
AdaDelay: Delay Adaptive Distributed Stochastic Optimization

Suvrit Sra
MIT

Adams Wei Yu
CMU

Mu Li
CMU

Alexander J. Smola
CMU

Abstract

We develop distributed stochastic convex optimization algorithms under a delayed gradient model in which server nodes update parameters and worker nodes compute stochastic (sub)gradients. Our setup is motivated by the behavior of real-world distributed computation systems; in particular, we analyze a setting wherein worker nodes can be differently slow at different times. In contrast to existing approaches, we do not impose a worst-case bound on the delays experienced but rather allow the updates to be sensitive to the actual delays experienced. This sensitivity allows use of larger stepsizes, which can help speed up initial convergence without having to wait too long for slower machines; the global convergence rate is still preserved. We experiment with different delay patterns, and obtain noticeable improvements for large-scale real datasets with billions of examples and features.

1 Introduction

We study the stochastic convex optimization problem

$$\min_{x \in \mathcal{X}} f(x) := \mathbb{E}[F(x; \xi)], \quad (1.1)$$

where $\mathcal{X} \subset \mathbb{R}^d$ is a compact convex set, $F(\cdot, \xi)$ is a convex loss for each $\xi \sim \mathbb{P}$, and \mathbb{P} is a (possibly unknown) probability distribution from which we can draw i.i.d. samples. Problem (1.1) is important throughout optimization and machine learning [7, 14, 19–21]. It should be distinguished from (the easier) finite-sum optimization problems [3, 18, 23, 24], for which sharper results on the empirical loss are possible but not on the generalization error.

Appearing in Proceedings of the 19th International Conference on Artificial Intelligence and Statistics (AISTATS) 2016, Cadiz, Spain. JMLR: W&CP volume 51. Copyright 2016 by the authors.

A classic approach to solve (1.1) is stochastic gradient descent (SGD) [17] (or stochastic approximation [14]). SGD iteratively computes $x_{t+1} \leftarrow \Pi_{\mathcal{X}}(x_t - \alpha_t g_t)$, where $\Pi_{\mathcal{X}}$ denotes orthogonal projection onto \mathcal{X} , while $\alpha_t \geq 0$ is a suitable stepsize and g_t is an unbiased stochastic gradient, i.e., $\mathbb{E}[g_t] \in \partial f(x_t)$.

Although much more scalable than gradient descent, SGD is a sequential method that does not immediately apply to huge-scale problems which need distributed optimization [2]. This setting is central to real-world machine learning and has attracted great research interest, a large part of which is dedicated to scaling up SGD [1, 6, 9, 11, 15, 16].

Motivation. We also focus on huge-scale problems, and develop a new distributed SGD method that incorporates (and benefits from) more precise models of real-world cloud computing networks. Indeed, the delay properties exhibited by machines in cloud computing settings are often quite different from what one may observe on small clusters owned by individuals or small labs. Cloud resources are shared across users who run variegated tasks. Consequently, the cloud environment is invariably more diverse in its availability of resources such as CPU, disk, or network bandwidth, in contrast to environments where resources are shared by a small number of individuals. Thus, being able to model network delays in a more fine-grained way, and to use them to guide the optimization procedure can be of great value to both providers and users of large-scale distributed computing.

We investigate a new delay sensitive asynchronous SGD algorithm that adapts to the actual delays experienced, rather than relying on pessimistic global worst-case “bounded delays”. One may envision the following practical scheme: In the beginning, the server updates parameters as soon as it receives a gradient from any machine, weighting it inversely proportional to the actual delay. Towards the end, the server takes larger update steps when it obtains gradients from infrequent contributors, and smaller ones with gradients from frequent contributors, to reduce the bias caused by the initial aggressive steps. This broad scheme partially motivates our approach.

Contributions. We introduce and analyze *AdaDelay* (**Adaptive Delay**), an asynchronous SGD algorithm that more closely follows the actual delays experienced during computation. Specifically, AdaDelay uses step sizes sensitive to the actual delays observed. While this allows us to use larger stepsizes, it requires somewhat more intricate analysis because: (i) step sizes are no longer guaranteed to be monotonically decreasing; and (ii) residuals that measure progress are not independent across time as they are coupled by the delay random variable.

We validate AdaDelay by experimenting with real-world large-scale datasets containing over a billion samples and features. The experiments reveal that the models that we introduce on network delays are a reasonable approximation to the actual observed delays, and that in the regime of large delays (e.g., when there are stragglers) using delay sensitive steps is very helpful for faster generalization, that is, for models that converge more quickly on *test accuracy*; this is revealed by experiments where using AdaDelay leads to significant improvements on the test error (AUC).

Related Work. Our work is built on [1]. However, the most important difference is that [1] uses worst case delays that can be overly pessimistic, while AdaDelay uses the actual delays experienced (albeit at the cost of more involved theoretical analysis).

The body of related work on distributed optimization is quite large, so we can hardly be comprehensive. We summarize below a few closely related works, breaking up our summary into two typical settings: synchronous and asynchronous approaches.

