

Learning from Data with Heterogeneous Noise using SGD

Shuang Song* Kamalika Chaudhuri† Anand D. Sarwate‡

Friday 19th December, 2014

Abstract

We consider learning from data of variable quality that may be obtained from different heterogeneous sources. Addressing learning from heterogeneous data in its full generality is a challenging problem. In this paper, we adopt instead a model in which data is observed through heterogeneous noise, where the noise level reflects the quality of the data source. We study how to use stochastic gradient algorithms to learn in this model. Our study is motivated by two concrete examples where this problem arises naturally: learning with local differential privacy based on data from multiple sources with different privacy requirements, and learning from data with labels of variable quality.

The main contribution of this paper is to identify how heterogeneous noise impacts performance. We show that given two datasets with heterogeneous noise, the order in which to use them in standard SGD depends on the learning rate. We propose a method for changing the learning rate as a function of the heterogeneity, and prove new regret bounds for our method in two cases of interest. Experiments on real data show that our method performs better than using a single learning rate and using only the less noisy of the two datasets when the noise level is low to moderate.

1 Introduction

Modern large-scale machine learning systems often integrate data from several different sources. In many cases, these sources provide data of a similar type (i.e. with the same features) but collected under different circumstances. For example, patient records from different studies of a particular drug may be combined to perform a more comprehensive analysis, or a collection of images with annotations from experts as well as non-experts may be combined to learn a predictor. In particular, data from different sources may be of varying *quality*. In this paper we adopt a model in which data is observed through heterogeneous noise, where the noise level reflects the quality of the data source. We study how to use stochastic gradient algorithms to learn from data of heterogeneous quality.

In full generality, learning from heterogeneous data is essentially the problem of domain adaptation – a challenge for which good and complete solutions are difficult to obtain. Instead, we focus on the special case of heterogeneous noise and show how to use information about the data quality to improve the performance of learning algorithms which ignore this information.

Two concrete instances of this problem motivate our study: locally differentially private learning from multiple sites, and classification with random label noise. Differential privacy (Dwork et al.,

*Computer Science and Engineering Dept., University of California, San Diego, shs037@eng.ucsd.edu

†Computer Science and Engineering Dept., University of California, San Diego, kamalika@cs.ucsd.edu

‡Electrical and Computer Engineering Dept., Rutgers University, asarwate@ece.rutgers.edu

2006b,a) is a privacy model that has received significant attention in machine-learning and data-mining applications. A variant of differential privacy is *local privacy* – the learner can only access the data via noisy estimates, where the noise guarantees privacy (Duchi et al., 2012, 2013). In many applications, we are required to learn from sensitive data collected from individuals with heterogeneous privacy preferences, or from multiple sites with different privacy requirements; this results in the heterogeneity of noise added to ensure privacy. Under random classification noise (RCN) (Kearns, 1998), labels are randomly flipped before being presented to the algorithm. The heterogeneity in the noise addition comes from combining labels of variable quality – such as labels assigned by domain experts with those assigned by a crowd.

To our knowledge, Crammer et al. (2006) were the first to provide a theoretical study of how to learn classifiers from data of variable quality. In their formulation, like ours, data is observed through heterogeneous noise. Given data with known noise levels, their study focuses on finding an optimal ordering of the data and a stopping rule without any constraint on the computational complexity. We instead shift our attention to studying *computationally efficient strategies* for learning classifiers from data of variable quality.

We propose a model for variable data quality which is natural in the context of large-scale learning using stochastic gradient descent (SGD) and its variants (Bottou, 2010; Bekkerman et al., 2011). We assume that the training data are accessed through an oracle which provides an unbiased but noisy estimate of the gradient of the objective. The noise comes from two sources: the random sampling of a data point, and additional noise due to the data quality. Our two motivating applications – learning with local differential privacy and learning from data of variable quality – can both be modeled as solving a regularized convex optimization problem using SGD. Learning from data with heterogeneous noise in this framework thus reduces to running SGD with noisy gradient estimates, where the magnitude of the added noise varies across iterations.

Main results. In this paper we study noisy stochastic gradient methods when learning from multiple data sets with different noise levels. For simplicity we consider the case where there are two data sets, which we call *Clean* and *Noisy*. We process these data sets sequentially using SGD with learning rate $\mathcal{O}(1/t)$. In a future full version of this work we also analyze averaged gradient descent (AGD) with learning rate $\mathcal{O}(1/\sqrt{t})$. We address some basic questions in this setup:

In what order should we process the data? Suppose we use standard SGD on the union of *Clean* and *Noisy*. We show theoretically and empirically that the order in which we should process the datasets to get good performance depends on the learning rate of the algorithm: in some cases we should use the order (*Clean*, *Noisy*) and in others (*Noisy*, *Clean*).

Can we use knowledge of the noise rates? We show that using separate learning rates that depend on the noise levels for the clean and noisy datasets improves the performance of SGD. We provide a heuristic for choosing these rates by optimizing an upper bound on the error for SGD that depends on the ratio of the noise levels. We analytically quantify the performance of our algorithm in two regimes of interest. For moderate noise levels, we demonstrate empirically that our algorithm outperforms using a single learning rate and using clean data only.

Does using noisy data always help? The work of Crammer et al. (2006) suggests that if the noise level of noisy data is above some threshold, then noisy data will not help. Moreover, when the noise levels are very high, our heuristic does not always empirically outperform simply using the clean data. On the other hand, our theoretical results suggest that changing the learning rate can make noisy data useful. How do we resolve this apparent contradiction?

We perform an empirical study to address this question. Our experiments demonstrate that very often, there exists a learning rate at which noisy data helps; however, because the actual noise level may be far from the upper bound used in our algorithm, our optimization may not choose the best learning rate for every data set. We demonstrate that by adjusting the learning rate we can

still take advantage of noisy data.

For simplicity we, like previous work [Crammer et al. \(2006\)](#), assume that the algorithms know the noise levels exactly. However, our algorithms can still be applied in the presence of approximate knowledge of the noise levels, and our result on the optimal data order only needs to know which dataset has more noise.

Related Work. There has been significant work on the convergence of SGD assuming analytic properties of the objective function, such as strong convexity and smoothness. When the objective function is λ -strongly convex, the learning rate used for SGD is $\mathcal{O}(1/\lambda t)$ ([Nemirovsky and Yudin, 1983](#); [Agarwal et al., 2009](#); [Rakhlin et al., 2012](#); [Moulines and Bach, 2011](#)), which leads to a regret of $\mathcal{O}(1/\lambda^2 t)$ for smooth objectives. For non-smooth objectives, SGD with learning rate $\mathcal{O}(1/\lambda t)$ followed by some form of averaging of the iterates achieves $\mathcal{O}(1/\lambda t)$ ([Nesterov and Vial, 2008](#); [Nemirovski et al., 2009](#); [Shalev-Shwartz et al., 2009](#); [Xiao, 2010](#); [Duchi and Singer, 2009](#)).

There is also a body of literature on differentially private classification by regularized convex optimization in the batch ([Chaudhuri et al., 2011](#); [Rubinstein et al., 2012](#); [Kifer et al., 2012](#)) as well as the online ([Jain et al., 2012](#)) setting. In this paper, we consider classification with *local differential privacy* ([Wasserman and Zhou, 2010](#); [Duchi et al., 2012](#)), a stronger form of privacy than ordinary differential privacy. [Duchi et al. \(2012\)](#) propose learning a classifier with local differential privacy using SGD, and [Song et al. \(2013\)](#) show empirically that using mini-batches significantly improves the performance of differentially private SGD. Recent work by [Bassily et al. \(2014\)](#) provides an improved privacy analysis for non-local privacy. Our work is an extension of these papers to heterogeneous privacy requirements.

[Crammer et al. \(2006\)](#) study classification when the labels in each data set are corrupted by RCN of different rates. Assuming the classifier minimizing the empirical 0/1 classification error can always be found, they propose a general theoretical procedure that processes the datasets in increasing order of noise, and determines when to stop using more data. In contrast, our noise model is more general and we provide a polynomial time algorithm for learning. Our results imply that in some cases the algorithm should process the noisy data first, and finally, our algorithm uses all the data.

2 The Model

We consider linear classification in the presence of noise. We are given T labelled examples $(x_1, y_1), \dots, (x_T, y_T)$, where $x_i \in \mathbb{R}^d$, and $y_i \in \{-1, 1\}$ and our goal is to find a hyperplane w that largely separates the examples labeled 1 from those labeled -1 . A standard solution is via the following regularized convex optimization problem:

$$w^* = \underset{w \in \mathcal{W}}{\operatorname{argmin}} f(w) := \frac{\lambda}{2} \|w\|^2 + \frac{1}{T} \sum_{i=1}^T \ell(w, x_i, y_i). \quad (1)$$

Here ℓ is a convex loss function, and $\frac{\lambda}{2} \|w\|^2$ is a regularization term. Popular choices for ℓ include the logistic loss $\ell(w, x, y) = \log(1 + e^{-yw^\top x})$ and the hinge loss $\ell(w, x, y) = \max(0, 1 - yw^\top x)$.

Stochastic Gradient Descent (SGD) is a popular approach to solving (1): starting with an initial w_1 , at step t , SGD updates w_{t+1} using the point (x_t, y_t) as follows:

$$w_{t+1} = \Pi_{\mathcal{W}}(w_t - \eta_t(\lambda w_t + \nabla \ell(w_t, x_t, y_t))). \quad (2)$$

Here Π is a projection operator onto the convex feasible set \mathcal{W} , typically set to $\{w : \|w\|_2 \leq 1/\lambda\}$ and η_t is a learning rate (or step size) which specifies how fast w_t changes. A common choice for the learning rate for the case when $\lambda > 0$ is c/t , where $c = \Theta(1/\lambda)$.

2.1 The Heterogeneous Noise Model

We propose an abstract model for heterogeneous noise that can be specialized to two important scenarios: differentially private learning, and random classification noise. By heterogeneous noise we mean that the distribution of the noise can depend on the data points themselves. More formally, we assume that the learning algorithm may only access the labeled data through an oracle \mathcal{G} which, given a $w \in \mathbb{R}^d$, draws a fresh independent sample (x, y) from the underlying data distribution, and returns an unbiased noisy gradient of the objective function $\nabla f(w)$, based on the example (x, y) :

$$\mathbb{E}[\mathcal{G}(w)] = \lambda w + \nabla \ell(w, x, y), \quad \mathbb{E}[\|\mathcal{G}(w)\|^2] \leq \Gamma^2. \quad (3)$$

The precise manner in which $\mathcal{G}(w)$ is generated depends on the application. Define the *noise level* for the oracle \mathcal{G} as the constant Γ in (3); larger Γ means more noisy data. Finally, to model finite training datasets, we assume that an oracle \mathcal{G} may be called only a limited number of times.

Observe that in this noise model, we can easily use the noisy gradient returned by \mathcal{G} to perform SGD. The update rule becomes:

$$w_{t+1} = \Pi_{\mathcal{W}}(w_t - \eta_t \mathcal{G}(w_t)). \quad (4)$$

The SGD estimate is w_{t+1} .

In practice, we can implement an oracle such as \mathcal{G} based on a finite labelled training set D as follows. We apply a random permutation on the samples in D , and at each invocation, compute a noisy gradient based on the next sample in the permutation. The number of calls to the oracle is limited to $|D|$. If the samples in D are drawn iid from the underlying data distribution, and if any extraneous noise added to the gradient at each iteration is unbiased and drawn independently, then this process will implement the oracle correctly.

To model heterogeneous noise, we assume that we have access to two oracles \mathcal{G}_1 and \mathcal{G}_2 implemented based on datasets D_1 and D_2 , which can be called at most $|D_1|$ and $|D_2|$ times respectively. For $j = 1, 2$, the noise level of oracle \mathcal{G}_j is Γ_j , and the values of Γ_1 and Γ_2 are known to the algorithm. In some practical situations, Γ_1 and Γ_2 will not be known exactly; however, our algorithm in Section 4 also applies when approximate noise levels are known, and our algorithm in Section 3 applies even when only the relative noise levels are known.

2.1.1 Local Differential Privacy

Local differential privacy (Wasserman and Zhou, 2010; Duchi et al., 2012; Kasiviswanathan et al., 2008) is a strong notion of privacy motivated by differential privacy (Dwork et al., 2006b). An untrusted algorithm is allowed to access a perturbed version of a sensitive dataset through a sanitization interface, and must use this perturbed data to perform some estimation. The amount of perturbation is controlled by a parameter ϵ , which measures the privacy risk.

Definition 1 (Local Differential Privacy). *Let $D = (X_1, \dots, X_n)$ be a sensitive dataset where each $X_i \in \mathcal{D}$ corresponds to data about individual i . A randomized sanitization mechanism M which outputs a disguised version (U_1, \dots, U_n) of D is said to provide ϵ -local differential privacy to individual i , if for all $x, x' \in \mathcal{D}$ and for all $S \subseteq \mathcal{S}$,*

$$\Pr(U_i \in S | X_i = x) \leq e^\epsilon \Pr(U_i \in S | X_i = x'). \quad (5)$$

Here the probability is taken over the randomization in the sanitization mechanism, and ϵ is a parameter that measures privacy risk where smaller ϵ means less privacy risk.

