

Big Batch SGD: Automated Inference using Adaptive Batch Sizes

Soham De, Abhay Yadav, David Jacobs and Tom Goldstein

Department of Computer Science, University of Maryland, College Park, USA
`{sohamde, ayadav, djacobs, tomg}@cs.umd.edu`

October 20, 2016

Abstract

Classical stochastic gradient methods for optimization rely on noisy gradient approximations that become progressively less accurate as iterates approach a solution. The large noise and small signal in the resulting gradients makes it difficult to use them for adaptive stepsize selection and automatic stopping. We propose alternative “big batch” SGD schemes that adaptively grow the batch size over time to maintain a nearly constant signal-to-noise ratio in the gradient approximation. The resulting methods have similar convergence rates to classical SGD methods without requiring convexity of the objective function. The high fidelity gradients enable automated learning rate selection and do not require stepsize decay. For this reason, big batch methods are easily automated and can run with little or no user oversight.

1 Introduction

We are interested in problems of the form

$$\min_{x \in \mathcal{X}} \ell(x) := \begin{cases} \mathbb{E}_{z \sim p}[f(x; z)], \\ \frac{1}{N} \sum_{i=1}^N f(x; z_i), \end{cases} \quad (1)$$

where $\{z_i\}$ is a collection of data drawn from a probability distribution p . We assume that ℓ and f are differentiable, but possibly non-convex, and domain \mathcal{X} is convex. In typical applications, each term $f(x, z)$ measures how well a model with parameters x fits one particular data observation z . The expectation over z measures how well the model fits the entire corpus of data on average.

When N is large (or even infinite), it becomes intractable to exactly evaluate $\ell(x)$ or its gradient, which makes classical gradient methods impossible. In such situations, the method of choice for minimizing (1) is stochastic gradient descent (SGD) [29]. On iteration t , the SGD method randomly selects a batch $\mathcal{B} \subset \{z_i\}$ of sample data, and then computes

$$x_{t+1} = x_t - \alpha_t \nabla_x \ell_{\mathcal{B}}(x_t), \quad (2)$$

where $\ell_{\mathcal{B}}(x) = \frac{1}{|\mathcal{B}|} \sum_{z \in \mathcal{B}} f(x; z).$

Note that $\mathbb{E}_{\mathcal{B}}[\nabla_x \ell_{\mathcal{B}}(x_t)] = \nabla_x \ell(x)$, and so the approximate gradient $\nabla_x \ell_{\mathcal{B}}(x_t)$ can be interpreted as a “noisy” approximation to the true gradient.

Because the gradient approximations are noisy, the stepsize α_t must vanish as $t \rightarrow \infty$. Intuitively, as the iterates approach the minimizer of the objective, we must take progressively smaller steps to allow the (mean-zero) noise in our gradient approximations to average out to a small error. Typical stepsize rules require the user to find the optimal decay rate schedule, which may require a brute-force parameter search (a.k.a. “fine tuning” in the neural network literature).

In this paper, we consider a “big batch” strategy for SGD. Rather than letting the stepsize vanish as the approximation becomes more accurate, we let the minibatch \mathcal{B} adaptively grow in size to maintain a constant signal-to-noise ratio of the gradient approximation and prevent the algorithm from getting overwhelmed with noise. Using this batching strategy, we could keep the step size constant, or let it adapt using a simple Armijo backtracking line search, making the method completely adaptive with no user-defined parameters. We also consider an adaptive step size method based on the Barzilai-Borwein curvature estimate [1], that fully automates the big batch method, while empirically enjoying a faster convergence rate than the Armijo backtracking line search.

Big batch methods have several potential advantages over conventional small-batch SGD:

- Big batch methods don’t require the user to choose stepsize decay parameters. Larger batch sizes with less noise enable easy estimation of the accuracy of the approximate gradient, making it straightforward to adaptively scale up the batch size and maintain fast convergence.
- Backtracking line search tends to work very well when combined with big batch methods. The use of line search makes big batch methods completely adaptive with no parameters. The constant signal-to-noise ratio also enables us to define an adaptive step size method based on the Barzilai-Borwein curvature estimate. This stepsize rule performs better empirically on a range of convex problems than the backtracking line search.
- High-order methods like stochastic L-BFGS typically require a large amount of work to correct each gradient direction before a descent step is taken. When using big batches, the overhead of more complex methods like L-BFGS can be amortized over more costly gradient approximations. Furthermore, better Hessian approximations can be computed using less noisy gradient terms.
- For a restricted class of non-convex problems (functions satisfying the Polyak-Łojasiewicz Inequality), the per-iteration complexity of big batch SGD is linear and the approximate gradients vanish as the method approaches a solution. This makes it easy to define automated stopping conditions for convex problems. In contrast, small batch SGD exhibits sub-linear convergence, and the magnitude of the (very) noisy gradients is not usable as a stopping criterion.

For the reasons above, big batch SGD is potentially much easier to automate and requires much less user oversight than classical small batch SGD.

1.1 Related work

In this paper, we focus on automating stochastic optimization methods by reducing the noise in SGD. We do this by adaptively growing the batch size to control the variance in the gradient estimates, maintaining a constant signal-to-noise ratio, leading to automated methods that do not require vanishing stepsize parameters. We are not the first to study methods with dynamically increasing batch size. In [10], the authors propose to increase the size of the batch by a constant factor on *every* iteration, and prove linear convergence in terms of the iterates of the algorithm. In [5], the authors propose an adaptive strategy for growing the batch size; however, the authors do not present a theoretical guarantee for this method, and instead prove linear convergence for a continuously growing batch, similar to [10].

Other methods have been proposed to control gradient noise and mitigate the effects of vanishing stepsizes. Variance reduction (VR) methods are variants of SGD where an error correction term is used to decrease the variance in the noisy gradient estimates. These methods preserve a linear convergence rate (on strongly convex problems) and have been shown to outperform SGD on convex problems [8, 13, 31, 9] and in the parallel [28] and distributed setting [7]. A caveat, however, is that these methods require either extra storage or full gradient computations, both limiting factors when the dataset is very large. The step size also requires tuning. In a recent paper [12], the authors propose a growing batch strategy for a VR method that enjoys the same linear convergence guarantees. Another conceptually related approach is importance sampling, i.e., choosing training points such that the variance in the gradient estimates is reduced, leading to speedups [4, 6, 21].

While there has been some work on adaptive step size methods for stochastic optimization [20, 30, 32, 15, 33], the methods are largely heuristic without any kind of theoretical guarantees on convergence rates. The work [32] was a first step towards provable automated stochastic methods, and we build towards this direction to show provable convergence rates for the automated big batch method.

2 Big Batch SGD

2.1 Preliminaries and motivation

Classical gradient methods thrive when the current iterate is far from optimal. In this case, a small amount of data is necessary to find a descent direction, and optimization progresses efficiently. As x_t starts approaching the true solution x^* , however, noisy gradient estimates frequently do not produce descent directions and do not reliably decrease the objective. By choosing larger batches with less noise, we may be able to maintain descent directions on each iteration and uphold fast convergence. This observation motivates the proposed “big batch” method. We now explore this idea more rigorously.

To simplify notation, we hereon use $\nabla \ell$ to denote $\nabla_x \ell$. We wish to show that a noisy gradient approximation produces a descent direction when the noise is comparable in magnitude to the true gradient.

Lemma 1. *A sufficient condition for $-\nabla \ell_{\mathcal{B}}(x)$ to be a descent direction is*

$$\|\nabla \ell_{\mathcal{B}}(x) - \nabla \ell(x)\|^2 < \|\nabla \ell_{\mathcal{B}}(x)\|^2.$$

This is a standard result in stochastic optimization (see the short proof in the supplemental) that shows that, if the error $\|\nabla \ell_{\mathcal{B}}(x) - \nabla \ell(x)\|^2$ is small relative to the gradient $\|\nabla \ell_{\mathcal{B}}(x)\|^2$, the stochastic approximation provides us with a descent direction. But how big is this error and how large does a batch need to be to guarantee this condition? By the weak law of large numbers¹

$$\begin{aligned} \mathbb{E}[\|\nabla \ell_{\mathcal{B}}(x) - \nabla \ell(x)\|^2] &= \frac{1}{|\mathcal{B}|} \mathbb{E}[\|\nabla f(x, z) - \nabla \ell(x)\|^2] \\ &= \frac{1}{|\mathcal{B}|} \text{Tr Var}_z \nabla f(x, z), \end{aligned}$$

and so we can estimate the error of a stochastic gradient if we have some knowledge of the variance of $\nabla f(x, z)$. In practice, this variance could be estimated using the sample variance of a batch $\{\nabla f(x, z)\}_{z \in \mathcal{B}}$. However, we would like some bounds on the magnitude of this gradient to show that it is well-behaved, and also to analyze worst-case convergence behavior. To this end, we make the following assumption.