Synchronous methods usually proceed in epochs, where a central node (parameter server) updates the global parameters, and waits until all the workers have finished their updates. For example, [22] and [8] leverage the finite sum structure of empirical risk minimization problems to derive a separable dual formulation for which a synchronous distributed coordinate ascent algorithm is adopted. Both [25] and [26] use an extreme strategy that only carries out a single round of communication at the end of the algorithm and then simply averages the results trained on the local workers. Performance of these algorithms, however, depends heavily on the slowest machine; commonly seen delay phenomena in commercial cloud computing systems involve multiple uncontrollable factors, and can contribute to a huge waste of resources due to waiting.

Asynchronous algorithms operate by letting each worker node run its local update without waiting for the others. Since network delays are inevitable, asynchronous methods ameliorate slow downs by avoiding waiting; thus, any updates computed by a local worker

can be immediately used to update the global parameter. The work [2] is a classic reference that introduces important asynchronous strategies; several more recent key works are [1, 9, 13, 19]. In [4, 15] the authors base their convergence on sparse data or variable settings; in comparison, our framework is more general, and thus capable of covering more applications. Of particular relevance to our paper is the recent work on delay adaptive gradient scaling in an AdaGrad [5] like framework [12]. The work [12] claims substantial improvements under specialized settings over [4]. Our experiments also confirm [12]’s claims that their best learning rate is insensitive to maximum delays. However, in our experience the method of [12] overly smooths the optimization path, which can have adverse effects on real-world data (see Section 4).

Finally, to our knowledge, previous works on asynchronous SGD (and its AdaGrad variants) assume monotonically decreasing step sizes. Our analysis involves non-monotonic steps to allow using actual delays instead of worst-case bounds; this proves to be quite beneficial in realistic settings. For instance, when there are stragglers that can slow down progress for all the machines in a worst-case delay model.

2 Problem Setup and Algorithm

We build on groundwork laid by [1, 13]; like them, we also optimize (1.1) under a delayed gradient model. We use the parameter-server computational framework [11], so that a central server¹ maintains the global parameters, while worker nodes compute stochastic gradients using their share of the data. The workers communicate their gradients back to the central server (using powerful communication saving techniques implemented in the work of [10]), which then updates the shared parameters and communicates them back.

To highlight our key ideas and avoid getting bogged down in unilluminating details, we consider only smooth stochastic optimization, i.e., $f \in C_L^1$, in this paper. Straightforward, (though tedious) extensions are possible to non-smooth problems, strongly convex costs, mirror descent versions, proximal splitting settings. We leave these details to the future.

As in [1, 7, 14], we make some standard assumptions:²

Assumption 2.1 (Lipschitz gradients). The function f has a locally L -Lipschitz gradients. That is,

$$\|\nabla f(x) - \nabla f(y)\| \leq L\|x - y\|, \quad \forall x, y \in \mathcal{X}.$$

¹This server is virtual; its physical realization may involve several machines, e.g., [10].

²These are easily satisfied for logistic-regression, least-squares, if the training data are bounded.

Assumption 2.2 (Bounded variance). There exists a constant $\sigma < \infty$ such that

$$\mathbb{E}_\xi[\|\nabla f(x) - \nabla F(x; \xi)\|^2] \leq \sigma^2, \quad \forall x \in \mathcal{X}.$$

Assumption 2.3 (Compact domain). Let $x^* \in \operatorname{argmin}_{x \in \mathcal{X}} f(x)$. Then,

$$\max_{x \in \mathcal{X}} \|x - x^*\| \leq R.$$

Finally, an additional assumption, also made in [1] is that of bounded gradients.

Assumption 2.4 (Bounded gradient). There exists a constant $G > 0$ so that

$$\|\nabla f(x)\| \leq G \quad \forall x \in \mathcal{X}.$$

These assumptions are typically satisfied in many machine learning problems, for instance, with logistic and least-squares losses, as long as the data samples ξ remain bounded, which is easy to satisfy.

Notation: Whenever a worker node returns an updated gradient at time t , we denote its associated random delay as τ_t , the delayed gradient as $g(t - \tau_t)$, and the step size as $\alpha(t, \tau_t)$. For a differentiable convex function h , the corresponding *Bregman divergence* is $D_h(x, y) := h(x) - h(y) - \langle \nabla h(y), x - y \rangle$. For simplicity, all norms are assumed to be Euclidean.

2.1 Delay model

Assumption 2.5 (Delay). We consider the following two practical delay models:

- (A) **Uniform:** Here $\tau_t \sim U(\{0, 2\bar{\tau}\})$. This model is a reasonable approximation to observed delays after an initial startup time of the network. We could make a more refined assumption that for $1 \leq t \leq T_1$ the delays are uniform on $\{0, \dots, T_1 - 1\}$. The analysis can be easily modified to handle this case; we omit it for brevity. Furthermore, the analysis also extends to delays having distributions with bounded support. Therefore, it indeed captures a wide spectrum of practical models.
- (B) **Scaled:** For each t , there is a $\theta_t \in (0, 1)$ such that $\tau_t < \theta_t t$. Moreover, assume that

$$\mathbb{E}[\tau_t] = \bar{\tau}_t, \quad \mathbb{E}[\tau_t^2] = B_t^2,$$

where $\bar{\tau}_t$ and B_t are constants that do not grow with t (the subscript only indicates that for each t , the random variable τ_t may have a different distribution). This model allows delay processes that are richer than uniform, as it no longer requires the support to be bounded. What it needs instead are bounded first and second moments.