Consider learning a linear classifier from a sensitive labelled dataset while ensuring local privacy of the participants. This problem can be expressed in our noise model by setting the sanitization mechanism as the oracle. Given a privacy risk ϵ , for $w \in \mathbb{R}^d$, the oracle \mathcal{G}^{DP} draws a random labelled sample (x, y) from the underlying data distribution, and returns the noisy gradient of the objective function at w computed based on (x, y) as

$$\mathcal{G}^{\text{DP}}(w) = \lambda w + \nabla \ell(w, x, y) + Z, \quad (6)$$

where Z is independent random noise drawn from the density: $\rho(z) \propto e^{-(\epsilon/2)\|z\|}$.

Duchi et al. (2012) showed that this mechanism provides ϵ -local privacy assuming analytic conditions on the loss function, bounded data, and that the oracle generates a fresh random sample at each invocation. The following result shows how to set the parameters to fit in our heterogeneous noise model. The full proof is provided in Appendix A.2.

Theorem 1. *If $\|\nabla \ell(w, x, y)\| \leq 1$ for all w and (x, y) , then $\mathcal{G}^{\text{DP}}(w)$ is ϵ -local differentially private. Moreover, for any w such that $\|w\| \leq \frac{1}{\lambda}$, $\mathbb{E}[\mathcal{G}^{\text{DP}}(w)] = \lambda w + \nabla \mathbb{E}_{(x,y)}[\ell(w, x, y)]$, and*

$$\mathbb{E}[\|\mathcal{G}^{\text{DP}}(w)\|^2] \leq 4 + \frac{4(d^2 + d)}{\epsilon^2}.$$

Proof. (Sketch) The term 4 comes from upper bounding $\mathbb{E}[\|\lambda w + \nabla \ell(w, x, y)\|^2]$ by $\max_{w,x,y} \|\lambda w + \nabla \ell(w, x, y)\|^2$ using $\|w\| \leq 1/\lambda$ and $\|\nabla \ell(w, x, y)\| \leq 1$. The term $4(d^2 + d)/\epsilon^2$ comes from properties of the noise distribution. \square

In practice, we may wish to learn classifiers from multiple sensitive datasets with different privacy parameters. For example, suppose we wish to learn a classifier from sensitive patient records in two different hospitals holding data sets D_1 and D_2 , respectively. The hospitals have different privacy policies, and thus different privacy parameters ϵ_1 and ϵ_2 . This corresponds to a heterogeneous noise model in which we have two sanitizing oracles – $\mathcal{G}_1^{\text{DP}}$ and $\mathcal{G}_2^{\text{DP}}$. For $j = 1, 2$, $\mathcal{G}_j^{\text{DP}}$ implements a differentially private oracle with privacy parameter ϵ_j based on dataset D_j and may be called at most $|D_j|$ times.

2.1.2 Random Classification Noise

In the random classification noise model of Kearns (1998), the learning algorithm is presented with labelled examples $(x_1, \tilde{y}_1), \dots, (x_T, \tilde{y}_T)$, where each $\tilde{y}_i \in \{-1, 1\}$ has been obtained by independently flipping the *true label* y_i with some probability σ . Natarajan et al. (2013) showed that solving

$$\underset{w}{\operatorname{argmin}} \frac{\lambda}{2} \|w\|^2 + \frac{1}{T} \sum_{i=1}^T \tilde{\ell}(w, x_i, \tilde{y}_i, \sigma) \quad (7)$$

yields a linear classifier from data with random classification noise, where $\tilde{\ell}$ is a surrogate loss function corresponding to a convex loss ℓ :

$$\tilde{\ell}(w, x, y, \sigma) = \frac{(1 - \sigma)\ell(w, x, y) - \sigma\ell(w, x, -y)}{1 - 2\sigma},$$

and σ is the probability that each label is flipped. This problem can be expressed in our noise model using an oracle \mathcal{G}^{RCN} which on input w draws a fresh labelled example (x, \tilde{y}) and returns

$$\mathcal{G}^{\text{RCN}}(w) = \lambda w + \nabla \tilde{\ell}(w, x, \tilde{y}, \sigma).$$

The SGD updates in (4) with respect to \mathcal{G}^{RCN} minimize (7). If $\|x\| \leq 1$ and $\|\nabla\ell(w, x, y)\| \leq 1$, we have $\mathbb{E}[\mathcal{G}^{\text{RCN}}(w)] = \lambda w + \nabla\ell(w, x, y)$ and $\mathbb{E}[\|\mathcal{G}^{\text{RCN}}(w)\|_2^2] \leq 3 + 1/(1 - 2\sigma)^2$, under the random classification noise assumption, so the oracle \mathcal{G}^{RCN} satisfies the conditions in (3) with $\Gamma^2 = 3 + 1/(1 - 2\sigma)^2$.

In practice, we may wish to learn classifiers from multiple datasets with different amounts of classification noise (Crammer et al., 2006); for example, we may have a small dataset D_1 labeled by domain experts, and a larger noisier dataset D_2 , labeled via crowdsourcing, with flip probabilities σ_1 and σ_2 . We model this scenario using two oracles – $\mathcal{G}_1^{\text{RCN}}$ and $\mathcal{G}_2^{\text{RCN}}$. For $j = 1, 2$, oracle $\mathcal{G}_j^{\text{RCN}}$ is implemented based on D_j and flip probability σ_j , and may be called at most $|D_j|$ times.

3 Data order depends on learning rate

Suppose we have two oracles \mathcal{G}_C (for “clean”) and \mathcal{G}_N (for “noisy”) implemented based on datasets D_C, D_N with noise levels Γ_C, Γ_N (where $\Gamma_C < \Gamma_N$) respectively. In which order should we query the oracle when using SGD? Perhaps surprisingly, it turns out that the answer depends on the learning rate. Below, we show a specific example of a convex optimization problem such that with $\eta_t = c/t$, the optimal ordering is to use \mathcal{G}_C first when $c \in (0, 1/\lambda)$, and the optimal ordering is to use \mathcal{G}_N first when $c > 1/\lambda$.

Let $|D_C| + |D_N| = T$ and consider the convex optimization problem:

$$\min_{w \in \mathcal{W}} \frac{\lambda}{2} \|w\|^2 - \frac{1}{T} \sum_{i=1}^T y_i w^\top x_i, \quad (8)$$

where the points $\{(x_i, y_i)\}$ are drawn from the underlying distribution by \mathcal{G}_C or \mathcal{G}_N . Suppose $\mathcal{G}(w) = \lambda w - yx + Z$ where Z is an independent noise vector such that $\mathbb{E}[Z] = 0$, $\mathbb{E}[\|Z\|^2] = V_C^2$ if \mathcal{G} is \mathcal{G}_C , and $\mathbb{E}[\|Z\|^2] = V_N^2$ if \mathcal{G} is \mathcal{G}_N with $V_N^2 \geq V_C^2$.

For our example, we consider the following three variants of SGD: CF and NF for “clean first” and “noisy first” and AO for an “arbitrary ordering”:

1. CF: For $t \leq |D_C|$, query \mathcal{G}_C in the SGD update (4). For $t > |D_C|$, query \mathcal{G}_N .
2. NF: For $t \leq |D_N|$, query \mathcal{G}_N in the SGD update (4). For $t > |D_N|$, query \mathcal{G}_C .
3. AO: Let S be an arbitrary sequence of length T consisting of $|D_C|$ C’s and $|D_N|$ N’s. In the SGD update (4) in round t , if the t -th element S_t of S is C, then query \mathcal{G}_C ; else, query \mathcal{G}_N .

In order to isolate the effect of the noise, we consider two additional oracles \mathcal{G}'_C and \mathcal{G}'_N ; the oracle \mathcal{G}'_C (resp. \mathcal{G}'_N) is implemented based on the dataset D_C (resp. D_N), and iterates over D_C (resp. D_N) in exactly the same order as \mathcal{G}_C (resp. \mathcal{G}_N); the only difference is that for \mathcal{G}'_C (resp. \mathcal{G}'_N), no extra noise is added to the gradient (that is, $Z = 0$). The main result of this section is stated in Theorem 2.

Theorem 2. *Let $\{w_t^{\text{CF}}\}$, $\{w_t^{\text{NF}}\}$ and $\{w_t^{\text{AO}}\}$ be the sequences of updates obtained by running SGD for objective function (8) under CF, NF and AO respectively, and let $\{v_t^{\text{CF}}\}$, $\{v_t^{\text{NF}}\}$ and $\{v_t^{\text{AO}}\}$ be the sequences of updates under CF, NF and AO with calls to \mathcal{G}_C and \mathcal{G}_N replaced by calls to \mathcal{G}'_C and \mathcal{G}'_N . Let $T = |D_C| + |D_N|$.*

1. *If the learning rate $\eta_t = c/t$ where $c \in (0, 1/\lambda)$, then*

$$\mathbb{E} \left[\|v_{T+1}^{\text{CF}} - w_{T+1}^{\text{CF}}\|^2 \right] \leq \mathbb{E} \left[\|v_{T+1}^{\text{AO}} - w_{T+1}^{\text{AO}}\|^2 \right].$$

2. If the learning rate $\eta_t = c/t$ where $c > 1/\lambda$, then

$$\mathbb{E} \left[\|v_{T+1}^{\text{NF}} - w_{T+1}^{\text{NF}}\|^2 \right] \leq \mathbb{E} \left[\|v_{T+1}^{\text{AO}} - w_{T+1}^{\text{AO}}\|^2 \right].$$

Proof. Let the superscripts CF, NF and AO indicate the iterates for the CF, NF and AO algorithms. Let w_1 denote the initial point of the optimization. Let $(x_t^{\text{O}}, y_t^{\text{O}})$ be the data used under order $\text{O} = \text{CF}, \text{NF}$ or AO to update w at time t , Z_t^{O} be the noise added to the exact gradient by \mathcal{G}_{C} or \mathcal{G}_{N} , depending on which oracle is used by O at t and w_t^{O} be the w obtained under order O at time t . Then by expanding the expression for w_t in terms of the gradients, we have

$$w_{T+1}^{\text{O}} = w_1 \prod_{i=1}^T (1 - \eta_i \lambda) - \sum_{t=1}^T \eta_t \left(\prod_{s=t+1}^T (1 - \eta_s \lambda) \right) (y_t^{\text{O}} x_t^{\text{O}} + Z_t^{\text{O}}). \quad (9)$$

Similarly, if $v_1 = w_1$, we have

$$v_{T+1}^{\text{O}} = w_1 \prod_{i=1}^T (1 - \eta_i \lambda) - \sum_{t=1}^T \eta_t \left(\prod_{s=t+1}^T (1 - \eta_s \lambda) \right) y_t^{\text{O}} x_t^{\text{O}}. \quad (10)$$

Define

$$\Delta_t = \eta_t \prod_{s=t+1}^{\top} (1 - \eta_s \lambda).$$

Taking the expected squared difference between (9) from (10), we obtain

$$\begin{aligned} \mathbb{E} \left[\|v_{T+1}^{\text{O}} - w_{T+1}^{\text{O}}\|^2 \right] &= \mathbb{E} \left[\left\| \sum_{t=1}^T \eta_t \left(\prod_{s=t+1}^T (1 - \eta_s \lambda) \right) Z_t^{\text{O}} \right\|^2 \right] \\ &= \mathbb{E} \left[\left\| \sum_{t=1}^T \Delta_t Z_t^{\text{O}} \right\|^2 \right] \\ &= \sum_{t=1}^T \Delta_t^2 \mathbb{E} \left[\|Z_t^{\text{O}}\|^2 \right], \end{aligned} \quad (11)$$

where the second step follows because the Z_i^{O} 's are independent.

If $\eta_t = c/t$, then

$$\Delta_t = \frac{c}{t} \prod_{s=t+1}^{\top} \left(1 - \frac{c\lambda}{s} \right).$$

Therefore

$$\frac{\Delta_{t+1}^2}{\Delta_t^2} = \left(\frac{\frac{c}{t+1} \prod_{s=t+2}^{\top} \left(1 - \frac{c\lambda}{s} \right)}{\frac{c}{t} \prod_{s=t+1}^{\top} \left(1 - \frac{c\lambda}{s} \right)} \right)^2 = \left(\frac{t}{(t+1) \left(1 - \frac{c\lambda}{t+1} \right)} \right)^2 = \left(\frac{1}{1 + \frac{1-c\lambda}{t}} \right)^2,$$

which is smaller than 1 if $c < 1/\lambda$, equal to 1 if $c = 1/\lambda$, and greater than 1 if $c > 1/\lambda$. Therefore Δ_t is decreasing if $c < 1/\lambda$ and is increasing if $c > 1/\lambda$.

