta}, \hat{x}) \in \mathcal{P} \times \mathcal{X}} d((\tilde{\theta}, \tilde{x}), (\hat{\theta}, \hat{x})),$$

363 we obtain the following estimate:

$$\begin{aligned}
\mathbb{E} \left[\int_{\mathcal{X}} Z_n(\cdot, x) \mu_n(dx) \right] &= \mathbb{E} \left[\int_{\mathcal{X}} Z_n(\cdot, x) 1_{Z_n(\cdot, x) \leq \tilde{\varepsilon}_n} + Z_n(\cdot, x) 1_{Z_n(\cdot, x) > \tilde{\varepsilon}_n} \mu_n(dx) \right] \\
364 \quad &\leq \tilde{\varepsilon}_n + D \mathbb{E} \left[\int_{\mathcal{X}} 1_{Z_n(\cdot, x) > \tilde{\varepsilon}_n} \mu_n(dx) \right] \\
&\leq \tilde{\varepsilon}_n + D \mathbb{E} \left[\int_{\mathcal{X}} \prod_{k=1}^n 1_{d((\Theta_k(\cdot), X_k(\cdot)), (\Theta_n(\cdot), x)) > \tilde{\varepsilon}_n} \mu_n(dx) \right].
\end{aligned}$$

365 Setting $i_0 := \lceil n - a_n + 1 \rceil$ as in the proof of Lemma 4.4 yields

$$366 \quad \mathbb{E} \left[\int_{\mathcal{X}} Z_n(\cdot, x) \mu_n(dx) \right] \leq \tilde{\varepsilon}_n + D \mathbb{E} \left[\int_{\mathcal{X}} \prod_{k=i_0}^n 1_{d((\Theta_k(\cdot), X_k(\cdot)), (\Theta_n(\cdot), x)) > \tilde{\varepsilon}_n} \mu_n(dx) \right].$$

367 Since μ_n is the empirical measure as defined in (4.3) and due to the linearity of \mathbb{E} , we obtain

$$\begin{aligned}
368 \quad \mathbb{E} \left[\int_{\mathcal{X}} Z_n(\cdot, x) \mu_n(dx) \right] &= \tilde{\varepsilon}_n + \frac{D}{n} \sum_{i=1}^n \mathbb{E} \left[\prod_{\substack{k=i_0 \\ k \neq i}}^n 1_{d((\Theta_k(\cdot), X_k(\cdot)), (\Theta_n(\cdot), X_i(\cdot))) > \tilde{\varepsilon}_n} \right] \\
&= \tilde{\varepsilon}_n + \frac{D}{n} \sum_{i=1}^n \prod_{\substack{k=i_0 \\ k \neq i}}^n \mathbb{P}(d((\Theta_k(\cdot), X_k(\cdot)), (\Theta_n(\cdot), X_i(\cdot))) > \tilde{\varepsilon}_n),
\end{aligned}$$

369 Where we used the independency of all $(X_i)_{i \in \mathbb{N}}$. Finally, applying Fubini's theorem results in

$$370 \quad (4.8) \quad \mathbb{E} \left[\int_{\mathcal{X}} Z_n(\cdot, x) \mu_n(dx) \right] = \tilde{\varepsilon}_n + \frac{D}{n} \sum_{i=1}^n \prod_{\substack{k=i_0 \\ k \neq i}}^n \int_{\mathcal{X}} \mathbb{P}(d((\Theta_k(\cdot), X_k(\cdot)), (\Theta_n(\cdot), x)) > \tilde{\varepsilon}_n) \mu(dx).$$

371 Let $\mathcal{X}_{\tilde{\varepsilon}_n} \subset \mathcal{X}$ be the set given in Assumption 2.7. Following the same argumentation as in the
 372 proof of Lemma 4.4, we obtain

$$\begin{aligned}
 & \int_{\mathcal{X}} \mathbb{P}(d((\Theta_k(\cdot), X_k(\cdot)), (\Theta_n(\cdot), x)) > \tilde{\varepsilon}_n) \mu(dx) \\
 &= \int_{\mathcal{X} \setminus \mathcal{X}_{\tilde{\varepsilon}_n}} \mathbb{P}(d((\Theta_k(\cdot), X_k(\cdot)), (\Theta_n(\cdot), x)) > \tilde{\varepsilon}_n) \mu(dx) \\
 373 &+ \int_{\mathcal{X}_{\tilde{\varepsilon}_n}} \mathbb{P}(d((\Theta_k(\cdot), X_k(\cdot)), (\Theta_n(\cdot), x)) > \tilde{\varepsilon}_n) \mu(dx) \\
 &\leq c' \tilde{\varepsilon}_n + \sup_{x \in \mathcal{X}_{\tilde{\varepsilon}_n}} \mathbb{P}(d((\Theta_k(\cdot), X_k(\cdot)), (\Theta_n(\cdot), x)) > \tilde{\varepsilon}_n) \leq c' \tilde{\varepsilon}_n + \left(1 - c \cdot \frac{n^{\frac{D}{2}}}{a_n}\right)^{a_n},
 \end{aligned}$$