Assumption 1. *Suppose that f has L_z -Lipschitz dependence on data z , i.e., given two data points $z_1, z_2 \sim p(z)$, we have*

$$\|\nabla f(x; z_1) - \nabla f(x; z_2)\| \leq L_z \|z_1 - z_2\|.$$

Under this assumption, we can bound the error of the stochastic gradient. The bound is uniform with respect to the parameter x , which makes it rather useful in analyzing the convergence rate for big batch methods.

Theorem 1. *Given the current iterate x , suppose Assumption 1 holds and that the data distribution p has bounded second moment. Then the estimated gradient $\nabla \ell_{\mathcal{B}}(x)$ has variance bounded by*

$$\begin{aligned} \mathbb{E}_{\mathcal{B}} \|\nabla \ell_{\mathcal{B}}(x) - \nabla \ell(x)\|^2 &:= \text{Tr Var}_{\mathcal{B}}(\nabla \ell_{\mathcal{B}}(x)) \\ &\leq \frac{4L_z^2 \text{Tr Var}_z(z)}{|\mathcal{B}|}, \end{aligned}$$

where $z \sim p(z)$. Note the bound is uniform in x .

¹We assume the random variable $\nabla f(x, z)$ is measurable and has bounded second moment. These conditions will be guaranteed by the hypothesis of Theorem 1.

We present the proof in the Supplementary Material. Note that, using a finite number of samples $\{z\}$, one can easily approximate the quantity $\text{Var}_z(z)$ that appears in our bound.

2.2 A template for big batch SGD

Theorem 1 and Lemma 1 together suggest that we should expect $d = -\nabla\ell_{\mathcal{B}}$ to be a descent direction reasonably often provided

$$\begin{aligned} \theta^2 \|\nabla\ell_{\mathcal{B}}(x)\|^2 &\geq \frac{1}{|\mathcal{B}|} [\text{Tr Var}_z(\nabla f(x; z_i))], \\ \text{or } \theta^2 \|\nabla\ell_{\mathcal{B}}(x)\|^2 &\geq \frac{4L_z^2 \text{Tr Var}_z(z)}{|\mathcal{B}|}, \end{aligned} \quad (3)$$

for some $\theta < 1$. Big batch methods capitalize on this observation.

On each iteration t , starting from a point x_t , the big batch method performs the following steps:

1. Estimate the variance $\text{Tr Var}_z[\nabla f(x, z)]$, and a batch size K large enough that

$$\begin{aligned} \theta^2 \mathbb{E} \|\nabla\ell_{\mathcal{B}}(x)\|^2 &\geq \mathbb{E} \|\nabla\ell_{\mathcal{B}}(x) - \nabla\ell(x)\|^2 \\ &= \frac{1}{K} \text{Tr Var}_z f(x, z), \end{aligned} \quad (4)$$

where $\theta \in (0, 1)$ and $|\mathcal{B}| = K$.

2. Choose a stepsize α_t
3. Perform the update $x_{t+1} = x_t - \alpha_t \nabla\ell_{\mathcal{B}_t}(x_t)$.

Clearly, we have a lot of latitude in how to implement these steps using different variance estimators and different stepsize strategies. In the following section, we show that, if condition (4) holds, then linear convergence can be achieved using an appropriate constant stepsize. In subsequent sections, we address the issue of how to build practical big batch implementations using automated variance and stepsize estimators that require no user oversight.

3 Convergence Analysis

We now present convergence bounds for big batch SGD methods (5). We study stochastic gradient updates of the form

$$x_{t+1} = x_t - \alpha \nabla\ell_{\mathcal{B}_t}(x_t) = x_t - \alpha(\nabla\ell(x_t) + e_t), \quad (5)$$

where $e_t = \nabla\ell_{\mathcal{B}_t}(x_t) - \nabla\ell(x_t)$, and $\mathbb{E}_{\mathcal{B}}[e_t] = 0$. Let us also define $g_t = \nabla\ell(x_t) + e_t$.

Before we present our results, we first state two assumptions about the loss function $\ell(x)$ that are fairly standard in the optimization literature.

Assumption 2. *We assume that the objective function ℓ has L -Lipschitz gradients:*

$$\ell(x) \leq \ell(y) + \nabla\ell(y)^T(x - y) + \frac{L}{2} \|x - y\|^2.$$

Note that a consequence of Assumption 2 is the property: $\|\nabla\ell(x) - \nabla\ell(y)\| \leq L\|x - y\|$.

Assumption 3. *We also assume that the objective function ℓ satisfies the Polyak-Łojasiewicz Inequality:*

$$\|\nabla\ell(x)\|^2 \geq 2\mu(\ell(x) - \ell(x^*)).$$

Note that this inequality does *not* require ℓ to be convex. It does, however, imply that every stationary point is a global minimizer [14, 24].

We now present a result that establishes an upper bound on the objective value in terms of the error in the gradient of the sampled batch. We present all the proofs in the Supplementary Material.

Lemma 2. *Suppose we apply an update of the form (5) where the batch \mathcal{B}_t is uniformly sampled from the distribution p on each iteration t . If the objective ℓ satisfies Assumption 2, then we have:*

$$\mathbb{E}[\ell(x_{t+1}) - \ell(x^*)] \leq \mathbb{E}[\ell(x_t) - \ell(x^*)] - \left(\alpha - \frac{L\alpha^2}{2}\right) \mathbb{E}\|\nabla\ell(x_t)\|^2 + \frac{L\alpha^2}{2} \mathbb{E}\|e_t\|^2.$$

Further, if the objective ℓ satisfies the PL Inequality (Assumption 3), we have:

$$\mathbb{E}[\ell(x_{t+1}) - \ell(x^*)] \leq \left(1 - 2\mu\left(\alpha - \frac{L\alpha^2}{2}\right)\right) \mathbb{E}[\ell(x_t) - \ell(x^*)] + \frac{L\alpha^2}{2} \mathbb{E}\|e_t\|^2.$$

Using Lemma 2, we now provide convergence rates for big batch SGD.

Theorem 2. *Suppose ℓ satisfies Assumptions 2 and 3. Suppose further that on each iteration the batch size is large enough to satisfy (4) for $\theta \in (0, 1)$. If $0 \leq \alpha < \frac{2}{L\beta}$, where $\beta = \frac{\theta^2 + (1-\theta)^2}{(1-\theta)^2}$, then we get the following linear convergence bound for big batch SGD using updates of the form 5:*

$$\mathbb{E}[\ell(x_{t+1}) - \ell(x^*)] \leq \gamma \cdot \mathbb{E}[\ell(x_t) - \ell(x^*)],$$

where $\gamma = \left(1 - 2\mu\left(\alpha - \frac{L\alpha^2}{2}\right)\right)$. Choosing the optimal step size of $\alpha = \frac{1}{\beta\kappa}$, we get

$$\mathbb{E}[\ell(x_{t+1}) - \ell(x^*)] \leq \left(1 - \frac{\mu}{\beta L}\right) \cdot \mathbb{E}[\ell(x_t) - \ell(x^*)].$$

Note that the above linear convergence rate bound holds without requiring convexity. Comparing it with the convergence rate of deterministic gradient descent under similar assumptions, we see that big batch SGD suffers a slowdown by a factor β , due to the noise in the estimation of the gradients. We now present a result proving a $\mathcal{O}(1/t)$ convergence rate for general smooth convex functions.

Theorem 3. *Suppose ℓ satisfies Assumptions 2, is convex, and the condition (4) is satisfied on each iteration. Then we get sub-linear convergence of the form:*

$$\mathbb{E}[\ell(x_t) - \ell(x^*)] \leq \frac{\|x_0 - x^*\|^2}{(2\alpha - 2L\alpha^2\beta)(t+1)} = \mathcal{O}(1/t),$$

where $\beta = \frac{\theta^2 + (1-\theta)^2}{(1-\theta)^2}$ and $\alpha < \frac{1}{L\beta}$. Choosing the optimal step size of $\alpha = \frac{1}{2L\beta}$, we get

$$\mathbb{E}[\ell(x_t) - \ell(x^*)] \leq \frac{2L\beta\|x_0 - x^*\|^2}{t+1} = \mathcal{O}(1/t).$$

3.1 Comparison to classical SGD

Conventional small batch SGD methods can attain only $\mathcal{O}(1/t)$ convergence for strongly convex problems, thus requiring $\mathcal{O}(1/\epsilon)$ gradient evaluations to achieve an optimality gap less than ϵ , and this has been shown to be *optimal* in the online setting (or infinite data setting) [26]. In the previous section, however, we have shown that big batch SGD methods converge linearly in the number of iterations, under assumptions weaker than strong convexity, in the online setting. Unfortunately, per-iteration converge rates are not a fair comparison between these methods because the cost of a big batch iteration grows with the iteration count,

whereas the cost of classical SGD does not. For this reason, it is interesting to study the converge rate of big batch SGD as a function of *gradient evaluations*.

