The Parallel Knowledge Gradient Method for Batch Bayesian Optimization

Jian Wu, Peter I. Frazier

Cornell University Ithaca, NY, 14853 {jw926, pf98}@cornell.edu

Abstract

In many applications of black-box optimization, one can evaluate multiple points simultaneously, e.g. when evaluating the performances of several different neural networks in a parallel computing environment. In this paper, we develop a novel batch Bayesian optimization algorithm — the parallel knowledge gradient method. By construction, this method provides the one-step Bayes-optimal batch of points to sample. We provide an efficient strategy for computing this Bayes-optimal batch of points, and we demonstrate that the parallel knowledge gradient method finds global optima significantly faster than previous batch Bayesian optimization algorithms on both synthetic test functions and when tuning hyperparameters of practical machine learning algorithms, especially when function evaluations are noisy.

1 Introduction

In Bayesian optimization [20] (BO), we wish to optimize a derivative-free expensive-to-evaluate function f with feasible domain $\mathbb{A} \subseteq \mathbb{R}^d$,

$$\min_{\boldsymbol{x}\in\mathbb{A}}f(\boldsymbol{x}),$$

with as few function evaluations as possible. In this paper, we assume that membership in the domain \mathbb{A} is easy to evaluate and we can evaluate f only at points in \mathbb{A} . We assume that evaluations of f are either noise-free, or have additive independent normally distributed noise. We consider the parallel setting, in which we perform more than one simultaneous evaluation of f.

BO typically puts a Gaussian process prior distribution on the function f, updating this prior distribution with each new observation of f, and choosing the next point or points to evaluate by maximizing an acquisition function that quantifies the benefit of evaluating the objective as a function of where it is evaluated. In comparison with other global optimization algorithms, BO often finds "near optimal" function values with fewer evaluations [20]. As a consequence, BO is useful when function evaluation is time-consuming, such as when training and testing complex machine learning algorithms (e.g. deep neural networks) or tuning algorithms on large-scale dataset (e.g. ImageNet) [4]. Recently, BO has become popular in machine learning as it is highly effective in tuning hyperparameters of machine learning algorithms [8, 9, 20, 23].

Most previous work in BO assumes that we evaluate the objective function sequentially [13], though a few recent papers have considered parallel evaluations [3, 5, 18, 27]. While in practice, we can often evaluate several different choices in parallel, such as multiple machines can simultaneously train the machine learning algorithm with different sets of hyperparameters. In this paper, we assume that we can access $q \ge 1$ evaluations simultaneously at each iteration. Then we develop a new parallel acquisition function to guide where to evaluate next based on the decision-theoretical analysis.

Our Contributions. We propose a novel batch BO method which measures the information gain of evaluating q points via a new acquisition function, the parallel knowledge gradient (q-KG). This method is derived using a decision-theoretic analysis that chooses the set of points to evaluate next that is optimal in the average-case with respect to the posterior when there is only one batch of points remaining. Naively maximizing q-KG would be extremely computationally intensive, especially when q is large, and so, in this paper, we develop a method based on infinitesimal perturbation analysis (IPA) [14] to evaluate q-KG's gradient efficiently, allowing its efficient optimization. In our experiments on both synthetic functions and tuning practical machine learning algorithms, q-KG consistently finds better function values than other parallel BO algorithms, such as parallel EI [2, 20, 27], batch UCB [5] and parallel UCB with exploration [3]. q-KG provides especially large value when function evaluations are noisy. The code in this paper is available at https://github.com/wujian16/Cornel1-MOE.

The rest of the paper is organized as follows. Section 2 reviews related work. Section 3 gives background on Gaussian processes and defines notation used later. Section 4 proposes our new acquisition function q-KG for batch BO. Section 5 provides our computationally efficient approach to maximizing q-KG. Section 6 presents the empirical performance of q-KG and several benchmarks on synthetic functions and real problems. Finally, Section 7 concludes the paper.

2 Related work

Within the past several years, the machine learning community has revisited BO [8, 9, 18, 20, 21, 23] due to its huge success in tuning hyperparameters of complex machine learning algorithms. BO algorithms consist of two components: a statistical model describing the function and an acquisition function guiding evaluations. In practice, Gaussian Process (GP) [16] is the mostly widely used statistical model due to its flexibility and tractability. Much of the literature in BO focuses on designing good acquisition functions that reach optima with as few evaluations as possible. Maximizing this acquisition function usually provides a single point to evaluate next, with common acquisition functions for sequential Bayesian optimization including probability of improvement (PI) [25], expected improvement (EI) [13], upper confidence bound (UCB) [22], entropy search (ES) [11], and knowledge gradient (KG) [17].

Recently, a few papers have extended BO to the parallel setting, aiming to choose a batch of points to evaluate next in each iteration, rather than just a single point. [10, 20] suggests parallelizing EI by iteratively constructing a batch, in each iteration adding the point with maximal single-evaluation EI averaged over the posterior distribution of previously selected points. [10] also proposes an algorithm called "constant liar", which iteratively constructs a batch of points to sample by maximizing single-evaluation while pretending that points previously added to the batch have already returned values.

There are also work extending UCB to the parallel setting. [5] proposes the GP-BUCB policy, which selects points sequentially by a UCB criterion until filling the batch. Each time one point is selected, the algorithm updates the kernel function while keeping the mean function fixed. [3] proposes an algorithm combining UCB with pure exploration, called GP-UCB-PE. In this algorithm, the first point is selected according to a UCB criterion; then the remaining points are selected to encourage the diversity of the batch. These two algorithms extending UCB do not require Monte Carlo sampling, making them fast and scalable. However, UCB criteria are usually designed to minimize cumulative regret rather than immediate regret, causing these methods to underperform in BO, where we wish to minimize simple regret.

The parallel methods above construct the batch of points in an iterative greedy fashion, optimizing some single-evaluation acquisition function while holding the other points in the batch fixed. The acquisition function we propose considers the batch of points collectively, and we choose the batch to jointly optimize this acquisition function. Other recent papers that value points collectively include [2] which optimizes the parallel EI by a closed-form formula, [15, 27], in which gradient-based methods are proposed to jointly optimize a parallel EI criterion, and [18], which proposes a parallel version of the ES algorithm and uses Monte Carlo sampling to optimize the parallel ES acquisition function.

We compare against methods from a number of these previous papers in our numerical experiments, and demonstrate that we provide an improvement, especially in problems with noisy evaluations.