374 with a_n defined as in Lemma 4.4 and $c > 0$. Utilizing $\log(1 - x) \leq -x$ for all $x < 1$ shows

$$\begin{aligned}
 375 \left(1 - c \cdot \frac{n^{\frac{D}{2}}}{a_n}\right)^{a_n} &= \exp\left(a_n \log\left(1 - c \cdot n^{\frac{1}{\max\{d_{\text{par}}, 2\}} - 1}\right)\right) \leq \exp\left(-ca_n \cdot n^{\frac{D}{2}}\right) \\
 &= \exp\left(-c \cdot n^{1+D - \frac{1}{\max\{d_{\text{par}}, 2\}}}\right) \leq \exp\left(-c \cdot n^D\right) \leq \tilde{\varepsilon}_n
 \end{aligned}$$

376 for n large enough. Therefore, we have

$$377 \int_{\mathcal{X}} \mathbb{P}(d((\Theta_k(\cdot), X_k(\cdot)), (\Theta_n(\cdot), x)) > \tilde{\varepsilon}_n) \mu(dx) \leq (c' + 1) \tilde{\varepsilon}_n.$$

378 Inserting into (4.8) yields

$$379 (4.9) \quad \mathbb{E} \left[\int_{\mathcal{X}} Z_n(\cdot, x) \mu_n(dx) \right] \leq \tilde{\varepsilon}_n + D((c' + 1) \tilde{\varepsilon}_n)^{n-i_0+1} \leq \bar{c} \cdot \tilde{\varepsilon}_n$$

380 for n large enough and some $\bar{c} > 0$.

381 Now, in order to bound the second term in (4.7), recall from (4.5) that

$$382 \mathbb{E} \left[\left\| \int_{\mathcal{X}} \boldsymbol{\delta}(\theta_n, x) \mu_n(dx) - \boldsymbol{\delta}(\theta_n, x) \mu(dx) \right\|_{\mathcal{P}} \right] \leq \mathbb{L}_{\delta} \cdot \mathbb{E} [d_W(\mu_n, \mu)],$$

383 where d_W is the Wasserstein distance. By [9, Thm. 1 for $q = 3$ and $p = 1$], for all $d_{\text{par}} \geq 1$
 384 there exists $\hat{C}(d_{\text{par}}) \in \mathbb{R}_{>0}$ s.t.:

$$\begin{aligned}
 385 \mathbb{E} [d_W(\mu_n, \mu)] &\leq \hat{C}(d_{\text{par}}) \cdot M_3 \begin{cases} n^{-\frac{1}{2}}, & d_{\text{par}} = 1, \\ n^{-\frac{1}{2}} \log(1 + n), & d_{\text{par}} = 2, \\ n^{-\frac{1}{d_{\text{par}}}}, & d_{\text{par}} \geq 3, \end{cases} \\
 386 (4.10) &\leq \hat{C}(d_{\text{par}}) \cdot M_3 \cdot n^{-\frac{1}{\max\{d_{\text{par}}, 2\}}} \log(1 + n),
 \end{aligned}$$

388 with $M_3 := \left(\int_{\mathcal{X}} \|x\|_{\mathcal{X}}^3 \mu(dx)\right)^{1/3}$.

389 Substituting (4.9) and (4.10) into (4.7) yields

(4.11)

$$390 \quad \sum_{n=N}^{\infty} \tau_n \mathbb{E} \left[\|\hat{G}_n - \mathbb{E}[\delta(\theta, X)]\|_{\mathcal{P}} \right] \leq \bar{c} \sum_{n=N}^{\infty} \tau_n \bar{\varepsilon}_n + \hat{C}(d_{\text{par}}) M_3 \sum_{n=N}^{\infty} \tau_n n^{-\frac{1}{\max\{d_{\text{par}}, 2\}}} \log(1+n)$$

391 for $N \in \mathbb{N}$ large enough. By Assumption 4.2, we have

$$392 \quad \tau_n \leq \bar{S} n^{-1-D+\frac{1}{\max\{d_{\text{par}}, 2\}}},$$

393 which, when inserted into (4.11), gives

$$394 \quad \sum_{n=N}^{\infty} \tau_n \mathbb{E} \left[\|\hat{G}_n - \mathbb{E}[\delta(\theta, X)]\|_{\mathcal{P}} \right] \leq \bar{c} \bar{S} \sum_{n=N}^{\infty} n^{-1-\frac{D}{2}} + \hat{C}(d_{\text{par}}) M_3 \bar{S} \sum_{n=N}^{\infty} n^{-1-D} \log(1+n),$$

395 showing that

$$\sum_{n=1}^{\infty} \tau_n \mathbb{E} \left[\|\hat{G}_n - \mathbb{E}[\delta(\theta, X)]\|_{\mathcal{P}} \right] \leq \infty.$$