From Lemma 2, we see that we should not expect to achieve an optimality gap less than ϵ until we have: $\frac{L\alpha^2}{2}\mathbb{E}_{\mathcal{B}_t}\|e_t\|^2 < \epsilon$. In the worst case, under Assumption 1, Theorem 1 suggest that this requires $\frac{L\alpha^2}{2} \frac{4L_z^2 \text{Tr Var}_z(z)}{|\mathcal{B}|} < \epsilon$, or $|\mathcal{B}| \geq O(1/\epsilon)$ gradient evaluations. Note that in the online or infinite data sample, this is an optimal bound, and matches that of other SGD methods.

We choose to study the infinite sample case since the finite sample case is fairly trivial with a growing batch size: asymptotically, the batch size becomes the whole dataset, at which point we get the same asymptotic behavior as deterministic gradient descent, achieving linear convergence rates. Note, this same convergence rate is also achieved in the *finite* (but not the infinite) sample setting by variance reduction methods.

4 Practical Implementation with Backtracking Line Search

While one could implement a big batch method using analytical bounds on the gradient and its variance (such as that provided by Theorem 1), the purpose of big batch methods is to enable automated adaptive estimation of algorithm parameters. Furthermore, the stepsize bounds provided by our convergence analysis, like the stepsize bounds for classical SGD, are fairly conservative and more aggressive stepsize choices are likely to be more effective.

The framework outlined in Section 2.2 requires two ingredients: estimating the batch size and estimating the stepsize. Estimating the batch size needed to achieve (4) is fairly straight forward. We start with an initial batch size K , and draw a random batch with $|\mathcal{B}| = K$. We then compute the stochastic gradient estimate $\nabla\ell_{\mathcal{B}}(x_t)$ and the sample variance

$$\begin{aligned} V_{\mathcal{B}} &:= \frac{1}{|\mathcal{B}| - 1} \sum_{z \in \mathcal{B}} \|f(x_t, z) - \nabla\ell_{\mathcal{B}}(x_t)\|^2 \\ &\approx \text{Tr Var}_{z \in \mathcal{B}}(\nabla f(x; z_i)). \end{aligned} \tag{6}$$

We then test whether $\|\nabla\ell_{\mathcal{B}}(x_t)\|^2 > V_{\mathcal{B}}/|\mathcal{B}|$ as a proxy for (4). If this condition holds, then we proceed with a gradient step. If not, we increase the batch size $K \leftarrow K + \delta_K$, and check our condition again. For efficiency, instead of resampling the batch each time the condition fails, one could also simply add data to the current batch, and update the variance using an online scheme like Welford's method [16]. The fixed stepsize big batch method is listed in Algorithm 1. Notice that, while we introduced the parameter θ for our theory, we set the highest possible noise tolerance in all our implementations ($\theta = 1$).

We also consider a backtracking variant of SGD that adaptively tunes the stepsize. This method selects batch sizes using the same criteria (6) as in the constant stepsize case. However, after a batch has been selected, a backtracking Armijo line search is used to select a stepsize. In the Armijo line search, we keep decreasing the step size (by a constant factor) until the following condition is satisfied on each iteration:

$$\ell_{\mathcal{B}}(x_{t+1}) \leq \ell_{\mathcal{B}}(x_t) - c\alpha_t \|\nabla\ell_{\mathcal{B}}(x_t)\|^2, \tag{7}$$

where c is a parameter of the line search usually set to $0 < c \leq 0.5$. We now present a convergence result of big batch SGD using the Armijo line search.

Theorem 4. *Suppose that ℓ satisfies Assumptions 2 and 3 and on each iteration, and the batch size is large enough to satisfy (4) for $\theta \in (0, 1)$. If an Armijo line search, given by (7), is used, and the step size is decreased by a factor of 2 failing (7), then we get the following linear convergence bound for big batch SGD using updates of the form 5:*

$$\mathbb{E}[\ell(x_{t+1}) - \ell(x^*)] \leq \gamma \cdot \mathbb{E}[\ell(x_t) - \ell(x^*)],$$

where $\gamma = \left(1 - 2c\mu \min\left(\alpha_0, \frac{1}{2\beta L}\right)\right)$ and $0 < c \leq 0.5$. If, further the initial step size α_0 is set large enough such that $\alpha_0 \geq \frac{1}{2\beta L}$, then we get:

$$\mathbb{E}[\ell(x_{t+1}) - \ell(x^*)] \leq \left(1 - \frac{c\mu}{\beta L}\right) \mathbb{E}[\ell(x_t) - \ell(x^*)].$$

In practice, on iterations where the batch size increases, we double the stepsize before running the line search to prevent the stepsizes from always decreasing monotonically.

Algorithm 1 Big batch SGD: fixed stepsize

```

1: initialize starting pt.  $x_0$ , stepsize  $\alpha$ , initial batch size  $K > 1$ , batch size increment  $\delta_k$ 
2: while not converged do
3:   Draw random batch with size  $|\mathcal{B}| = K$ 
4:   Calculate  $V_{\mathcal{B}}$  and  $\nabla \ell_{\mathcal{B}}(x_t)$  using (6)
5:   while  $\|\nabla \ell_{\mathcal{B}}(x_t)\|^2 \leq V_{\mathcal{B}}/|\mathcal{B}|$  do
6:     Increase batch size  $K \leftarrow K + \delta_K$ 
7:     Sample more gradients
8:     Update  $V_{\mathcal{B}}$  and  $\nabla \ell_{\mathcal{B}}(x_t)$ 
9:   end while
10:   $x_{t+1} = x_t - \alpha \nabla \ell_{\mathcal{B}}(x_t)$ 
11: end while

```

5 Adaptive Step Sizes using the Barzilai-Borwein Estimate

While the Armijo backtracking line search leads to an automated big batch method, the step size sequence is monotonic (neglecting the heuristic mentioned in the previous section). In this section, we derive a non-monotonic step size scheme that uses curvature estimates to propose new stepsize choices.

Our derivation follows the classical adaptive Barzilai and Borwein (BB) method [1]. The BB method fits a quadratic model to the objective on each iteration, and a stepsize is proposed that is optimal for the local quadratic model [11]. To derive the analog of the BB method for stochastic problems, we consider quadratic approximations of the form $\ell(x) = \mathbb{E}_{\theta} f(x, \theta)$, where $f(x, \theta) = \frac{\nu}{2} \|x - \theta\|^2$ and $\theta \sim \mathcal{N}(x^*, \sigma^2 I)$. We now derive the optimal step size on each iteration for this quadratic approximation (for details refer to the Supplementary). We have:

$$\ell(x) = \mathbb{E}_{\theta} f(x, \theta) = \frac{\nu}{2} (\|x - x^*\|^2 + d\sigma^2),$$

Further, notice that:

$$\begin{aligned} \mathbb{E}_{\theta} [\nabla \ell(x)] &= \nu(x - x^*), \quad \text{and} \\ \text{Tr Var}_{\theta} [\nabla \ell(x)] &= d\nu^2 \sigma^2. \end{aligned}$$

Now, we can rewrite the big batch SGD update as:

$$\begin{aligned} x_{t+1} &= x_t - \alpha_t \frac{1}{|\mathcal{B}|} \sum_{i \in \mathcal{B}} \nu(x_t - \theta_i) \\ &= (1 - \nu\alpha_t)x_t + \nu\alpha_t x^* + \frac{\nu\sigma\alpha_t}{|\mathcal{B}|} \sum_{i \in \mathcal{B}} \xi_i, \end{aligned}$$

where we write $\theta_i = x^* + \sigma\xi_i$ with $\xi_i \sim \mathcal{N}(0, 1)$. Thus, the expected value of the function is:

$$\mathbb{E}[\ell(x_{t+1})] = \frac{\nu}{2} (\|(1 - \nu\alpha_t)(x_t - x^*)\|^2 + (1 + \frac{\nu^2\alpha_t^2}{|\mathcal{B}|})d\sigma^2).$$

Minimizing $\mathbb{E}[\ell(x_{t+1})]$ w.r.t. α_t we get:

$$\begin{aligned}\alpha_t &= \frac{1}{\nu} \cdot \frac{\|\mathbb{E}[\nabla \ell_{\mathcal{B}_t}(x_t)]\|^2}{\|\mathbb{E}[\nabla \ell_{\mathcal{B}_t}(x_t)]\|^2 + \frac{1}{|\mathcal{B}_t|} \text{Tr Var}[\nabla f(x_t)]} \\ &= \frac{1}{\nu} \cdot \left(1 - \frac{\frac{1}{|\mathcal{B}_t|} \text{Tr Var}[\nabla f(x_t)]}{\mathbb{E}\|\nabla \ell_{\mathcal{B}_t}(x_t)\|^2}\right) \\ &\geq \frac{1 - \theta^2}{\nu}.\end{aligned}\tag{8}$$

Here ν denotes the curvature of the quadratic approximation. Note that, in the case of *deterministic* gradient descent, the optimal step size is simply $1/\nu$ [1, 11]. Thus, the step size for big batch SGD is scaled down by at most $1 - \theta^2$. Note that the constant signal-to-noise ratio maintained by the big batch method, enables us to provide this lower bound.