Our method is also closely related to the knowledge gradient (KG) method [7, 17] for the non-batch (sequential) setting, which chooses the Bayes-optimal point to evaluate if only one iteration is left [17], and the final solution that we choose is not restricted to be one of the points we evaluate. (Expected improvement is Bayes-optimal if the solution is restricted to be one of the points we evaluate.) We go beyond this previous work in two aspects. First, we generalize to the parallel setting. Second, while the sequential setting allows evaluating the KG acquisition function exactly, evaluation requires Monte Carlo in the parallel setting, and so we develop more sophisticated computational techniques to optimize our acquisition function. Recently, [28] studies a nested batch knowledge gradient policy. However, they optimize over a finite discrete feasible set, where the gradient of KG does not exist. As a result, their computation of KG is much less efficient than ours. Moreover, they focus on a nesting structure from materials science not present in our setting.

3 Background on Gaussian processes

In this section, we state our prior on f, briefly discuss well known results about Gaussian processes (GP), and introduce notation used later. We put a Gaussian process prior over the function $f: \mathbb{A} \to \mathbb{R}$, which is specified by its mean function $\mu(x): \mathbb{A} \to \mathbb{R}$ and kernel function $K(x_1, x_2): \mathbb{A} \times \mathbb{A} \to \mathbb{R}$. We assume either exact or independent normally distributed measurement errors, i.e. the evaluation $y(x_i)$ at point x_i satisfies

$$\boldsymbol{y}(\boldsymbol{x}_i) \mid f(\boldsymbol{x}_i) \sim \mathcal{N}(f(\boldsymbol{x}_i), \sigma^2(\boldsymbol{x}_i)),$$

where $\sigma^2: \mathbb{A} \to \mathbb{R}^+$ is a known function describing the variance of the measurement errors. If σ^2 is not known, we can also estimate it as we do in Section 6.

Supposing we have measured f at n points $\boldsymbol{x}^{(1:n)} := \{\boldsymbol{x}^{(1)}, \boldsymbol{x}^{(2)}, \cdots, \boldsymbol{x}^{(n)}\}$ and obtained corresponding measurements $y^{(1:n)}$, we can then combine these observed function values with our prior to obtain a posterior distribution on f. This posterior distribution is still a Gaussian process with the mean function $\boldsymbol{\mu}^{(n)}$ and the kernel function $K^{(n)}$ as follows

$$\mu^{(n)}(\boldsymbol{x}) = \mu(\boldsymbol{x})$$

$$+ K(\boldsymbol{x}, \boldsymbol{x}^{(1:n)}) \left(K(\boldsymbol{x}^{(1:n)}, \boldsymbol{x}^{(1:n)}) + \text{diag} \{ \sigma^{2}(\boldsymbol{x}^{(1)}), \cdots, \sigma^{2}(\boldsymbol{x}^{(n)}) \} \right)^{-1} (\boldsymbol{y}^{(1:n)} - \mu(\boldsymbol{x}^{(1:n)})),$$

$$K^{(n)}(\boldsymbol{x}_{1}, \boldsymbol{x}_{2}) = K(\boldsymbol{x}_{1}, \boldsymbol{x}_{2})$$

$$- K(\boldsymbol{x}_{1}, \boldsymbol{x}^{(1:n)}) \left(K(\boldsymbol{x}^{(1:n)}, \boldsymbol{x}^{(1:n)}) + \text{diag} \{ \sigma^{2}(\boldsymbol{x}^{(1)}), \cdots, \sigma^{2}(\boldsymbol{x}^{(n)}) \} \right)^{-1} K(\boldsymbol{x}^{(1:n)}, \boldsymbol{x}_{2}).$$

$$(3.1)$$

4 Parallel knowledge gradient (q-KG)

In this section, we propose a novel parallel Bayesian optimization algorithm by generalizing the concept of the knowledge gradient from [7] to the parallel setting. The knowledge gradient policy in [7] for discrete A chooses the next sampling decision by maximizing the expected incremental value of a measurement, without assuming (as expected improvement does) that the point returned as the optimum must be a previously sampled point.

We now show how to compute this expected incremental value of an additional iteration in the parallel setting. Suppose that we have observed n function values. If we were to stop measuring now, $\min_{x\in\mathbb{A}}\boldsymbol{\mu}^{(n)}(x)$ would be the minimum of the predictor of the GP. If instead we took one more batch of samples, $\min_{x\in\mathbb{A}}\boldsymbol{\mu}^{(n+q)}(x)$ would be the minimum of the predictor of the GP. The difference between these quantities, $\min_{x\in\mathbb{A}}\boldsymbol{\mu}^{(n)}(x)-\min_{x\in\mathbb{A}}\boldsymbol{\mu}^{(n+q)}(x)$, is the increment in expected solution quality (given the posterior after n+q samples) that results from the additional batch of samples.

This increment in solution quality is random given the posterior after n samples, because $\min_{x\in\mathbb{A}} \mu^{(n+q)}(x)$ is itself a random vector due to its dependence on the outcome of the samples. We can compute the probability distribution of this difference (with more details given below), and the q-KG algorithm values the sampling decision $z^{(1:q)} := \{z_1, z_2, \cdots, z_q\}$ according to its expected value, which we call the parallel knowledge gradient factor, and indicate it using the notation

q-KG. Formally, we define the q-KG factor for a set of candidate points to sample $z^{(1:q)}$ as

$$q\text{-}KG(\boldsymbol{z}^{(1:q)}, \mathbb{A}) = \min_{x \in \mathbb{A}} \boldsymbol{\mu}^{(n)}(x) - \mathbb{E}_n \left[\min_{x \in \mathbb{A}} \boldsymbol{\mu}^{(n+q)}(x) | \boldsymbol{z}^{(1:q)} \right], \tag{4.1}$$

where $\mathbb{E}_n\left[\cdot\right]:=\mathbb{E}\left[\cdot|x^{(1:n)},y^{(1:n)}\right]$ is the expectation taken with respect to the posterior distribution after n evaluations. Then we choose to evaluate the next batch of q points that maximizes the parallel knowledge gradient,

$$\max_{\boldsymbol{z}^{(1:q)} \subset \mathbb{A}} q \text{-} KG(\boldsymbol{z}^{(1:q)}, \mathbb{A}). \tag{4.2}$$

By construction, the parallel knowledge gradient policy is Bayes-optimal for minimizing the minimum of the predictor of the GP if only one decision is remaining. The *q-KG* algorithm will reduce to the parallel EI algorithm if function evaluations are noise-free and the final recommendation is restricted to the previous sampling decisions. Because under the two conditions above, the increment in expected solution quality will become

$$\min_{x \in \boldsymbol{x}^{(1:n)}} \boldsymbol{\mu}^{(n)}(x) - \min_{x \in \boldsymbol{x}^{(1:n)} \cup \boldsymbol{z}^{(1:q)}} \boldsymbol{\mu}^{(n+q)}(x) = \min y^{(1:n)} - \min \left\{ y^{(1:n)}, \min_{x \in \boldsymbol{z}^{(1:q)}} \boldsymbol{\mu}^{(n+q)}(x) \right\} \\
= \left(\min y^{(1:n)} - \min_{x \in \boldsymbol{z}^{(1:q)}} \boldsymbol{\mu}^{(n+q)}(x) \right)^+,$$

which is exactly the parallel EI acquisition function. However, computing q-KG and its gradient is very expensive. We will address the computational issues in Section 5. The full description of the q-KG algorithm is summarized as follows.