396

397 Before we can present our main result, we collect a few auxiliary results.

398 **Lemma 4.10 (Collection of auxiliary results).**

399

400 (a) *The objective functional value in iteration $n \in \mathbb{N}$ satisfies*

$$401 \quad J_{n+1} - J_n \leq -\frac{1}{\tau_n} \|\theta_{n+1} - \theta_n\|_{\mathcal{P}}^2 + \phi_n,$$

402 where $\phi_n := \tau_n \|\nabla J_n - \hat{G}_n\|_{\mathcal{P}} \|\hat{G}_n\|_{\mathcal{P}} + \tau_n^2 C \|\hat{G}_n\|_{\mathcal{P}}^2$ and $C \geq 0$ denotes a constant de-
403 pending only on the Lipschitz constants and suprema of the involved functions.

404 (b) *For ϕ_n as defined above, it holds $\sum_{n=1}^{\infty} \mathbb{E}[\phi_n] < \infty$.*

405 (c) *For all $t \geq 0$, we have*

$$406 \quad \|\text{Proj}_{\mathcal{P}}(\theta_n - t\hat{G}_n) - \theta_n\|_{\mathcal{P}} \leq \frac{t}{\tau_n} \|\theta_{n+1} - \theta_n\|_{\mathcal{P}}.$$

407 *Proof.* Assertions (a), (b) and (c) correspond to Lemma 16, Corollary 17 and Lemma 18
408 in [16]. Note that, by Theorem 4.9, the proofs given therein can be carried over to our setting
409 as well. ■

410 **Theorem 4.11 (Main theorem).** *Let $(\theta_n)_{n \in \mathbb{N}}$ be generated by Algorithm 3.1 with weights
411 calculated by one of the methods mentioned in Section 3.1. Then there exists a sub-sequence
412 $(\theta_{n_k})_{k \in \mathbb{N}}$ converging to a stationary point, i.e.,*

$$413 \quad \liminf_{n \rightarrow \infty} \mathbb{E} \left[\|\text{Proj}_{\mathcal{P}}(\theta_n - t\nabla J_n) - \theta_n\|_{\mathcal{P}}^2 \right] = 0 \quad \text{for all } t \geq 0.$$

414 *On the other hand, assume the time-step series $(\tau_n)_{n \in \mathbb{N}}$ satisfies $\tau_n \geq \tau$ for all $n \in \mathbb{N}$ and
415 some $\tau > 0$. Let further $(x_n)_{n \in \mathbb{N}}$ be dense in \mathcal{X} and assume $(\theta_n)_{n \in \mathbb{N}}$ converges to $\theta^* \in \mathcal{P}$.
416 Then θ^* is a stationary point of J , i.e.*

$$417 \quad \|\text{Proj}_{\mathcal{P}}(\theta^* - t\nabla J(\theta^*)) - \theta^*\|_{\mathcal{P}}^2 = 0 \quad \text{for all } t \geq 0.$$

418 *Proof.* To prove the first part, we show

$$419 \quad (4.12) \quad \sum_{n=1}^{\infty} \tau_n \mathbb{E} \left[\|\text{Proj}_{\mathcal{P}}(\theta_n - t\nabla J_n) - \theta_n\|_{\mathcal{P}}^2 \right] < \infty \quad \text{for all } t \geq 0.$$

420 By the assumed compactness of \mathcal{P} and regularity of J , we have

$$421 \quad J_{\text{inf}} := \inf_{\theta \in \mathcal{P}} J(\theta) > -\infty.$$

422 For arbitrary $N \in \mathbb{N}$, Lemma 4.10 (a) gives

$$423 \quad J_{\text{inf}} - J_1 \leq \mathbb{E}[J_{N+1}] - J_1 = \sum_{n=1}^N \mathbb{E}[J_{n+1} - J_n] \leq \sum_{n=1}^N \left(-\frac{1}{\tau_n} \mathbb{E} \left[\|\theta_{n+1} - \theta_n\|_{\mathcal{P}}^2 \right] + \mathbb{E}[\phi_n] \right).$$

424 Rearranging terms and utilizing Lemma 4.10 (b) yields

$$425 \quad (4.13) \quad \sum_{n=1}^{\infty} \frac{1}{\tau_n} \mathbb{E} \left[\|\theta_{n+1} - \theta_n\|_{\mathcal{P}}^2 \right] \leq J_1 - J_{\text{inf}} + \sum_{n=1}^{\infty} \mathbb{E}[\phi_n] < \infty.$$