We estimate the curvature ν_t on each iteration using the BB least-squares rule [1, 11, 32] as follows

$$\nu_t = \frac{\langle x_t - x_{t-1}, \nabla \ell_{\mathcal{B}_t}(x_t) - \nabla \ell_{\mathcal{B}_t}(x_{t-1}) \rangle}{\|x_t - x_{t-1}\|^2}.\tag{9}$$

Thus, each time we sample a batch \mathcal{B}_t on the t -th iteration, we calculate the gradient on that batch in the previous iterate, i.e., we calculate $\nabla \ell_{\mathcal{B}_t}(x_{t-1})$. This gives us an approximate curvature estimate, with which we derive the step size α_t using (8).

5.1 Convergence Proof

Here we present convergence proof for the adaptive step size method described above. For the convergence proof, we first state two assumptions:

Assumption 4. *Each f has L -Lipschitz gradients:*

$$f(x) \leq f(y) + \nabla f(y)^T(x - y) + \frac{L}{2}\|x - y\|^2.$$

Assumption 5. *Each f satisfies the Restricted Secant Inequality (RSI):*

$$\langle \nabla f(x) - \nabla f(y), x - y \rangle \geq \mu\|x - y\|^2.$$

Note that both assumptions are slightly stronger than Assumptions 2 and 3, i.e., Assumption 4 implies 2 and Assumption 5 implies 3 [14]. Further, note that Assumption 5 does not require convexity.

Further, for simplicity we assume that the step size is set to be: $\alpha_t = (1 - \theta^2)/\nu_t$. Using this and Assumption 4 and 5, we can upper and lower bound the curvature parameter, and thus the step size, as follows:

$$\mu \leq \nu_t \leq L \quad \implies \quad \frac{1 - \theta^2}{L} \leq \alpha_t \leq \frac{1 - \theta^2}{\mu}.$$

From Theorem 2, we see that we have linear convergence with the adaptive step size method when:

$$\begin{aligned}1 - 2\mu\left(\alpha - \frac{L\alpha^2\beta}{2}\right) &\leq 1 - \frac{2(1 - \theta^2)}{\kappa} + \beta(1 - \theta^2)^2\kappa < 1 \\ \implies \left(\frac{\mu}{L}\right)^2 &< \frac{2}{\beta(1 - \theta^2)}.\end{aligned}$$

5.2 Practical Implementation

From the convergence proof above, we see that the adaptive step size method enjoys a linear convergence rate when the problem is well-conditioned. To achieve robustness for poorly conditioned problems, we include a backtracking line search step after calculating (8), which ensures that the stepsize does not blow up.

Further, instead of calculating two gradients on each iteration t ($\nabla \ell_{\mathcal{B}_t}(x_t)$ and $\nabla \ell_{\mathcal{B}_t}(x_{t-1})$) our practical implementation uses the same batch (as well as each calculated step size) on two consecutive iterations so that no gradients are wasted.

We found the step size calculated from (8) tends to be noisy when the batch is small. While this did not affect long-term performance, we perform a smoothing operation to even out the step sizes and make performance more predictable. Let $\tilde{\alpha}_t$ denote the step size calculated from (8). Then, the step size on each iteration is given by

$$\alpha_t = \left(1 - \frac{|\mathcal{B}|}{N}\right) \alpha_{t-1} + \frac{|\mathcal{B}|}{N} \tilde{\alpha}_t.$$

This ensures that the update is proportional to how accurate the estimate on each iteration is. Further, when $|\mathcal{B}_t| = N$, we just use $\alpha_t = 1/\nu_t$, since in this case there is no noise in the algorithm.

6 Experiments

In this section, we present our experimental results. We explore big batch methods with both convex and non-convex (neural network) experiments on large and high-dimensional datasets.

6.1 Convex Experiments

For performing convex experiments, we test big batch SGD on a binary classification problem with logistic regression, and a linear regression problem:

$$\begin{aligned} \min_x \frac{1}{n} \sum_{i=1}^n \log(1 + \exp(b_i a_i^T x)) \quad \text{and} \\ \min_x \frac{1}{n} \sum_{i=1}^n (a_i^T x - b_i)^2. \end{aligned}$$

Figure 1 presents the results of our convex experiments. We perform experiments on three standard real world datasets: IJCNN1 [25] and COVERTYPE [3] for logistic regression, and MILLIONSONG [2] for linear regression. As a preprocessing step, we normalize the features for each dataset. We compare standard SGD with a decreasing learning rate given by $\alpha_t = a/(b+t)$ to: big batch SGD using a fixed step size (BBS Fixed), with backtracking line search (BBS Armijo) and with the adaptive step size (8) (BBS BB). Further, we also compare with deterministic gradient descent (GD) as well as the growing batch method described in [10] (which we denote as SF). While [10] proposes a quasi-Newton method using the growing batch strategy, we adapt their algorithm to a first-order method with a growing batch. For all algorithms (except the automated methods), we fine-tuned the stepsize parameters for fastest convergence using a grid search. BBS Armijo and BBS BB require no parameter tuning and are completely adaptive.

We see that across all three problems, the big batch methods outperform the other algorithms. Further, we notice that both fully automated big batch methods significantly outperform big batch SGD with a fixed step size. Further, notice that the automated methods increase the batch size more slowly than BBS Fixed and FS. Thus, the automated methods can take more steps with smaller batches, leveraging the advantages of stochastic methods for longer. The step sizes derived by the automated methods are very close to the optimal fixed step size rate.