Algorithm 1 The *q-KG* algorithm

Require: the number of initial stage samples I, and the number of main stage sampling iterations N.

- 1: Initial Stage: draw I initial samples from a latin hypercube design in \mathbb{A} , $x^{(i)}$ for $i=1,\ldots,I$.
- 2: Main Stange:
- 3: **for** s = 1 to N **do**
- 4: Solve (4.2), i.e. get $(z_1^*, z_2^*, \dots, z_q^*) = \operatorname{argmax}_{\mathbf{z}^{(1:q)} \subset \mathbb{A}} q \text{-} KG(\mathbf{z}^{(1:q)}, \mathbb{A})$
- 5: Sample these points $(z_1^*, z_2^*, \dots, z_q^*)$, re-train the hyperparameters of the GP by MLE, and update the posterior distribution of f.
- 6: end for
- 7: **return** $x^* = \operatorname{argmin}_{x \in \mathbb{A}} \mu^{(I+Nq)}(x)$.

5 Computation of q-KG

In this section, we provide the strategy to maximize q-KG by a gradient-based optimizer. In Section 5.1 and Section 5.2, we describe how to compute q-KG and its gradient when $\mathbb A$ is finite in (4.1). Section 5.3 describes an effective way to discretize $\mathbb A$ in (4.1). The readers should note that there are two $\mathbb A$ s here, one is in (4.1) which is used to compute the q-KG factor given a sampling decision $z^{(1:q)}$. The other is the feasible domain in (4.2) ($z^{(1:q)} \subset \mathbb A$) that we optimize over. We are discretizing the first $\mathbb A$.

5.1 Estimating q-KG when \mathbb{A} is finite in (4.1)

Following [7], we express $\mu^{(n+q)}(x)$ as

$$\begin{split} \pmb{\mu}^{(n+q)}(x) &= \pmb{\mu}^{(n)}(x) + K^{(n)}(x, \pmb{z}^{(1:q)}) \left(K^{(n)}(\pmb{z}^{(1:q)}, \pmb{z}^{(1:q)}) \right. \\ &+ \mathrm{diag} \{ \sigma^2(\pmb{z}^{(1)}), \cdots, \sigma^2(\pmb{z}^{(q)}) \} \right)^{-1} \left(y(\pmb{z}^{(1:q)}) - \pmb{\mu}^{(n)}(\pmb{z}^{(1:q)}) \right). \end{split}$$

Because $y(\boldsymbol{z}^{(1:q)}) - \boldsymbol{\mu}^{(n)}(\boldsymbol{z}^{(1:q)})$ is normally distributed with zero mean and covariance matrix $K^{(n)}(\boldsymbol{z}^{(1:q)}, \boldsymbol{z}^{(1:q)}) + \operatorname{diag}\{\sigma^2(\boldsymbol{z}^{(1)}), \cdots, \sigma^2(\boldsymbol{z}^{(q)})\}$ with respect to the posterior after n observations, we can rewrite $\boldsymbol{\mu}^{(n+q)}(x)$ as

$$\mu^{(n+q)}(x) = \mu^{(n)}(x) + \tilde{\sigma}_n(x, z^{(1:q)}) Z_q,$$
 (5.1)

where Z_q is a standard q-dimensional normal random vector, and

$$\tilde{\sigma}_n(x, \boldsymbol{z}^{(1:q)}) = K^{(n)}(x, \boldsymbol{z}^{(1:q)})(D^{(n)}(\boldsymbol{z}^{(1:q)})^T)^{-1},$$

where $D^{(n)}(\boldsymbol{z}^{(1:q)})$ is the Cholesky factor of the covariance matrix $K^{(n)}(\boldsymbol{z}^{(1:q)}, \boldsymbol{z}^{(1:q)}) + \operatorname{diag}\{\sigma^2(\boldsymbol{z}^{(1)}), \cdots, \sigma^2(\boldsymbol{z}^{(q)})\}$. Now we can compute the q-KG factor using Monte Carlo sampling when $\mathbb A$ is finite: we can sample Z_q , compute (5.1), then plug in (4.1), repeat many times and take average.

5.2 Estimating the gradient of q-KG when \mathbb{A} is finite in (4.1)

In this section, we propose an unbiased estimator of the gradient of q-KG using IPA [14] when \mathbb{A} is finite. Accessing a stochastic gradient makes optimization much easier. By (5.1), we express q-KG as

$$q\text{-}KG(\boldsymbol{z}^{(1:q)}, \mathbb{A}) = \mathbb{E}_n\left(g(\boldsymbol{z}^{(1:q)}, \mathbb{A}, Z_q)\right), \tag{5.2}$$

where $g = \min_{x \in \mathbb{A}} \boldsymbol{\mu}^{(n)}(x) - \min_{x \in \mathbb{A}} \left(\boldsymbol{\mu}^{(n)}(x) + \tilde{\sigma}_n(x, \boldsymbol{z}^{(1:q)}) Z_q \right)$. Under the condition that μ and K are continuously differentiable, one can show that (please see details in the supplementary materials)

$$\frac{\partial}{\partial z_{ij}} q \text{-} KG(\boldsymbol{z}^{(1:q)}, \mathbb{A}) = \mathbb{E}_n \left(\frac{\partial}{\partial z_{ij}} g(\boldsymbol{z}^{(1:q)}, \mathbb{A}, Z_q) \right), \tag{5.3}$$

where z_{ij} is the jth dimension of the ith point in $z^{(1:q)}$. By the formula of g,

$$\frac{\partial}{\partial z_{ij}} g(\mathbf{z}^{(1:q)}, \mathbb{A}, Z_q) = \frac{\partial}{\partial z_{ij}} \boldsymbol{\mu}^{(n)}(x^*(\text{before})) - \frac{\partial}{\partial z_{ij}} \boldsymbol{\mu}^{(n)}(x^*(\text{after})) - \frac{\partial}{\partial z_{ij}} \tilde{\sigma}_n(x^*(\text{after}), \mathbf{z}^{(1:q)}) Z_q$$