426 By Lemma 3.2 (c) and Lemma 4.10 (c), we obtain

$$\begin{aligned} & \|\text{Proj}_{\mathcal{P}}(\theta_n - t\nabla J_n) - \theta_n\|_{\mathcal{P}}^2 \\ 427 \quad & \leq \left(\|\text{Proj}_{\mathcal{P}}(\theta_n - t\hat{G}_n) - \theta_n\|_{\mathcal{P}} + \|\text{Proj}_{\mathcal{P}}(\theta_n - t\nabla J_n) - \text{Proj}_{\mathcal{P}}(\theta_n - t\hat{G}_n)\|_{\mathcal{P}} \right)^2 \\ & \leq \left(\frac{t}{\tau_n} \|\theta_{n+1} - \theta_n\|_{\mathcal{P}} + t\|\hat{G}_n - \nabla J_n\|_{\mathcal{P}} \right)^2 \leq \frac{2t^2}{\tau_n^2} \|\theta_{n+1} - \theta_n\|_{\mathcal{P}}^2 + 2t^2 \|\hat{G}_n - \nabla J_n\|_{\mathcal{P}}^2, \end{aligned}$$

428 where we used Young's inequality in the last line. Therefore, it holds

$$429 \quad \sum_{n=1}^{\infty} \tau_n \mathbb{E} \left[\|\text{Proj}_{\mathcal{P}}(\theta_n - t\nabla J_n) - \theta_n\|_{\mathcal{P}}^2 \right] \leq 2t^2 \sum_{n=1}^{\infty} \frac{1}{\tau_n} \mathbb{E} \left[\|\theta_{n+1} - \theta_n\|_{\mathcal{P}}^2 \right] + 2t^2 \sum_{n=1}^{\infty} \tau_n \mathbb{E} \left[\|\hat{G}_n - \nabla J_n\|_{\mathcal{P}}^2 \right].$$

430 (4.12) now follows from (4.13) and Theorem 4.9.

431 For the second part, observe that convergence of $(\theta_n)_{n \in \mathbb{N}}$ and density of $(x_n)_{n \in \mathbb{N}}$ in \mathcal{X} yield

$$432 \quad Z_n(x) \rightarrow 0 \quad \text{for all } x \in \mathcal{X}.$$

433 Therefore, by similar steps as performed in the proof of Theorem 4.7, it holds

$$434 \quad \|\hat{G}_n - \nabla J_n\|_{\mathcal{P}} \rightarrow 0,$$

435 where ∇J_n denotes $\nabla J(\theta_n)$. Hence, for all $t \geq 0$ we obtain

$$\begin{aligned} & \|\text{Proj}_{\mathcal{P}}(\theta^* - t\nabla J(\theta^*)) - \theta^*\|_{\mathcal{P}}^2 \\ & = \lim_{n \rightarrow \infty} \|\text{Proj}_{\mathcal{P}}(\theta_n - t\nabla J(\theta_n)) - \theta_n\|_{\mathcal{P}} \\ 436 \quad & \leq \lim_{n \rightarrow \infty} \left(\|\text{Proj}_{\mathcal{P}}(\theta_n - t\hat{G}_n) - \theta_n\|_{\mathcal{P}} + \|\text{Proj}_{\mathcal{P}}(\theta_n - t\hat{G}_n) - \text{Proj}_{\mathcal{P}}(\theta_n - t\nabla J_n)\|_{\mathcal{P}} \right) \\ & \leq \lim_{n \rightarrow \infty} \left(\frac{t}{\tau_n} \|\theta_{n+1} - \theta_n\|_{\mathcal{P}} + \|\text{Proj}_{\mathcal{P}}(\theta_n - t\hat{G}_n) - \text{Proj}_{\mathcal{P}}(\theta_n - t\nabla J_n)\|_{\mathcal{P}} \right) \\ & \leq \lim_{n \rightarrow \infty} \frac{t}{\tau} \|\theta_{n+1} - \theta_n\|_{\mathcal{P}} + \lim_{n \rightarrow \infty} t \|\hat{G}_n - \nabla J_n\|_{\mathcal{P}} = 0, \end{aligned}$$

437 where we used Lemma 4.10 (c) for the second inequality. ■

438 **5. Numerical Results.** In this section, we consider three different settings in which we
 439 compare the CSG methods to suiting algorithms from the literature. The comparison is based
 440 on the number of gradient evaluations, since these represent the time-consuming computations
 441 in complex optimization tasks.

442 **5.1. Comparison with SG.** To start our numerical analysis, we consider the problem

$$443 \quad \min_{\theta \in \mathcal{P}} \frac{1}{2} \int_{\mathcal{X}} (x - \theta)^2 dx,$$

444 where $\mathcal{P} = \mathcal{X} = [-\frac{1}{2}, \frac{1}{2}]$.

445 To study the behavior of the algorithms, we choose four different stepsizes (n^{-1} , $n^{-2/3}$,
 446 $n^{-1/3}$ and a constant stepsize of 1) and track the absolute error in each iteration $|\theta_n - \theta^*|$.
 447 In order to obtain meaningful results, the 10000 starting points were chosen randomly in \mathcal{P} .
 448 For a comparison, we do the same for the ordinary stochastic gradient descent method (SG),
 since it is one of the most commonly used techniques for problems like our example.