If Δ_t is decreasing, then (11) is minimized if $\mathbb{E} [\|Z_t^{\text{O}}\|^2]$ is increasing; if Δ_t is increasing, then (11)

is minimized if $\mathbb{E} [\|Z_t^{\text{O}}\|^2]$ is decreasing; and if Δ_t is constant, then (11) is the same under any order of $\mathbb{E} [\|Z_t^{\text{O}}\|^2]$.

Therefore for $c < 1/\lambda$,

$$\mathbb{E} \left[\left\| v_{T+1}^{\text{CF}} - w_{T+1}^{\text{CF}} \right\|^2 \right] \leq \mathbb{E} \left[\left\| v_{T+1}^{\text{AO}} - w_{T+1}^{\text{AO}} \right\|^2 \right] \leq \mathbb{E} \left[\left\| v_{T+1}^{\text{NF}} - w_{T+1}^{\text{NF}} \right\|^2 \right].$$

For $c = 1/\lambda$,

$$\mathbb{E} \left[\left\| v_{T+1}^{\text{CF}} - w_{T+1}^{\text{CF}} \right\|^2 \right] = \mathbb{E} \left[\left\| v_{T+1}^{\text{AO}} - w_{T+1}^{\text{AO}} \right\|^2 \right] = \mathbb{E} \left[\left\| v_{T+1}^{\text{NF}} - w_{T+1}^{\text{NF}} \right\|^2 \right].$$

For $c > 1/\lambda$,

$$\mathbb{E} \left[\left\| v_{T+1}^{\text{CF}} - w_{T+1}^{\text{CF}} \right\|^2 \right] \geq \mathbb{E} \left[\left\| v_{T+1}^{\text{AO}} - w_{T+1}^{\text{AO}} \right\|^2 \right] \geq \mathbb{E} \left[\left\| v_{T+1}^{\text{NF}} - w_{T+1}^{\text{NF}} \right\|^2 \right].$$

□

This result says that arbitrary ordering of the data is worse than sequentially processing one data set after the other except in the special case where $c = 1/\lambda$. If the step size is small ($c < 1/\lambda$), the SGD should use the clean data first to more aggressively proceed towards the optimum. If the step size is larger ($c > 1/\lambda$), then SGD should reserve the clean data for refining the initial estimates given by processing the noisy data.

4 Adapting the learning rate to the noise level

We now investigate whether the performance of SGD can be improved by using different learning rates for oracles with different noise levels. Suppose we have oracles \mathcal{G}_1 and \mathcal{G}_2 with noise levels Γ_1 and Γ_2 that are implemented based on two datasets D_1 and D_2 . Unlike the previous section, we do not assume any relation between Γ_1 and Γ_2 – we analyze the error for using oracle \mathcal{G}_1 followed by \mathcal{G}_2 in terms of Γ_1 and Γ_2 to choose a data order. Let $T = |D_1| + |D_2|$. Let $\beta_1 = \frac{|D_1|}{T}$ and $\beta_2 = 1 - \beta_1 = \frac{|D_2|}{T}$ be the fraction of the data coming from \mathcal{G}_1 and \mathcal{G}_2 , respectively. We adapt the gradient updates in (4) to heterogeneous noise by choosing the learning rate η_t as a function of the noise level. Algorithm 1 shows a modified SGD for heterogeneous learning rates.

Algorithm 1 SGD with varying learning rate

- 1: **Inputs:** Oracles $\mathcal{G}_1, \mathcal{G}_2$ implemented by data sets D_1, D_2 . Learning rates c_1 and c_2 .
 - 2: Set $w_1 = 0$.
 - 3: **for** $t = 1, 2, \dots, |D_1|$ **do**
 - 4: $w_{t+1} = \Pi_{\mathcal{W}} \left(w_t - \frac{c_1}{t} \mathcal{G}_1(w_t) \right)$
 - 5: **end for**
 - 6: **for** $t = |D_1| + 1, |D_1| + 2, \dots, |D_1| + |D_2|$ **do**
 - 7: $w_{t+1} = \Pi_{\mathcal{W}} \left(w_t - \frac{c_2}{t} \mathcal{G}_2(w_t) \right)$
 - 8: **end for**
 - 9: **return** $w_{|D_1|+|D_2|+1}$.
-

Consider SGD with learning rate $\eta_t = c_1/t$ while querying \mathcal{G}_1 and with $\eta_t = c_2/t$ while querying \mathcal{G}_2 in the update (4). We must choose an order in which to query \mathcal{G}_1 and \mathcal{G}_2 as well as the constants

c_1 and c_2 to get the best performance. We do this by minimizing an upper bound on the distance between the final iterate w_{T+1} and the optimal solution w^* to $\mathbb{E}[f(w)]$ where f is defined in (1), and the expectation is with respect to the data distribution and the gradient noise; the upper bound we choose is based on [Rakhlin et al. \(2012\)](#). Note that for smooth functions f , a bound on the distance $\|w_{T+1} - w^*\|$ automatically translates to a bound on the regret $f(w_{T+1}) - f(w^*)$.

Theorem 3 generalizes the results of [Rakhlin et al. \(2012\)](#) to our heterogeneous noise setting; the proof is in the supplement.

Theorem 3. *If $2\lambda c_1 > 1$ and if $2\lambda c_2 \neq 1$, and if we query \mathcal{G}_1 before \mathcal{G}_2 with learning rates c_1/t and c_2/t respectively, then the SGD algorithm satisfies*

$$\begin{aligned} \mathbb{E} [\|w_{T+1} - w^*\|^2] &\leq \frac{4\Gamma_1^2}{T} \cdot \frac{\beta_1^{2\lambda c_2 - 1} c_1^2}{2\lambda c_1 - 1} \\ &+ \frac{4\Gamma_2^2}{T} \cdot \frac{(1 - \beta_1^{2\lambda c_2 - 1}) c_2^2}{2\lambda c_2 - 1} + \mathcal{O}\left(\frac{1}{T^{\min(2, 2\lambda c_1)}}\right). \end{aligned} \quad (12)$$

Proof. (Sketch) Let $g(w)$ be the true gradient $\nabla f(w)$ and $\hat{g}(w)$ be the unbiased noisy gradient provided by the oracle \mathcal{G}_1 or \mathcal{G}_2 , whichever is queried.

By strong convexity of f , we have

$$\mathbb{E}_{1, \dots, t} [\|w_{t+1} - w^*\|^2] \leq (1 - 2\lambda\eta_t) \mathbb{E}_{1, \dots, t} [\|w_t - w^*\|^2] + \eta_t^2 \gamma_t^2.$$

Solving it inductively, with $\gamma_t = \Gamma_1, \eta_t = c_1/t$ for $t \leq \beta_1 T$ and $\gamma_t = \Gamma_2, \eta_t = c_2/t$ for $t > \beta_1 T$, we have

$$\begin{aligned} \mathbb{E}_{1, \dots, T} [\|w_{T+1} - w^*\|^2] &\leq \prod_{i=i_0}^{\beta_1 T} \left(1 - \frac{2\lambda c_1}{i}\right) \prod_{i=\beta_1 T+1}^T \left(1 - \frac{2\lambda c_2}{i}\right) \mathbb{E}_{1, \dots, T} [\|w_{i_0} - w^*\|^2] \\ &+ \Gamma_1^2 \prod_{i=\beta_1 T+1}^T \left(1 - \frac{2\lambda c_2}{i}\right) \sum_{i=i_0}^{\beta_1 T} \frac{c_1^2}{i^2} \prod_{j=i+1}^{\beta_1 T} \left(1 - \frac{2\lambda c_1}{j}\right) \\ &+ \Gamma_2^2 \sum_{i=\beta_1 T+1}^T \frac{c_2^2}{i^2} \prod_{j=i+1}^T \left(1 - \frac{2\lambda c_2}{j}\right), \end{aligned}$$

where i_0 is the smallest positive integer such that $2\lambda\eta_{i_0} < 1$, i.e., $i_0 = \lceil 2c_1\lambda \rceil$.

Then using $1 - x \leq e^{-x}$ and upper bounding each term using integrals we get the (12). \square

Two remarks are in order. First, observe that the first two terms in the right hand side dominate the other term. Second, our proof techniques for Theorem 3, adapted from [Rakhlin et al. \(2012\)](#), require that $2\lambda c_1 > 1$ in order to get a $O(1/T)$ rate of convergence; without this condition, the dependence on T is $\Omega(1/T)$.

4.1 Algorithm description

Our algorithm for selecting c_1 and c_2 is motivated by Theorem 3. We propose an algorithm that selects c_1 and c_2 by minimizing the quantity $B(c_1, c_2)$ which represents the highest order terms in Theorem 3:

$$B(c_1, c_2) = \frac{4\Gamma_1^2 \beta_1^{2\lambda c_2 - 1} c_1^2}{T(2\lambda c_1 - 1)} + \frac{4\Gamma_2^2 (1 - \beta_1^{2\lambda c_2 - 1}) c_2^2}{T(2\lambda c_2 - 1)}. \quad (13)$$

Given $\lambda, \Gamma_1, \Gamma_2$ and β_1 , we use c_1^* and c_2^* to denote the values of c_1 and c_2 that minimize $B(c_1, c_2)$. We can optimize for fixed c_2 with respect to c_1 by minimizing $\frac{c_1^2}{2\lambda c_1 - 1}$; this gives $c_1^* = 1/\lambda$, and $\frac{c_1^{*2}}{2\lambda c_1^* - 1} = 1/\lambda^2$, which is independent of β_1 or the noise levels Γ_1 and Γ_2 . Minimizing $B(c_1^*, c_2)$ with respect to c_2 can be now performed numerically to yield $c_2^* = \mathbf{argmin}_{c_2} B(c_1^*, c_2)$. This yields optimal values of c_1 and c_2 .

Now suppose we have two oracles $\mathcal{G}_C, \mathcal{G}_N$ with noise levels Γ_C and Γ_N that are implemented based on datasets D_C and D_N respectively. Let $\Gamma_C < \Gamma_N$, and let $\beta_C = \frac{|D_C|}{|D_C|+|D_N|}$ and $\beta_N = \frac{|D_N|}{|D_C|+|D_N|}$ be the fraction of the total data in each data set. Define the following functions:

$$H_{CN}(c) = \frac{4\Gamma_C^2 \beta_C^{2\lambda c - 1}}{\lambda^2} + \frac{4\Gamma_N^2 (1 - \beta_C^{2\lambda c - 1}) c^2}{2\lambda c - 1},$$

$$H_{NC}(c) = \frac{4\Gamma_N^2 \beta_N^{2\lambda c - 1}}{\lambda^2} + \frac{4\Gamma_C^2 (1 - \beta_N^{2\lambda c - 1}) c^2}{2\lambda c - 1}.$$

These represent the constant of the leading term in the upper bound in Theorem 3 for $(\mathcal{G}_1, \mathcal{G}_2) = (\mathcal{G}_C, \mathcal{G}_N)$ and $(\mathcal{G}_1, \mathcal{G}_2) = (\mathcal{G}_N, \mathcal{G}_C)$, respectively. Algorithm 2 repeats the process of choosing optimal c_1, c_2 with two orderings of the data – \mathcal{G}_C first and \mathcal{G}_N first – and selects the solution which provides the best bounds (according to the higher order terms of Theorem 3).

Algorithm 2 Selecting the Learning Rates

- 1: **Inputs:** Data sets D_C and D_N accessed through oracles \mathcal{G}_C and \mathcal{G}_N with noise levels Γ_C and Γ_N .
 - 2: Let $\beta_C = \frac{|D_C|}{|D_C|+|D_N|}$ and $\beta_N = \frac{|D_N|}{|D_C|+|D_N|}$.
 - 3: Calculate $c_{CN} = \mathbf{argmin}_c H_{CN}(c)$ and $c_{NC} = \mathbf{argmin}_c H_{NC}(c)$.
 - 4: **if** $H_{CN}(c_{CN}) \leq H_{NC}(c_{NC})$ **then**
 - 5: Run Algorithm 1 using oracles $(\mathcal{G}_C, \mathcal{G}_N)$, learning rates $c_1 = \frac{1}{\lambda}$ and $c_2 = c_{CN}$.
 - 6: **else**
 - 7: Run Algorithm 1 using oracles $(\mathcal{G}_N, \mathcal{G}_C)$, learning rates $c_1 = \frac{1}{\lambda}$ and $c_2 = c_{NC}$.
 - 8: **end if**
-

4.2 Regret Bounds

To provide a regret bound on the performance of SGD with two learning rates, we need to plug the optimal values of c_1 and c_2 into the right hand side of (13). Observe that as $c_1 = c_2$ and $c_2 = 0$ are feasible inputs to (13), our algorithm by construction has a superior regret bound than using a single learning rate only, or using clean data only.