where $x^*(\text{before}) = \operatorname{argmin}_{x \in \mathbb{A}} \boldsymbol{\mu}^{(n)}(x), \ x^*(\text{after}) = \operatorname{argmin}_{x \in \mathbb{A}} \left(\boldsymbol{\mu}^{(n)}(x) + \tilde{\sigma}_n(x, \boldsymbol{z}^{(1:q)}) Z_q\right),$ and

$$\begin{split} \frac{\partial}{\partial z_{ij}} \tilde{\sigma}_n(x^*(\text{after}), \pmb{z}^{(1:q)}) &= \left(\frac{\partial}{\partial z_{ij}} K^{(n)}(x^*(\text{after}), \pmb{z}^{(1:q)})\right) (D^{(n)}(\pmb{z}^{(1:q)})^T)^{-1} \\ &- K^{(n)}(x^*(\text{after}), \pmb{z}^{(1:q)}) (D^{(n)}(\pmb{z}^{(1:q)})^T)^{-1} \\ &\left(\frac{\partial}{\partial z_{ij}} D^{(n)}(\pmb{z}^{(1:q)})^T\right) (D^{(n)}(\pmb{z}^{(1:q)})^T)^{-1}. \end{split}$$

Now we can sample many times and take average to estimate the gradient of q-KG via (5.3). This technique is called infinitesimal perturbation analysis (IPA) in gradient estimation [14]. Since we can estimate the gradient of q-KG efficiently when \mathbb{A} is finite, we will apply some standard gradient-based optimization algorithms, such as multi-start stochastic gradient ascent to maximize q-KG.

5.3 Approximating q-KG when \mathbb{A} is infinite in (4.1) through discretization

We have specified how to maximize q-KG when \mathbb{A} is finite in (4.1), but usually \mathbb{A} is infinite. In this case, we will discretize \mathbb{A} to approximate q-KG, and then maximize over the approximate q-KG. The discretization itself is an interesting research topic [17].

In this paper, the discrete set \mathbb{A}_n is not chosen statically, but evolves over time: specifically, we suggest drawing M samples from the global optima of the posterior distribution of the Gaussian process (please refer to [11, 18] for a description of this technique). This sample set, denoted by \mathbb{A}_n^M , is then extended by the locations of previously sampled points $\boldsymbol{x}^{(1:n)}$ and the set of candidate points $\boldsymbol{z}^{(1:q)}$. Then (4.1) can be restated as

$$q\text{-}KG(\boldsymbol{z}^{(1:q)}, \mathbb{A}_n) = \min_{x \in \mathbb{A}_n} \mu^{(n)}(x) - \mathbb{E}_n \left[\min_{x \in \mathbb{A}_n} \mu^{(n+q)}(x) | \boldsymbol{z}^{(1:q)} \right], \tag{5.4}$$

where $\mathbb{A}_n = \mathbb{A}_n^M \cup \boldsymbol{x}^{(1:n)} \cup \boldsymbol{z}^{(1:q)}$. For the experimental evaluation we recompute \mathbb{A}_n^M in every iteration after updating the posterior of the Gaussian process.

6 Numerical experiments

We conduct experiments in two different settings: the noise-free setting and the noisy setting. In both settings, we test the algorithms on well-known synthetic functions chosen from [1] and practical problems. Following previous literature [20], we use a constant mean prior and the ARD Matérn 5/2 kernel. In the noisy setting, we assume that $\sigma^2(x)$ is constant across the domain \mathbb{A} , and we estimate it together with other hyperparameters in the GP using maximum likelihood estimation (MLE). We set M=1000 to discretize the domain following the strategy in Section 5.3. In general, the q-KG algorithm performs as well or better than state-of-the-art benchmark algorithms on both synthetic and real problems. It performs especially well in the noisy setting.

Before describing the details of the empirical results, we highlight the implementation details of our method and the open-source implementations of the benchmark methods. Our implementation inherits the open-source implementation of parallel EI from the Metrics Optimization Engine [26], which is fully implemented in C++ with a python interface. We reuse their GP regression and GP hyperparameter fitting methods and implement the *q-KG* method in C++. The code in this paper is available at https://github.com/wujian16/Cornell-MOE. Besides comparing to parallel EI in [26], we also compare our method to a well-known heuristic parallel EI implemented in Spearmint [12], the parallel UCB algorithm (GP-BUCB) and parallel UCB with pure exploration (GP-UCB-PE) both implemented in Gpoptimization [6].

6.1 Noise-free problems

In this section, we focus our attention on the noise-free setting, in which we can evaluate the objective exactly. We show that parallel knowledge gradient outperforms or is competitive with state-of-the-art benchmarks on several well-known test functions and tuning practical machine learning algorithms.

6.1.1 Synthetic functions

First, we test our algorithm along with the benchmarks on 4 well-known synthetic test functions: Branin2 on the domain $[-15, 15]^2$, Rosenbrock3 on the domain $[-2, 2]^3$, Ackley5 on the domain $[-2, 2]^5$, and Hartmann6 on the domain $[0, 1]^6$. We initiate our algorithms by randomly sampling 2d + 2 points from a Latin hypercube design, where d is the dimension of the problem. Figure 1 reports the mean and the standard deviation of the base 10 logarithm of the immediate regret by running 100 random initializations with batch size q = 4.

The results show that q-KG is significantly better on Rosenbrock3, Ackley5 and Hartmann6, and is slightly worse than the best of the other benchmarks on Branin2. Especially on Rosenbrock3 and Ackley5, q-KG makes dramatic progress in early iterations.

6.1.2 Tuning logistic regression and convolutional neural networks (CNN)

In this section, we test the algorithms on two practical problems: tuning logistic regression on the MNIST dataset and tuning CNN on the CIFAR10 dataset. We set the batch size to q = 4.

First, we tune logistic regression on the MNIST dataset. This task is to classify handwritten digits from images, and is a 10-class classification problem. We train logistic regression on a training set with 60000 instances with a given set of hyperparameters and test it on a test set with 10000 instances. We tune 4 hyperparameters: mini batch size from 10 to 2000, training iterations from 100 to 10000, the $\ell 2$ regularization parameter from 0 to 1, and learning rate from 0 to 1. We report the mean and standard deviation of the test error for 20 independent runs. From the results, one can see that both algorithms are making progress at the initial stage while q-KG can maintain this progress for longer and results in a better algorithm configuration in general.