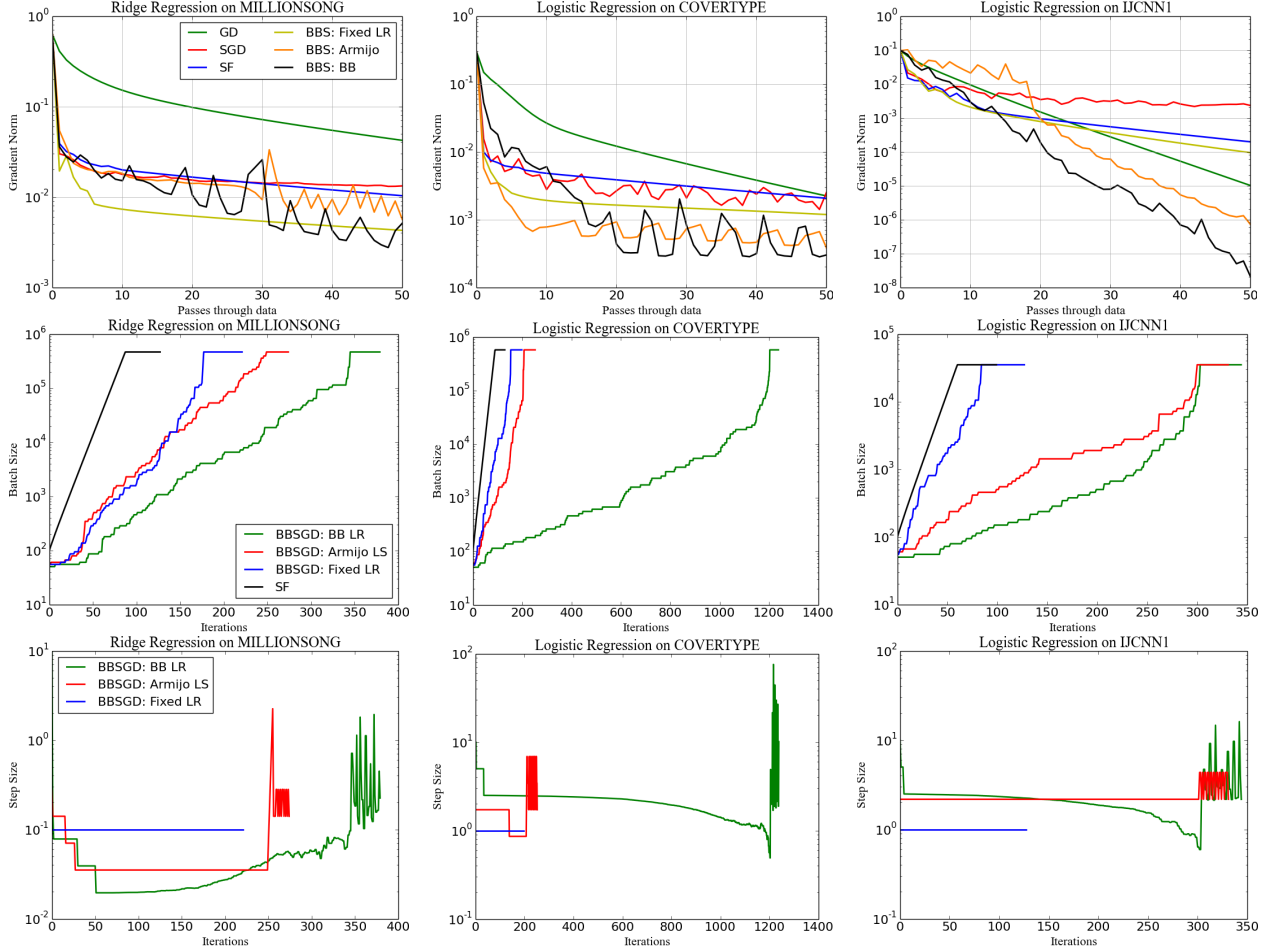


Figure 1: Convex experiments. Left to right: Ridge regression on MILLIONSONG; Logistic regression on COVERTYPE; Logistic regression on IJCNN1. The top row shows how the norm of the gradient decreases with the number of epochs, the middle and bottom rows show the batch sizes and step sizes used on each iteration by the big batch methods.

6.2 Neural Network Experiments

To demonstrate the versatility of the big batch SGD framework, we also present results on neural network experiments. We compare big batch SGD against SGD (Finely tuned LR schedules), SGD (Fixed LR) and Adadelta [33] (which is an adaptive method proposed for neural networks and does not require manual tuning of a learning rate). We also combine the big batch method with the AdaDelta update rule (BB+AdaDelta) to show that more complex SGD variants can benefit from growing batch sizes.

We train a convolutional neural network [19] (ConvNet) to classify three benchmark image datasets: CIFAR-10 [17], SVHN [23] and MNIST [19]. The ConvNet used in our experiments is composed of 4 layers, excluding the input layer. The inputs to the ConvNet are 32×32 pixel images. The first layer of the ConvNet contains $16 \times 3 \times 3$ filters, while the second layer contains $256 \times 3 \times 3$ filters. The third and fourth layers are fully connected [19] with 256 and 10 outputs respectively. Each layer except the last one is followed by a ReLu non-linearity [18] and a max pooling stage [27] of size 2×2 . This ConvNet has over 4.3 million weights.

To compare against fine-tuned SGD, we used a grid search on the learning rate schedule to identify optimal learning and decay parameters (up to a factor of 2 accuracy). For CIFAR10, the learning rate starts from 0.5 and is divided by 2 every 5 epochs with 0 learning rate decay. For SVHN, the learning rate starts from 0.5 and is divided by 2 every 5 epochs with $1e-5$ learning rate decay. For MNIST, the learning rate starts

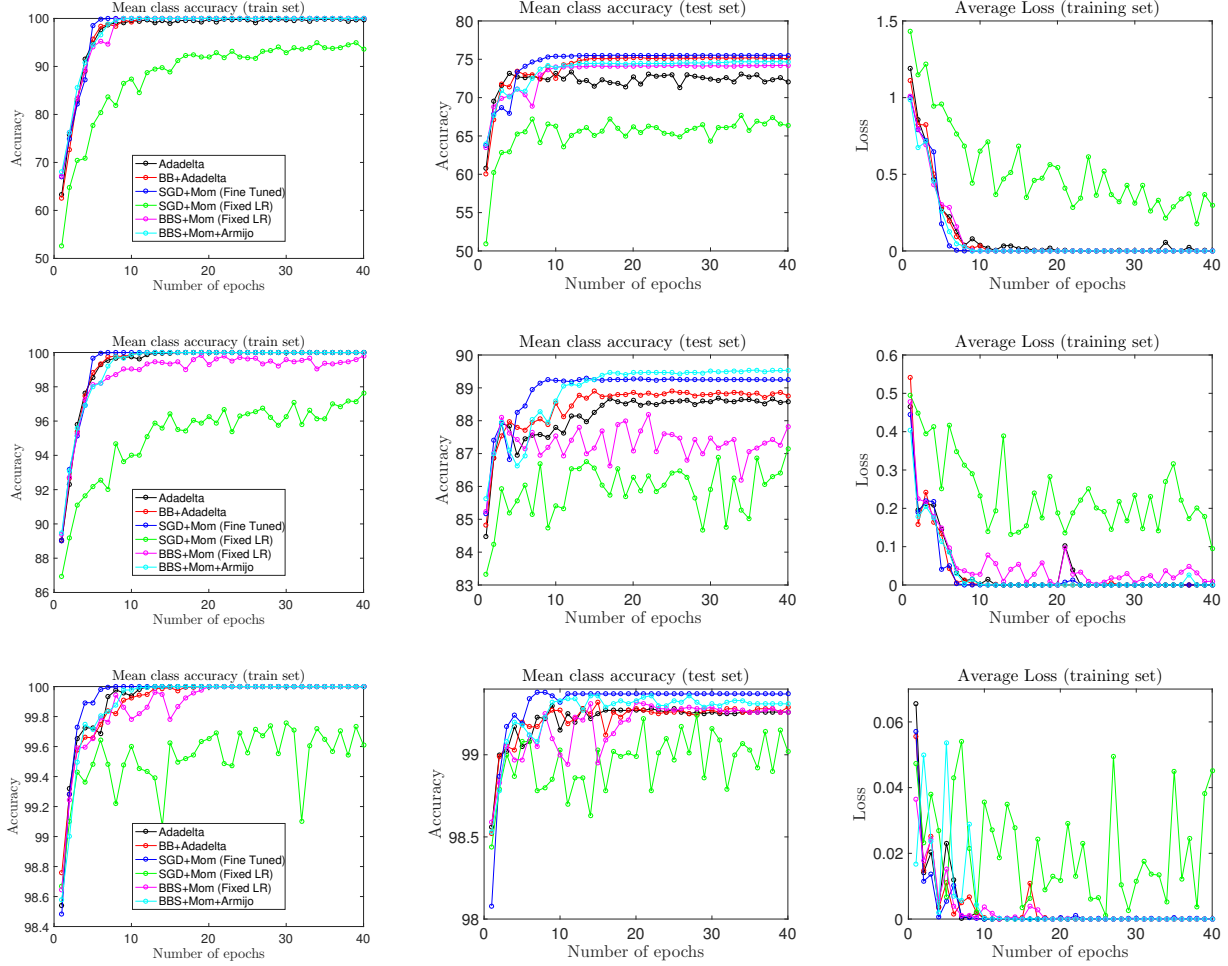


Figure 2: Neural Network Experiments. Top row presents results for CIFAR-10, middle row for SVHN, and bottom row for MNIST. The first column presents classification accuracies on the training set, the middle column presents classification accuracies on the test set, and the last column shows the change in the loss function.

from 1 and is divided by 2 every 3 epochs with 0 learning rate decay. All datasets use a momentum of 0.9 and mini-batches of size 128. Fixed learning rate methods use the default stepsize decay rule of the *Torch* library, which is $\alpha_t = \alpha_0 / (1 + 10^{-7}t)$, where α_0 was chosen to be the same initial LR used in the fine tuned experiments. We also tune the hyper-parameter ρ in the Adadelta and found 0.9, 0.9 and 0.8 to be best performing for CIFAR10, SVHN and MNIST respectively.

We plot the accuracy on the train and test set vs the number of epochs (full passes through the dataset) in Figure 2. We notice that the big batch SGD with backtracking performs better than both Adadelta and SGD (Fixed LR) in terms of both train and test error. big batch SGD even performs comparably to fine tuned SGD but without the trouble of fine tuning. Finally, we note that the big batch AdaDelta performs consistently better than plain AdaDelta on both large scale problems (SVHN and CIFAR-10), and performance is nearly identical on the small-scale MNIST problem.

7 Conclusion

In this work, we analyzed and studied the behavior of alternative SGD methods in which the batch size increases over time. Unlike classical SGD methods, in which stochastic gradients quickly become swamped with noise, these “big batch” methods maintain a nearly constant signal to noise ratio of the approximate gradient. As a result, big batch methods are able to adaptively adjust batch sizes without user oversight. The proposed automated methods are shown to be empirically comparable or better performing than other standard methods, but without requiring an expert user to choose learning rates and decay parameters.