Unfortunately, the value of c_2 that minimizes (13) does not have a closed form solution, and as such it is difficult to provide a general simplified regret bound that holds for all Γ_1, Γ_2 and β_1 . In this section, we consider two cases of interest, and derive simplified versions of the regret bound for SGD with two learning rates for these cases.

We consider the two data orders $(\Gamma_1, \Gamma_2) = (\Gamma_N, \Gamma_C)$ and $(\Gamma_1, \Gamma_2) = (\Gamma_C, \Gamma_N)$ in a scenario where $\Gamma_N/\Gamma_C \gg 1$ and both β_N and β_C are bounded away from 0 and 1. That is, the noisy data is much noisier. The following two lemmas provide upper and lower bounds on $B(c_1^*, c_2^*)$ in this setting.

Lemma 1. Suppose $(\Gamma_1, \Gamma_2) = (\Gamma_N, \Gamma_C)$ and $0 < \beta_N < 1$. Then for sufficiently large Γ_N/Γ_C , the optimal solution c_2^* to (13) satisfies

$$2c_2^*\lambda \in \left[1 + \frac{2\log(\Gamma_N/\Gamma_C) + \log\log(1/\beta_N)}{\log(1/\beta_N)}, 1 + \frac{2\log(4\Gamma_N/\Gamma_C) + \log\log(1/\beta_N)}{\log(1/\beta_N)} \right].$$

Moreover, $B(c_1^*, c_2^*)$ satisfies:

$$B(c_1^*, c_2^*) \geq \frac{4\Gamma_C^2(\log(\frac{\Gamma_N}{\Gamma_C}) + \frac{1}{2}\log\log(\frac{1}{\beta_N}))}{\lambda^2 T \log(\frac{1}{\beta_N})}$$

$$B(c_1^*, c_2^*) \leq \frac{4\Gamma_C^2}{\lambda^2 T} \left(4 + \frac{4 + 2\log(\frac{\Gamma_N}{\Gamma_C}) + \log\log(\frac{1}{\beta_N})}{\log(\frac{1}{\beta_N})} \right).$$

Proof. (Sketch) We prove that for any $2\lambda c_2 \leq 1 + \frac{2\log(\Gamma_N/\Gamma_C) + \log\log(1/\beta_N)}{\log(1/\beta_N)}$, $B(c_1^*, c_2)$ is decreasing with respect to c_2 when Γ_N/Γ_C is sufficiently large; and for any $2\lambda c_2 \geq 1 + \frac{2\log(4\Gamma_N/\Gamma_C) + \log\log(1/\beta_N)}{\log(1/\beta_N)}$, $B(c_1^*, c_2)$ is increasing when Γ_N/Γ_C is sufficiently large. Therefore the minimum of $B(c_1^*, c_2)$ is achieved when $2\lambda c_2$ is in between. \square

Observe that the regret bound grows logarithmically with Γ_N/Γ_C . Moreover, if we only used the cleaner data, then the regret bound would be $\frac{4\Gamma_C^2}{\lambda^2 \beta_C T}$, which is better, especially for large Γ_N/Γ_C . This means that using two learning rates with the noisy data first gives poor results at high noise levels.

Our second bound takes the opposite data order, processing the clean data first.

Lemma 2. Suppose $(\Gamma_1, \Gamma_2) = (\Gamma_C, \Gamma_N)$ and $0 < \beta_C < 1$. Let $\sigma = (\Gamma_N/\Gamma_C)^{-2}$. Then for sufficiently large Γ_N/Γ_C , the optimal solution c_2^* to (13) satisfies: $2c_2^*\lambda \in \left[\sigma, \frac{8}{\beta_C} \sigma \right]$. Moreover, $B(c_1^*, c_2^*)$ satisfies:

$$B(c_1^*, c_2^*) \geq \frac{4\Gamma_C^2}{\lambda^2 \beta_C T} \beta_C^{8\sigma/\beta_C}$$

$$B(c_1^*, c_2^*) \leq \frac{4\Gamma_C^2}{\lambda^2 \beta_C T} \beta_C^\sigma \left(1 + \sigma \frac{\log(1/\beta_C)}{4} \right).$$

Proof. (Sketch) Similar as the proof of Lemma 1, we prove that for any $2\lambda c_2 \leq \sigma$, $B(c_1^*, c_2)$ is decreasing with respect to c_2 ; and for any $2\lambda c_2 \geq \frac{8}{\beta_C} \sigma$, $B(c_1^*, c_2)$ is increasing when Γ_N/Γ_C is sufficiently large. Therefore the minimum of $B(c_1^*, c_2)$ is achieved when $2\lambda c_2$ is in between. \square

If we only used the clean dataset, then the regret bound would be $\frac{4\Gamma_C^2}{\lambda^2 \beta_C T}$, so Lemma 2 yields an improvement by a factor of $\beta_C^{(\Gamma_N/\Gamma_C)^{-2}} \left(1 + \left(\frac{\Gamma_N}{\Gamma_C} \right)^{-2} \frac{\log(1/\beta_C)}{4} \right)$. As $\beta_C < 1$, observe that this factor is always less than 1, and tends to 1 as Γ_N/Γ_C tends to infinity; therefore the difference between the regret bounds narrows as the noisy data grows noisier. We conclude that using two learning rates with clean data first gives a better regret bound than using only clean data or using two learning rates with noisy data first.

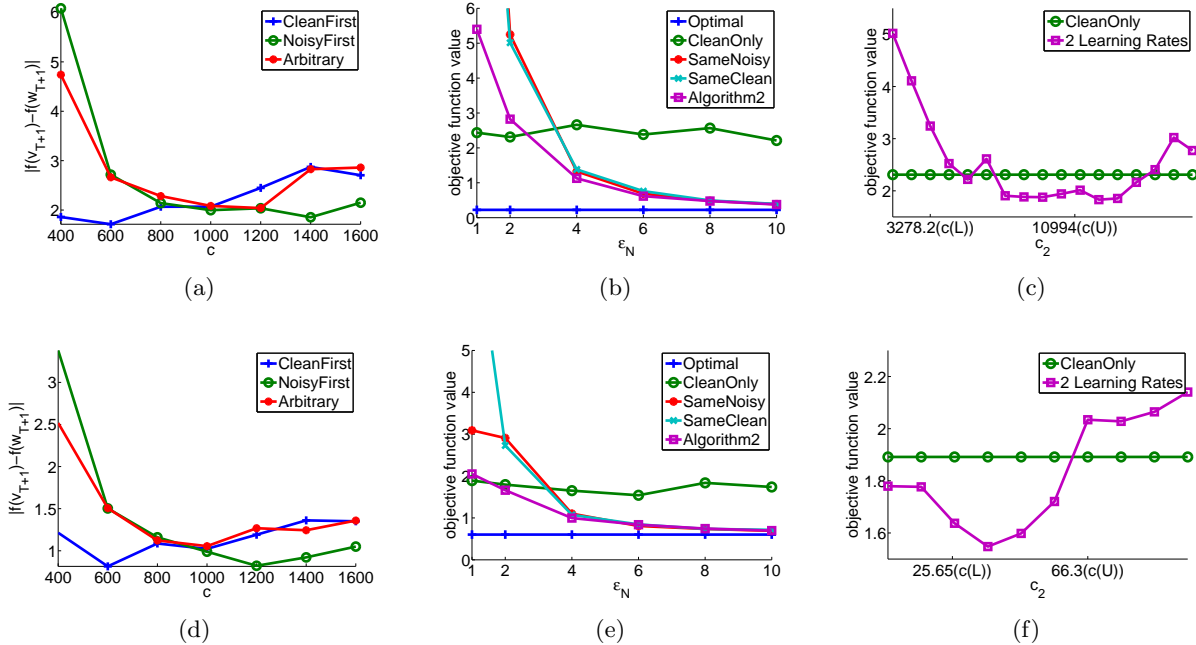


Figure 1: Column 1 plots $|f(w_{T+1}) - f(v_{T+1})|$ vs. constant c for $\lambda = 0.001$. Column 2 plots final objective function value vs. ϵ_N for $\epsilon_C = 10$. Column 3 plots final objective function value vs. c_2 for $\epsilon_N = 2$ (top) and $\epsilon_N = 1$ (bottom). Top row shows figures for MNIST and bottom row for Covertypes.

5 Experiments

We next illustrate our theoretical results through experiments on real data. We consider the task of training a regularized logistic regression classifier for binary classification under local differential privacy. For our experiments, we consider two real datasets – MNIST (with the task 1 vs. Rest) and Covertypes (Type 2 vs. Rest). The former consists of 60,000 samples in 784 dimensions, while the latter consists of 500,000 samples in 54-dimensions. We reduce the dimension of the MNIST dataset to 25 via random projections.

To investigate the effect of heterogeneous noise, we divide the training data into subsets (D_C, D_N) to be accessed through oracles $(\mathcal{G}_C, \mathcal{G}_N)$ with privacy parameters (ϵ_C, ϵ_N) respectively. We pick $\epsilon_C > \epsilon_N$, so \mathcal{G}_N is noisier than \mathcal{G}_C . To simulate typical practical situations where cleaner data is rare, we set the size of D_C to be $\beta_C = 10\%$ of the total data size. We set the regularization parameter $\lambda = 10^{-3}$, Γ_C and Γ_N according to Theorem 1 and use SGD with mini-batching (batch size 50).

Does Data Order Change Performance? Our first task is to investigate the effect of data order on performance. For this purpose, we compare three methods – CleanFirst, where all of D_C is used before D_N , NoisyFirst, where all of D_N is used before D_C , and Arbitrary, where data from $D_N \cup D_C$ is presented to the algorithm in a random order.

The results are in Figures 1(a) and 1(d). We use $\epsilon_C = 10, \epsilon_N = 3$. For each algorithm, we plot $|f(w_{T+1}) - f(v_{T+1})|$ as a function of the constant c in the learning rate. Here $f(w_{T+1})$ is the function value obtained after T rounds of SGD, and $f(v_{T+1})$ is the function value obtained after T rounds of SGD if we iterate over the data in the same order, but add no extra noise to the gradient. (See Theorem 2 for more details.) As predicted by Theorem 2, the results show that for $c < \frac{1}{\lambda}$,

CleanFirst has the best performance, while for $c > \frac{1}{\lambda}$, NoisyFirst performs best. Arbitrary performs close to NoisyFirst for a range of values of c , which we expect as only 10% of the data belongs to D_C .

Are Two Learning Rates Better than One? We next investigate whether using two learning rates in SGD can improve performance. We compare five approaches. **Optimal** is the gold standard where we access the raw data without any intervening noisy oracle. **CleanOnly** uses only D_C with learning rate with the optimal value of c obtained from Section 4. **SameClean** and **SameNoisy** use a single value of the constant c in the learning rate for $D_N \cup D_C$, where c is obtained by optimizing (13)¹ under the constraint that $c_1 = c_2$. **SameClean** uses all of D_C before using D_N , while **SameNoisy** uses all of D_N before using D_C . In **Algorithm2**, we use **Algorithm 2** to set the two learning rates and the data order (D_C first or D_N first). In each case, we set $\epsilon_C = 10$, vary ϵ_N from 1 to 10, and plot the function value obtained at the end of the optimization.

The results are plotted in Figures 1(b) and 1(e). Each plotted point is an average of 100 runs. It is clear that **Algorithm2**, which uses two learning rates, performs better than both **SameNoisy** and **SameClean**. As expected, the performance difference diminishes as ϵ_N increases (that is, the noisy data gets cleaner). For moderate and high ϵ_N , **Algorithm2** performs best, while for low ϵ_N (very noisy D_N), **CleanOnly** has slightly better performance. We therefore conclude that using two learning rates is better than using a single learning rate with both datasets, and that **Algorithm2** performs best for moderate to low noise levels.

Does Noisy Data Always Help? A natural question to ask is whether using noisy data always helps performance, or if there is some threshold noise level beyond which we should not use noisy data. Lemma 2 shows that in theory, we obtain a better upper bound on performance when we use noisy data; in contrast, Figures 1(b) and 1(e) show that for low ϵ_N (high noise), **Algorithm2** performs worse than **CleanOnly**. How do we explain this apparent contradiction?

To understand this effect, in Figures 1(c) and 1(f) we plot the performance of SGD using two learning rates (with $c_1 = \frac{1}{\lambda}$) against **CleanOnly** as a function of the second learning rate c_2 . The figures show that the best performance is attained at a value of c_2 which is different from the value predicted by **Algorithm2**, and *this best performance is better than CleanOnly*. Thus, noisy data always improves performance; however, the improvement may not be achieved at the learning rate predicted by our algorithm.