Note. Our analysis is flexible, and can actually cover many other delay distributions by combining the above two delay models. For example, with Gaussian delays (where τ_t obeys a Gaussian distribution but its support is truncated, since $t > 0$) may be seen as a combination of the following: 1) When $t \geq C$ (a suitable constant), the Gaussian assumption indicates $\tau_t < \theta t$, which falls under our second delay model; 2) When $0 \leq t \leq C$, our proof technique with bounded support (same as uniform model) applies. Of course, more refined analysis for specific delay models may help tighten constants.

2.2 Algorithm

Under the above delay models, we consider the following projected stochastic gradient iteration:

$$x_{t+1} \leftarrow \operatorname{argmin}_{x \in \mathcal{X}} \left[\langle g(t - \tau_t), x \rangle + \frac{1}{2\alpha(t, \tau_t)} \|x - x_t\|^2 \right], \tag{2.1}$$

where the stepsize $\alpha(t, \tau_t)$ is sensitive to the actual delay observed. Whenever a worker transmits a delayed gradient $g(t - \tau_t)$ at time t , the parameter server conducts an update of (2.1). Thus, (2.1) generates a sequence of iterates $\{x_t\}_{t \geq 1}$; the server also maintains the averaged iterate

$$\bar{x}_T := \frac{1}{T} \sum_{t=1}^T x_{t+1}; \tag{2.2}$$

our convergence analysis is stated for this iterate.

3 Convergence analysis

We use stepsizes of the form $\alpha(t, \tau_t) = (L + \eta(t, \tau_t))^{-1}$, where the step offsets $\eta(t, \tau_t)$ are chosen to be sensitive to the actual delay of the incoming gradients. We typically use

$$\eta(t, \tau_t) = c\sqrt{t + \tau_t}, \tag{3.1}$$

for some constant c (to be chosen later). We can also consider time-varying c_t multipliers in (3.1) (see Corollary 3.4), but initially for clarity of presentation we let c be independent of t . If there are no delays, then $\tau_t = 0$ and iteration (2.1) reduces to the usual synchronous SGD. The constant c is used to trade off contributions to the error bound from the noise variance σ , the feasible set radius R , and the bounds on gradient norms.

Our convergence analysis builds on the groundwork of [1]. But the key difference is that our step sizes $\alpha(t, \tau_t)$ depend on the actual delay τ_t experienced, rather than on a fixed worst-case bounds on the maximum possible delay. These delay sensitive step sizes necessitate a more intricate analysis. There are two

main reasons for this: (i) the stepsize $\alpha(t, \tau_t)$ is no longer independent of the actual delay τ_t ; whereby (ii) the $\alpha(t, \tau_t)$ values are *no longer* monotonically decreasing, a property that the analysis of [1] relies on (and usual SGD convergence analysis also uses). We highlight our main theoretical results below; to streamline presentation, auxiliary technical lemmas are available in the supplement.

Theorem 3.1. *Let x_t be generated according to (2.1). Under Assumption 2.5 (A) (uniform delay) we have*

$$\begin{aligned} & \mathbb{E} \left[\sum_{t=1}^T (f(x_{t+1}) - f(x^*)) \right] \\ & \leq \left(\sqrt{2}cR^2\bar{\tau} + \frac{\sigma^2}{c} \right) \sqrt{T} + \frac{LG^2(4\bar{\tau}+3)(\bar{\tau}+1)}{6c^2} \log T \\ & \quad + \frac{1}{2}(L+c)R^2 + \bar{\tau}GR + \frac{LG^2\bar{\tau}(\bar{\tau}+1)(2\bar{\tau}+1)^2}{6(L^2+c^2)}, \end{aligned}$$

while under Assumption 2.5 (B) (scaled delay) we have

$$\begin{aligned} & \mathbb{E} \left[\sum_{t=1}^T (f(x_{t+1}) - f(x^*)) \right] \\ & \leq \frac{\sigma^2}{c} \sqrt{T} + \frac{1}{2}cR^2 \sum_{t=2}^T \frac{\bar{\tau}_t + 1}{\sqrt{2t-1}} + GR \left[1 + \sum_{t=1}^{T-1} \frac{B_t^2}{(T-t)^2} \right] \\ & \quad + G^2 \sum_{t=1}^T \frac{B_t^2 + 1 + \bar{\tau}_t}{L^2 + c^2(1-\theta_t)t} + \frac{1}{2}R^2(L+c). \end{aligned}$$

Proof Sketch. The proof begins by analyzing the difference $f(x_{t+1}) - f(x^*)$; Lemma A.2 (provided in the supplement) bounds this difference, ultimately leading to an inequality of the form:

$$\mathbb{E} \left[\sum_{t=1}^T (f(x_{t+1}) - f(x^*)) \right] \leq \mathbb{E} \left[\sum_{t=1}^T \Delta(t) + \Gamma(t) + \Sigma(t) \right].$$

The random variables $\Delta(t)$, $\Gamma(t)$, $\Sigma(t)$ are defined as

$$\Delta(t) := \frac{1}{2\alpha(t, \tau_t)} \left[\|x^* - x_t\|^2 - \|x^* - x_{t+1}\|^2 \right]; \quad (3.2)$$

$$\Gamma(t) := \langle \nabla f(x_t) - \nabla f(x_{t-\tau_t}), x_{t+1} - x^* \rangle; \quad (3.3)$$

$$\Sigma(t) := \frac{1}{2\eta(t, \tau_t)} \|\nabla f(x_{t-\tau_t}) - g(t - \tau_t)\|^2. \quad (3.4)$$