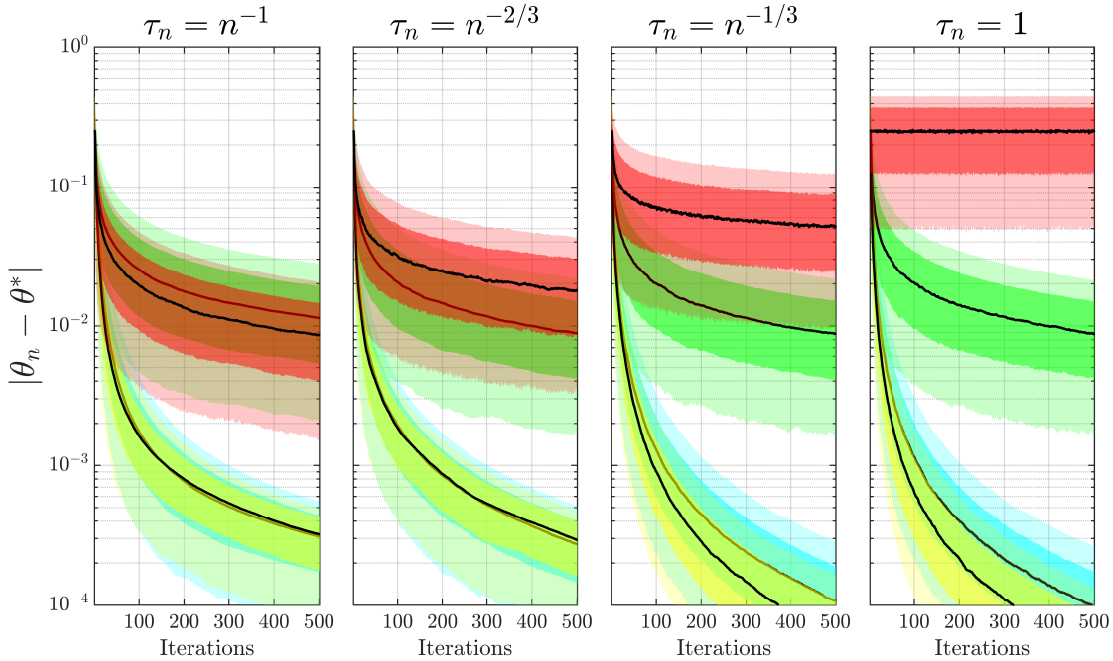


Figure 2. Comparison of the absolute error $|\theta_n - \theta^*|$ for SG (red), CSG with empirical weights (green), exact hybrid CSG (cyan) and exact CSG (yellow).

449 Notice that, in contrast to SG, a larger stepsize does not worsen the performance of the
 450 CSG algorithms for our example. Instead, a constant stepsize leads to a faster convergence
 451 for the hybrid and exact CSG method, whereas SG fails to solve the problem.
 452

453 **5.2. Comparison with SCGD.** As mentioned in Remark 2.3, the vanishing error in inner
 454 function value approximations allows us to solve optimization problems in which the cost

455 function depends non-linearly on a suiting expectation value. For instance, we may solve the
 456 problem

$$457 \quad (5.1) \quad \min_{\theta \in \mathcal{P}} \frac{1}{20} \int_{\mathcal{Y}} \left(2y + 5 \int_{\mathcal{X}} \cos\left(\frac{\theta - x}{\pi}\right) dx \right)^2 dy,$$

458 where $\mathcal{P} = [0, 10]$, $\mathcal{X} = [-1, 1]$ and $\mathcal{Y} = [-3, 3]$. The optimal solution $\theta^* = \frac{\pi^2}{2}$ to this example
 459 can be found analytically. Setting

$$460 \quad f_y(t) := \frac{3}{10}(2y + t)^2 \quad \text{and} \quad g_x(t) := 10 \cos\left(\frac{\theta - x}{\pi}\right),$$

461 problem (5.1) can be reformulated as

$$462 \quad (5.2) \quad \min_{\theta \in \mathcal{P}} \mathbb{E}_{\mathcal{Y}}[f_y(\mathbb{E}_{\mathcal{X}}[g_x(\theta)])].$$

463 Since f_y is non-linear, the SG algorithm can not be used to solve (5.1). Therefore, we compare
 464 our results with the so called stochastic compositional gradient descent (SCGD) method (see
 465 [21]), which is specifically designed for problems of the form (5.2).