Why does our algorithm perform suboptimally? We believe this happens because the values of Γ_N and Γ_C used by our algorithm are fairly loose upper bounds. For local differential privacy, an easy lower bound on Γ is $\sqrt{\frac{4(d^2+d)}{\epsilon^2 b}}$, where b is the mini-batch size; let $c_2(L)$ (resp. $c_2(U)$) be the value of c_2 obtained by plugging in these lower bounds (resp. upper bounds from Theorem 1) to **Algorithm 1**. Our experiments show that the optimal value of c_2 always lies between $c_2(L)$ and $c_2(U)$, which indicates that the suboptimal performance may be due to the looseness in the bounds.

We thus find that even in these high noise cases, theoretical analysis often allows us to identify *an interval* containing the optimal value of c_2 . In practice, we recommend running **Algorithm 2** twice – once with upper, and once with lower bounds to obtain an interval containing c_2 , and then performing a line search to find the optimal c_2 .

¹Note that we plug in separate noise rates for \mathcal{G}_C and \mathcal{G}_N in the learning rate calculations.

6 Conclusion

In this paper we propose a model for learning from heterogeneous noise that is appropriate for studying stochastic gradient approaches to learning. In our model, data from different sites are accessed through different oracles which provide noisy versions of the gradient. Learning under local differential privacy and random classification noise are both instances of our model. We show that for two sites with different noise levels, processing data from one site followed by the other is better than randomly sampling the data, and the optimal data order depends on the learning rate. We then provide a method for choosing learning rates that depends on the noise levels and showed that these choices achieve lower regret than using a common learning rate. We validate these findings through experiments on two standard data sets and show that our method for choosing learning rates often yields improvements when the noise levels are moderate. In the case where one data set is much noisier than the other, we provide a different heuristic to choose a learning rate that improves the regret.

There are several different directions towards generalizing the work here. Firstly, extending the results to multiple sites and multiple noise levels will give more insights as to how to leverage large numbers of data sources. This leads naturally to cost and budgeting questions: how much should we pay for additional noisy data? Our results for data order do not depend on the actual noise levels, but rather their relative level. However, we use the noise levels to tune the learning rates for different sites. If bounds on the noise levels are available, we can still apply our heuristic. Adaptive approaches for estimating the noise levels while learning are also an interesting approach for future study.

Acknowledgements. The work of K. Chaudhuri and S. Song was sponsored by NIH under U54 HL108460 and the NSF under IIS 1253942.

A Appendix

A.1 Mathematical miscellany

In many cases we would like to bound a summation using an integral.

Lemma 3. *For $x \geq 0$, we have*

$$\sum_{i=a}^b i^x \leq \int_a^{b+1} i^x di = \frac{(b+1)^{x+1} - a^{x+1}}{x+1} \quad (14)$$

$$\sum_{i=a}^b i^x \geq \int_{a-1}^b i^x di = \frac{b^{x+1} - (a-1)^{x+1}}{x+1} \quad (15)$$

For $x < 0$ and $x \neq -1$, we have

$$\sum_{i=a}^b i^x \leq \int_{a-1}^b i^x di = \frac{b^{x+1} - (a-1)^{x+1}}{x+1} \quad (16)$$

$$\sum_{i=a}^b i^x \geq \int_a^{b+1} i^x di = \frac{(b+1)^{x+1} - a^{x+1}}{x+1} \quad (17)$$

For $x = -1$, we have

$$\sum_{i=a}^b i^x \leq \int_{a-1}^b i^x di = \log \frac{b}{a-1} \quad (18)$$

$$\sum_{i=a}^b i^x \geq \int_a^{b+1} i^x di = \log \frac{b+1}{a} \quad (19)$$

The sequence $\{i^x\}$ is increasing when $x > 0$ and is decreasing when $x < 0$. The proof follows directly from applying standard technique of bounding summation with integral.

A.2 Details from Section 2

Proof. (Of Theorem 1) Consider an oracle \mathcal{G} implemented based on a dataset D of size T . Given any sequence w_1, w_2, \dots, w_T , the *disguised version* of D output by \mathcal{G} is the sequence of gradients $\mathcal{G}(w_1), \dots, \mathcal{G}(w_T)$. Suppose that the oracle accesses the data in a (random) order specified by a permutation π ; for any t , any $x, x' \in \mathcal{X}$, $y, y' \in \{1, -1\}$, we have

$$\begin{aligned} \frac{\rho(\mathcal{G}(w_t) = g | (x_{\pi(t)}, y_{\pi(t)}) = (x, y))}{\rho(\mathcal{G}(w_t) = g | (x_{\pi(t)}, y_{\pi(t)}) = (x', y'))} &= \frac{\rho(Z_t = g - \lambda w - \nabla \ell(w, x, y))}{\rho(Z_t = g - \lambda w - \nabla \ell(w, x', y'))} \\ &= \frac{e^{-(\epsilon/2)\|g - \lambda w - \nabla \ell(w, x, y)\|}}{e^{-(\epsilon/2)\|g - \lambda w - \nabla \ell(w, x', y')\|}} \\ &\leq \exp((\epsilon/2)(\|\nabla \ell(w, x, y)\| + \|\nabla \ell(w, x', y')\|)) \\ &\leq \exp(\epsilon). \end{aligned}$$

The first inequality follows from the triangle inequality, and the last step follows from the fact that $\|\nabla \ell(w, x, y)\| \leq 1$. The privacy proof follows.

For the rest of the theorem, we consider a slightly generalized version of SGD that includes mini-batch updates. Suppose the batch size is b ; for standard SGD, $b = 1$. For a given t , we call $\mathcal{G}(w_t)$ b successive times to obtain noisy gradient estimates $g_1(w_t), \dots, g_b(w_t)$; these are gradient estimates at w_t but are based on separate (private) samples. The SGD update rule is:

$$w_{t+1} = \Pi_{\mathcal{W}} \left(w_t - \frac{\eta_t}{b} (g_1(w_t) + \dots + g_b(w_t)) \right).$$

For any i , $\mathbb{E}[g_i(w_t)] = \lambda w + \mathbb{E}[\nabla \ell(w, x, y)]$, where the first expectation is with respect to the data distribution and the noise, and the second is with respect to the data distribution; the unbiasedness result follows.

We now bound the norm of the noisy gradient calculated from a batch. Suppose that the oracle accesses the dataset D in an order π . Then, $g_i(w_t) = \lambda w + \nabla \ell(w_t, x_{\pi((t-1)b+i)}, y_{\pi((t-1)b+i)}) + Z_{(t-1)b+i}$. Expanding on the expression for the expected squared norm of the gradient, we have

$$\begin{aligned} \mathbb{E} \left[\left\| \frac{1}{b} (g_1(w_t) + \dots + g_b(w_t)) \right\|^2 \right] &= \mathbb{E} \left[\left\| \lambda w + \frac{1}{b} \sum_{i=1}^b \nabla \ell(w_t, x_{\pi((t-1)b+i)}, y_{\pi((t-1)b+i)}) \right\|^2 \right] \\ &\quad + \frac{2}{b} \mathbb{E} \left[\left(\lambda w + \frac{1}{b} \sum_{i=1}^b \nabla \ell(w_t, x_{\pi((t-1)b+i)}, y_{\pi((t-1)b+i)}) \right) \cdot \left(\sum_{i=1}^b Z_{(t-1)b+i} \right) \right] \\ &\quad + \frac{1}{b^2} \mathbb{E} \left[\left\| \sum_{i=1}^b Z_{(t-1)b+i} \right\|^2 \right] \quad (20) \end{aligned}$$

We now look at the three terms in (20) separately. The first term can be further expanded to:

$$\begin{aligned} \mathbb{E} \left[\|\lambda w\|^2 \right] + \mathbb{E} \left[\left\| \frac{1}{b^2} \sum_{i=1}^b \nabla \ell(w_t, x_{\pi((t-1)b+i)}, y_{\pi((t-1)b+i)}) \right\|^2 \right] \\ + 2\lambda w \cdot \left(\sum_{i=1}^b \mathbb{E} \left[\nabla \ell(w_t, x_{\pi((t-1)b+i)}, y_{\pi((t-1)b+i)}) \right] \right) \end{aligned} \quad (21)$$

The first term in (21) is at most $\lambda^2 \max_{w \in \mathcal{W}} \|w\|^2$, which is at most 1. The second term is at most $\max_w \lambda \|w\| \cdot \max_{w,x,y} \|\nabla \ell(w, x, y)\| \leq 1$, and the third term is at most 2. Thus, the first term in (20) is at most 4. Notice that this upper bound can be pretty loose compare to the average $\left\| \lambda w + \frac{1}{b} \sum_{i=1}^b \nabla \ell(w_t, x_{\pi((t-1)b+i)}, y_{\pi((t-1)b+i)}) \right\|^2$ values seen in experiment. This leads to a loose estimation of the noise level for oracle \mathcal{G}^{DP} .

To bound the second term in (20), observe that for all i , $Z_{(t-1)b+i}$ is independent of any $Z_{(t-1)b+i'}$ when $i \neq i'$, as well as of the dataset. Combining this with the fact that $\mathbb{E}[Z_\tau] = 0$ for any τ , we get that this term is 0.

To bound the third term in (20), we have:

$$\begin{aligned} \frac{1}{b^2} \mathbb{E} \left[\left\| \sum_{t \in B} Z_t \right\|_2^2 \right] &= \frac{1}{b^2} \mathbb{E} \left[\sum_{t \in B} \|Z_t\|_2^2 + \sum_{t \in B, s \in B, t \neq s} Z_t \cdot Z_s \right] \\ &= \frac{1}{b^2} \sum_{t \in B} \mathbb{E} \left[\|Z_t\|_2^2 \right] + \frac{1}{b^2} \sum_{t \in B, s \in B, t \neq s} \mathbb{E} [Z_t] \cdot \mathbb{E} [Z_s] \\ &= \frac{1}{b^2} \sum_{t \in B} \mathbb{E} \left[\|Z_t\|_2^2 \right], \end{aligned}$$

where the first equality is from the linearity of expectation and the last two equalities is from the fact that Z_i is independently drawn zeros mean vector. Because Z_t follows $\rho(Z_t = z) \propto e^{-(\epsilon/2)\|z\|}$, we have

$$\rho(\|Z_t\| = x) \propto x^{d-1} e^{-(\epsilon/2)x},$$

which is a Gamma distribution. For $X \sim \text{Gamma}(k, \theta)$, $\mathbb{E}[X] = k\theta$ and $\text{Var}(X) = k\theta^2$. Also, by property of expectation, $\mathbb{E}[X^2] = (\mathbb{E}[X])^2 + \text{Var}(X)$. We then have $\mathbb{E}[\|Z_t\|_2^2] = \frac{4(d^2 + d)}{\epsilon^2}$ and the whole term equals to $\frac{4(d^2 + d)}{\epsilon^2 b}$.

Combining the three bounds together, we have a final bound of $4 + \frac{4(d^2 + d)}{\epsilon^2 b}$. The lemma follows. \square

A.3 Proofs from Section 4

Recall that we have oracles $\mathcal{G}_1, \mathcal{G}_2$ based on data sets D_1 and D_2 . The fractions of data in each data set are $\beta_1 = \frac{|D_1|}{|D_1| + |D_2|}$ and $\beta_2 = \frac{|D_2|}{|D_1| + |D_2|}$, respectively.

A.3.1 Proof of Theorem 3

Theorem 3 is a corollary of the following Lemma.

Lemma 4. *Consider the SGD algorithm that follows Algorithm 1. Suppose the objective function is λ -strongly convex, and define $\mathcal{W} = \{w : \|w\| \leq B\}$. If $2\lambda c_1 > 1$ and $i_0 = \lceil 2c_1\lambda \rceil$, then we have the following two cases:*

1. If $2\lambda c_2 \neq 1$,

$$\mathbb{E} [\|w_{t+1} - w^*\|^2] \leq \left(4\Gamma_1^2 \frac{\beta_1^{2\lambda c_2 - 1} c_1^2}{2\lambda c_1 - 1} + 4\Gamma_2^2 \frac{c_2^2 (1 - \beta_1^{2\lambda c_2 - 1})}{2\lambda c_2 - 1} \right) \cdot \frac{1}{T} + \mathcal{O} \left(\frac{1}{T^{\min(2\lambda c_1, 2)}} \right)$$

2. If $2\lambda c_2 = 1$,

$$\mathbb{E} [\|w_{t+1} - w^*\|^2] \leq \left(4\Gamma_1^2 \frac{\beta_1^{2\lambda c_2 - 1} c_1^2}{2\lambda c_1 - 1} + 4\Gamma_2^2 c_2^2 \log \frac{1}{\beta_1} \right) \cdot \frac{1}{T} + \mathcal{O} \left(\frac{1}{T^{\min(2\lambda c_1, 2)}} \right)$$

We first begin with a lemma which follows from arguments very similar to those made in [Rakhlin et al. \(2012\)](#).

Lemma 5. *Let w^* be the optimal solution to $\mathbb{E}[f(w)]$. Then,*

$$\mathbb{E}_{1, \dots, t} [\|w_{t+1} - w^*\|^2] \leq (1 - 2\lambda\eta_t) \mathbb{E}_{1, \dots, t} [\|w_t - w^*\|^2] + \eta_t^2 \gamma_t^2.$$

where the expectation is taken wrt the oracle as well as sampling from the data distribution.