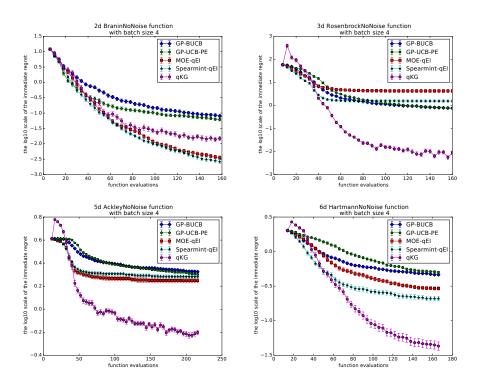


Figure 1: Performances on noise-free synthetic functions with q=4. We report the mean and the standard deviation of the log10 scale of the immediate regret vs. the number of function evaluations.

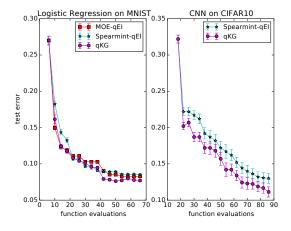


Figure 2: Performances on tuning machine learning algorithms with q=4

In the second experiment, we tune a CNN on CIFAR10 dataset. This is also a 10-class classification problem. We train the CNN on the 50000 training data with certain hyperparameters and test it on the test set with 10000 instances. For the network architecture, we choose the one in tensorflow tutorial. It consists of 2 convolutional layers, 2 fully connected layers, and on top of them is a softmax layer for final classification. We tune totally 8 hyperparameters: the mini batch size from 10 to 1000, training epoch from 1 to 10, the $\ell 2$ regularization parameter from 0 to 1, learning rate from 0 to 1, the kernel size from 0 to 1000, the number of channels in convolutional layers from 0 to 1000, the number of hidden units in fully connected layers from 100 to 1000, and the dropout rate from 0 to 1. We report the mean and standard deviation of the test error for 0 to 1000 in the initial stage and

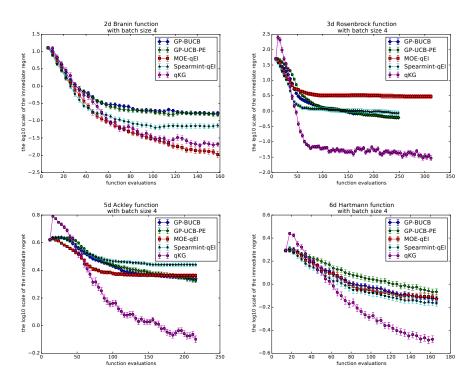


Figure 3: Performances on noisy synthetic functions with q=4. We report the mean and the standard deviation of the log 10 scale of the immediate regret vs. the number of function evaluations.

maintain this advantage to the end. This architecture has been carefully tuned by the human expert, and achieve a test error around 14%, and our automatic algorithm improves it to around 11%.

6.2 Noisy problems

In this section, we study problems with noisy function evaluations. Our results show that the performance gains over benchmark algorithms from q-KG evident in the noise-free setting are even larger in the noisy setting.

6.2.1 Noisy synthetic functions

We test on the same 4 synthetic functions from the noise-free setting, and add independent gaussian noise with standard deviation $\sigma=0.5$ to the function evaluation. The algorithms are not given this standard deviation, and must learn it from data.

The results in Figure 3 show that q-KG is consistently better than or at least competitive with all competing methods. Also observe that the performance advantage of q-KG is larger than for noise-free problems.

6.2.2 Noisy logistic regression with small test sets

Testing on a large test set such as ImageNet is slow, especially when we must test many times for different hyperparameters. To speed up hyperparameter tuning, we may instead test the algorithm on a subset of the testing data to approximate the test error on the full set. We study the performance of our algorithm and benchmarks in this scenario, focusing on tuning logistic regression on MNIST. We train logistic regression on the full training set of 60,000, but we test the algorithm by testing on 1,000 randomly selected samples from the test set, which provides a noisy approximation of the test error on the full test set.

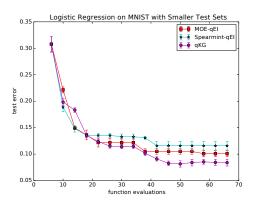


Figure 4: Tuning logistic regression on smaller test sets with q=4

We report the mean and standard deviation of the test error on the full set using the hyperparameters recommended by each parallel BO algorithm for 20 independent runs. The result shows that q-KG is better than both versions of parallel EI, and its final test error is close to the noise-free test error (which is substantially more expensive to obtain). As we saw with synthetic test functions, q-KG's performance advantage in the noisy setting is wider than in the noise-free setting.

Acknowledgments

The authors were partially supported by NSF CAREER CMMI-1254298, NSF CMMI-1536895, NSF IIS-1247696, AFOSR FA9550-12-1-0200, AFOSR FA9550-15-1-0038, and AFOSR FA9550-16-1-0046.

7 Conclusions

In this paper, we introduce a novel batch Bayesian optimization method q-KG, derived from a decision-theoretical perspective, and develop a computational method to implement it efficiently. We show that q-KG outperforms or is competitive with the state-of-the-art benchmark algorithms on several synthetic functions and in tuning practical machine learning algorithms.

References

- [1] D. Bingham. Optimization test problems. http://www.sfu.ca/~ssurjano/optimization.html, 2015.
- [2] C. Chevalier and D. Ginsbourger. Fast computation of the multi-points expected improvement with applications in batch selection. In *Learning and Intelligent Optimization*, pages 59–69. Springer, 2013.
- [3] E. Contal, D. Buffoni, A. Robicquet, and N. Vayatis. Parallel gaussian process optimization with upper confidence bound and pure exploration. In *Machine Learning and Knowledge Discovery in Databases*, pages 225–240. Springer, 2013.
- [4] L. Deng, D. Yu, et al. Deep learning: methods and applications. *Foundations and Trends in Signal Processing*, 7(3–4):197–387, 2014.
- [5] T. Desautels, A. Krause, and J. W. Burdick. Parallelizing exploration-exploitation tradeoffs in gaussian process bandit optimization. *The Journal of Machine Learning Research*, 15(1): 3873–3923, 2014.
- [6] E. C. etal. Gpoptimization. https://reine.cmla.ens-cachan.fr/e.contal/gpoptimization, 2015.
- [7] P. Frazier, W. Powell, and S. Dayanik. The knowledge-gradient policy for correlated normal beliefs. *INFORMS journal on Computing*, 21(4):599–613, 2009.

