

## ENHANCED DERIVATIVE-FREE OPTIMIZATION USING ADAPTIVE CORRELATION-INDUCED FINITE DIFFERENCE ESTIMATORS

Guo Liang<sup>1</sup>, Guangwu Liu<sup>2</sup>, and Kun Zhang<sup>1</sup>

<sup>1</sup>Institute of Statistics and Big Data, Renmin University of China, Beijing, CHINA

<sup>2</sup>Dept. of Decision Analytics and Operations, College of Business, City University of Hong Kong,  
Kowloon, Hong Kong SAR, CHINA

### ABSTRACT

Gradient-based methods are well-suited for derivative-free optimization (DFO), where finite-difference (FD) estimates are commonly used as gradient surrogates. Traditional stochastic approximation methods, such as Kiefer-Wolfowitz (KW) and simultaneous perturbation stochastic approximation (SPSA), typically utilize only two samples per iteration, resulting in imprecise gradient estimates and necessitating diminishing step sizes for convergence. In this paper, we combine a batch-based FD estimate and an adaptive sampling strategy, developing an algorithm designed to enhance DFO in terms of both gradient estimation efficiency and sample efficiency. Furthermore, we establish the consistency of our proposed algorithm and demonstrate that, despite using a batch of samples per iteration, it achieves the same sample complexity as the KW and SPSA methods. Additionally, we propose a novel stochastic line search technique to adaptively tune the step size in practice. Finally, comprehensive numerical experiments confirm the superior empirical performance of the proposed algorithm.

### 1 INTRODUCTION

Stochastic optimization aims to find the minimization (or maximization) of a function in the presence of noise. Specifically, in this paper, we consider solving the following unconstrained stochastic optimization problem:

$$\min_{\mathbf{x} \in \mathcal{X}} f(\mathbf{x}) = \mathbb{E}[F(\mathbf{x})], \quad (1)$$

where  $\mathcal{X} \subseteq \mathbb{R}^d$  is a convex set,  $f : \mathcal{X} \rightarrow \mathbb{R}$  is the true performance, and  $F$  is the response function. This problem has a wide range of application, including simulation optimization (Chang et al. 2013; Hu and Fu 2025) and reinforcement learning (Fazel et al. 2018). Among problem (1), a difficult but important case lies in the lack of the closed form of  $f(\mathbf{x})$ , and only estimates of the output function are available. That is, for any  $\mathbf{x} \in \mathcal{X}$ , we can only get an unbiased but noisy estimate of  $f(\mathbf{x})$ , i.e.,  $F(\mathbf{x})$ . Such problem is derivative-free optimization (DFO, sometimes referred to as black-box optimization). As the problem becomes complex, the DFO will become increasingly important, and Golovin et al. (2017) mention that “any sufficiently complex system acts as a blackbox when it becomes easier to experiment with than to understand”.

Much literature has discussed the methodology development of DFO. The first category of algorithms are heuristic methods. For instance, the Nelder-Mead simplex algorithm, a direct-search-based method, is widely applied in practical scenarios (Barton and Ivey Jr 1996). A key limitation of these algorithms is the lack of theoretical convergence guarantees (Spall 2005). Another line involves transforming stochastic problems into deterministic ones, taking advantages of deterministic optimization. For examples, sample average approximation (Kim et al. 2015) generates many sample paths and uses the sample mean to estimate the unknown expectation; metamodels such as response surface methodology and Gaussian process are also

used to fit the unknown function (Hong and Zhang 2021). Recently, model-based trust region methods, which integrate model fitting and trust region techniques, have been developed to address simulation optimization problems. Notable examples include the interpolation-based trust-region approach (NEWUOA) (Powell 2006), stochastic trust-region response-surface method (STRONG) (Chang et al. 2013) and adaptive sampling trust-region algorithm (ASTRO) (Shashaani et al. 2018). Interested readers may refer to Audet and Hare (2017) and Larson et al. (2019) for a survey of these methods.

Although our focus is on derivative-free approaches, Audet and Hare (2017) comment that “*if gradient information is available, reliable and obtainable at reasonable cost, then gradient-based methods should be used.*” On the other hand, Shi et al. (2023) and Wang et al. (2025) perform lots of experiments, showing the efficiency of gradient-based methods with finite-difference (FD) gradient surrogates. Therefore, we would like to consider gradient-based stochastic search algorithm to solve (1) in this paper. Note that under some smoothness conditions, (1) can be solved by the following iteration

$$\mathbf{x}_{k+1} = \Pi_{\mathcal{X}} \left( \mathbf{x}_k - a_k \widehat{\nabla} f(\mathbf{x}_k) \right), \quad (2)$$

where  $\mathbf{x}_k$  and  $\mathbf{x}_{k+1}$  are the current solution and next prediction, respectively,  $\widehat{\nabla} f(\mathbf{x}_k)$  is the estimate of the true gradient  $\nabla f(\mathbf{x}_k)$  at  $k$ -th iteration,  $a_k > 0$  is the step size, and  $\Pi_{\mathcal{X}}(\cdot)$  is the projection operator onto the feasible region  $\mathcal{X}$ .

Under the black-box setting, the gradient estimate in (2) is usually substituted by FD gradient and in this case, (2) can date back to the Kiefer-Wolfowitz stochastic approximation (KWSA) algorithm (Kiefer and Wolfowitz 1952). Due to the variance of the FD estimate, the step size should tend to 0 to ensure the convergence. A similar method in high dimension is the simultaneous perturbation stochastic approximation (SPSA) method (Spall 1992; Spall 1997). In practice, the initial step size is crucial and an inappropriate step size will lead to poor performance. Note that when the simulation error is sufficiently small, the performance of KW-type methods is satisfactory (Shi et al. 2023). Therefore, a possible improvement is using more samples to increase the accuracy of the FD estimate. The above idea is inspired by the mini-batch method, which serves as an improvement over stochastic gradient descent or the Robbins-Monro (RM) algorithm (Robbins and Monro 1951). However, the “batch-based” idea is rarely used in DFO and complexity analysis is still open (Shashaani 2024). Most existing methods consider only a fixed perturbation scheme (Bollapragada et al. 2024), which leads to biased gradient estimates. As a result, the complexity analysis typically focuses on convergence to a neighborhood of the optimum, rather than to the optimum itself. In this paper, we fill the gap and strengthen the convergence results.

The first challenge, when using the FD gradient, lies in the construction of an accurate batch-based FD estimate. Fox and Glynn (1989) and Zazanis and Suri (1993) study the convergence of the standard batch-based FD estimator and provide the theoretically optimal perturbation size by minimizing the mean squared error (MSE). However, the optimal perturbation is related to the structure information of the blackbox and we do not know it when using the FD estimator. To overcome this issue, Li and Lam (2020) propose a two-stage method, estimating the perturbation in the first stage and then generating remaining samples at the estimated perturbation to give the expectation-minimization FD (EM-FD) estimator. Recently, Liang et al. (2024) propose a correlation-induced FD (Cor-FD) estimator, using all samples in a batch to estimate the perturbation and then recycling them to estimate the gradient. Cor-FD is available when the batch size is small and it is shown that Cor-FD possesses a reduced variance, and in some cases a reduced bias, compared to the optimal FD estimator.