Proof. (Of Lemma 5) By strong convexity of f , we have

$$f(w') \geq f(w) + g(w)^\top (w' - w) + \frac{\lambda}{2} \|w - w'\|^2. \quad (22)$$

Then by taking $w = w_t$, $w' = w^*$ we have

$$g(w_t)^\top (w_t - w^*) \geq f(w_t) - f(w^*) + \frac{\lambda}{2} \|w_t - w^*\|^2. \quad (23)$$

And similarly by taking $w' = w_t$, $w = w^*$, we have

$$f(w_t) - f(w^*) \geq \frac{\lambda}{2} \|w_t - w^*\|^2. \quad (24)$$

By the update rule and convexity of \mathcal{W} , we have

$$\begin{aligned} \mathbb{E}_{1, \dots, t} [\|w_{t+1} - w^*\|^2] &= \mathbb{E}_{1, \dots, t} [\|\Pi_{\mathcal{W}}(w_t - \eta_t \hat{g}(w_t)) - w^*\|^2] \\ &\leq \mathbb{E}_{1, \dots, t} [\|w_t - \eta_t \hat{g}(w_t) - w^*\|^2] \\ &= \mathbb{E}_{1, \dots, t} [\|w_t - w^*\|^2] - 2\eta_t \mathbb{E}_{1, \dots, t} [\hat{g}(w_t)^\top (w_t - w^*)] + \eta_t^2 \mathbb{E}_{1, \dots, t} [\|\hat{g}(w_t)\|^2]. \end{aligned}$$

Consider the term $\mathbb{E}_{1, \dots, t} [\hat{g}(w_t)^\top (w_t - w^*)]$, where the expectation is taken over the randomness from time 1 to t . Since w_t is a function of the samples used from time 1 to $t - 1$, it is independent

of the sample used at t . So we have

$$\begin{aligned}
\mathbb{E}_{1,\dots,t} [\|w_{t+1} - w^*\|^2] &\leq \mathbb{E}_{1,\dots,t} [\hat{g}(w_t)^\top (w_t - w^*)] \\
&= \mathbb{E}_{1,\dots,t-1} \left[\mathbb{E}_t [\hat{g}(w_t)^\top (w_t - w^*) | w_t] \right] \\
&= \mathbb{E}_{1,\dots,t-1} \left[\mathbb{E}_t [\hat{g}(w_t)^\top | w_t] (w_t - w^*) \right] \\
&= \mathbb{E}_{1,\dots,t-1} [g(w_t)^\top (w_t - w^*)].
\end{aligned}$$

We have the following upper bound:

$$\begin{aligned}
\mathbb{E}_{1,\dots,t} [\|w_{t+1} - w^*\|^2] &\leq \mathbb{E}_{1,\dots,t} [\|w_t - w^*\|^2] - 2\eta_t \mathbb{E}_{1,\dots,t-1} [g(w_t)^\top (w_t - w^*)] \\
&\quad + \eta_t^2 \mathbb{E}_{1,\dots,t} [\|\hat{g}(w_t)\|^2].
\end{aligned}$$

By (23) and the bound $\mathbb{E} [\|\hat{g}(w_t)\|^2] \leq \gamma_t^2$, we have

$$\mathbb{E}_{1,\dots,t} [\|w_{t+1} - w^*\|^2] \leq \mathbb{E}_{1,\dots,t} [\|w_t - w^*\|^2] - 2\eta_t \mathbb{E}_{1,\dots,t-1} \left[f(w_t) - f(w^*) + \frac{\lambda}{2} \|w_t - w^*\|^2 \right] + \eta_t^2 \gamma_t^2.$$

Then by (24) and the fact that w_t is independent of the sample used in time t , we have the following recursion:

$$\mathbb{E}_{1,\dots,t} [\|w_{t+1} - w^*\|^2] \leq (1 - 2\lambda\eta_t) \mathbb{E}_{1,\dots,t} [\|w_t - w^*\|^2] + \eta_t^2 \gamma_t^2.$$

□

Proof. (Of Lemma 4) Let $g(w)$ be the true gradient $\nabla f(w)$ and $\hat{g}(w)$ be the unbiased noisy gradient provided by the oracle \mathcal{G}_1 or \mathcal{G}_2 , whichever is queried. From Lemma 5, we have the following recursion:

$$\mathbb{E}_{1,\dots,t} [\|w_{t+1} - w^*\|^2] \leq (1 - 2\lambda\eta_t) \mathbb{E}_{1,\dots,t} [\|w_t - w^*\|^2] + \eta_t^2 \gamma_t^2.$$

Let i_0 be the smallest positive integer such that $2\lambda\eta_{i_0} < 1$, i.e. $i_0 = \lceil 2c_1\lambda \rceil$. Notice that for fixed step size constant c and λ , i_0 would be a fixed constant. Therefore we assume that $i_0 < \beta T$. Using the above inequality inductively, and substituting $\gamma_t = \Gamma_1$ for $t \leq \beta_1 T$ and $\gamma_t = \Gamma_2$ for $t > \beta_1 T$, we have

$$\begin{aligned}
\mathbb{E}_{1,\dots,T} [\|w_{T+1} - w^*\|^2] &\leq \prod_{i=i_0}^{\beta_1 T} (1 - 2\lambda\eta_i) \prod_{i=\beta_1 T+1}^T (1 - 2\lambda\eta_i) \mathbb{E}_{1,\dots,T} [\|w_{i_0} - w^*\|^2] \\
&\quad + \Gamma_1^2 \prod_{i=\beta_1 T+1}^T (1 - 2\lambda\eta_i) \sum_{i=i_0}^{\beta_1 T} \eta_i^2 \prod_{j=i+1}^{\beta_1 T} (1 - 2\lambda\eta_j) \\
&\quad + \Gamma_2^2 \sum_{i=\beta_1 T+1}^T \eta_i^2 \prod_{j=i+1}^T (1 - 2\lambda\eta_j).
\end{aligned}$$

By substituting $\eta_t = \frac{c_1}{t}$ for D_1 and $\eta_t = \frac{c_2}{t}$ for D_2 , we have

$$\begin{aligned} \mathbb{E}_{1,\dots,T} [\|w_{T+1} - w^*\|^2] &\leq \prod_{i=i_0}^{\beta_1 T} \left(1 - \frac{2\lambda c_1}{i}\right) \prod_{i=\beta_1 T+1}^T \left(1 - \frac{2\lambda c_2}{i}\right) \mathbb{E}_{1,\dots,T} [\|w_{i_0} - w^*\|^2] \\ &\quad + \Gamma_1^2 \prod_{i=\beta_1 T+1}^T \left(1 - \frac{2\lambda c_2}{i}\right) \sum_{i=i_0}^{\beta_1 T} \frac{c_1^2}{i^2} \prod_{j=i+1}^{\beta_1 T} \left(1 - \frac{2\lambda c_1}{j}\right) \\ &\quad + \Gamma_2^2 \sum_{i=\beta_1 T+1}^T \frac{c_2^2}{i^2} \prod_{j=i+1}^T \left(1 - \frac{2\lambda c_2}{j}\right). \end{aligned}$$

Applying the inequality $1 - x \leq e^{-x}$ to each of the terms in the products, and simplifying, we get:

$$\begin{aligned} \mathbb{E}_{1,\dots,T} [\|w_{T+1} - w^*\|^2] &\leq e^{-2\lambda c_1 \sum_{i=i_0}^{\beta_1 T} \frac{1}{i}} e^{-2\lambda c_2 \sum_{i=\beta_1 T+1}^T \frac{1}{i}} \mathbb{E}_{1,\dots,T} [\|w_{i_0} - w^*\|^2] \\ &\quad + \Gamma_1^2 e^{-2\lambda c_2 \sum_{i=\beta_1 T+1}^T \frac{1}{i}} \sum_{i=i_0}^{\beta_1 T} \frac{c_1^2}{i^2} e^{-2\lambda c_1 \sum_{j=i+1}^{\beta_1 T} \frac{1}{j}} \\ &\quad + \Gamma_2^2 \sum_{i=\beta_1 T+1}^T \frac{c_2^2}{i^2} e^{-2\lambda c_2 \sum_{j=i+1}^T \frac{1}{j}}. \end{aligned} \tag{25}$$

We would like to bound (25) term by term.

A bound we will use later is:

$$e^{2\lambda c_2 / \beta_1 T} = 1 + \frac{2\lambda c_2}{\beta_1 T} e^{2\lambda c_2 / \beta_1 T'} \leq 1 + \frac{2\lambda c_2}{\beta_1 T} e^{2\lambda c_2 / \beta_1}, \tag{26}$$

where the equality is obtained using Taylor's theorem, and the inequality follows because T' is in the range $[1, \infty)$. Now we can bound the three terms in (25) separately.

The first term in (25): We bound this as follows:

$$\begin{aligned} &e^{-2\lambda c_1 \sum_{i=i_0}^{\beta_1 T} \frac{1}{i}} e^{-2\lambda c_2 \sum_{i=\beta_1 T+1}^T \frac{1}{i}} \mathbb{E}_{1,\dots,T} [\|w_{i_0} - w^*\|^2] \\ &\leq e^{-2\lambda c_1 \log \frac{\beta_1 T}{i_0}} e^{-2\lambda c_2 (\log \frac{1}{\beta_1} - \frac{1}{\beta_1 T})} \mathbb{E}_{1,\dots,T} [\|w_{i_0} - w^*\|^2] \\ &\leq \left(\frac{i_0}{T}\right)^{2\lambda c_1} \beta_1^{2\lambda(c_2 - c_1)} e^{2\lambda c_2 / \beta_1 T} (4B^2) \\ &\leq \left(\frac{i_0}{T}\right)^{2\lambda c_1} \beta_1^{2\lambda(c_2 - c_1)} \left(1 + \frac{2\lambda c_2}{\beta_1 T} e^{2\lambda c_2 / \beta_1}\right) 4B^2 \\ &= 4B^2 i_0^{2\lambda c_1} \beta_1^{2\lambda(c_2 - c_1)} \frac{1}{T^{2\lambda c_1}} + \mathcal{O}\left(\frac{1}{T^{2\lambda c_1 + 1}}\right), \end{aligned}$$

where the first equality follows from (17). The second inequality follows from $\|w\| \leq B$, $\|w - w'\| \leq \|w\| + \|w'\| \leq 2B$, and bounding expectation using maximum. The third follows from (26).

The second term in (25): We bound this as follows:

$$\begin{aligned}
\Gamma_1^2 e^{-2\lambda c_2 \sum_{i=\beta_1 T+1}^{\top} \frac{1}{i}} \sum_{i=i_0}^{\beta_1 T} \frac{c_1^2}{i^2} e^{-2\lambda c_1 \sum_{j=i+1}^{\beta_1 T} \frac{1}{j}} &\leq \Gamma_1^2 e^{-2\lambda c_2 (\log \frac{1}{\beta_1} - \frac{1}{\beta_1 T})} \sum_{i=i_0}^{\beta_1 T} \frac{c_1^2}{i^2} e^{-2\lambda c_1 \log \frac{\beta_1 T}{i+1}} \\
&= \Gamma_1^2 \beta_1^{2\lambda c_2} e^{2\lambda c_2 / \beta_1 T} \sum_{i=i_0}^{\beta_1 T} \frac{c_1^2}{i^2} \left(\frac{i+1}{\beta_1 T} \right)^{2\lambda c_1} \\
&= \Gamma_1^2 \beta_1^{2\lambda(c_2-c_1)} e^{2\lambda c_2 / \beta_1 T} c_1^2 T^{-2\lambda c_1} \sum_{i=i_0}^{\beta_1 T} \frac{(i+1)^{2\lambda c_1}}{i^2} \\
&\leq \Gamma_1^2 \beta_1^{2\lambda(c_2-c_1)} e^{2\lambda c_2 / \beta_1 T} c_1^2 T^{-2\lambda c_1} \sum_{i=i_0}^{\beta_1 T} 4(i+1)^{2\lambda c_1-2} \\
&\leq 4\Gamma_1^2 \beta_1^{2\lambda(c_2-c_1)} \left(1 + \frac{2\lambda c_2}{\beta_1 T} e^{2\lambda c_2 / \beta_1} \right) c_1^2 T^{-2\lambda c_1} \sum_{i=i_0+1}^{\beta_1 T+1} i^{2\lambda c_1-2}, \tag{27}
\end{aligned}$$

where the first inequality follows from (17), the second inequality follows from $(1 + \frac{1}{i})^2 \leq (1 + \frac{1}{1})^2 = 4$, and the last inequality follows from (26).