Thus, all that remains to do is bound each of these random variables and combine the bounds to obtain the claim of the theorem. Lemma A.3 bounds $\Delta(t)$ under Assumption 2.5(A), while Lemma A.4 bounds it under Assumption 2.5(B). Similarly, Lemmas A.5 and Lemma A.6 bound (3.3), while Lemma A.7 bounds (3.4). \square

Theorem 3.1 has several implications, which we now present as corollaries. Corollaries 3.2 and 3.3 show that both our delay models share a similar convergence rate of $\mathcal{O}(\frac{1}{\sqrt{T}})$. Corollary 3.4 shows that such

results continue to hold even if we replace the constant c with a bounded (away from zero, and from above) sequence $\{c_t\}$, a setting of great practical value. Finally, Corollary 3.5 gives the convergence of a more general choice of step sizes by considering $\eta_t = c_t(t + \tau_t)^\beta$ for $0 < \beta < 1$. It also highlights the known fact that for $\beta = 0.5$, the algorithm achieves the best theoretical convergence.

Corollary 3.2. *Let τ_t satisfy Assumption 2.5 (A). Then we have the following bound on \bar{x}_T :*

$$\mathbb{E}[f(\bar{x}_T) - f^*] = \mathcal{O} \left(D_1 \frac{\sqrt{T}}{T} + D_2 \frac{\log T}{T} + D_3 \frac{1}{T} \right),$$

where

$$\begin{aligned} D_1 &= \sqrt{2}cR^2\bar{\tau} + \frac{\sigma^2}{c}, D_2 = \frac{LG^2(4\bar{\tau}+3)(\bar{\tau}+1)}{6c^2}, \\ D_3 &= \frac{1}{2}(L+c)R^2 + \bar{\tau}GR + \frac{LG^2\bar{\tau}(\bar{\tau}+1)(2\bar{\tau}+1)^2}{6(L^2+c^2)}. \end{aligned}$$

The constant D_1 captures the variance due to stochastic gradients as well as the added variance due to the nonmonotone step sizes. The terms D_2 and D_3 capture the contribution to convergence based on delays and the Lipschitz smoothness of the gradients. We believe that it may be possible to get rid of the extra $\log T$ factor in the bound by a more refined analysis.

The following corollary follows easily from adapting the proof of Theorem 3.1, and combining it with Lemma A.6 which bounds $\Gamma(t)$ under the stated assumptions on the first two moments of the delay random variables.

Corollary 3.3. *Let τ_t satisfy Assumption 2.5 (B); let $\bar{\tau}_t = \tau$, $\theta_t = \theta$, and $B_t = B$ for all t . Then,*

$$\mathbb{E}[f(\bar{x}_T) - f^*] = \mathcal{O} \left(\frac{D_4}{\sqrt{T}} + D_5 \frac{\log(1 + \frac{c^2(1-\theta)T}{L^2})}{T} + \frac{D_6}{T} \right).$$

where

$$\begin{aligned} D_4 &= \left[\frac{1}{\sqrt{2}}cR^2(\bar{\tau}+1) + \frac{\sigma^2}{c} \right], D_5 = \frac{G^2(B^2 + \tau + 1)}{c^2(1-\theta)}, \\ D_6 &= \frac{1}{2}(L+c)R^2 + GR \left(1 + \frac{\pi^2 B^2}{6} \right). \end{aligned}$$

Here, the constant D_4 describes the contribution due to the variance introduced by stochastic gradients as well as the non-monotone step sizes. The role of D_5 and D_6 is similar to D_2 and D_3 from Corollary 3.2. As the reader may notice, the impact of the noise model on delay is on its contribution to D_5 and D_6 (which shrink as $O(\log T/T)$ and $O(1/T)$ respectively), and using properties of more specialized noise models one may be able to obtain more refined estimates of these constants.

The next corollary is of great value empirically. Theoretically, it states the impact on convergence rates for time varying choice of c_t .

Corollary 3.4. *If $\eta_t = c_t\sqrt{t + \tau_t}$ with $0 < M_1 \leq c_t \leq M_2$, then the conclusion of Theorem 3.1, Corollary 3.2 and 3.3 still hold, except that the term c is replaced by M_2 and $\frac{1}{c}$ by $\frac{1}{M_1}$.*

Finally, if one wishes to use step size offsets $\eta_t = c_t(t + \tau_t)^\beta$ where $0 < \beta < 1$, we obtain a convergence bound of the form stated below (we report only the asymptotically worst term, as this result is of limited importance).

Corollary 3.5. *Let $\eta_t = c_t(t + \tau_t)^\beta$ with $0 < M_1 \leq c_t \leq M_2$ and $0 < \beta < 1$. Then, there exists a constant D_τ such that*

$$\mathbb{E}[f(\bar{x}_T) - f^*] = \mathcal{O}\left(\frac{D_\tau}{T^{\min(\beta, 1-\beta)}}\right).$$

4 Experiments

We now evaluate the efficiency of AdaDelay in a distributed environment using large-scale real datasets.