Given that Cor-FD is efficient when the budget is limited, it is suitable for batch-based DFO algorithm. However, selecting an appropriate batch size in each iteration is crucial. Specifically, if the batch size is too small, the descent direction may lack sufficient accuracy, limiting adjustments to minor corrections along this direction. Conversely, if the batch size is too large, samples may be wasted, as the descent direction does not require excessive precision. In fact, the most efficient algorithms should employ a progressive batching approach in which the batch size is initially small, and increases as the iteration progresses (Bollapragada et al. 2018). For this purpose, the adaptive sampling condition called the *norm condition* (Bollapragada

et al. 2024) has been proposed, which sets the batch size based on the *signal-to-noise ratio* (i.e., the ratio between the true gradient and the gradient estimate error).

In addition to estimate the gradient, another challenge lies in selecting an appropriate step size: too small leads to slow convergence, while too large may cause divergence. Existing stochastic line search methods based on a relaxed Armijo condition (Shi et al. 2023) often yield suboptimal step sizes due to the presence of noise in function evaluations, and can only guarantee convergence to a neighborhood of the optimum (Berahas et al. 2019). To address this, we propose increasing the number of simulations when evaluating candidate step sizes, enabling a more accurate assessment under the relaxed Armijo condition.

The rest of the paper is organized as follows. In Section 2, we give some backgrounds about the gradient-based stochastic optimization. Section 3 presents the adaptive sampling condition, the stochastic line search and the complete algorithms based on constant step sizes. In Section 4, we present the main results about the algorithm in Section 3. Sections 5 applies our algorithm to solve DFO problems, followed by conclusions in Section 6.

## 2 PRELIMINARIES

### 2.1 Gradient-Based Stochastic Search

To find the optimal solution  $\mathbf{x}^*$  of (1), the classical method is gradient-based stochastic search (also known as stochastic approximation). Specifically,  $\mathbf{x}^*$  can be obtained by the recursion (2). Without loss of generality, we assume that  $\mathcal{X} = \mathbb{R}^d$  and remove the projection operator. Then, the recursion is

$$\mathbf{x}_{k+1} = \mathbf{x}_k - a_k \widehat{\nabla} f(\mathbf{x}_k). \quad (3)$$

There are two key elements in (3): the step size  $a_k$  and the gradient estimate  $\widehat{\nabla} f(\mathbf{x}_k)$ . If the unbiased gradient estimate  $\nabla F(\mathbf{x}_k)$  can be obtained, then (3) is the RM algorithm, which is the origin of the stochastic gradient descent. In this paper, we assume that one can only get the zeroth-order information with noise and the first-order information is unavailable. In this case, the gradient estimate  $\widehat{\nabla} f(\mathbf{x}_k)$  is usually substituted by the FD method. Such methods include the KW and SPSA methods, which obtain  $\widehat{\nabla} f(\mathbf{x}_k)$  with only 2 samples at each iteration. Consequently, the variance of  $\widehat{\nabla} f(\mathbf{x}_k)$  is large and the step size  $a_k$  should tend to 0 to ensure the recursion goes towards the optimal solution  $\mathbf{x}^*$ . Although it is shown that the convergence rate of the KW and SPSA algorithms can reach  $\mathcal{O}(1/k^{2/3})$ , where  $k$  is the iteration number, the convergence rate may be unattainable in practice. Even when there is no noise, using the diminishing step size may lead to degeneration, as can be seen in Example 2.1.

**Example 2.1** (Broadie et al. (2011)) Consider finding the infimum point of the deterministic function  $f(x) = 0.001x^2$  with the KW method:

$$x_{k+1} = x_k - a_k \frac{f(x_k + h_k) - f(x_k - h_k)}{2h_k}.$$

If we set  $a_k = \theta_a/k$  and  $h_k = \theta_h/k^{1/4}$  with  $\theta_a = 1$  and  $\theta_h = 1$ , respectively, then the KW algorithm becomes  $x_{k+1} = x_k(1 - 1/(500k))$ . Starting with  $x_1 = 1$ , we have

$$x_{k+1} = \prod_{i=1}^k \left(1 - \frac{1}{500i}\right) = \exp\left(\sum_{i=1}^k \log\left(1 - \frac{1}{500i}\right)\right) \geq \exp\left(-\sum_{i=1}^k \frac{1}{500i}\right) \geq \mathcal{O}\left(\frac{1}{k^{0.002}}\right),$$

where the first inequality is because for any  $x \in (0, 1)$ ,  $\log(1 - x) \geq -x$ . Therefore, the MSE cannot converge faster than  $\mathcal{O}(k^{-0.004})$ , which is significantly slower than the theoretically optimal rate of the KW algorithm  $\mathcal{O}(k^{-2/3})$ .

To address the problem in Example 2.1, one approach is to carefully adjust the initial value of the step size,  $\theta_a$ , to prevent degeneration in the algorithm's convergence rate. Along this line, Broadie et al. (2011) propose a scaled-and-shifted KW (SSKW) method, which adaptively tunes both the step size and perturbation by introducing 9 additional hyperparameters. Another approach is to use a constant (or non-diminishing) step size, which demands a highly accurate gradient estimate, requiring more samples to compute  $\widehat{\nabla}f(\mathbf{x}_k)$ . In this paper, we focus on the second method and use the correlation-induced central FD (Cor-CFD) estimate as the surrogate of  $\widehat{\nabla}f(\mathbf{x}_k)$  in (3). The benefit of this surrogate is that Cor-CFD possesses a variance reduction property and its performance is close to (or even better than) that of the optimal CFD estimate. Consequently, this surrogate is efficient when the sample size is limited and is suitable for the initial stage of the optimization. For the completeness, we briefly outline the construction of the Cor-CFD estimator below. For more details, interested readers may refer to Liang et al. (2024).

## 2.2 Correlation-Induced Central Finite Difference

In this section, we introduce the Cor-CFD estimate. For simplicity, we assume  $d = 1$  because when  $d > 1$ , we can apply Cor-CFD across all the coordinates or randomized directions. Consider the  $k$ -th iteration and denote  $x_k$  by the point of interest and  $n_k$  by the total sample pairs used for gradient estimation (i.e.,  $2n_k$  function evaluations). The conventional CFD estimate is

$$g_{n_k, h_k}(x_k) = \frac{1}{n_k} \sum_{i=1}^{n_k} \frac{F_i(x_k + h_k) - F_i(x_k - h_k)}{2h_k},$$

where  $h_k$  is the perturbation and  $(F_i(x_k + h_k), F_i(x_k - h_k))$  is a sample pair. Under mild conditions, the expectation and variance of  $g_{n_k, h_k}(x_k)$  can be easily calculated:

$$\mathbb{E}[g_{n_k, h_k}(x_k)] = \nabla f(x_k) + B_k h_k^2 + o(h_k^2), \quad \text{Var}[g_{n_k, h_k}(x_k)] = \frac{\sigma_k^2}{2n_k h_k^2} + o\left(\frac{1}{h_k^2}\right), \quad (4)$$

where  $B_k = \nabla^{(3)}f(x_k)/6$  with  $\nabla^{(3)}f(\cdot)$  being the third derivative of  $f(\cdot)$  and  $\sigma_k^2 = \text{Var}[F(x_k)]$ . To minimize the MSE ( $= \text{Bias}^2 + \text{Variance}$ ), the optimal perturbation is  $h_k^* = (\sigma_k^2/(4n_k B_k^2))^{1/6}$ , which relies on the unknown constants  $B_k, \sigma_k$  and is more challenge than the estimation of  $\nabla f(x_k)$ .