References

- [1] Jonathan Barzilai and Jonathan M Borwein. Two-point step size gradient methods. *IMA Journal of Numerical Analysis*, 8(1):141–148, 1988.
- [2] Thierry Bertin-Mahieux, Daniel P.W. Ellis, Brian Whitman, and Paul Lamere. The million song dataset. In *Proceedings of the 12th International Conference on Music Information Retrieval (ISMIR 2011)*, 2011.
- [3] Jock A Blackard and Denis J Dean. Comparative accuracies of artificial neural networks and discriminant analysis in predicting forest cover types from cartographic variables. *Computers and electronics in agriculture*, 24(3):131–151, 1999.
- [4] Guillaume Bouchard, Théo Trouillon, Julien Perez, and Adrien Gaidon. Accelerating stochastic gradient descent via online learning to sample. *arXiv preprint arXiv:1506.09016*, 2015.
- [5] Richard H Byrd, Gillian M Chin, Jorge Nocedal, and Yuchen Wu. Sample size selection in optimization methods for machine learning. *Mathematical programming*, 134(1):127–155, 2012.
- [6] Dominik Csiba and Peter Richtárik. Importance sampling for minibatches. *arXiv preprint arXiv:1602.02283*, 2016.
- [7] Soham De and Tom Goldstein. Efficient distributed SGD with variance reduction. In *2016 IEEE International Conference on Data Mining*. IEEE, 2016.
- [8] Aaron Defazio, Francis Bach, and Simon Lacoste-Julien. Saga: A fast incremental gradient method with support for non-strongly convex composite objectives. In *Advances in Neural Information Processing Systems*, pages 1646–1654, 2014.
- [9] Aaron J Defazio, Tibério S Caetano, and Justin Domke. Finito: A faster, permutable incremental gradient method for big data problems. *arXiv preprint arXiv:1407.2710*, 2014.
- [10] Michael P Friedlander and Mark Schmidt. Hybrid deterministic-stochastic methods for data fitting. *SIAM Journal on Scientific Computing*, 34(3):A1380–A1405, 2012.
- [11] Tom Goldstein, Christoph Studer, and Richard Baraniuk. A field guide to forward-backward splitting with a FASTA implementation. *arXiv eprint, abs/1411.3406*, 2014.
- [12] Reza Harikandeh, Mohamed Osama Ahmed, Alim Virani, Mark Schmidt, Jakub Konečný, and Scott Sallinen. Stopwasting my gradients: Practical svrg. In *Advances in Neural Information Processing Systems*, pages 2251–2259, 2015.
- [13] Rie Johnson and Tong Zhang. Accelerating stochastic gradient descent using predictive variance reduction. In *Advances in Neural Information Processing Systems*, pages 315–323, 2013.
- [14] Hamed Karimi, Julie Nutini, and Mark Schmidt. Linear convergence of gradient and proximal-gradient methods under the polyak-łojasiewicz condition. In *Joint European Conference on Machine Learning and Knowledge Discovery in Databases*, pages 795–811. Springer, 2016.
- [15] Diederik Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014.

- [16] Donald Ervin Knuth. *The art of computer programming: sorting and searching*, volume 2. Pearson Education, 1998.
- [17] Alex Krizhevsky and Geoffrey Hinton. Learning multiple layers of features from tiny images, 2009.
- [18] Alex Krizhevsky, Ilya Sutskever, and Geoffrey E Hinton. Imagenet classification with deep convolutional neural networks. In *Advances in neural information processing systems*, pages 1097–1105, 2012.
- [19] Yann LeCun, Léon Bottou, Yoshua Bengio, and Patrick Haffner. Gradient-based learning applied to document recognition. *Proceedings of the IEEE*, 86(11):2278–2324, 1998.
- [20] Maren Mahsereci and Philipp Hennig. Probabilistic line searches for stochastic optimization. In *Advances In Neural Information Processing Systems*, pages 181–189, 2015.
- [21] Deanna Needell, Rachel Ward, and Nati Srebro. Stochastic gradient descent, weighted sampling, and the randomized kaczmarz algorithm. In *Advances in Neural Information Processing Systems*, pages 1017–1025, 2014.
- [22] Yurii Nesterov. *Introductory lectures on convex optimization: A basic course*, volume 87. Springer Science & Business Media, 2013.
- [23] Yuval Netzer, Tao Wang, Adam Coates, Alessandro Bissacco, Bo Wu, and Andrew Y Ng. Reading digits in natural images with unsupervised feature learning. In *NIPS workshop on deep learning and unsupervised feature learning*, volume 2011, page 4. Granada, Spain, 2011.
- [24] Boris Teodorovich Polyak. Gradient methods for minimizing functionals. *Zhurnal Vychislitel’noi Matematiki i Matematicheskoi Fiziki*, 3(4):643–653, 1963.
- [25] Danil Prokhorov. Ijcnn 2001 neural network competition. *Slide presentation in IJCNN*, 1, 2001.
- [26] Alexander Rakhlin, Ohad Shamir, and Karthik Sridharan. Making gradient descent optimal for strongly convex stochastic optimization. *arXiv preprint arXiv:1109.5647*, 2011.
- [27] Marc Aurelio Ranzato, Fu Jie Huang, Y-Lan Boureau, and Yann LeCun. Unsupervised learning of invariant feature hierarchies with applications to object recognition. In *Computer Vision and Pattern Recognition, 2007. CVPR’07. IEEE Conference on*, pages 1–8. IEEE, 2007.
- [28] Sashank J Reddi, Ahmed Hefny, Suvrit Sra, Barnabas Poczos, and Alex J Smola. On variance reduction in stochastic gradient descent and its asynchronous variants. In *Advances in Neural Information Processing Systems*, pages 2647–2655, 2015.
- [29] Herbert Robbins and Sutton Monro. A stochastic approximation method. *The annals of mathematical statistics*, pages 400–407, 1951.
- [30] Tom Schaul, Sixin Zhang, and Yann LeCun. No more pesky learning rates. In *Proceedings of The 30th International Conference on Machine Learning*, pages 343–351, 2013.
- [31] Mark Schmidt, Nicolas Le Roux, and Francis Bach. Minimizing finite sums with the stochastic average gradient. *arXiv preprint arXiv:1309.2388*, 2013.
- [32] Conghui Tan, Shiqian Ma, Yu-Hong Dai, and Yuqiu Qian. Barzilai-borwein step size for stochastic gradient descent. *arXiv preprint arXiv:1605.04131*, 2016.
- [33] Matthew D Zeiler. Adadelta: an adaptive learning rate method. *arXiv preprint arXiv:1212.5701*, 2012.

Supplementary Material

A Proof of Lemma 1

Proof. We know that $-\nabla\ell_{\mathcal{B}}(x)$ is a descent direction iff the following condition holds:

$$\nabla\ell_{\mathcal{B}}(x)^T \nabla\ell(x) > 0. \quad (10)$$

Expanding $\|\nabla\ell_{\mathcal{B}}(x) - \nabla\ell(x)\|^2$ we get

$$\begin{aligned} \|\nabla\ell_{\mathcal{B}}(x)\|^2 + \|\nabla\ell(x)\|^2 - 2\nabla\ell_{\mathcal{B}}(x)^T \nabla\ell(x) &< \|\nabla\ell_{\mathcal{B}}(x)\|^2, \\ \implies -2\nabla\ell_{\mathcal{B}}(x)^T \nabla\ell(x) &< -\|\nabla\ell(x)\|_2^2 \leq 0, \end{aligned}$$

which is always true for a descent direction (10). \square

B Proof of Theorem 1

Proof. Let $\bar{z} = \mathbb{E}[z]$ be the mean of z . Given the current iterate x , we assume that the batch \mathcal{B} is sampled uniformly with replacement from p . We then have the following bound:

$$\begin{aligned} \|\nabla f(x; z) - \nabla\ell(x)\|^2 &\leq 2\|\nabla f(x; z) - \nabla f(x, \bar{z})\|^2 + 2\|\nabla f(x, \bar{z}) - \nabla\ell(x)\|^2 \\ &\leq 2L_z^2\|z - \bar{z}\|^2 + 2\|\nabla f(x, \bar{z}) - \nabla\ell(x)\|^2 \\ &= 2L_z^2\|z - \bar{z}\|^2 + 2\|\mathbb{E}_z[\nabla f(x, \bar{z}) - \nabla f(x, z)]\|^2 \\ &\leq 2L_z^2\|z - \bar{z}\|^2 + 2\mathbb{E}_z\|\nabla f(x, \bar{z}) - \nabla f(x, z)\|^2 \\ &\leq 2L_z^2\|z - \bar{z}\|^2 + 2L_z^2\mathbb{E}_z\|\bar{z} - z\|^2 \\ &= 2L_z^2\|z - \bar{z}\|^2 + 2L_z^2 \text{Tr Var}_z(z), \end{aligned}$$

where the first inequality uses the property $\|a + b\|^2 \leq 2\|a\|^2 + 2\|b\|^2$, the second and fourth inequalities use Assumption 1, and the third line uses Jensen's inequality. This bound is *uniform* in x . We then have

$$\begin{aligned} \mathbb{E}_z\|\nabla f(x; z) - \nabla\ell(x)\|^2 &\leq 2L_z^2\mathbb{E}_z\|z - \bar{z}\|^2 + 2L_z^2 \text{Tr Var}_z(z) \\ &= 4L_z^2 \text{Tr Var}_z(z) \end{aligned}$$

uniformly for all x . The result follows from the observation that

$$\mathbb{E}_{\mathcal{B}}\|\nabla f_{\mathcal{B}}(x) - \nabla\ell(x)\|^2 = \frac{1}{|\mathcal{B}|} \mathbb{E}_z\|\nabla f(x; z) - \nabla\ell(x)\|^2.$$