4.1 Datasets and setup

We collected two click-through rate datasets for evaluation, which are shown in Table 1. One is the Criteo dataset³, where the first 8 days are used for training while the following 2 days are used for validation. We applied one-hot encoding for category and string features. The other dataset, named CTR2, is collected from a large Internet company. We sampled 100 million examples from three weeks for training, and 20 millions examples from the next week for validation. We extracted 2 billion unique features using the on-production feature extraction module. These two datasets have comparable size, but different example-feature ratios. We adopt Logistic Regression as our classification model.

	# train	# test	# features	nnz
Criteo	1.5B	400M	360M	58B
CTR2	110M	20M	1.9B	13B

Table 1: CTR datasets. M denotes millions; B billions.

All experiments were carried on a cluster with 20 machines. Most machines are equipped with dual Intel Xeon 2.40GHz CPUs, 32 GB memory and 1 Gbit/s Ethernet.

³<http://labs.criteo.com/downloads/download-terabyte-click-logs/>

4.2 Algorithms

We compare AdaDelay with two related methods AsyncAdaGrad [1] and AdaptiveRevision [12]. Their main difference lies in the choice of the learning rate at time t : $\alpha(t, \tau_t) = (L + \eta(t, \tau_t))^{-1}$. Denote by $\eta_j(t, \tau_t)$ the j -th element of $\eta(t, \tau_t)$, and similarly $g_j(t - \tau_t)$ the delayed gradient on feature j . AsyncAdaGrad adopts a scaled learning rate $\eta_j(t, \tau_t) = \sqrt{\sum_{i=1}^t g_j^2(i, \tau_i)}$. AdaptiveRevision takes into account actual delays by considering $g_j^{\text{bak}}(t, \tau_t) = \sum_{i=t-\tau}^{t-1} g_j(i, \tau_i)$. It uses a non-decreasing learning rate based on $\sqrt{\sum_{i=1}^t g_j^2(i, \tau_i) + 2g_j(t, \tau_t)g_j^{\text{bak}}(t, \tau_t)}$.

Similar to AsyncAdaGrad and AdaptiveRevision, we use a scaled learning rate in AdaDelay to better model the nonuniform sparsity of the dataset (this step size choice falls within the purview of Corollary 3.4). In other words, we set $\eta_j(t, \tau_t) = c_j\sqrt{t + \tau_t}$, where $c_j = \sqrt{\frac{1}{t} \sum_{i=1}^t \frac{i}{i+\tau_i} g_j^2(i - \tau_i)}$ averages the weighted delayed gradients on feature j . We follow the common practice of fixing L to 1 while choosing the best $\alpha(t, \tau_t) = \alpha_0(L + \eta(t, \tau_t))^{-1}$ by a grid search over α_0 .

We set the minibatch size to 10^5 and 10^4 for Criteo and CTR2, respectively, to reduce the communication frequency for better system performance⁴. We search for α_0 in the range $[10^{-4}, 1]$ and report the best results for each compared algorithm.

4.3 Implementation

We implemented these three methods in the parameter server framework [11], which is a high-performance asynchronous communication library supporting various data consistency models. There are two groups of nodes in this framework: workers and servers. Worker nodes run independently from each other. At each time, a worker first reads a minibatch of data from a distributed filesystem, and then pulls the relevant recent working set of parameters, namely the weights of the features that appear in this minibatch, from the server nodes. It next computes the gradients and then pushes these gradients to the server nodes.

The server nodes maintain the weights. For each feature, both AsyncAdaGrad and AdaDelay store the weight and the accumulated gradient which is used to compute the scaled learning rate. While AdaptiveRevision needs two more entries for each feature.

To compute the actual delay τ for AdaDelay, we let the server nodes record the time $t(w, i)$ when worker

⁴Probably due to the scale and the sparsity of the datasets, we observed no significant improvement when decreasing the minibatch size.

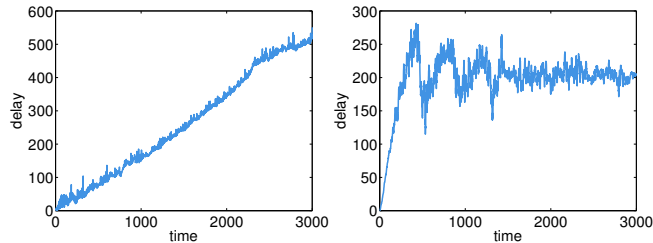


Figure 1: The first 3,000 observed delays at server nodes. Left: Criteo dataset with 1,600 workers; Right: CTR2 dataset with 400 workers.

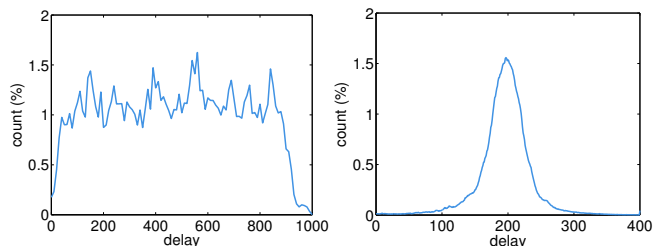


Figure 2: Histogram of all delays (left: Criteo with 1,600 workers; right: CTR2 with 400 workers).