To set an appropriate perturbation in practice, the Cor-CFD method generates  $R$  perturbations  $h_{k,1}, \dots, h_{k,R}$  randomly from a pilot distribution  $\mathcal{P}_0$ . Without loss of generality, assume  $n_k = b_k R$ . Then, for each perturbation  $h_{k,r}$  ( $r = 1, \dots, R$ ),  $b_k$  sample pairs  $(F_i(x_k + h_{k,r}), F_i(x_k - h_{k,r}))$  ( $i = 1, \dots, b_k$ ) are generated and all of the CFD estimates are shown in Figure 1.

Figure 1: All of the CFD estimates generated by the Cor-CFD method.

By resampling the CFD estimates  $b_k$  times with replacement at each perturbation  $h_{k,r}$ , we obtain the estimates of  $\mathbb{E}[g_{b_k, h_{k,r}}(x_k)]$  and  $\text{Var}[g_{b_k, h_{k,r}}(x_k)]$ , denoted by  $\mathbb{E}_*[g_{b_k, h_{k,r}}(x_k)]$  and  $\text{Var}_*[g_{b_k, h_{k,r}}(x_k)]$ , respectively. It follows from (4) that the bias and variance of  $g_{b_k, h_k}(x_k)$  are linear with respect to  $h_k^2$  and  $1/(b_k h_k^2)$ .

Therefore, regressing  $[\mathbb{E}_*[g_{b_k, h_{k,1}}(x_k)], \dots, \mathbb{E}_*[g_{b_k, h_{k,R}}(x_k)]]$  on  $[1, \dots, 1]$  and  $[h_{k,1}^2, \dots, h_{k,R}^2]$  gives the estimates  $[\hat{\nabla}' f(x_k), \hat{B}_k]$ . Regressing  $[\text{Var}_*[g_{b_k, h_{k,1}}(x_k)], \dots, \text{Var}_*[g_{b_k, h_{k,R}}(x_k)]]$  on  $[1/(2b_k h_{k,1}^2), \dots, 1/(2b_k h_{k,R}^2)]$  gives the estimate  $\hat{\sigma}_k^2$ . Then, we can set the estimated optimal perturbation  $\hat{h}_k = \left( \hat{\sigma}_k^2 / (4n_k \hat{B}_k^2) \right)^{1/6}$ .

After choosing the perturbation, the Cor-CFD method reuses all of the CFD estimates shown in Figure 1 by adjusting their location and scale. Our goal is that the expectation and variance of the transformed CFD estimates are similar to those of the optimal CFD estimates. Specifically, for  $i = 1, \dots, b_k$  and  $r = 1, \dots, R$ , we transform the CFD estimates in Figure 1 to

$$\frac{h_{k,r}}{\hat{h}_k} \left[ \frac{F_i(x_k + h_{k,r}) - F_i(x_k - h_{k,r})}{2h_{k,r}} - \hat{\nabla}' f(x_k) - \hat{B}_k h_{k,r}^2 \right] + \hat{\nabla}' f(x_k) + \hat{B}_k \hat{h}_k^2. \quad (5)$$

Finally, the Cor-CFD estimate of  $\nabla f(x_k)$  is defined as  $(1/n_k) \sum_{i=1}^{b_k} \sum_{r=1}^R (5)$ .

### 3 PROPOSED ALGORITHM

Here and after, at  $k$ -th iteration, we denote  $n_k$  by the number of sample pairs at each coordinate and  $g_k(\mathbf{x}_k)$  by the corresponding Cor-CFD estimate. Back to (2), although we have selected an appropriate method to surrogate the gradient, there are still two questions should be addressed. The first question is how to set  $n_k$  at  $k$ -th iteration and the second question is how to set the step size  $a_k$ . In Section 3.1, we consider the first question and propose an algorithm with constant step sizes. In Section 3.2, we consider the second question and propose a heuristic line search technique.

#### 3.1 Adaptive Sampling

For the first question, note that we aim to ensure that the estimated gradient aligns with the descending direction, meaning the angle between  $g_k(\mathbf{x}_k)$  and  $\nabla f(\mathbf{x}_k)$  is acute. However, this is not fully achievable due to the inherent uncertainty in gradient estimation. To address this, a straightforward method is increasing  $n_k$  to reduce the uncertainty. Equivalently, we can incorporate the uncertainty by defining a confidence region for  $\nabla f(\mathbf{x}_k)$ , ensuring that all  $d$ -dimensional vectors within this region align with the descending direction. Specifically, let  $\mathcal{F}_k = \sigma\{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_k\}$  be the  $\sigma$ -field generated by  $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_k$ . Denote  $\mathbf{b}_k = \mathbb{E}[g_k(\mathbf{x}_k)|\mathcal{F}_k] - \nabla f(\mathbf{x}_k)$  and  $\boldsymbol{\epsilon}_k = g_k(\mathbf{x}_k) - \mathbb{E}[g_k(\mathbf{x}_k)|\mathcal{F}_k]$ , which are the bias and noise terms, respectively. Then, consider the confidence region

$$[\mathbf{l}, \mathbf{u}] := \left[ g_k(\mathbf{x}_k) - \mathbf{b}_k - \sqrt{\mathbb{E}[\boldsymbol{\epsilon}_k \circ \boldsymbol{\epsilon}_k | \mathcal{F}_k]} / \theta, g_k(\mathbf{x}_k) - \mathbf{b}_k + \sqrt{\mathbb{E}[\boldsymbol{\epsilon}_k \circ \boldsymbol{\epsilon}_k | \mathcal{F}_k]} / \theta \right],$$

where  $\circ$  denotes the element-wise product,  $\sqrt{\mathbb{E}[\boldsymbol{\epsilon}_k \circ \boldsymbol{\epsilon}_k | \mathcal{F}_k]}$  is a  $d$ -dimensional vector with each element denoting the standard deviation of the corresponding element in  $g_k(\mathbf{x}_k)$ , and  $\theta$  is a user-specified hyperparameter. To ensure all  $d$ -dimensional vectors in the confidence region are the descent direction, we increase  $n_k$  until  $\min \{ \mathbb{E}[\mathbf{l} | \mathcal{F}_k]^\top \nabla f(\mathbf{x}_k), \mathbb{E}[\mathbf{u} | \mathcal{F}_k]^\top \nabla f(\mathbf{x}_k) \} \geq 0$ . Equivalently, this condition holds when both  $(\nabla f(\mathbf{x}_k) \pm \sqrt{\mathbb{E}[\boldsymbol{\epsilon}_k \circ \boldsymbol{\epsilon}_k | \mathcal{F}_k]} / \theta)^\top \nabla f(\mathbf{x}_k) \geq 0$ , which can be derived from

$$\mathbb{E}[||\boldsymbol{\epsilon}_k||^2 | \mathcal{F}_k] \leq \theta^2 ||\nabla f(\mathbf{x}_k)||^2 \quad (6)$$

using the Cauchy-Schwarz inequality.