- [8] J. R. Gardner, M. J. Kusner, Z. E. Xu, K. Q. Weinberger, and J. Cunningham. Bayesian optimization with inequality constraints. In *Proceedings of the 31st International Conference on Machine Learning (ICML-14)*, pages 937–945, 2014.
- [9] M. Gelbart, J. Snoek, and R. Adams. Bayesian optimization with unknown constraints. In *Proceedings of the Thirtieth Conference Annual Conference on Uncertainty in Artificial Intelligence (UAI-14)*, pages 250–259, Corvallis, Oregon, 2014. AUAI Press.
- [10] D. Ginsbourger, R. Le Riche, and L. Carraro. Kriging is well-suited to parallelize optimization. In Computational Intelligence in Expensive Optimization Problems, pages 131–162. Springer, 2010.
- [11] J. M. Hernández-Lobato, M. W. Hoffman, and Z. Ghahramani. Predictive entropy search for efficient global optimization of black-box functions. In *Advances in Neural Information Processing Systems*, pages 918–926, 2014.
- [12] R. A. Jasper Snoek, Hugo Larochelle and other Spearmint developers. Spearmint. https://github.com/HIPS/Spearmint, 2015.
- [13] D. R. Jones, M. Schonlau, and W. J. Welch. Efficient global optimization of expensive black-box functions. *Journal of Global optimization*, 13(4):455–492, 1998.
- [14] P. L'Ecuyer. A unified view of the IPA, SF, and LR gradient estimation techniques. *Management Science*, 36(11):1364–1383, 1990.
- [15] S. Marmin, C. Chevalier, and D. Ginsbourger. Differentiating the multipoint expected improvement for optimal batch design. In *International Workshop on Machine Learning, Optimization and Big Data*, pages 37–48. Springer, 2015.
- [16] C. E. Rasmussen. Gaussian processes for machine learning. MIT Press, 2006.
- [17] W. Scott, P. Frazier, and W. Powell. The correlated knowledge gradient for simulation optimization of continuous parameters using gaussian process regression. SIAM Journal on Optimization, 21(3):996–1026, 2011.
- [18] A. Shah and Z. Ghahramani. Parallel predictive entropy search for batch global optimization of expensive objective functions. In *Advances in Neural Information Processing Systems*, pages 3312–3320, 2015.
- [19] S. P. Smith. Differentiation of the cholesky algorithm. *Journal of Computational and Graphical Statistics*, 4(2):134–147, 1995.
- [20] J. Snoek, H. Larochelle, and R. P. Adams. Practical bayesian optimization of machine learning algorithms. In *Advances in neural information processing systems*, pages 2951–2959, 2012.
- [21] J. Snoek, K. Swersky, R. Zemel, and R. Adams. Input warping for bayesian optimization of non-stationary functions. In *Proceedings of the 31st International Conference on Machine Learning (ICML-14)*, pages 1674–1682, 2014.
- [22] N. Srinivas, A. Krause, M. Seeger, and S. M. Kakade. Gaussian process optimization in the bandit setting: No regret and experimental design. In *Proceedings of the 27th International Conference on Machine Learning (ICML-10)*, pages 1015–1022, 2010.
- [23] K. Swersky, J. Snoek, and R. P. Adams. Multi-task bayesian optimization. In Advances in neural information processing systems, pages 2004–2012, 2013.
- [24] B. S. Thomson, J. B. Bruckner, and A. M. Bruckner. *Elementary real analysis*. ClassicalReal-Analysis. com, 2008.
- [25] A. Törn and A. Zilinskas. Global optimization, volume 350 of lecture notes in computer science, 1989.
- [26] J. Wang, S. C. Clark, E. Liu, and P. I. Frazier. Metrics optimization engine. http://yelp.github.io/MOE/, 2014. Last accessed on 2016-01-21.

- [27] J. Wang, S. C. Clark, E. Liu, and P. I. Frazier. Parallel bayesian global optimization of expensive functions. *Operations Research*, 2015.
- [28] Y. Wang, K. G. Reyes, K. A. Brown, C. A. Mirkin, and W. B. Powell. Nested-batch-mode learning and stochastic optimization with an application to sequential multistage testing in materials science. *SIAM Journal on Scientific Computing*, 37(3):B361–B381, 2015.

Supplementary Material

A Asynchronous q-KG Optimization

The (A.1) corresponds to the synchronous q-KG optimization, in which we wait for all q points from our previous batch to finish before searching for a new batch of q points. However, in some applications, we may wish to generate a new batch of points to evaluate next while p(<q) points are still being evaluated, before we have their values. This is common in training machine learning algorithms, where different machine learning models do not necessarily finish at the same time.

$$\max_{\boldsymbol{z}^{(1:q)} \subset \mathbb{A}} q \text{-}KG(\boldsymbol{z}^{(1:q)}, \mathbb{A}). \tag{A.1}$$

We can generalize (A.1) to the asynchronous q-KG optimization. Given that p points are still under evaluation, now we would like to recommend a batch of q points to evaluate. As we did for the synchronous q-KG optimization above, now we estimate the gradient of the q-KG of the combined q+p points only with respect to the q points that we need to recommend. Then we proceed the same way via gradient-based algorithms.

B Speed-up analysis

Next, we compare q-KG at different levels of parallelism against the fully sequential KG algorithm. We test the algorithms with different batch sizes on two noisy synthetic functions Branin2 and Hartmann6, whose standard deviation of the noise is $\sigma=0.5$. From the results, our parallel knowledge gradient method does provide a speed-up as q goes up.

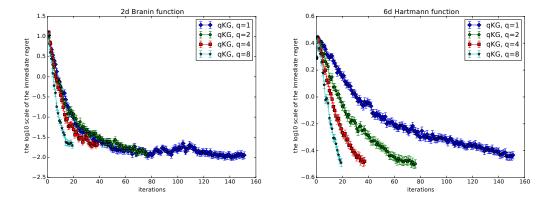


Figure 5: The performances of q-KG with different batch sizes. We report the mean and the standard deviation of the log 10 scale of the immediate regret vs. the number of iterations. Iteration 0 is the initial designs. For each iteration later, we evaluate q points recommended by the q-KG algorithm.

C Unbiasedness of the stochastic gradient estimator

Recall that in Section 5 of the main document, we have expressed the q-KG factor as follows,

$$q\text{-}KG(\boldsymbol{z}^{(1:q)}, \mathbb{A}) = \mathbb{E}_n\left(g(\boldsymbol{z}^{(1:q)}, \mathbb{A}, Z_q)\right)$$
 (C.1)

where the expectation is taken over Z_q and

$$g(\boldsymbol{z}^{(1:q)}, \mathbb{A}, Z_q) = \min_{x \in \mathbb{A}} \boldsymbol{\mu}^{(n)}(x) - \min_{x \in \mathbb{A}} \left(\boldsymbol{\mu}^{(n)}(x) + \tilde{\sigma}_n(x, \boldsymbol{z}^{(1:q)}) Z_q \right),$$

$$\tilde{\sigma}_n(x, \boldsymbol{z}^{(1:q)}) = K^{(n)}(x, \boldsymbol{z}^{(1:q)}) (D^{(n)}(\boldsymbol{z}^{(1:q)})^T)^{-1}.$$

The main purpose of this section is to prove the following proposition.