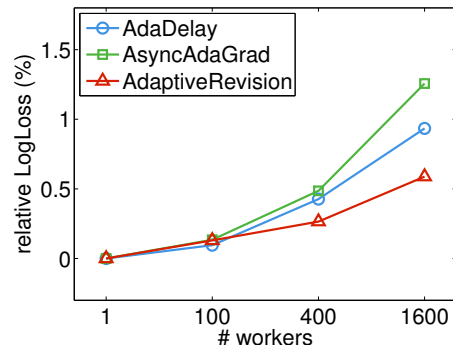
w is pulling the weight for minibatch i . Denote by $t'(w, i)$ the time when the server nodes are updating the weight by using the gradients of this minibatch. Then the actual delay of this minibatch can be obtained by $t'(w, i) - t(w, i)$.

AdaptiveRevision needs gradient components g_j^{bck} for each feature j to calculate its learning rate. If we send g_j^{bck} over the network by following [12], we increase the total network communication by 50%, which greatly harms system performance due to the limited network bandwidth. Instead, we store g_j^{bck} at the server node during while processing this minibatch. This incurs no extra network overhead, however, it increases the memory consumption of the server nodes.

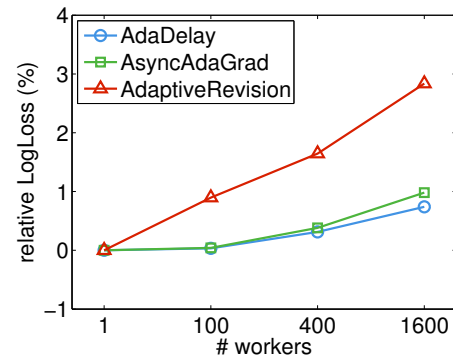
The parameter server implements a node using an operating system process, which has its own communication and computation threads. In order to run thousands of workers on our limited hardware, we may combine server workers into a single process to reduce the system overhead.

4.4 Results

Delays. We first visualize the actual delays observed at server nodes. As noted from Figure 1, delay τ_t is around θt at the early stage of the training, where the scaling constant θ varies for different tasks. For example, it is close to 0.2 when training the Criteo dataset with 1,600 workers, while it increases to 1 for the CTR2 dataset with 400 workers. After the delay



(a) Criteo



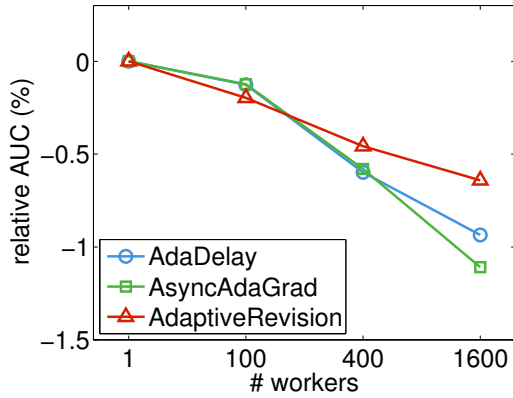
(b) CTR2

Figure 3: Relative (% worsening) of online LogLoss as function of maximal delays (lower is better).

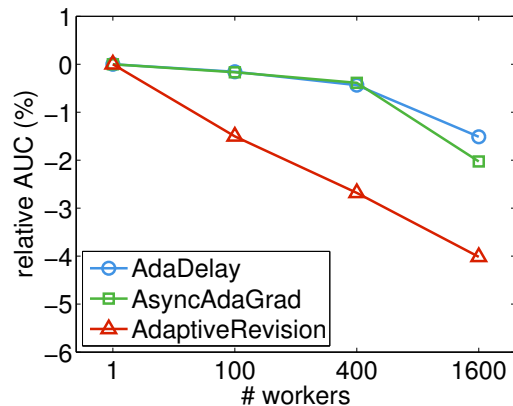
hits the value u , which is often half of the number of worker nodes, it behaves like a Gaussian distribution with mean u , which are shown in Figure 2.

Loss and AUC. Next, we present the comparison results of these three algorithms by varying the number of workers. Following [12] we use the online LogLoss as the criterion. That is, given example with feature vector d and label $y \in \{+1, -1\}$, we calculate the loss function $f(x, (d, y)) = \log(1 + \exp(-y\langle d, x \rangle))$ before updating x using (y, d) . Similar to [12], we report the average LogLoss over the second half of the training data to ignore the possible large values when starting training.

Figure 3 reports the relative change in online LogLoss for the three algorithms compared (smaller value is better). It is seen that on the Criteo dataset, AdaDelay performs better than AsyncAdaGrad, though AdaptiveRevision is slightly better than both. However, for the larger CTR2 dataset, both AdaDelay and AsyncAdaGrad are substantially better than AdaptiveRevision. The reason why it differs from [12] is probably due to the datasets we used are 1000 times larger than the ones reported by [12], and we evaluated the algorithms in a distributed environment rather than a simulated setting where a large minibatch size



(a) Criteo



(b) CTR2

Figure 4: Relative test AUC (higher is better) as function of maximal delays.

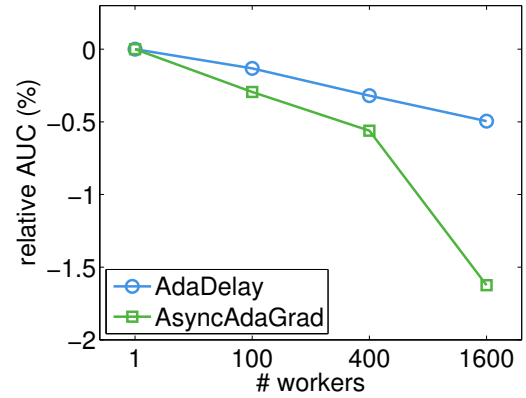
is necessary for the former. However, as reported [12], we also observed that AdaptiveRevision’s best learning rate is insensitive to the number of workers.