In fact, (6) represents an efficient sampling condition, indicating that the variance must be sufficiently small. This condition controls the *noise-to-signal ratio* in gradient estimation, thereby improving its reliability. (6) is called the *norm condition* and has been considered by Bollapragada et al. (2018) and Bollapragada et al. (2024). To apply this condition, we need to identify surrogates for  $\mathbb{E}[||\boldsymbol{\epsilon}_k||^2 | \mathcal{F}_k]$  and  $||\nabla f(\mathbf{x}_k)||$ . Note that  $\mathbb{E}[||\boldsymbol{\epsilon}_k||^2 | \mathcal{F}_k]$  represents the sum of variances across all coordinates. We can employ

sample variance to estimate the variance of each component of  $g_k(\mathbf{x}_k)$  and subsequently  $\mathbb{E}[\|\epsilon_k\|^2 | \mathcal{F}_k]$ . Despite the existence of the correlation, it is efficient to use the sample variance because it has been shown that the correlation in Cor-CFD tends to reduce the variance (Liang et al. 2024). For  $\|\nabla f(\mathbf{x}_k)\|$ ,  $\|g_k(\mathbf{x}_k)\|$  can be chosen as an appropriate surrogate. Specifically, the estimated version of (6) is

$$\frac{\sum_{i=1}^d \hat{\sigma}_i^2}{n_k} \leq \theta^2 \|g_k(\mathbf{x}_k)\|^2, \quad (7)$$

where  $\hat{\sigma}_i^2$  is the sample variance of (5) at  $i$ -th coordinate in the  $k$ -th iteration, which is an estimated upper bound of the true variance. Then the algorithm with constant step size is proposed (see Algorithm 1).

### 3.2 Stochastic Line Search

In practice, it is difficult to set an appropriate step size because the optimization problem is a blackbox. To address this issue, a stochastic line search method has been proposed to adjust the step size (Berahas et al. 2019; Shi et al. 2023), inspired by the backtracking line search in deterministic optimization.

Specifically, the classical stochastic line search in Berahas et al. (2019) begins with an initial step size  $a_k = \tilde{a}$  at  $k$ -th iteration and determine whether

$$F(\mathbf{x}_k - a_k g_k(\mathbf{x}_k)) > F(\mathbf{x}_k) - l_1 a_k \|g_k(\mathbf{x}_k)\|^2 + 2\sigma_F \quad (8)$$

holds, where  $0 < l_1 < 1$  is a parameter. If (8) holds, then the function value at next predicted step is significantly larger than that at the current step and the step size should not be chosen. In this case, we shrink  $a_k \rightarrow l_2 a_k$ , where  $l_1 < l_2 < 1$ . The procedure will stop until (8) does not hold. Note that the simulation noise  $\sigma_F$  can be substituted by the upper bound. For example, if it is very large, then (8) will never hold and the stochastic line search outputs the constant step size  $\tilde{a}$ .

After the stochastic line search (8), we get a step size which is not too “bad” (a “bad” step size means that the next predicted value is much larger than the current value), but is not guaranteed to be “good” (a “good” step size means that the next predicted value is not larger than the current value). In other words,

---

#### Algorithm 1: Cor-CFD-based DFO Algorithm with Constant Step Size

---

**Input:** Total number of function evaluations  $\mathcal{S}$ , initial sample pairs  $n_0$ , starting point  $\mathbf{x}_0$ , adaptive sampling threshold  $\theta$  and step length  $a > 0$ .

**Output:** The ultimate estimate  $\mathbf{x}_k$ .

**Initialization:** Set  $k \leftarrow 0$ ,  $s \leftarrow 0$ .

**while**  $s < \mathcal{S}$  **do**

**foreach** coordinate  $i = 1, \dots, d$  **do**

        Compute gradient estimate  $g_{k,i}(\mathbf{x}_k)$  using Cor-CFD with  $n_k$  sample pairs, where  $g_{k,i}(\mathbf{x}_k)$  denotes the  $i$ -th component of  $g_k(\mathbf{x}_k)$ .

        Compute the sample variance  $\hat{\sigma}_i^2$ .

**if** (7) does not hold **then**

        Increase  $n_k$  to  $\lfloor \sum_{i=1}^d \hat{\sigma}_i^2 / (\theta^2 \|g_k(\mathbf{x}_k)\|^2) \rfloor + 1$ , where  $\lfloor \cdot \rfloor$  represents the largest integer that does not exceed  $\cdot$ .

        Add  $\lfloor \sum_{i=1}^d \hat{\sigma}_i^2 / (\theta^2 \|g_k(\mathbf{x}_k)\|^2) \rfloor + 1 - n_k$  sample pairs to update gradient estimate  $g_k(\mathbf{x}_k)$ .

        Set  $n_k = \lfloor \sum_{i=1}^d \hat{\sigma}_i^2 / (\theta^2 \|g_k(\mathbf{x}_k)\|^2) \rfloor + 1$ .

    Update  $\mathbf{x}_{k+1} \leftarrow \mathbf{x}_k - a g_k(\mathbf{x}_k)$ .

    Update  $s \leftarrow s + 2dn_k$  and  $k \leftarrow k + 1$ .

**return**  $\mathbf{x}_k$

---

(8) is conservative and only used to exclude the “bad” cases. Similar to (8), we present a new criterion:

$$F(\mathbf{x}_k - a_k g_k(\mathbf{x}_k)) \leq F(\mathbf{x}_k) - l_1 a_k \|g_k(\mathbf{x}_k)\|^2 - 2\sigma_F. \quad (9)$$

The intuition of (9) is that it indicates that the function value at the next predicted step is significantly smaller than that at the current step. Note that condition (9) may not always hold, particularly when the current solution is near the optimal value. However, increasing simulation replications to reduce the black-box function’s uncertainty will eventually satisfy condition (9). Specifically, there exists a number  $N$  such that

$$\frac{1}{N} \sum_{i=1}^N F_i(\mathbf{x}_k - a_k g_k(\mathbf{x}_k)) \leq \frac{1}{N} \sum_{i=1}^N F_i(\mathbf{x}_k) - l_1 a_k \|g_k(\mathbf{x}_k)\|^2 - 2 \frac{\sigma_F}{\sqrt{N}} \quad (10)$$

because when  $N \rightarrow \infty$ , it is identical to  $f(\mathbf{x}_k - a_k g_k(\mathbf{x}_k)) \leq f(\mathbf{x}_k) - l_1 a_k \|g_k(\mathbf{x}_k)\|^2$  which closes to the standard line search criterion and holds under mild conditions. Conditions (9) and (10) require using many samples to assess the appropriateness of the step size. This effort is justified because, near the optimal solution, maintaining algorithmic progress requires highly accurate step-size selection. Note that when  $N \rightarrow \infty$  and  $a_k \rightarrow 0$ , (10) always tends to hold, and if a step size is small enough, it is also a “good” step size (see Theorem 1 for theoretical evidences). Thus, in practice, it is enough to evaluate condition (9) for a finite number of times (i.e.,  $N \leq N_0$ , where  $N_0$  is a user-specified parameter). If it remains unsatisfied, we shrink the step size. Additionally, we enforce a lower bound on the step size to avoid it becoming too small.