Proposition 1. When \mathbb{A} is finite, under the condition that μ and K are continuous differentiable,

$$\frac{\partial}{\partial z_{ij}} q\text{-}KG(\boldsymbol{z}^{(1:q)}, \mathbb{A}) \bigg|_{\boldsymbol{z}^{(1:q)} = \boldsymbol{\theta}^{(1:q)}} = \mathbb{E}_n \left(\frac{\partial}{\partial z_{ij}} g(\boldsymbol{z}^{(1:q)}, \mathbb{A}, Z_q) \bigg|_{\boldsymbol{z}^{(1:q)} = \boldsymbol{\theta}^{(1:q)}} \right), \quad (C.2)$$

where $1 \le i \le q$, $1 \le j \le d$, z_{ij} is the jth dimension of the ith point in $z^{(1:q)}$ and $\theta^{(1:q)} \in$ the interior of \mathbb{A}^q .

Without loss of generality, we assume that (1) i and j are fixed in advance and (2) $\mathbb{A} = [0,1]^d$, we would like to prove that (C.2) is correct. Before proceeding, we define one more notation $f_{\mathbb{A},Z_q}(z_{ij}) := g(\mathbf{z}^{(1:q)},\mathbb{A},Z_q)$ where $\mathbf{z}^{(1:q)}$ equals to $\theta^{(1:q)}$ component-wise except for z_{ij} . To prove it, we cite Theorem 1 in [14], which requires three conditions to make (C.2) valid: there exists an open neighborhood $\Theta \subset [0,1]$ of θ_{ij} where θ_{ij} is the jth dimension of ith point in $\theta^{(1:q)}$ such that

- (i) $f_{\mathbb{A},Z_q}(z_{ij})$ is continuous in Θ for any fixed \mathbb{A} and Z_q ;
- (ii) $f_{\mathbb{A},Z_q}(z_{ij})$ is differentiable except on a denumerable set in Θ for any given \mathbb{A} and Z_q ;
- (iii) the derivative of $f_{\mathbb{A},Z_q}(z_{ij})$ (when it exists) is uniformly bounded by $\Gamma(Z_q)$ for all $z_{ij} \in \Theta$, and the expectation of $\Gamma(Z_q)$ is finite.

C.1 Proof of condition (i)

Under the condition that the mean function μ and the kernel function K are continuous differentiable, we see that for any given x, $\tilde{\sigma}_n(x, \boldsymbol{z}^{(1:q)})$ is continuous differentiable in $\boldsymbol{z}^{(1:q)}$ by the result that the multiplication, the inverse (when the inverse exists) and the Cholesky operators preserve continuous differentiability Smith [19]. When A is finite, we see that $g(\boldsymbol{z}^{(1:q)}, \mathbb{A}, Z_q) = \min_{x \in \mathbb{A}} \boldsymbol{\mu}^{(n)}(x) - \min_{x \in \mathbb{A}} \left(\boldsymbol{\mu}^{(n)}(x) + \tilde{\sigma}_n(x, \boldsymbol{z}^{(1:q)}) Z_q\right)$ is continuous in $\boldsymbol{z}^{(1:q)}$. Then $f_{\mathbb{A}, Z_q}(z_{ij})$ is also continuous in z_{ij} by the definition of the function $f_{\mathbb{A}, Z_q}(z_{ij})$.

C.2 Proof of condition (ii)

By the expression that $f_{\mathbb{A},Z_q}(z_{ij}) = \min_{x \in \mathbb{A}} \boldsymbol{\mu}^{(n)}(x) - \min_{x \in \mathbb{A}} \left(\boldsymbol{\mu}^{(n)}(x) + \tilde{\sigma}_n(x,\boldsymbol{z}^{(1:q)})Z_q\right)$, if both $\operatorname{argmin}_{x \in \mathbb{A}} \boldsymbol{\mu}^{(n)}(x)$ and $\operatorname{argmin}_{x \in \mathbb{A}} \left(\boldsymbol{\mu}^{(n)}(x) + \tilde{\sigma}_n(x,\boldsymbol{z}^{(1:q)})Z_q\right)$ are unique, then $f_{\mathbb{A},Z_q}(z_{ij})$ is differentiable at z_{ij} . We define $D(\mathbb{A}) \subset \Theta$ to be the set that $f_{\mathbb{A},Z_q}(z_{ij})$ is not differentiable, then we see that

$$D(\mathbb{A}) \subset \bigcup_{x,x'\in\mathbb{A}} \left\{ z_{ij} \in \Theta : \boldsymbol{\mu}^{(n)}(x) = \boldsymbol{\mu}^{(n)}(x'), \frac{d\boldsymbol{\mu}^{(n)}(x)}{dz_{ij}} \neq \frac{d\boldsymbol{\mu}^{(n)}(x')}{dz_{ij}} \right\} \cup \bigcup_{x,x'\in\mathbb{A}} \left\{ z_{ij} \in \Theta : h_x(z_{ij}) = h_{x'}(z_{ij}), \frac{dh_x(z_{ij})}{dz_{ij}} \neq \frac{dh_{x'}(z_{ij})}{dz_{ij}} \right\}$$

where $h_x(z_{ij}) := \boldsymbol{\mu}^{(n)}(x) + \tilde{\sigma}_n(x, \boldsymbol{z}^{(1:q)}) Z_q$. $\boldsymbol{\mu}^{(n)}(x) \left(\boldsymbol{\mu}^{(n)}(x')\right)$ depend on z_{ij} if $x = z_i$ ($x' = z_i$) where z_i is the ith point of $\boldsymbol{z}^{(1:q)}$. As \mathbb{A} is finite, we only need to show that

$$\left\{z_{ij} \in \Theta : \boldsymbol{\mu}^{(n)}(x) = \boldsymbol{\mu}^{(n)}(x'), \frac{d\boldsymbol{\mu}^{(n)}(x)}{dz_{ij}} \neq \frac{d\boldsymbol{\mu}^{(n)}(x')}{dz_{ij}}\right\}$$

and

$$\left\{z_{ij} \in \Theta : h_x(z_{ij}) = h_{x'}(z_{ij}), \frac{dh_x(z_{ij})}{dz_{ij}} \neq \frac{dh_{x'}(z_{ij})}{dz_{ij}}\right\}$$

is denumerable.