AdaDelay seems to have a tiny edge over AsyncAdaGrad, and as predicted by our theory, this edge grows much bigger when there are large delays (e.g., due to stragglers)—we report on this in greater detail in Section 4.4.1.

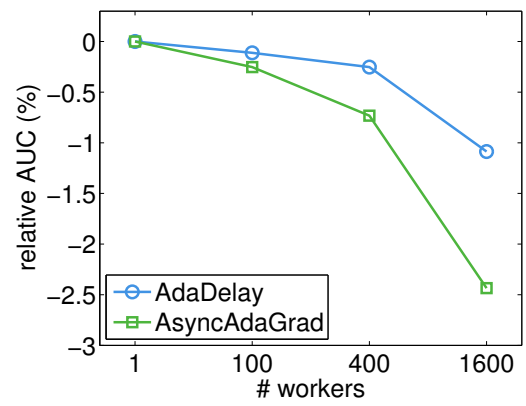
Besides the LogLoss, AUC is another important merit for computational advertising, which measures the ranking ability of the model and often 1% difference is significant for click-through rate estimation. We made a separate validation dataset for calculating the AUC, and shown the results on Figure 4. As can be seen, the test AUC results are consistent with the online LogLoss.

4.4.1 Stragglers

Previous experiments indicate that AdaDelay improves upon AsyncAdaGrad when a large number of workers (greater than 400) is used, which means the



(a) Criteo



(b) CTR2

Figure 5: Relative test AUC (higher is better) as function of maximal delays with the existence of stragglers.

delay adaptive learning rate takes effect when the delay can be large. To further investigate this phenomenon, we simulated an overloaded cluster where several stragglers may produce large delays; we do this by slowing down half of the workers by a random factor in [1, 4] when computing gradients. The test AUC are shown in Figure 5⁵. As can be seen, AdaDelay consistently outperforms AsyncAdaGrad, which shows that adaptive modeling of the actual delay is better than using a constant worst case delay when the variance of the delays is large.

4.4.2 Scalability

Finally we report the system performance. We first present the speedup from 1 machine to 16 machines, where each machine runs 100 workers. We observed a near linear speedup of AdaDelay, which is shown in Figure 6. The main reason is due to the asynchronous updating which removes the dependencies between worker nodes. In addition, using multiple

⁵As before, the results on online LogLoss are similar to the test AUC and therefore omitted.

	AdaDelay	AsyncAdaGrad	AdaptiveRevision
Criteo	24 GB	24 GB	55 GB
CTR2	97 GB	97 GB	200 GB

Table 2: Total memory used by server nodes.

workers within a machine can fully utilize the computational resources by hiding the overhead of reading data and communicating the parameters. The results of AsyncAdaGrad and AdaptiveRevision are similar to AdaDelay because their computational workloads are identical except for parameter updating, which affects the overall system performance little.

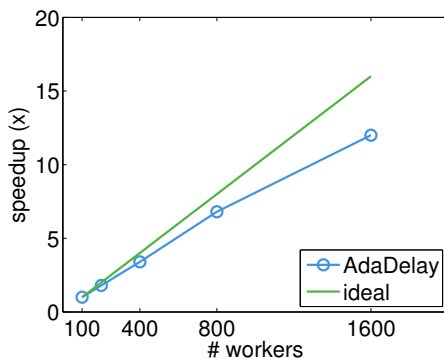


Figure 6: The speedup of AdaDelay. The results of AsyncAdaGrad and AdaptiveRevision are almost identical to AdaDelay and therefore omitted.

In the parameter server framework, worker nodes only need to cache one or a few data minibatches. Most memory is used by the server nodes to store the model. We summarize the server memory usage for the three algorithms compared in Table 2.

As expected, AdaDelay and AsyncAdaGrad have similar memory consumption because the extra storage needed by AdaDelay to track and compute the incurred delays τ_t is tiny. However AdaptiveRevision doubles memory usage, because of the extra entries that it needs for each feature, and because of the cached delayed gradient g^{bak} .

5 Conclusions

In real distributed computing environment, there are multiple factors contributing to delay, such as the CPU speed, I/O of disk, and network throughput. With the inevitable and sometimes unpredictable phenomenon of delay, we considered distributed convex optimization by developing and analyzing *AdaDelay*, an asynchronous SGD method that tolerates stale gradients.

A key component of our work that differs from existing approaches is the use of (server-side) updates sensitive

to the actual delay observed in the network. This allows us to use larger stepsizes initially, which can lead to more rapid initial convergence, and stronger ability to adapt to the environment. We discussed details of two different realistic delay models: (i) uniform (more generally, bounded support) delays, and (ii) scaled delays with constant first and second moments but not necessarily bounded support. Under both models, we obtain theoretically optimal convergence rates.