**Remark 3.1** Stochastic line search is similar to the “shifted” procedure in Broadie et al. (2011), as both methods reduce a “too large” step size until it becomes appropriate. However, these two methods are fundamentally different. The “shifted” procedure in Broadie et al. (2011) is based on KWSA, with the step size sequence defined as  $\theta_a/k^\gamma$ , where  $\gamma > 0$  is a hyperparameter. In contrast, stochastic line search is based on the standard line search used in deterministic optimization. In stochastic line search, the step size does not necessarily decrease at a fixed rate and it may decrease, not change or even increase. The only requirement is that the predicted value of the next iteration is smaller than the current function value to some extent. As a result, stochastic line search is more flexible.

## 4 THEORETICAL RESULTS

In this section, we present the convergence results of Algorithm 1. Firstly, we state some assumptions, which suppose that the objective function  $F(\mathbf{x})$  is smooth and strongly convex, satisfying the following regularity conditions.

**Assumption 4.1** (Differentiability) The function  $f(\mathbf{x})$  is fifth continuously differentiable on  $\mathcal{X}$  and  $\nabla^5 f(\mathbf{x}) \neq 0$  for any  $\mathbf{x} \in \mathcal{X}$ .

**Assumption 4.2** (Lipschitz smoothness) The gradient of  $f(\mathbf{x})$  is  $M$ -Lipschitz, i.e., there exists a constant  $M > 0$  such that  $\|\nabla f(\mathbf{x}_1) - \nabla f(\mathbf{x}_2)\| \leq M\|\mathbf{x}_1 - \mathbf{x}_2\|$  for all  $\mathbf{x}_1, \mathbf{x}_2 \in \mathcal{X}$ , where  $\|\cdot\|$  denotes the Euclidean norm.

**Assumption 4.3** (Strongly convex) There exists a constant  $m > 0$  such that  $\lambda(\mathbf{x}) \geq m$  for all  $\mathbf{x} \in \mathcal{X}$ , where  $\lambda(\mathbf{x})$  denotes the smallest eigenvalue of the Hessian matrix  $H(\mathbf{x}) := \nabla^2 f(\mathbf{x})$ .

Assumption 4.1 is used for the theoretical completeness of the Cor-CFD method. Note that when  $\mathcal{X}$  is a compact set, Assumption 4.2 holds automatically if Assumption 4.1 holds. Assumption 4.3 ensures that (1) has a unique solution  $\mathbf{x}^*$ . Assumptions 4.2 and 4.3 are standard conditions when studying the convergence results of the gradient-based method (Scheinberg 2022; Hu and Fu 2025). Specifically, these two assumptions mean that  $M$  and  $m$  are the upper and lower bounds for all eigenvalues of the Hessian matrix.

In the following, we present Theorem 1 to show the convergence result of Algorithm 1. Note that in Theorem 1, we use the condition  $\max\{\mathbb{E}[\|\mathbf{b}_k\|^2|\mathcal{F}_k], \mathbb{E}[\|\epsilon_k\|^2|\mathcal{F}_k]\} \leq \theta^2 \|\nabla f(\mathbf{x}_k)\|^2$  which is different from (6). This condition is introduced solely for the convenience of the proof. It is reasonable because if we use Cor-CFD gradient estimate,  $\mathbb{E}[\|\mathbf{b}_k\|^2|\mathcal{F}_k]$  and  $\mathbb{E}[\|\epsilon_k\|^2|\mathcal{F}_k]$  have the same order. While the conventional CFD estimate can also be used without affecting the convergence result, its performance depends on the manual choice of perturbation, which may reduce the algorithm's stability and adaptability in black-box settings. Conversely, the Cor-CFD method automatically yields a high-quality gradient estimate with a variance reduction effect.

**Theorem 1** (Iteration complexity) Suppose that Assumptions 4.1, 4.2 and 4.3 hold. Let  $\mathbf{x}_0$  be the initial point and at  $k$ -th iteration, let  $\max\{\mathbb{E}[\|\mathbf{b}_k\|^2|\mathcal{F}_k], \mathbb{E}[\|\epsilon_k\|^2|\mathcal{F}_k]\} \leq \theta^2 \|\nabla f(\mathbf{x}_k)\|^2$ , where  $0 < \theta < m/(2M)$  is a threshold. If  $0 < a_k = a \leq 1/((2\theta^2 + 2\theta + 1)M)$  for any  $k \geq 0$ , then we have

$$\mathbb{E}[\|\mathbf{x}_k - \mathbf{x}^*\|^2] \leq (1 - (m - 2\theta M)a)^k \|\mathbf{x}_0 - \mathbf{x}^*\|^2. \quad (11)$$

The proof of Theorem 1 can be found in Liang et al. (2025). Observe that  $0 < (m - 2\theta M)a < 1$  when  $0 < \theta < m/(2M)$  and  $a \leq 1/((2\theta^2 + 2\theta + 1)M)$ . Under these conditions,  $\mathbf{x}_k$  converges linearly in expectation to the minimum point  $\mathbf{x}^*$ . Theorem 1 generalizes the convergence result of standard gradient descent. Notably, as  $\theta \rightarrow 0$ , the gradient estimate in each iteration approaches the true gradient at the current solution. Consequently, the right hand side (RHS) on (11) converges to  $(1 - ma)^k \|\mathbf{x}_0 - \mathbf{x}^*\|^2$ , where  $0 < a < 1/M$ . The optimal case occurs as  $a \rightarrow 1/M$ , with  $\mathbb{E}[\|\mathbf{x}_k - \mathbf{x}^*\|^2]$  converging to 0 at a rate comparable to a geometric series, featuring an exponent approaching  $m/M$ . This result aligns perfectly with deterministic gradient descent.

It is important to note that focusing solely on the iteration complexity is insufficient. For instance, when  $\theta \rightarrow 0$ , the gradient descent result is recovered, but this requires the batch size to be infinity at each iteration to satisfy the adaptive sampling condition. Therefore, it is crucial to consider both the iteration complexity and the associated sample complexity, i.e., the total stochastic function evaluations required to get an  $\epsilon$ -accurate solution. We employ the metric that  $\mathbf{x}_k$  is said to be an  $\epsilon$ -accurate solution if  $\mathbb{E}[\|\mathbf{x}_k - \mathbf{x}^*\|^2] \leq \epsilon$ . Note that the number of stochastic function evaluations at any iteration  $k$  is  $S_k = 2dn_k$ , where  $n_k$  is the number of sample pairs at each coordinate. In the following, we present Theorem 2 to analyze the sample complexity of Algorithm 1.