Bounding summation using integral following (16) and (14) of Lemma 3, if $2\lambda c_1 > 1$, the term on the right hand side would be in the order of $\mathcal{O}(1/T)$; if $2\lambda c_1 = 1$, it would be $\mathcal{O}(\log T/T)$; if $2\lambda c_1 < 1$, it would be $\mathcal{O}(1/T^{2\lambda c_1})$. Therefore to minimize the bound in terms of order, we would choose c_1 such that $2\lambda c_1 > 1$. To get an upper bound of the summation in (27), using (16) of Lemma 3, for $2\lambda c_1 < 2$,

$$\sum_{j=i_0+1}^{\beta_1 T+1} i^{2\lambda c_1-2} = \sum_{j=i_0+1}^{\beta_1 T} i^{2\lambda c_1-2} + (\beta_1 T + 1)^{2\lambda c_1-2} \leq \frac{(\beta_1 T)^{2\lambda c_1-1}}{2\lambda c_1 - 1} + \mathcal{O}(T^{2\lambda c_1-2}).$$

For $2\lambda c_1 > 2$, using (14) of Lemma 3,

$$\sum_{j=i_0+1}^{\beta_1 T+1} i^{2\lambda c_1-2} = \sum_{j=i_0+1}^{\beta_1 T-1} i^{2\lambda c_1-2} + (\beta_1 T)^{2\lambda c_1-2} + (\beta_1 T + 1)^{2\lambda c_1-2} \leq \frac{(\beta_1 T)^{2\lambda c_1-1}}{2\lambda c_1 - 1} + \mathcal{O}(T^{2\lambda c_1-2}).$$

Finally, for $2\lambda c_1 = 2$,

$$\sum_{j=i_0+1}^{\beta_1 T+1} i^{2\lambda c_1-2} = (\beta_1 T + 1) - (i_0 + 1) + 1 = \beta_1 T + \mathcal{O}(1).$$

Combining the three cases together, we have

$$\sum_{j=i_0+1}^{\beta_1 T+1} i^{2\lambda c_1-2} \leq \frac{(\beta_1 T)^{2\lambda c_1-1}}{2\lambda c_1 - 1} + \mathcal{O}(T^{2\lambda c_1-2}).$$

This allows us to further upper bound (27):

$$\begin{aligned}
& 4\Gamma_1^2 \beta_1^{2\lambda(c_2-c_1)} \left(1 + \frac{2\lambda c_2}{\beta_1 T} e^{2\lambda c_2/\beta_1}\right) c_1^2 T^{-2\lambda c_1} \sum_{i=i_0+1}^{\beta_1 T+1} i^{2\lambda c_1-2} \\
& \leq 4\Gamma_1^2 \beta_1^{2\lambda(c_2-c_1)} \left(1 + \frac{2\lambda c_2}{\beta_1 T} e^{2\lambda c_2/\beta_1}\right) c_1^2 T^{-2\lambda c_1} \left(\frac{(\beta_1 T)^{2\lambda c_1-1}}{2\lambda c_1-1} + \mathcal{O}\left(T^{2\lambda c_1-2}\right)\right) \\
& = \frac{4\Gamma_1^2 c_1^2 \beta_1^{2\lambda c_2-1}}{2\lambda c_1-1} \cdot \frac{1}{T} + \mathcal{O}\left(\frac{1}{T^2}\right) + \mathcal{O}\left(\frac{1}{T^3}\right).
\end{aligned}$$

The last term in (25): We bound this as follows:

$$\begin{aligned}
& \Gamma_2^2 \sum_{i=\beta_1 T+1}^{\top} \frac{c_2^2}{i^2} e^{-2\lambda c_2 \sum_{j=i+1}^{\top} \frac{1}{j}} \leq \Gamma_2^2 \sum_{i=\beta_1 T+1}^{\top} \frac{c_2^2}{i^2} e^{-2\lambda c_2 \log \frac{T}{i+1}} \\
& = \Gamma_2^2 c_2^2 T^{-2\lambda c_2} \sum_{i=\beta_1 T+1}^{\top} \frac{(i+1)^{2\lambda c_2}}{i^2} \leq 4\Gamma_2^2 c_2^2 T^{-2\lambda c_2} \sum_{i=\beta_1 T+1}^{\top} \frac{(i+1)^{2\lambda c_2}}{(i+1)^2} \\
& = 4\Gamma_2^2 c_2^2 T^{-2\lambda c_2} \sum_{i=\beta_1 T+2}^{\top+1} i^{2\lambda c_2-2}, \tag{28}
\end{aligned}$$

where the first inequality follows from (17) and the last inequality from $(1 + \frac{1}{i})^2 \leq 4$.
If $2\lambda c_2 \neq 1$ and $2\lambda c_2 \leq 2$, using (16) from Lemma 3,

$$\sum_{j=\beta_1 T+2}^{T+1} i^{2\lambda c_2-2} \leq \frac{1 - \beta_1^{2\lambda c_2-1}}{2\lambda c_2 - 1} T^{2\lambda c_2-1}.$$

If $2\lambda c_2 > 2$, using (14) from Lemma 3,

$$\begin{aligned}
\sum_{j=\beta_1 T+2}^{T+1} i^{2\lambda c_2-2} &= \sum_{j=\beta_1 T}^{T-1} i^{2\lambda c_2-2} + T^{2\lambda c_2-2} + (T+1)^{2\lambda c_2-2} - (\beta_1 T+1)^{2\lambda c_2-2} - (\beta_1 T)^{2\lambda c_2-2} \\
&= \frac{1 - \beta_1^{2\lambda c_2-1}}{2\lambda c_2 - 1} T^{2\lambda c_2-1} + \mathcal{O}\left(T^{2\lambda c_2-2}\right).
\end{aligned}$$

If $2\lambda c_2 = 2$,

$$\sum_{j=\beta_1 T+2}^{T+1} i^{2\lambda c_2-2} = \sum_{j=\beta_1 T+2}^{T+1} 1 = (1 - \beta_1)T.$$

In all three cases we have

$$\sum_{j=\beta_1 T+2}^{T+1} i^{2\lambda c_2-2} \leq \frac{1 - \beta_1^{2\lambda c_2-1}}{2\lambda c_2 - 1} T^{2\lambda c_2-1} + \mathcal{O}\left(T^{2\lambda c_2-2}\right).$$

Then (28) can be further upper bounded for $2\lambda c_2 \neq 1$

$$4\Gamma_2^2 c_2^2 T^{-2\lambda c_2} \sum_{i=\beta_1 T+2}^{\top+1} i^{2\lambda c_2-2} \leq 4\Gamma_2^2 \frac{c_2^2 (1 - \beta_1^{2\lambda c_2-1})}{2\lambda c_2 - 1} \cdot \frac{1}{T} + \mathcal{O}\left(\frac{1}{T^2}\right). \tag{29}$$

If $2\lambda c_2 = 1$, we have

$$\sum_{j=\beta_1 T+2}^{T+1} i^{2\lambda c_2-2} = \sum_{j=\beta_1 T+1}^T i^{-1} - (\beta_1 T + 1)^{-1} + (T + 1)^{-1} \leq \log \frac{1}{\beta_1},$$

and then

$$4\Gamma_2^2 c_2^2 T^{-2\lambda c_2} \sum_{i=\beta_1 T+2}^{T+1} i^{2\lambda c_2-2} \leq 4\Gamma_2^2 c_2^2 \log \frac{1}{\beta_1} \cdot \frac{1}{T}.$$

which is basically taking the limit as $2\lambda c_2 \rightarrow 1$ of the highest order term of (29).

Therefore the summation of the three terms is of order $\mathcal{O}(\frac{1}{T})$ (from the second and third terms), and the constant in the front of the highest order term takes on one of two values:

1. If $2\lambda c_2 \neq 1$,

$$4\Gamma_1^2 \frac{c_1^2 \beta_1^{2\lambda c_2-1}}{2\lambda c_1 - 1} + 4\Gamma_2^2 \frac{c_2^2 (1 - \beta_1^{2\lambda c_2-1})}{2\lambda c_2 - 1}.$$

2. If $2\lambda c_2 = 1$,

$$4\Gamma_1^2 \frac{c_1^2 \beta_1^{2\lambda c_2-1}}{2\lambda c_1 - 1} + 4\Gamma_2^2 c_2^2 \log \frac{1}{\beta_1}.$$

□

A.3.2 Proof of Lemma 1

Proof. (Of Lemma 1) Omitting the constant terms and setting $k_1 = 2\lambda c_1, k_2 = 2\lambda c_2$, we can re-write (13) as $1/T$ times

$$Q(k_1, k_2) = \Gamma_1^2 \frac{\beta_1^{k_2-1} k_1^2}{k_1 - 1} + \Gamma_2^2 \frac{(1 - \beta_1^{k_2-1}) k_2^2}{k_2 - 1}, \quad (30)$$

with $k_1^* = 2\lambda c_1^* = 2$.

Observe that in this case, $k_2^* \geq 2$. Let $x = k_2 - 1$; then $x \geq 1$. Plugging in $k_1^* = 2$, we can re-write (30) as

$$Q(x) = 4\Gamma_1^2 \beta_1^x + \Gamma_2^2 (1 - \beta_1^x) \left(x + \frac{1}{x} + 2 \right). \quad (31)$$

Taking the derivative, we see that

$$Q'(x) = -4\Gamma_1^2 \beta_1^x \log(1/\beta_1) + \Gamma_2^2 (1 - \beta_1^x) \left(1 - \frac{1}{x^2} \right) + \Gamma_2^2 \left(x + \frac{1}{x} + 2 \right) \beta_1^x \log(1/\beta_1). \quad (32)$$

Suppose

$$l = \frac{2 \log(\Gamma_1/\Gamma_2) + \log \log(1/\beta_1)}{\log(1/\beta_1)}.$$

Observe that $\beta_1^l \log(1/\beta_1) = \frac{\Gamma_2^2}{\Gamma_1^2}$. Plugging $x = l$ in to (32), the first term is $-4\Gamma_2^2$, the second term is at most Γ_2^2 , and the third term is at most $\frac{\Gamma_2^4}{\Gamma_1^2} (l + \frac{1}{l} + 2)$. Observe that for any fixed β_1 , for large enough Γ_1/Γ_2 , $l \geq 1$. Thus, the right hand side of (32) is at most: $-4\Gamma_2^2 + \Gamma_2^2 + \frac{\Gamma_2^4}{\Gamma_1^2} (l + 3)$. For fixed β_1 , l grows logarithmically in Γ_1/Γ_2 , and hence, for large enough Γ_1/Γ_2 , $\frac{\Gamma_2^2(l+3)}{\Gamma_1^2}$ will become

arbitrarily small. Therefore, for large enough Γ_1/Γ_2 , $Q'(l) < 0$.

Suppose

$$u = \frac{2 \log(4\Gamma_1/\Gamma_2) + \log \log(1/\beta_1)}{\log(1/\beta_1)}.$$

Observe that $\beta_1^u \log(1/\beta_1) = \frac{\Gamma_2^2}{16\Gamma_1^2}$. Plugging in $x = u$ to (32), the first term reduces to $-\frac{1}{4}\Gamma_2^2$, the second term is $\Gamma_2^2(1 - \beta_1^u)(1 - \frac{1}{u^2})$, and the third term is ≥ 0 . Observe that as $\Gamma_1/\Gamma_2 \rightarrow \infty$ with β_1 fixed, $\beta_1^u \rightarrow 0$ and $1/u^2 \rightarrow 0$. Thus, for large enough Γ_1/Γ_2 , $\Gamma_2^2(1 - \beta_1^u)(1 - \frac{1}{u^2}) \rightarrow \Gamma_2^2$, and therefore $Q'(u) > 0$. Thus, $Q'(x) = 0$ somewhere between l and u and the first part of the lemma follows.