466 Again, the 1000 starting points are randomly generated. This time however, we draw the
 467 starting points only from the interval $[\frac{11}{2}, \frac{19}{2}]$ instead of $\mathcal{P} = [0, 10]$. The reason for this is that
 468 the optimal solution $\frac{\pi^2}{2} \approx 4.935$ would otherwise be very close to the median starting point,
 469 resulting in artificially small absolute errors for all methods. Since the objective function in
 470 (5.1) is strongly convex in a neighborhood of the optimal solution, the accelerated SCGD
 471 method (see [21]) performed better than the standard version. Therefore, we compared our
 472 results to the aSCGD algorithm and chose the optimal stepsizes for aSCGD according to
 473 Theorem 7 in [21]. For the hybrid, inexact hybrid and empirical CSG algorithm, we chose a
 474 constant stepsize of $\frac{1}{30}$, which is a rough approximation to the inverse of the Lipschitz constant
 475 $L_{\nabla J}$. The resulting graphs are shown in Figure 3.

476 From a practical viewpoint, one is mainly interested in how many iterations it takes the
 477 error to fall below a desired tolerance. For this purpose, we analyzed the number of steps after
 478 which the different methods achieved a given absolute error with 90% certainty. The results
 479 can be seen in Figure 4.

480 **5.3. Chance constraint problems.** As a prototype example for chance constraint prob-
 481 lems, we consider

$$482 \quad \begin{aligned} & \max_{\theta \in [0, \frac{3}{4}]} \theta \\ & \text{s.t.} \quad \mathbb{P}(\theta - X^2 \leq 0) \geq \frac{1}{2}, \quad X \sim \mathcal{U}_{[-1,1]} \end{aligned}$$

483 with optimal solution $\theta^* = \frac{1}{4}$. By introducing the characteristic function $\chi_{[0, \infty)}$ and trans-
 484 forming the constraint to a penalty term, we arrive at

$$485 \quad \max_{\theta \in [0, \frac{3}{4}]} \theta - \lambda \max \left\{ 0, \frac{1}{2} \int_{-1}^1 \chi_{[0, \infty)}(\theta - x^2) dx - \frac{1}{2} \right\}.$$

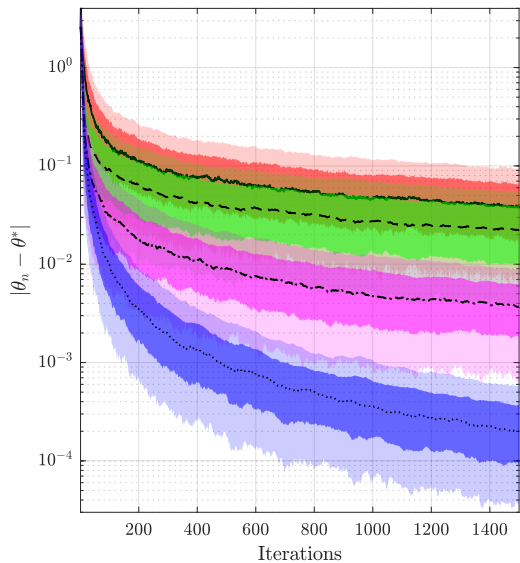


Figure 3. Comparison of the absolute error $|\theta_n - \theta^*|$. From top to bottom: aSCGD (red/solid), CSG with empirical weights (green/dashed), inexact hybrid CSG with $\beta = 1.5$ (magenta/dash-dotted) and hybrid CSG (blue/dotted). The shaded areas indicate the quantiles $P_{0.1,0.9}$ (light) and $P_{0.25,0.75}$ (dark).

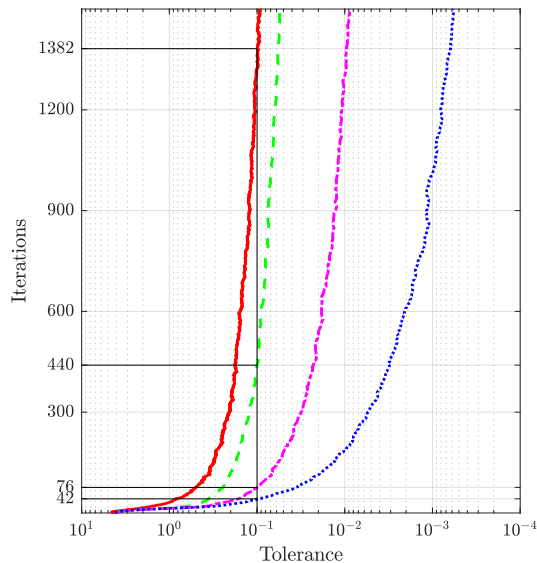


Figure 4. Minimum number of steps needed for aSCGD (red/solid), CSG with empirical weights (green/dashed), inexact hybrid CSG with $\beta = 1.5$ (magenta/dash-dotted) and hybrid CSG (blue/dotted) such that at least 90% of the runs achieved an absolute error smaller than the given tolerance.