**Theorem 2** (Sample complexity) Under the same conditions as those in Theorem 1. Let  $d = \mathcal{O}(1)$ . Denote  $\mathcal{S}(\epsilon)$  by the total stochastic function evaluations to get an  $\epsilon$ -accurate solution. If the third derivative  $\nabla^3 f(\mathbf{x})$  and the function noise  $\sigma(\mathbf{x})$  are bounded below away from 0, then we have  $\mathbb{E}[\mathcal{S}(\epsilon)] \geq \mathcal{C}_1 \epsilon^{-3/2} + \mathcal{C}_2$ , where  $\mathcal{C}_1$  and  $\mathcal{C}_2$  are constants that depends on the threshold  $\theta$ , step size  $a$ , problem dimension  $d$ , unknown function and the simulation error. Specifically,

$$\mathcal{C}_1 = \frac{4d (\mathcal{C}_\theta M^2)^{-3/2}}{3\bar{a}^{-3/2} \log(1/\bar{a})}, \quad \mathcal{C}_2 = \frac{4d (\mathcal{C}_\theta M^2 \|\mathbf{x}_0 - \mathbf{x}^*\|^2)^{-3/2}}{3 \log(1/\bar{a})},$$

where  $\bar{a} = 1 - (m - 2\theta M)a$ ,  $\mathcal{C}_\theta = \theta^2/\mathcal{C}$  and  $\mathcal{C}$  depends on the function  $f(\cdot)$  and simulation noise.

*Proof.* Let  $\boldsymbol{\eta}_k := \mathbf{x}_k - \mathbf{x}^*$ . Denote  $\mathcal{K}(\epsilon)$  by the number of iteration to get an  $\epsilon$ -accurate solution. From Theorem 1, we should let  $(1 - (m - 2\theta M)a)^{\mathcal{K}(\epsilon)} \|\boldsymbol{\eta}_0\|^2 \leq \epsilon$ , which is equivalent to

$$\mathcal{K}(\epsilon) \geq \frac{\log \epsilon - \log \|\boldsymbol{\eta}_0\|^2}{\log \bar{a}}, \quad (12)$$

where we denote  $1 - (m - 2\theta M)a$  by  $\bar{a}$  for the sake of convenience. Then,  $\mathcal{S}(\epsilon) = \sum_{k=0}^{\mathcal{K}(\epsilon)-1} S_k = 2d \sum_{k=0}^{\mathcal{K}(\epsilon)-1} n_k$ . According to the properties of Cor-CFD,  $\max\{\mathbb{E}[\|\mathbf{b}_k\|^2|\mathcal{F}_k], \mathbb{E}[\|\epsilon_k\|^2|\mathcal{F}_k]\} \leq \mathcal{C}_k n_k^{-2/3}$ ,



where  $\mathcal{C}_k$  is a constant depending on the third derivative  $\nabla^3 f(\mathbf{x}_k)$ , the function noise  $\sigma(\mathbf{x}_k)$  and the dimension  $d$ . To satisfy the adaptive sampling condition, we should let

$$\mathcal{C}n_k^{-2/3} \leq \mathcal{C}_k n_k^{-2/3} \leq \theta^2 \|\nabla f(\mathbf{x}_k)\|^2, \quad (13)$$

where the first inequality is because that the third derivative  $\nabla^3 f(\mathbf{x})$  and the function noise  $\sigma(\mathbf{x})$  are bounded below away from 0 (see Liang et al. (2024) for more details). Taking expectation on both sides of (13) gives  $\mathbb{E}[n_k^{-2/3}] \leq \mathcal{C}_\theta \mathbb{E}[\|\nabla f(\mathbf{x}_k)\|^2]$ , where  $\mathcal{C}_\theta = \theta^2/\mathcal{C}$ . It follows Taylor expansion that

$$\nabla f(\mathbf{x}_k) = \nabla f(\mathbf{x}^*) + H(\mathbf{x}_k^\dagger)\boldsymbol{\eta}_k = H(\mathbf{x}_k^\dagger)\boldsymbol{\eta}_k,$$

where  $\mathbf{x}_k^\dagger$  lies on the line segment between  $\mathbf{x}_k$  and  $\mathbf{x}^*$ , and the second equality is because  $\nabla f(\mathbf{x}^*) = 0$ . Then we have

$$\mathbb{E}[\|\nabla f(\mathbf{x}_k)\|^2] = \mathbb{E}[\boldsymbol{\eta}_k^\top H^\top(\mathbf{x}_k^\dagger)H(\mathbf{x}_k^\dagger)\boldsymbol{\eta}_k] \leq M^2 \mathbb{E}[\|\boldsymbol{\eta}_k\|^2] \leq M^2 \bar{a}^k \|\boldsymbol{\eta}_0\|^2.$$

Therefore,  $\mathbb{E}[n_k^{-2/3}] \leq \mathcal{C}_\theta \mathbb{E}[\|\nabla f(\mathbf{x}_k)\|^2] \leq \mathcal{C}_\theta M^2 \bar{a}^k \|\boldsymbol{\eta}_0\|^2$ . By Jensen's inequality, we have  $\mathbb{E}[n_k]^{-2/3} \leq \mathbb{E}[n_k^{-2/3}]$ , which gives that  $\mathbb{E}[n_k] \geq (\mathcal{C}_\theta M^2 \bar{a}^k \|\boldsymbol{\eta}_0\|^2)^{-3/2}$ . Because  $\mathcal{S}(\epsilon) = 2d \sum_{k=0}^{\mathcal{K}(\epsilon)-1} n_k$ , we have

$$\begin{aligned} \mathbb{E}[\mathcal{S}(\epsilon)] &\geq 2d (\mathcal{C}_\theta M^2 \|\boldsymbol{\eta}_0\|^2)^{-3/2} \sum_{k=0}^{\mathcal{K}(\epsilon)-1} \bar{a}^{-3k/2} \geq 2d (\mathcal{C}_\theta M^2 \|\boldsymbol{\eta}_0\|^2)^{-3/2} \int_0^{\mathcal{K}(\epsilon)-1} \bar{a}^{-3u/2} du \\ &\geq 2d (\mathcal{C}_\theta M^2 \|\boldsymbol{\eta}_0\|^2)^{-3/2} \int_0^{\frac{\log \epsilon - \log \|\boldsymbol{\eta}_0\|^2}{\log \bar{a}} - 1} \bar{a}^{-3u/2} du = \mathcal{C}_1 \epsilon^{-3/2} + \mathcal{C}_2, \end{aligned}$$

where the second inequality is due to the relationship between the summation and integral, and

$$\mathcal{C}_1 = \frac{4d (\mathcal{C}_\theta M^2)^{-3/2}}{3\bar{a}^{-3/2} \log(1/\bar{a})}, \quad \mathcal{C}_2 = \frac{4d (\mathcal{C}_\theta M^2 \|\boldsymbol{\eta}_0\|^2)^{-3/2}}{3 \log(1/\bar{a})}.$$

□

According to Theorem 2, achieving an  $\epsilon$ -accurate solution requires at least  $\mathcal{O}(\epsilon^{-3/2})$  function evaluations. In other words, for a given total budget  $\mathcal{S}$ , the MSE of our algorithm is  $\mathcal{O}(\mathcal{S}^{-2/3})$ , matching the optimal performance of the KW algorithm (Hu and Fu 2025). This result stems from combining the Cor-CFD estimate with the adaptive sampling condition, which ensures that an appropriate number of samples are generated at each iteration, maximizing the use of sample information. For a fixed total budget, reliable gradient estimation allows the algorithm to maintain consistently sufficient descent, even with a reduced number of iterations. Furthermore, because the step size does not need to approach 0, the algorithm circumvents the degeneration scenario illustrated in Broadie et al. (2011).