Defining $\eta(z_{ij}):=h_{x_1}(z_{ij})-h_{x_2}(z_{ij})$ on Θ , one can see that $\eta(z_{ij})$ is continuous differentiable on Θ . We would like to show that $E:=\left\{z_{ij}\in\Theta:\eta(z_{ij})=0,\frac{d\eta(z_{ij})}{dz_{ij}}\neq0\right\}$ is denumerable. To prove it, we will show that E contains only isolated points. Then one can use a theorem in real analysis: any set of isolated points in $\mathbb R$ is denumerable (see the proof of statement 4.2.25 on page 165 in [24]). To prove that E only contains isolated points, we use the definition of an isolated point: $y\in E$ is an isolated point of E. We will prove by contradiction, suppose that $y\in E$ is a limit point of E, then it means that there exists a sequence

of points y_1, y_2, \cdots all belong to E such that $\lim_{n\to\infty} y_n = z_{ij}$. However, by the definition of derivative and

$$\begin{split} \eta(y_n) &= \eta(z_{ij}) = 0 \\ 0 &\neq \frac{d\eta(y)}{dy}\big|_{y=z_{ij}} = \lim_{n \to \infty} \frac{\eta(y_n) - \eta(z_{ij})}{y_n - z_{ij}} = \lim_{n \to \infty} 0 = 0, \end{split}$$

a contradiction. So we conclude that E only contains isolated points, so is denumerable.

Defining $\delta(z_{ij}) := \boldsymbol{\mu}^{(n)}(x_1) - \boldsymbol{\mu}^{(n)}(x_2)$ on Θ , $\delta(z_{ij})$ is also continuous differentiable on Θ , then one can similarly prove that $\left\{z_{ij} \in \Theta : \delta(z_{ij}) = 0, \frac{d\delta(z_{ij})}{dz_{ij}} \neq 0\right\}$ is denumerable.

C.3 Proof of condition (iii)

Recall from Section 5 of the main document,

$$\frac{d}{dz_{ij}} f(z_{ij}, \mathbb{A}, Z_q) = \frac{\partial}{\partial z_{ij}} g(\mathbf{z}^{(1:q)}, \mathbb{A}, Z_q)$$

$$= \frac{\partial}{\partial z_{ij}} \boldsymbol{\mu}^{(n)}(x^*(\text{before})) - \frac{\partial}{\partial z_{ij}} \boldsymbol{\mu}^{(n)}(x^*(\text{after}))$$

$$- \frac{\partial}{\partial z_{ij}} \tilde{\sigma}_n(\mathbf{z}^{(1:q)}, x^*(\text{after})) Z_q,$$

where $x^*(\text{before}) = \operatorname{argmin}_{x \in \mathbb{A}} \boldsymbol{\mu}^{(n)}(x)$, $x^*(\text{after}) = \operatorname{argmin}_{x \in \mathbb{A}} \left(\boldsymbol{\mu}^{(n)}(x) + \tilde{\sigma}_n(x, \boldsymbol{z}^{(1:q)}) Z_q\right)$, and

$$\begin{split} \frac{\partial}{\partial z_{ij}} \tilde{\sigma}_n(\pmb{z}^{(1:q)}, x^*(\text{after})) &= \left(\frac{\partial}{\partial z_{ij}} K^{(n)}(x^*(\text{after}), \pmb{z}^{(1:q)})\right) (D^{(n)}(\pmb{z}^{(1:q)})^T)^{-1} \\ &- D^{(n)}(x^*(\text{after}), \pmb{z}^{(1:q)}) (D^{(n)}(\pmb{z}^{(1:q)})^T)^{-1} \\ &\left(\frac{\partial}{\partial z_{ij}} D^{(n)}(\pmb{z}^{(1:q)})^T\right) (D^{(n)}(\pmb{z}^{(1:q)})^T)^{-1}. \end{split}$$

We can calculate the $\frac{\partial}{\partial z_{ij}} \boldsymbol{\mu}^{(n)}(x)$ as follows

$$\frac{\partial}{\partial z_{ij}}\boldsymbol{\mu}^{(n)}(x) = \begin{cases} \frac{\partial}{\partial z_{ij}}\boldsymbol{\mu}^{(n)}(z_i) & \text{if } x = z_i, \text{ i.e. the } i\text{th point of } \boldsymbol{z}^{(1:q)} \\ 0 & \text{otherwise.} \end{cases}$$

Using the fact that μ is continuously differentiable and \mathbb{A} is compact, then $\frac{\partial}{\partial z_{ij}}\mu^{(n)}(x)$ is bounded by some B>0. By the result that $\frac{\partial}{\partial z_{ij}}\tilde{\sigma}_n(\boldsymbol{z}^{(1:q)},x^*(\text{after}))$ is continuous, it is bounded by a vector $0 \leq \Lambda < \infty$ as \mathbb{A} is compact. Then

$$\left| \frac{d}{dz_{ij}} f(z_{ij}, \mathbb{A}, Z_q) \right| \le 2B + \sum_{i=1}^q \Lambda_i |z_i|,$$

where $Z_q = (z_1, \cdots, z_q)^T$. And

$$\mathbb{E}\left(\sum_{i=1}^{q} \Lambda_i |z_i|\right) = \sqrt{2/\pi} \sum_{i=1}^{q} \Lambda_i < \infty.$$

D Convergence of stochastic gradient ascent

In this section, we will prove that stochastic gradient ascent (SGA) converges to a stationary point. We follow the same idea of proving the Theorem 2 in [27].

First, it requires the step size γ_t satisfying $\gamma_t \to 0$ as $t \to \infty$, $\sum_{t=0}^\infty \gamma_t = \infty$ and $\sum_{t=0}^\infty \gamma_t^2 < \infty$. Second, it requires the second moment of the gradient estimator is finite. In the above section 1.3, we have shown that $|\frac{\partial}{\partial z_{ij}}g(\boldsymbol{z}^{(1:q)},\mathbb{A},Z_q)| \leq 2B + \sum_{i=1}^q \Lambda_i |z_i|$, then

$$\mathbb{E}\left(\frac{\partial}{\partial z_{ij}}g(\boldsymbol{z}^{(1:q)},\mathbb{A},Z_q)\right)^2 \leq 4B^2 + \sum_{i=1}^q \Lambda_i^2 + 4B\sqrt{2/\pi}\sum_{i=1}^q \Lambda_i + \frac{4}{\pi}\sum_{i\neq j} \Lambda_i \Lambda_j < \infty.$$