486 Since the penalized objective function is no longer continuously differentiable, we can not
 487 guarantee the existence of a gradient and will have to work with subgradients instead, cf. [2].
 488 Notice that the proofs provided above also hold true for a subgradient method, if the stepsize
 489 is chosen accordingly. While the computation of a (sub-)gradient of $\max\{0, \cdot\}$ is not an issue,
 490 $\chi_{[0,\infty)}$ needs to be regularized further. The final problem then reads as follows:

$$491 \quad (5.3) \quad \max_{\theta \in [0, \frac{3}{4}]} \theta - \lambda \max \left\{ 0, \frac{1}{4} \int_{-1}^1 ((\tanh(\alpha(\theta - x^2)) + 1) dx - \frac{1}{2} \right\}.$$

492 Due to the non-linearity of $\max\{0, \cdot\}$, we again choose the SCGD method for comparison.
 493 This time, the objective function is not strongly convex in a neighborhood of θ_{opt} . Therefore,
 494 the stepsizes for the standard SCGD method are chosen according to Theorem 6 in [21], i.e.,
 495 optimal for this setting. For the CSG algorithms, we choose $\tau_n = \frac{1}{n}$. Lastly, we fix $\lambda = 3$ and
 496 $\alpha = 25$. The optimal solution θ_{opt} to (5.3) then satisfies $|\theta^* - \theta_{opt}| < 1.5 \cdot 10^{-3}$. The results
 497 of 1000 runs with random starting points in $[0, \frac{3}{4}]$ are presented in Figure 5.

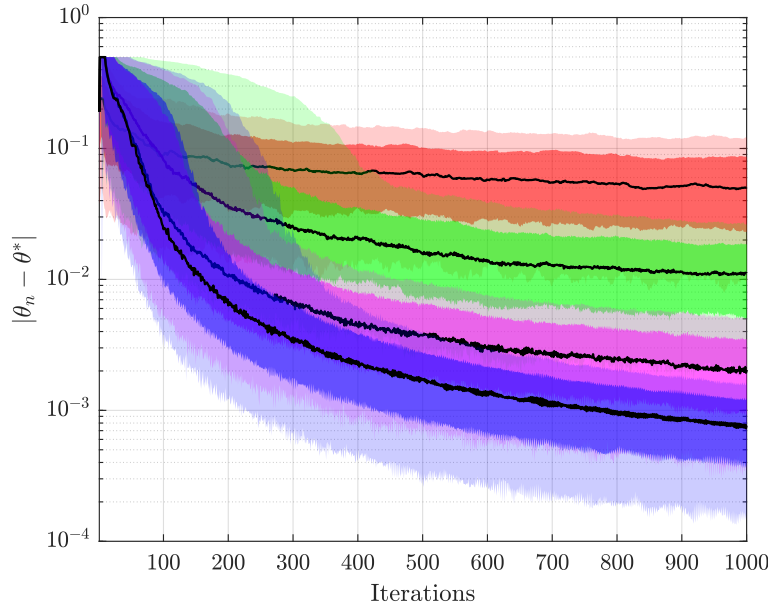


Figure 5. Comparison of the absolute error $|\theta_n - \theta_{opt}|$. From top to bottom: SCGD (red), CSG with empirical weights (green), inexact hybrid CSG with $\beta = 1.5$ (magenta) and hybrid CSG (blue). The shaded areas indicate the quantiles $P_{0.1,0.9}$ (light) and $P_{0.25,0.75}$ (dark).

498 **6. Conclusion and Outlook.** In this article a more flexible way to compute design de-
 499 pendent integration weights for the efficient approximation of the full cost function and its
 500 gradient when applying the CSG method to a class of stochastic optimization problems was
 501 introduced. While this significantly widened the scope of the CSG method, there are still
 502 a number of research questions, which would be very interesting to be investigated in the
 503 future. First, as a consequence of the strong convergence properties shown in this paper, the
 504 CSG method – in the course of the optimization iterations – behaves more and more like a
 505 fully deterministic descent method. This calls for more elaborate techniques to calculate the
 506 step length such as linesearch or trust region strategies. Another interesting question is, if
 507 convergence of the iterates generated by the CSG method can be shown for a constant choice
 508 of the step size. Indeed the numerical examples we have presented in this paper suggest that
 509 this should be possible. And finally, exploiting specific structures of the given probability
 510 distributions, one could come up with even more efficient integration techniques allowing to
 511 solve problems with high dimensional distributions more efficiently than using the empirical
 512 weight strategy presented in this article.

513 **Acknowledgements.** Funded by the Deutsche Forschungsgemeinschaft (DFG, German
 514 Research Foundation) through project D05 in the CRC 1411 (Project-ID 416229255) and
 515 subproject B06 in TRR 154 (Project-ID 239904186).