## 5 NUMERICAL EXPERIMENTS

In this section, we test the performance of our algorithm on two types of problems which are both constructed from deterministic optimization with added noise. The first type of problem is a simple power function  $f(x) = 0.001x^2$  with a noise  $\mathcal{N}(0, \sigma^2(x))$ , where  $\sigma(x) = 0.001$ . This problem is used to illustrate the convergence rate of our method. The second type of problem is the Rosenbrock function defined as

$$f(x_1, x_2) = 100(x_2 - x_1^2)^2 + (x_1 - 1)^2.$$

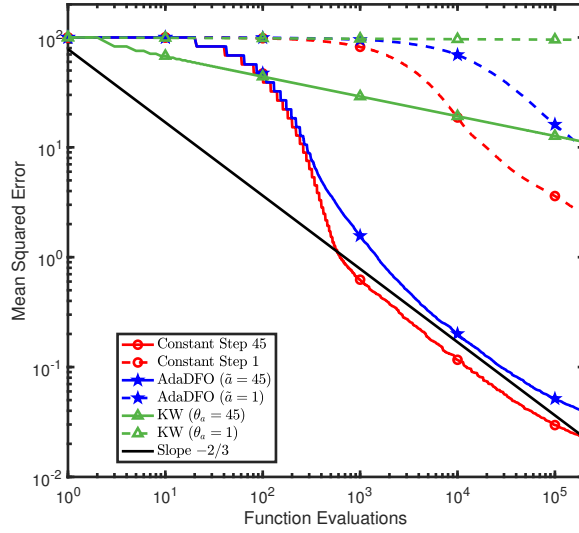


Figure 2: Comparison of the convergence rate of our algorithm and KWSA for the first type of problem.

The unique minimum  $(1, 1)$  lies in a narrow, parabolic valley  $(x_1, x_1^2)$ . Although the valley is easy to find, it is difficult to find the minimum even when there is no noise. In this problem, we set  $F(x_1, x_2) = f(x_1, x_2) + \mathcal{N}(0, 1)$ . Note that these two problems do not satisfy Assumption 4.1 because the fifth order derivative is equal to 0, but they serve to “stress test” our algorithm.

Our Algorithm is marked as “AdaDFO” across all the experiments. For the first type of problem, we set  $x_0 = 10$ . The results are shown in Figure 2, where “Constant Step 45” and “Constant Step 1” represent Algorithm 1 with step sizes 45 and 1, respectively, and the KWSA algorithm is applied without considering the noise. Figure 2 shows the number of function evaluations on the horizontal axis and the MSE of the current solution on the vertical axis. A line with a slope of  $-2/3$  is included for visualization. As shown in the figure, when adaptive sampling is coupled with a constant step size, the observed convergence rate is consistent with Theorem 2. Even in the absence of noise, however, the KWSA algorithm fails to achieve its theoretical optimum convergence rate. Furthermore, we find that incorporating a stochastic line search prevents degradation in the convergence rate of the AdaDFO algorithm, which recovers the  $\mathcal{O}(S^{-2/3})$ , where  $S$  denotes the number of function evaluations. Although AdaDFO performs slightly worse than the constant step-size baseline (due to smaller step sizes returned by stochastic line search), its performance closely matches that of constant-step versions, supporting the effectiveness of stochastic line search.

Table 1 reports the optimality gaps (OGs) of the second type of problem, defined as  $f(\mathbf{x}_k^*) - f(\mathbf{x}^*)$ , where  $\mathbf{x}^*$  is the optimal solution and  $\mathbf{x}_k^*$  is the final solution before termination. SD denotes the standard deviation, and SR denotes the success rate, defined as the proportion of runs where  $f(\mathbf{x}_k^*) < f(\mathbf{x}_0)$  over 1000 replications. Table 1 demonstrates the effectiveness of AdaDFO under noise. It consistently achieves 100% success rate across all budgets, with its OG steadily decreasing from 2.85 to 0.32, highlighting both accuracy and robustness. Some key points from the table to note:

- NoAdaDFO denotes the method without adaptive sampling and use  $20k$  samples at  $k$ -th iteration. As shown in the table, its performance is inferior to that of AdaDFO.
- NM stagnates early, with average OG being around 6.59, indicating difficulty in handling noisy Rosenbrock function.
- STRONG and ASTRO show improvement with increased evaluations. However, their convergence remains slow. Even at the highest budget, their OGs are much larger than that of AdaDFO.
- NEWUOA outperforms others at  $2 \times 10^3$  evaluations but converges slowly. When evaluations are  $2 \times 10^4$  and  $2 \times 10^5$ , its performance is inferior to that of AdaDFO.

Table 1: Optimality gap of different methods for the second type of problem.

Function Evaluations	$2 \times 10^3$			$2 \times 10^4$			$2 \times 10^5$		
Method	Mean	SD	SR	Mean	SD	SR	Mean	SD	SR
AdaDFO	<b>2.85</b>	<b>2.69</b>	<b>100%</b>	<b>0.51</b>	<b>0.99</b>	<b>100%</b>	<b>0.32</b>	<b>0.34</b>	<b>100%</b>
NoAdaDFO	4.91	2.17	100%	1.39	2.06	100%	0.34	0.71	100%
NM	6.59	0.25	100%	6.58	0.22	100%	6.59	0.22	100%
STRONG	16.74	0.02	100%	4.99	0.55	100%	4.91	0.71	100%
ASTRO	6.48	0.22	100%	6.42	0.21	100%	6.25	0.82	100%
NEUWOA	1.98	0.70	100%	0.67	0.57	100%	0.61	0.52	100%
SPSA	3.55	–	99.2%	1.08	–	98.4%	0.13	–	98.7%

- The mean OG of SPSA is calculated by averaging OGs of successful trails. Although SPSA can attain a small OG, its gradient estimates are unstable due to using only two evaluations per iteration, occasionally leading to algorithm failure (success rate  $< 100\%$ ).

## 6 CONCLUSIONS

In this paper, we propose a batch-based DFO algorithm by combining the Cor-FD gradient estimate with an adaptive sampling condition. This combination allows us to obtain an appropriate gradient surrogate for KW-type stochastic approximation method. We prove that, under mild conditions, the use of a properly chosen constant step size ensures convergence. Additionally, we derive the sample complexity of our method, which demonstrates that its convergence rate does not deteriorate compared to the KWSA method. In the black-box scenario, we introduce a new stochastic line search technique to adaptively tune the step size. Numerical experiments confirm the effectiveness of our proposed algorithm, showing that it is suitable for solving DFO problems.

## ACKNOWLEDGMENTS

The research of the third author was supported by National Natural Science Foundation of China (NNSFC) grants 72471232. The research of the second author is supported partially by the National Natural Science Foundation of China and the Research Grants Council of Hong Kong (RGC-HK), under the RGC-HK General Research Fund Projects 11508620 and 11508523 and Theme-based Research Project T32-615/24-R, and NSFC/RGC-HK Joint Research Scheme under project N\_CityU 105/21.