\square

C Proof of Lemma 2

Proof. From (5) and Assumption 2 we get

$$\ell(x_{t+1}) \leq \ell(x_t) - \alpha g_t^T \nabla\ell(x_t) + \frac{L\alpha^2}{2} \|g_t\|^2.$$

Taking expectation w.r.t. the batch \mathcal{B}_t and conditioning on x_t , we get

$$\mathbb{E}[\ell(x_{t+1}) - \ell(x^*)] \leq \ell(x_t) - \ell(x^*) - \alpha \mathbb{E}[g_t]^T \nabla\ell(x_t) + \frac{L\alpha^2}{2} \mathbb{E}\|g_t\|^2$$

$$\begin{aligned}
&= \ell(x_t) - \ell(x^*) - \alpha \|\nabla \ell(x_t)\|^2 + \frac{L\alpha^2}{2} (\|\nabla \ell(x_t)\|^2 + \mathbb{E}\|e_t\|^2 + \mathbb{E}[e_t]^T \nabla \ell(x_t)) \\
&= \ell(x_t) - \ell(x^*) - \left(\alpha - \frac{L\alpha^2}{2}\right) \|\nabla \ell(x_t)\|^2 + \frac{L\alpha^2}{2} \mathbb{E}\|e_t\|^2 \\
&\leq \left(1 - 2\mu\left(\alpha - \frac{L\alpha^2}{2}\right)\right) (\ell(x_t) - \ell(x^*)) + \frac{L\alpha^2}{2} \mathbb{E}\|e_t\|^2,
\end{aligned}$$

where the second inequality follows from Assumption 3. Taking expectation, the result follows. \square

D Proof of Theorem 2

Proof. We begin by applying the reverse triangle inequality to (4) to get

$$(1 - \theta) \mathbb{E} \|\nabla \ell_{\mathcal{B}}(x)\| \leq \mathbb{E} \|\nabla \ell(x)\|$$

which applied to (4) yields

$$\frac{\theta^2}{(1 - \theta)^2} \mathbb{E} \|\nabla \ell(x_t)\|^2 \geq \mathbb{E} \|\nabla \ell_{\mathcal{B}}(x_t) - \nabla \ell(x_t)\|^2 = \mathbb{E} \|e_t\|^2. \quad (11)$$

Now, we apply (11) to the result in Lemma 2 to get

$$\mathbb{E}[\ell(x_{t+1}) - \ell(x^*)] \leq \mathbb{E}[\ell(x_t) - \ell(x^*)] - \left(\alpha - \frac{L\alpha^2\beta}{2}\right) \mathbb{E} \|\nabla \ell(x_t)\|^2,$$

where $\beta = \frac{\theta^2 + (1-\theta)^2}{(1-\theta)^2} \geq 1$. Assuming $\alpha - \frac{L\alpha^2\beta}{2} \geq 0$, we can apply Assumption 3 to write

$$\mathbb{E}[\ell(x_{t+1}) - \ell(x^*)] \leq \left(1 - 2\mu\left(\alpha - \frac{L\alpha^2\beta}{2}\right)\right) \mathbb{E}[\ell(x_t) - \ell(x^*)],$$

which proves the theorem. Note that $\max_{\alpha} \{\alpha - \frac{L\alpha^2\beta}{2}\} = \frac{1}{2L\beta}$, and $\mu \leq L$. It follows that

$$0 \leq \left(1 - 2\mu\left(\alpha - \frac{L\alpha^2\beta}{2}\right)\right) < 1.$$

The second result follows immediately. \square

E Proof of Theorem 3

Proof. Applying the reverse triangle inequality to (4) and using Lemma 2 we get, as in Theorem 2:

$$\mathbb{E}[\ell(x_{t+1})] \leq \mathbb{E}[\ell(x_t)] - \left(\alpha - \frac{L\alpha^2\beta}{2}\right) \mathbb{E} \|\nabla \ell(x_t)\|^2, \quad (12)$$

where $\beta = \frac{\theta^2 + (1-\theta)^2}{(1-\theta)^2} \geq 1$. Note that $\alpha - \frac{L\alpha^2\beta}{2} > 0$ if $\alpha < \frac{2}{L\beta}$.

From (5), taking norm on both sides and taking expectation, conditioned on all x_k , with $k = 0, 1, \dots, t$, we get

$$\begin{aligned}
\mathbb{E} \|x_{t+1} - x^*\|^2 &= \|x_t - x^*\|^2 - 2\alpha \mathbb{E} \langle x_t - x^*, \nabla \ell(x_t) + \epsilon_t \rangle + \alpha^2 \mathbb{E} \|\nabla \ell(x_t) + \epsilon_t\|^2 \\
&= \|x_t - x^*\|^2 - 2\alpha \langle x_t - x^*, \nabla \ell(x_t) \rangle + \alpha^2 \|\nabla \ell(x_t)\|^2 + \alpha^2 \mathbb{E} \|\epsilon_t\|^2 \\
&\leq \|x_t - x^*\|^2 - 2\alpha \langle x_t - x^*, \nabla \ell(x_t) \rangle + \alpha^2 \|\nabla \ell(x_t)\|^2 + \alpha^2 \frac{\theta^2}{(1 - \theta)^2} \|\nabla \ell(x_t)\|^2
\end{aligned}$$

$$\begin{aligned}
&= \|x_t - x^*\|^2 - 2\alpha \langle x_t - x^*, \nabla \ell(x_t) \rangle + \alpha^2 \beta \|\nabla \ell(x_t)\|^2 \\
&\leq \|x_t - x^*\|^2 - 2\alpha(\ell(x_t) - \ell(x^*)) + \alpha^2 \beta \|\nabla \ell(x_t)\|^2 \\
&\leq \|x_t - x^*\|^2 - 2\alpha(\ell(x_t) - \ell(x^*)) + 2L\alpha^2 \beta (\ell(x_t) - \ell(x^*)) \\
&= \|x_t - x^*\|^2 - (2\alpha - 2L\alpha^2 \beta)(\ell(x_t) - \ell(x^*)),
\end{aligned}$$

where we use the property that $\mathbb{E}[\epsilon_t] = 0$, and the properties $\ell(x) \leq \ell(x^*) + \langle x - x^*, \nabla \ell(x) \rangle$ (which follows from the convexity of ℓ) and $\|\nabla \ell(x)\|^2 \leq 2L(\ell(x) - \ell(x^*))$ (a proof for this identity can be found in [22]).