516

REFERENCES

- 517 [1] R. A. ADAMS AND J. J. F. FOURNIER, *Sobolev spaces*, vol. 140 of Pure and Applied Mathematics
518 (Amsterdam), Elsevier/Academic Press, Amsterdam, second ed., 2003.
- 519 [2] Y. I. ALBER, A. N. IUSEM, AND M. V. SOLODOV, *On the projected subgradient method for nonsmooth*
520 *convex optimization in a Hilbert space*, Math. Programming, 81 (1998), pp. 23–35.
- 521 [3] J.-P. AUBIN, *Applied functional analysis*, Pure and Applied Mathematics (New York), Wiley-Interscience,
522 New York, second ed., 2000. With exercises by Bernard Cornet and Jean-Michel Lasry, Translated
523 from the French by Carole Labrousse.
- 524 [4] K. BEYER, J. GOLDSTEIN, R. RAMAKRISHNAN, AND U. SHAFT, *When is "nearest neighbor" meaningful?*,
525 ICDT 1999. LNCS, 1540 (1997).
- 526 [5] P. BILLINGSLEY, *Convergence of probability measures*, Wiley Series in Probability and Statistics: Prob-
527 ability and Statistics, John Wiley & Sons, Inc., New York, second ed., 1999. A Wiley-Interscience
528 Publication.
- 529 [6] L. BOTTOU, F. E. CURTIS, AND J. NOCEDAL, *Optimization methods for large-scale machine learning*,
530 SIAM Rev., 60 (2018), pp. 223–311.
- 531 [7] F. E. CURTIS, K. SCHEINBERG, AND R. SHI, *A stochastic trust region algorithm based on careful step*
532 *normalization*, INFORMS J. Optim., 1 (2019), pp. 200–220.
- 533 [8] G. B. FOLLAND, *A guide to advanced real analysis*, vol. 37 of The Dolciani Mathematical Expositions,
534 Mathematical Association of America, Washington, DC, 2009. MAA Guides, 2.
- 535 [9] N. FOURNIER AND A. GUILLIN, *On the rate of convergence in wasserstein distance of the empirical*
536 *measure*, Probability Theory and Related Fields, 162 (2015), pp. 707–738.
- 537 [10] A. L. GIBBS AND F. E. SU, *On choosing and bounding probability metrics*, International statistical review,
538 70 (2002), pp. 419–435.
- 539 [11] A. KLENKE, *Probability theory*, Universitext, Springer-Verlag London, Ltd., London, 2008. A compre-
540 hensive course, Translated from the 2006 German original.
- 541 [12] Q. LIN, S. NADARAJAH, N. SOHEILI, AND T. YANG, *A data efficient and feasible level set method for*
542 *stochastic convex optimization with expectation constraints*, J. Mach. Learn. Res., 21 (2020), pp. Paper
543 No. 143, 45.
- 544 [13] L. M. NGUYEN, K. SCHEINBERG, AND M. TAKÁVČ, *Inexact SARAH algorithm for stochastic optimization*,
545 Optim. Methods Softw., 36 (2021), pp. 237–258.
- 546 [14] C. PAQUETTE AND K. SCHEINBERG, *A stochastic line search method with expected complexity analysis*,
547 SIAM J. Optim., 30 (2020), pp. 349–376.
- 548 [15] V. PATEL, *Kalman-based stochastic gradient method with stop condition and insensitivity to conditioning*,
549 SIAM J. Optim., 26 (2016), pp. 2620–2648.
- 550 [16] L. PFLUG, N. BERNHARDT, M. GRIESHAMMER, AND M. STINGL, *CSG: a new stochastic gradient method*
551 *for the efficient solution of structural optimization problems with infinitely many states*, Struct. Mul-
552 tidiscip. Optim., 61 (2020), pp. 2595–2611.
- 553 [17] H. ROBBINS AND S. MONRO, *A stochastic approximation method*, Ann. Math. Statistics, 22 (1951),
554 pp. 400–407.
- 555 [18] M. SCHMIDT, N. LE ROUX, AND F. BACH, *Minimizing finite sums with the stochastic average gradient*,
556 Math. Program., 162 (2017), pp. 83–112.
- 557 [19] C. TAN, S. MA, Y.-H. DAI, AND Y. QIAN, *Barzilai-borwein step size for stochastic gradient descent*,
558 2016, <https://arxiv.org/abs/1605.04131>.
- 559 [20] V. S. VARADARAJAN, *On the convergence of sample probability distributions*, Sankhyā, 19 (1958), pp. 23–
560 26.
- 561 [21] M. WANG, E. X. FANG, AND H. LIU, *Stochastic compositional gradient descent: algorithms for minimizing*
562 *compositions of expected-value functions*, Math. Program., 161 (2017), pp. 419–449.
- 563 [22] X. WANG AND H. ZHANG, *Inexact proximal stochastic second-order methods for nonconvex composite*
564 *optimization*, Optim. Methods Softw., 35 (2020), pp. 808–835.
- 565 [23] S. YAKOWITZ, J. E. KRIMMEL, AND F. SZIDAROVSKY, *Weighted Monte Carlo integration*, SIAM J.
566 Numer. Anal., 15 (1978), pp. 1289–1300.
- 567 [24] Y.-N. ZHU AND X. ZHANG, *Stochastic primal dual fixed point method for composite optimization*, J. Sci.
568 Comput., 84 (2020), pp. Paper No. 16, 25.