Adapting more closely to observed delays and incorporating server-side delay sensitive gradient aggregation that combines the benefits of the adaptive revision framework [12] with our delayed gradient methods is an important future direction. Extension of our analysis to handle constrained convex optimization problems without projection oracles is an important part of future work. Finally, how to apply our techniques to large-scale nonconvex problems such as matrix factorization and deep neural networks is an important direction worth studying.

Acknowledgments

SS and AS were partially supported by NSF grant IIS-1409802.

References

- [1] A. Agarwal and J. C. Duchi. Distributed delayed stochastic optimization. In *Advances in Neural Information Processing Systems*, pages 873–881, 2011.
- [2] D. Bertsekas and J. Tsitsiklis. *Parallel and Distributed Computation: Numerical Methods*. Prentice-Hall, 1989.
- [3] D. P. Bertsekas. Incremental gradient, subgradient, and proximal methods for convex optimization: A survey. *Optimization for Machine Learning*, 2010:1–38, 2011.
- [4] J. Duchi, M. I. Jordan, and B. McMahan. Estimation, optimization, and parallelism when data is sparse. In *NIPS 26*, pages 2832–2840, 2013.
- [5] J. C. Duchi, E. Hazan, and Y. Singer. Adaptive subgradient methods for online learning and stochastic optimization. *Journal of Machine Learning Research*, 12:2121–2159, 2011.
- [6] J. C. Duchi, A. Agarwal, and M. J. Wainwright. Dual averaging for distributed optimization: convergence analysis and network scaling. *Automatic Control, IEEE Transactions on*, 57(3):592–606, 2012.
- [7] S. Ghadimi and G. Lan. Optimal stochastic approximation algorithms for strongly convex stochastic composite optimization i: A generic

- algorithmic framework. *SIAM Journal on Optimization*, 22(4):1469–1492, 2012.
- [8] M. Jaggi, V. Smith, M. Takáč, J. Terhorst, S. Krishnan, T. Hofmann, and M. I. Jordan. Communication-efficient distributed dual coordinate ascent. In *NIPS*, pages 3068–3076, 2014.
- [9] J. Langford, A. J. Smola, and M. Zinkevich. Slow learners are fast. In *Neural Information Processing Systems*, 2009. URL <http://arxiv.org/abs/0911.0491>.
- [10] M. Li, D. G. Andersen, J. W. Park, A. J. Smola, A. Ahmed, V. Josifovski, J. Long, E. J. Shekita, and B.-Y. Su. Scaling distributed machine learning with the parameter server. In *Operating Systems Design and Implementation (OSDI)*, 2014.
- [11] M. Li, D. G. Andersen, A. J. Smola, and K. Yu. Communication efficient distributed machine learning with the parameter server. In *NIPS 27*, pages 19–27, 2014.
- [12] B. McMahan and M. Streeter. Delay-tolerant algorithms for asynchronous distributed online learning. In *NIPS 27*, pages 2915–2923, 2014.
- [13] A. Nedić, D. P. Bertsekas, and V. S. Borkar. Distributed asynchronous incremental subgradient methods. *Studies in Computational Mathematics*, 8:381–407, 2001.
- [14] A. Nemirovski, A. Juditsky, G. Lan, and A. Shapiro. Robust stochastic approximation approach to stochastic programming. *SIAM Journal on Optimization*, 19(4):1574–1609, 2009.
- [15] F. Niu, B. Recht, C. Re, and S. J. Wright. Hogwild: A lock-free approach to parallelizing stochastic gradient descent. In *NIPS*, pages 693–701, 2011.
- [16] S. S. Ram, A. Nedić, and V. V. Veeravalli. Distributed stochastic subgradient projection algorithms for convex optimization. *Journal of optimization theory and applications*, 147(3):516–545, 2010.
- [17] H. Robbins and S. Monro. A stochastic approximation method. *Annals of Mathematical Statistics*, 22:400–407, 1951.
- [18] M. W. Schmidt, N. L. Roux, and F. R. Bach. Minimizing finite sums with the stochastic average gradient. *CoRR*, abs/1309.2388, 2013.
- [19] O. Shamir and N. Srebro. Distributed stochastic optimization and learning. In *Proceedings of the 52nd Annual Allerton Conference on Communication, Control, and Computing*, 2014.
- [20] A. Shapiro, D. Dentcheva, and A. Ruszczyński. *Lectures on stochastic programming: modeling and theory*, volume 16. SIAM, 2014.
- [21] N. Srebro and A. Tewari. Stochastic Optimization for Machine Learning. ICML 2010 Tutorial, 2010.
- [22] T. Yang. Trading computation for communication: Distributed stochastic dual coordinate ascent. In *NIPS*, pages 629–637, 2013.
- [23] A. W. Yu, Q. Lin, and T. Yang. Doubly stochastic primal-dual coordinate method for regularized empirical risk minimization with factorized data. *CoRR*, abs/1508.03390, 2015.
- [24] Y. Zhang and L. Xiao. Stochastic primal-dual coordinate method for regularized empirical risk minimization. In *ICML*, pages 353–361, 2015.
- [25] Y. Zhang, J. C. Duchi, and M. J. Wainwright. Communication-efficient algorithms for statistical optimization. *Journal of Machine Learning Research*, 14(1):3321–3363, 2013.
- [26] M. Zinkevich, M. Weimer, A. J. Smola, and L. Li. Parallelized stochastic gradient descent. In *NIPS*, pages 2595–2603, 2010.