Consider

$$x = \frac{2 \log(m\Gamma_1/\Gamma_2) + \log \log(1/\beta_1)}{\log(1/\beta_1)}$$

with $1 \leq m \leq 4$. The first term of (31) is always positive. As for the second term, $x + \frac{1}{x} + 2 \geq x$ for positive x and $\beta_1^x = \frac{\Gamma_2^2}{m^2\Gamma_1^2} \frac{1}{\log(1/\beta_1)}$ is small when Γ_1/Γ_2 is sufficiently large. Therefore for sufficiently large Γ_1/Γ_2 , we have $\Gamma_2^2(1 - \beta_1^x)(x + \frac{1}{x} + 2) \geq \frac{\Gamma_2^2}{2}x$, and thus $Q(x) \geq \frac{\Gamma_2^2}{2}x$, which gives the lower bound. And plugging in $x = l$ gives the upper bound. \square

A.3.3 Proof of Lemma 2

Proof. (Of Lemma 2) Let $k_2 = \epsilon$; then $\epsilon \geq 0$. Plugging in $k_1^* = 2$, we can re-write (30) as

$$Q(\epsilon) = 4\Gamma_1^2\beta_1^{\epsilon-1} + \Gamma_2^2(1 - \beta_1^{\epsilon-1}) \left(-1 + \epsilon + \frac{1}{-1 + \epsilon} + 2 \right). \quad (33)$$

Taking the derivative, we obtain the following:

$$\begin{aligned} Q'(\epsilon) &= -4\Gamma_1^2\beta_1^{\epsilon-1} \log(1/\beta_1) + \Gamma_2^2(1 - \beta_1^{\epsilon-1}) \left(1 - \frac{1}{(1 - \epsilon)^2} \right) - \frac{\Gamma_2^2\epsilon^2}{1 - \epsilon} \beta_1^{\epsilon-1} \log(1/\beta_1) \\ &= -\beta_1^{\epsilon-1} \log(1/\beta_1) \left(4\Gamma_1^2 + \frac{\Gamma_2^2\epsilon^2}{1 - \epsilon} \right) + \Gamma_2^2(\beta_1^{\epsilon-1} - 1) \left(\frac{1}{(1 - \epsilon)^2} - 1 \right) \\ &= -\beta_1^{\epsilon-1} \log(1/\beta_1) \left(4\Gamma_1^2 + \frac{\Gamma_2^2\epsilon^2}{1 - \epsilon} \right) + \Gamma_2^2(\beta_1^{\epsilon-1} - 1) \frac{\epsilon(2 - \epsilon)}{(1 - \epsilon)^2}. \end{aligned} \quad (34)$$

For $\epsilon = \frac{\Gamma_1^2}{\Gamma_2^2} \leq 1$, using $1 - \beta_1^{1-\epsilon} \leq (1 - \epsilon) \log(1/\beta_1)$ and $\beta_1^{\epsilon-1} - 1 = (1 - \beta_1^{1-\epsilon})\beta_1^{\epsilon-1}$, this is at most:

$$-\beta_1^{\epsilon-1} \log(1/\beta_1) \left(4\Gamma_1^2 + \frac{\Gamma_2^2\epsilon^2}{1 - \epsilon} - \frac{\Gamma_2^2\epsilon(2 - \epsilon)}{1 - \epsilon} \right) = -2\Gamma_1^2\beta_1^{\epsilon-1} \log(1/\beta_1).$$

Thus, at $l = \frac{\Gamma_1^2}{\Gamma_2^2}$, $Q'(l) < 0$.

Moreover, for $\epsilon \in [0, \frac{1}{2}]$, $1 - \beta_1^{1-\epsilon} \geq \beta_1(1 - \epsilon) \log(1/\beta_1)$. Therefore, $Q'(\epsilon)$ is at least:

$$\begin{aligned} Q'(\epsilon) &\geq -\beta_1^{\epsilon-1} \log(1/\beta_1) \left(4\Gamma_1^2 + \frac{\Gamma_2^2\epsilon^2}{1 - \epsilon} \right) + \Gamma_2^2\beta_1^\epsilon \log(1/\beta_1) \frac{\epsilon(2 - \epsilon)}{1 - \epsilon} \\ &\geq \beta_1^{\epsilon-1} \log(1/\beta_1) \left(\frac{\Gamma_2^2\beta_1\epsilon(2 - \epsilon)}{1 - \epsilon} - 4\Gamma_1^2 - \frac{\Gamma_2^2\epsilon^2}{1 - \epsilon} \right). \end{aligned}$$

Let $u = \frac{8\Gamma_2^2}{\beta_1\Gamma_1^2}$; suppose that Γ_2/Γ_1 is large enough such that $u \leq \beta_1/4$. Then, $u(2-u)\beta_1 - u^2 \geq \frac{15u\beta_1}{16}$, and

$$\frac{\Gamma_2^2(u(2-u)\beta_1 - u^2)}{1-u} \geq \frac{15\Gamma_2^2u\beta_1}{16(1-\beta_1)} \geq \frac{15\Gamma_1^2}{2(1-\beta_1)} \geq 5\Gamma_1^2.$$

Therefore, $Q'(u) > 0$, and thus $Q(\epsilon)$ is minimized at some $\epsilon \in [l, u]$.

For the second part of the lemma, the upper bound is obtained by plugging in $\epsilon = \frac{\Gamma_1}{\Gamma_2}$. For the lower bound, observe that for any $\epsilon \in [l, u]$, $Q(\epsilon) \geq 4\Gamma_1^2\beta_1^{u-1} \geq 4\Gamma_1^2\beta_1^{\Gamma_2^2/\beta_1\Gamma_1^2-1}$. \square

References

- A. Agarwal, P. L. Bartlett, P. Ravikumar, and M. J. Wainwright. Information-theoretic lower bounds on the oracle complexity of convex optimization. In Y. Bengio, D. Schuurmans, J. Lafferty, C. K. I. Williams, and A. Culotta, editors, *Advances in Neural Information Processing Systems 22*. MIT Press, 2009. URL <http://papers.nips.cc/paper/3689-information-theoretic-lower-bounds-on-the-oracle-complexity-of-convex-optimization>.
- R. Bassily, A. Thakurta, and A. Smith. Private empirical risk minimization, revisited. In *Proceedings of the IEEE Symposium on Foundations of Computer Science (FOCS 2014)*, Philadelphia, PA, USA, October 2014.
- R. Bekkerman, M. Bilenko, and J. Langford. *Scaling up Machine Learning, Parallel and Distributed Approaches*. Cambridge University Press, 2011.
- L. Bottou. Large-scale machine learning with stochastic gradient descent. In Yves Lechevalier and Gilbert Saporta, editors, *Proceedings of COMPSTAT'2010*, pages 177–186, Paris, France, August 2010. doi: 10.1007/978-3-7908-2604-3_16. URL http://dx.doi.org/10.1007/978-3-7908-2604-3_16.
- K. Chaudhuri, C. Monteleoni, and A. D. Sarwate. Differentially private empirical risk minimization. *Journal of Machine Learning Research*, 12:1069–1109, March 2011. URL <http://jmlr.csail.mit.edu/papers/v12/chaudhuri11a.html>.
- K. Crammer, M. Kearns, and J. Wortman. Learning from data of variable quality. In Y. Weiss, B. Schölkopf, and J.C. Platt., editors, *Advances in Neural Information Processing Systems 18*, pages 219–226. MIT Press, 2006. URL <http://papers.nips.cc/paper/2920-learning-from-data-of-variable-quality>.
- J. Duchi and Y. Singer. Efficient online and batch learning using forward backward splitting. *Journal of Machine Learning Research*, 10:2899–2934, December 2009. URL <http://www.jmlr.org/papers/v10/duchi09a.html>.
- J. Duchi, M. Jordan, and M. Wainwright. Privacy aware learning. In P. Bartlett, F. C. N. Pereira, C. J. C. Burges, L. Bottou, and K. Q. Weinberger, editors, *Advances in Neural Information Processing Systems 25*, pages 1439–1447. MIT Press, 2012. URL <http://papers.nips.cc/paper/4505-privacy-aware-learning>.
- J. Duchi, M. J. Wainwright, and M. Jordan. Local privacy and minimax bounds: Sharp rates for probability estimation. In C.J.C. Burges, L. Bottou, M. Welling, Z. Ghahramani, and K.Q. Weinberger, editors, *Advances in Neural Information Processing Systems*

- 26, pages 1529–1537. Curran Associates, Inc., 2013. URL <http://papers.nips.cc/paper/5013-local-privacy-and-minimax-bounds-sharp-rates-for-probability-estimation.pdf>.
- C. Dwork, K. Kenthapadi, F. McSherry, I. Mironov, and M. Naor. Our data, ourselves: Privacy via distributed noise generation. In Serge Vaudenay, editor, *Advances in Cryptology - EUROCRYPT 2006*, volume 4004 of *Lecture Notes in Computer Science*, pages 486–503, Berlin, Heidelberg, 2006a. Springer-Verlag. doi: 10.1007/11761679_29. URL http://dx.doi.org/10.1007/11761679_29.
- C. Dwork, F. McSherry, K. Nissim, and A. Smith. Calibrating noise to sensitivity in private data analysis. In Shai Halevi and Tal Rabin, editors, *Theory of Cryptography*, volume 3876 of *Lecture Notes in Computer Science*, pages 265–284, Berlin, Heidelberg, March 4–7 2006b. Springer. doi: 10.1007/11681878_14. URL http://dx.doi.org/10.1007/11681878_14.
- P. Jain, P. Kothari, and A. Thakurta. Differentially private online learning. In *Proceedings of the 25th Annual Conference on Learning Theory (COLT '12)*, volume 23 of *JMLR Workshop and Conference Proceedings*, pages 24.1–24.34, Edinburgh, Scotland, June 2012. URL <http://www.jmlr.org/proceedings/papers/v23/jain12/jain12.pdf>.
- S. A. Kasiviswanathan, H. K. Lee, K. Nissim, S. Raskhodnikova, and A. Smith. What can we learn privately? In *IEEE 49th Annual IEEE Symposium on Foundations of Computer Science (FOCS '08)*, pages 531–540, 2008. doi: 10.1109/FOCS.2008.27. URL <http://dx.doi.org/10.1109/FOCS.2008.27>.
- M. Kearns. Efficient noise-tolerant learning from statistical queries. *Journal of the ACM*, 45(6): 983–1006, November 1998. doi: 10.1145/293347.293351. URL <http://dx.doi.org/10.1145/293347.293351>.
- D. Kifer, A. Smith, and A. Thakurta. Private convex empirical risk minimization and high-dimensional regression. In Shie Mannor, Nathan Srebro, and Robert C. Williamson, editors, *Proceedings of the 25th Annual Conference on Learning Theory (COLT '12)*, volume 23 of *JMLR Workshop and Conference Proceedings*, pages 25.1–25.40, Edinburgh, Scotland, June 2012. URL <http://jmlr.csail.mit.edu/proceedings/papers/v23/kifer12/kifer12.pdf>.
- E. Moulines and Francis R. Bach. Non-asymptotic analysis of stochastic approximation algorithms for machine learning. In J. Shawe-Taylor, R.S. Zemel, P.L. Bartlett, F. Pereira, and K.Q. Weinberger, editors, *Advances in Neural Information Processing Systems 24*, pages 451–459. Curran Associates, Inc., 2011. URL <http://papers.nips.cc/paper/4316-non-asymptotic-analysis-of-stochastic-approximation-algorithms-for-machine-learning.pdf>.
- N. Natarajan, I. Dhillon, P. Ravikumar, and A. Tewari. Learning with noisy labels. In C.J.C. Burges, L. Bottou, M. Welling, Z. Ghahramani, and K.Q. Weinberger, editors, *Advances in Neural Information Processing Systems 26*, pages 1196–1204. Curran Associates, Inc., 2013. URL <http://papers.nips.cc/paper/5073-learning-with-noisy-labels>.
- 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. doi: 10.1137/070704277. URL <http://dx.doi.org/10.1137/070704277>.

- A.S. Nemirovsky and D.B. Yudin. *Problem complexity and method efficiency in optimization*. Wiley & Sons, 1983.
- Y. Nesterov and J. Vial. Confidence level solutions for stochastic programming. *Automatica*, 44(6):1559–1568, June 2008. doi: 10.1016/j.automatica.2008.01.017. URL <http://dx.doi.org/10.1016/j.automatica.2008.01.017>.
- A. Rakhlin, O. Shamir, and K. Sridharan. Making gradient descent optimal for strongly convex stochastic optimization. Technical Report arXiv:1109.5647 [cs.LG], ArXiv, 2012. URL <http://arxiv.org/abs/1109.5647>.
- B. I. P. Rubinstein, P. L. Bartlett, L. Huang, and N. Taft. Learning in a large function space: Privacy-preserving mechanisms for SVM learning. *Journal of Privacy and Confidentiality*, 4(1):65–100, 2012. URL <http://repository.cmu.edu/jpc/vol4/iss1/4/>.
- S. Shalev-Shwartz, O. Shamir, N. Srebro, and K. Sridaran. Stochastic convex optimization. In *Proceedings of the Conference on Learning Theory (COLT)*, 2009. URL <http://www.cs.mcgill.ca/~colt2009/papers/018.pdf>.
- S. Song, K. Chaudhuri, and Anand D. Sarwate. Stochastic gradient descent with differentially private updates. In *2013 IEEE Global Conference on Signal and Information Processing (GlobalSIP)*, pages 245–248, 2013. doi: 10.1109/GlobalSIP.2013.6736861. URL <http://dx.doi.org/10.1109/GlobalSIP.2013.6736861>.
- L. Wasserman and S. Zhou. A statistical framework for differential privacy. *Journal of the American Statistical Association*, 105(489):375–389, 2010. doi: 10.1198/jasa.2009.tm08651. URL <http://dx.doi.org/10.1198/jasa.2009.tm08651>.
- L. Xiao. Dual averaging methods for regularized stochastic learning and online optimization. *Journal of Machine Learning Research*, 11:2543–2596, October 2010. URL <http://www.jmlr.org/papers/v11/xiao10a.html>.