Note that $2\alpha - 2L\alpha^2 \beta > 0$ when $\alpha < \frac{1}{L\beta}$. Taking expectation on all x , we get

$$\mathbb{E}[\ell(x_t) - \ell(x^*)] \leq \frac{1}{2\alpha(1 - L\alpha\beta)} (\mathbb{E}\|x_t - x^*\|^2 - \mathbb{E}\|x_{t+1} - x^*\|^2). \quad (13)$$

Summing (13) over all $t = 0, 1, \dots, T$, and using the telescoping sum in $\|x_t - x^*\|^2$, we get:

$$\begin{aligned}
\sum_{t=0}^T \mathbb{E}[\ell(x_t) - \ell(x^*)] &\leq \frac{1}{2\alpha(1 - L\alpha\beta)} (\mathbb{E}\|x_0 - x^*\|^2 - \mathbb{E}\|x_{T+1} - x^*\|^2) \\
&\leq \frac{1}{2\alpha(1 - L\alpha\beta)} \|x_0 - x^*\|^2.
\end{aligned} \quad (14)$$

From (12) we see that $\mathbb{E}[\ell(x_{t+1})] \leq \mathbb{E}[\ell(x_t)]$ when $\alpha < \frac{2}{L\beta}$. Thus we can write (14) as

$$\mathbb{E}[\ell(x_T) - \ell(x^*)] \leq \frac{1}{T+1} \sum_{t=0}^T \mathbb{E}[\ell(x_t) - \ell(x^*)] \leq \frac{\|x_0 - x^*\|^2}{(2\alpha - 2L\alpha^2 \beta)(T+1)}.$$

Choosing the optimal step size of $\alpha = \frac{1}{2L\beta}$, we get

$$\mathbb{E}[\ell(x_T) - \ell(x^*)] \leq \frac{2L\beta\|x_0 - x^*\|^2}{T+1}.$$

□

F Proof of Theorem 4

Proof. Applying the reverse triangle inequality to (4) and using Lemma 2 we get, as in Theorem 2:

$$\mathbb{E}[\ell(x_{t+1}) - \ell(x^*)] \leq \mathbb{E}[\ell(x_t) - \ell(x^*)] - \left(\alpha - \frac{L\alpha^2 \beta}{2}\right) \mathbb{E}\|\nabla \ell(x_t)\|^2, \quad (15)$$

where $\beta = \frac{\theta^2 + (1-\theta)^2}{(1-\theta)^2} \geq 1$.

We will show that the backtracking condition in (7) is satisfied whenever $0 < \alpha_t \leq \frac{1}{\beta L}$. First notice that:

$$0 < \alpha_t \leq \frac{1}{\beta L} \implies -\alpha_t + \frac{L\alpha_t^2 \beta}{2} \leq -\frac{\alpha_t}{2}.$$

Thus, we can rewrite (15) as:

$$\begin{aligned}
\mathbb{E}[\ell(x_{t+1}) - \ell(x^*)] &\leq \mathbb{E}[\ell(x_t) - \ell(x^*)] - \frac{\alpha_t}{2} \mathbb{E}\|\nabla \ell(x_t)\|^2 \\
&\leq \mathbb{E}[\ell(x_t) - \ell(x^*)] - c\alpha_t \mathbb{E}\|\nabla \ell(x_t)\|^2,
\end{aligned}$$

where $0 < c \leq 0.5$. Thus, the backtracking line search condition (7) is satisfied whenever $0 < \alpha_t \leq \frac{1}{L\beta}$.

Now we know that α_t is either α_0 (the initial step size) or $\geq \frac{1}{2\beta L}$, where the step size is decreased by a factor of 2 each time the backtracking condition fails. Thus, we can rewrite the above as:

$$\mathbb{E}[\ell(x_{t+1}) - \ell(x^*)] \leq \mathbb{E}[\ell(x_t) - \ell(x^*)] - c \min\left(\alpha_0, \frac{1}{2\beta L}\right) \mathbb{E}\|\nabla \ell(x_t)\|^2.$$

Using Assumption 3 we get:

$$\mathbb{E}[\ell(x_{t+1}) - \ell(x^*)] \leq \left(1 - 2c\mu \min\left(\alpha_0, \frac{1}{2\beta L}\right)\right) \mathbb{E}[\ell(x_t) - \ell(x^*)].$$

Assuming we start off the step size at a large value s.t. $\min(\alpha_0, \frac{1}{2\beta L}) = \frac{1}{2\beta L}$, we can rewrite this as:

$$\mathbb{E}[\ell(x_{t+1}) - \ell(x^*)] \leq \left(1 - \frac{c\mu}{\beta L}\right) \mathbb{E}[\ell(x_t) - \ell(x^*)].$$

□

G Derivation of Adaptive Step Size

Consider the quadratic approximation $\ell(x) = \mathbb{E}_\theta f(x, \theta)$, where we define $f(x, \theta) = \frac{\nu}{2}\|x - \theta\|^2$ with $\theta \sim \mathcal{N}(x^*, \sigma^2 I)$.

We derive the optimal step size for this. We can rewrite the quadratic approximation as:

$$\ell(x) = \mathbb{E}_\theta f(x, \theta) = \frac{\nu}{2} \mathbb{E}_\theta \|x - \theta\|^2 = \frac{\nu}{2} [x^T x - 2x^T x^* - \mathbb{E}(\theta^T \theta)] = \frac{\nu}{2} (\|x - x^*\|^2 + d\sigma^2),$$

since we can write:

$$\mathbb{E}(\theta^T \theta) = \mathbb{E} \sum_{i=1}^d \theta_i^2 = \sum_{i=1}^d \mathbb{E} \theta_i^2 = \sum_{i=1}^d (x_i^*)^2 + \sigma^2 = \|x^*\|^2 + d\sigma^2.$$

Further, notice that:

$$\begin{aligned} \mathbb{E}_\theta [\nabla \ell(x)] &= \mathbb{E}_\theta [\nu(x - \theta)] = \nu(x - x^*), \quad \text{and} \\ \text{Tr Var}_\theta [\nabla \ell(x)] &= \mathbb{E}_\theta [\nu^2 (x - \theta)^T (x - \theta)] - \nu^2 (x - x^*)^T (x - x^*) = d\nu^2 \sigma^2. \end{aligned}$$

Using the quadratic approximation, we can rewrite the update for big batch SGD as follows:

$$x_{t+1} = x_t - \alpha_t \frac{1}{|\mathcal{B}|} \sum_{i \in \mathcal{B}} \nu(x_t - \theta_i) = (1 - \nu\alpha_t)x_t + \frac{\nu\alpha_t}{|\mathcal{B}|} \sum_{i \in \mathcal{B}} \theta_i = (1 - \nu\alpha_t)x_t + \nu\alpha_t x^* + \frac{\nu\sigma\alpha_t}{|\mathcal{B}|} \sum_{i \in \mathcal{B}} \xi_i,$$

where we write $\theta_i = x^* + \sigma\xi_i$ with $\xi_i \sim \mathcal{N}(0, 1)$. Thus, the expected value of the function is:

$$\begin{aligned} \mathbb{E}[\ell(x_{t+1})] &= \mathbb{E}_\xi \left[\ell \left((1 - \nu\alpha_t)x_t + \nu\alpha_t x^* + \frac{\nu\sigma\alpha_t}{|\mathcal{B}|} \sum_{i \in \mathcal{B}} \xi_i \right) \right] \\ &= \frac{\nu}{2} \mathbb{E}_\xi \left[\left\| (1 - \nu\alpha_t)(x_t - x^*) + \frac{\nu\sigma\alpha_t}{|\mathcal{B}|} \sum_{i \in \mathcal{B}} \xi_i \right\|^2 + d\sigma^2 \right] \\ &= \frac{\nu}{2} \left(\|(1 - \nu\alpha_t)(x_t - x^*)\|^2 + \mathbb{E}_\xi \left\| \frac{\nu\sigma\alpha_t}{|\mathcal{B}|} \sum_{i \in \mathcal{B}} \xi_i \right\|^2 + d\sigma^2 \right) \\ &= \frac{\nu}{2} \left(\|(1 - \nu\alpha_t)(x_t - x^*)\|^2 + \left(1 + \frac{\nu^2\alpha_t^2}{|\mathcal{B}|}\right) d\sigma^2 \right). \end{aligned}$$

Minimizing $\mathbb{E}[\ell(x_{t+1})]$ w.r.t. α_t we get:

$$\begin{aligned}
\alpha_t &= \frac{1}{\nu} \cdot \frac{\|\mathbb{E}[\nabla \ell_{\mathcal{B}_t}(x_t)]\|^2}{\|\mathbb{E}[\nabla \ell_{\mathcal{B}_t}(x_t)]\|^2 + \frac{1}{|\mathcal{B}_t|} \text{Tr Var}[\nabla f(x_t)]} \\
&= \frac{1}{\nu} \cdot \frac{\mathbb{E}\|\nabla \ell_{\mathcal{B}_t}(x_t)\|^2 - \frac{1}{|\mathcal{B}_t|} \text{Tr Var}[\nabla f(x_t)]}{\mathbb{E}\|\nabla \ell_{\mathcal{B}_t}(x_t)\|^2} \\
&= \frac{1}{\nu} \cdot \left(1 - \frac{\frac{1}{|\mathcal{B}_t|} \text{Tr Var}[\nabla f(x_t)]}{\mathbb{E}\|\nabla \ell_{\mathcal{B}_t}(x_t)\|^2} \right) \\
&\geq \frac{1 - \theta^2}{\nu}.
\end{aligned}$$

Here ν denotes the curvature of the quadratic approximation. Thus, the optimal step size for big batch SGD is the optimal step size for deterministic gradient descent scaled down by at most $1 - \theta^2$.