## REFERENCES

- Audet, C., and W. Hare. 2017. *Derivative-Free and Blackbox Optimization*. Cham: Springer.
- Barton, R. R., and J. S. Ivey Jr. 1996. “Nelder-Mead Simplex Modifications for Simulation Optimization”. *Management Science* 42(7):954–973.
- Berahas, A. S., R. H. Byrd, and J. Nocedal. 2019. “Derivative-Free Optimization of Noisy Functions via Quasi-Newton Methods”. *SIAM Journal on Optimization* 29(2):965–993.
- Bollapragada, R., C. Karamanli, and S. M. Wild. 2024. “Derivative-Free Optimization via Adaptive Sampling Strategies”. *arXiv preprint arXiv:2404.11893*.
- Bollapragada, R., J. Nocedal, D. Mudigere, H.-J. Shi, and P. T. P. Tang. 2018. “A Progressive Batching L-BFGS Method for Machine Learning”. In *International Conference on Machine Learning*. July 10<sup>th</sup>-15<sup>th</sup>, Stockholm, Sweden, 620–629.
- Broadie, M., D. Cicek, and A. Zeevi. 2011. “General Bounds and Finite-Time Improvement for the Kiefer-Wolfowitz Stochastic Approximation Algorithm”. *Operations Research* 59(5):1211–1224.
- Chang, K.-H., L. J. Hong, and H. Wan. 2013. “Stochastic Trust-Region Response-Surface Method (STRONG)—A New Response-Surface Framework for Simulation Optimization”. *INFORMS Journal on Computing* 25(2):230–243.
- Fazel, M., R. Ge, S. Kakade, and M. Mesbahi. 2018. “Global Convergence of Policy Gradient Methods for the Linear Quadratic Regulator”. In *International Conference on Machine Learning*. July 10<sup>th</sup>-15<sup>th</sup>, Stockholm, Sweden, 1467–1476.

- Fox, B. L., and P. W. Glynn. 1989. "Replication Schemes for Limiting Expectations". *Probability in the Engineering and Informational Sciences* 3(3):299–318.
- Golovin, D., B. Solnik, S. Moitra, G. Kochanski, J. Karro, and D. Sculley. 2017. "Google Vizier: A Service for Black-Box Optimization". In *Proceedings of the 23rd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*. August 13<sup>th</sup>-17<sup>th</sup>, New York, NY, United States, 1487-1495.
- Hong, L. J., and X. Zhang. 2021. "Surrogate-Based Simulation Optimization". *INFORMS TutORials in Operations Research*:287–311.
- Hu, J., and M. C. Fu. 2025. "On the Convergence Rate of Stochastic Approximation for Gradient-Based Stochastic Optimization". *Operations research* 73(2):1143–1150.
- Kiefer, J., and J. Wolfowitz. 1952. "Stochastic Estimation of the Maximum of A Regression Function". *The Annals of Mathematical Statistics* 23(3):462–466.
- Kim, S., R. Pasupathy, and S. G. Henderson. 2015. "A Guide to Sample Average Approximation". In *Handbook of Simulation Optimization*, 207–243. New York, NY: Springer.
- Larson, J., M. Menickelly, and S. M. Wild. 2019. "Derivative-Free Optimization Methods". *Acta Numerica* 28:287–404.
- Li, H., and H. Lam. 2020. "Optimally Tuning Finite-Difference Estimators". In *2020 Winter Simulation Conference (WSC)*, 457–468. IEEE.
- Liang, G., G. Liu, and K. Zhang. 2024. "A Correlation-Induced Finite Difference Estimator". *arXiv preprint arXiv:2405.05638*.
- Liang, G., G. Liu, and K. Zhang. 2025. "Enhanced Derivative-Free Optimization Using Adaptive Correlation-Induced Finite Difference Estimators". *arXiv preprint arXiv:2502.20819*.
- Powell, M. J. D. 2006. "The NEWUOA Software for Unconstrained Optimization without Derivatives". In *Large-Scale Nonlinear Optimization*, 255–297. Boston, MA: Springer US.
- Robbins, H., and S. Monro. 1951. "A Stochastic Approximation Method". *The Annals of Mathematical Statistics* 22(3):400–407.
- Scheinberg, K. 2022. "Finite Difference Gradient Approximation: To Randomize or Not?". *INFORMS Journal on Computing* 34(5):2384–2388.
- Shashaani, S. 2024. "Simulation Optimization: An Introductory Tutorial on Methodology". In *2024 Winter Simulation Conference (WSC)*, 1308–1322. IEEE.
- Shashaani, S., F. S. Hashemi, and P. Raghu. 2018. "ASTRO-DF: A Class of Adaptive Sampling Trust-Region Algorithms for Derivative-Free Stochastic Optimization". *SIAM Journal on Optimization* 28(4):3145–3176.
- Shi, H.-J. M., Q. M. Xuan, F. Oztoprak, and J. Nocedal. 2023. "On the Numerical Performance of Finite-Difference-Based Methods for Derivative-Free Optimization". *Optimization Methods and Software* 38(2):289–311.
- Spall, J. C. 1992. "Multivariate Stochastic Approximation Using a Simultaneous Perturbation Gradient Approximation". *IEEE Transactions on Automatic Control* 37(3):332–341.
- Spall, J. C. 1997. "A One-Measurement Form of Simultaneous Perturbation Stochastic Approximation". *Automatica* 33(1):109–112.
- Spall, J. C. 2005. *Introduction to Stochastic Search and Optimization: Estimation, Simulation, and Control*. Hoboken: John Wiley & Sons.
- Wang, D.-Y., G. Liang, G. Liu, and K. Zhang. 2025. "Derivative-Free Optimization via Finite Difference Approximation: An Experimental Study". *Asia-Pacific Journal of Operational Research*:2540005.
- Zazanis, M. A., and R. Suri. 1993. "Convergence Rates of Finite-Difference Sensitivity Estimates for Stochastic Systems". *Operations Research* 41(4):694–703.

## AUTHOR BIOGRAPHIES

**GUO LIANG** is a Ph.D. Candidate in the Institute of Statistics and Big Data at Renmin University of China. His research interests include stochastic simulation, stochastic gradient estimation, financial engineering and risk management. His email address is [lianguo000221@ruc.edu.cn](mailto:lianguo000221@ruc.edu.cn).

**GUANGWU LIU** is a Professor in the Department of Decision Analytics and Operations, College of Business at City University of Hong Kong. His research interests include stochastic simulation, business analytics, financial engineering, and risk management. He serves as an Associate Editor of the Asia-Pacific Journal of Operational Research, and Naval Research Logistic. His e-mail address is [msgw.liu@cityu.edu.hk](mailto:msgw.liu@cityu.edu.hk) and his website is <https://www.cb.cityu.edu.hk/staff/guanliu>.

**KUN ZHANG** is an Assistant Professor in the Institute of Statistics and Big Data at Renmin University of China. He holds a PhD in management science from City University of Hong Kong. His research interests include simulation optimization, machine learning, business analytics, financial engineering and risk management. His email address is [kunzhang@ruc.edu.cn](mailto:kunzhang@ruc.edu.cn).