# RÉNYI DIFFERENTIALLY PRIVATE ADMM FOR NON-SMOOTH REGULARIZED OPTIMIZATION

**Chen Chen**
Department of Computer Science
University of Georgia
Athens, GA 30602

**Jaewoo Lee**
Department of Computer Science
University of Georgia
Athens, GA 30602

June 18, 2021

## ABSTRACT

In this paper we consider the problem of minimizing composite objective functions consisting of a convex differentiable loss function plus a non-smooth regularization term, such as $L_1$ norm or nuclear norm, under Rényi differential privacy (RDP). To solve the problem, we propose two stochastic alternating direction method of multipliers (ADMM) algorithms: ssADMM based on gradient perturbation and mpADMM based on output perturbation. Both algorithms decompose the original problem into sub-problems that have closed-form solutions. The first algorithm, ssADMM, applies the recent privacy amplification result for RDP to reduce the amount of noise to add. The second algorithm, mpADMM, numerically computes the sensitivity of ADMM variable updates and releases the updated parameter vector at the end of each epoch. We compare the performance of our algorithms with several baseline algorithms on both real and simulated datasets. Experimental results show that, in high privacy regimes (small $\epsilon$), ssADMM and mpADMM outperform other baseline algorithms in terms of classification and feature selection performance, respectively.

## 1 Introduction

Concerns on privacy of individuals in the data used for training machine learning models have led to extensive research on private model building techniques [1, 2, 3, 4, 5, 6, 7], especially in the context of Empirical Risk Minimization (ERM). Let $D = (d_1, d_2, \ldots, d_n)$ be a dataset, where $d_i \in \mathcal{D}$. Many machine learning problems can be formulated as regularized optimization problems of form:

$$\min_{x \in \mathbb{R}^p} F(x) := \frac{1}{n} \sum_{i=1}^{n} f(x, d_i) + \lambda h(x) \tag{1}$$

where $\lambda > 0$ is a regularization coefficient, $f : \mathbb{R}^p \times \mathcal{D} \to \mathbb{R}$ is a smooth convex loss function, and $h : \mathbb{R}^p \to \mathbb{R}$ is a simple convex *non-smooth* regularizer such as $L_1$-norm or nuclear norm. This formulation has received substantial attention as it arises in many interesting applications of machine learning such as generalized lasso [8], matrix recovery [9, 10], and a class of $L_1$ regularized problems. Despite recent advances in methods for differentially private ERM, many existing solutions are not directly applicable to the problem in (1) due to requirement for differentiability [3, 4, 5, 7] or strong convexity [1] of the regularization term $h(x)$. Alternating direction method of multipliers (ADMM) [11] has shown to be effective in solving optimization problems with complicated structure regularization.

In this paper, we propose two stochastic ADMM algorithms that satisfy Rényi Differential Privacy (RDP), namely subsampled stochastic ADMM (ssADMM) and model perturbation based ADMM (mpADMM). The first algorithm has the following key features. First, ssADMM is scalable and fast. The algorithm splits the composite objective function into differentiable and non-smooth terms, $\sum_i f(x, d_i)$ and $h(x)$, using the ADMM framework. The differentiable term is further approximated by the first order Taylor expansion and linearization as in [12]. This approximated augmented Lagrangian function has a simple analytical solution. For the non-smooth regularization term $h(x)$, ssADMM applies proximal mappings. For many non-smooth regularization function popularly used in machine learning, such as $L_1$-norm,

SCAD [13], and MCP [14], those proximal mappings yield closed form solutions. Therefore, both subproblems can be solved efficiently.

Second, ssADMM makes use of recently proposed *privacy amplification* lemma [15] to tightly bound the total privacy loss across many iterations. In the closed-form solution of the modified augmented Lagrangian function, the only data dependent term is the gradient $\nabla f(x^k)$, where $x^k$ denotes the value of $x$ at iteration $k$. The algorithm computes the gradient $\nabla f(x^k)$ using a randomly *subsampled* data and add Gaussian noise to ensure $(\alpha, \epsilon_k)$-RDP, which allows us to exploit the randomness in the subsampling and to introduce less noise to each iteration.

The second algorithm, mpADMM, takes the output perturbation approach but substantially differs from the original method. Unlike the original method which releases model parameters once only at the end, the proposed method releases the output after each epoch. For each epoch, we numerically compute the sensitivity of both primal and dual variable updates in ADMM and release the parameter vector using the Gaussian mechanism. The algorithm uses the released (noisy) output as the starting value for the next epoch.

Our contributions are summarized as follows:

- We propose two efficient Rényi differentially private algorithms, based on stochastic ADMM, for solving non-smooth convex optimization problems. In our proposed ssADMM, each subproblem is solved exactly in closed form.

- We apply the recent privacy amplification result for RDP to stochastic ADMM and show that the inherent randomness in subsampling process can be used to achieve stronger privacy protection.

- We empirically show the effectiveness of the proposed algorithms by performing extensive empirical evaluations on generalized linear models and comparing with other baseline algorithms. The results show that, in high privacy regimes (small $\epsilon$), ssADMM and mpADMM outperform other baseline algorithms in terms of classification and feature selection performance, respectively.

The rest of this paper are organized as follow: Section 2 summarizes related work. In Section 3, we provide background on Rényi differential privacy and ADMM. Section 4 introduces the proposed Rényi differentially private ADMM algorithms. Section 5 provides the performance evaluations on both synthetic and real datasets. Section 6 concludes the paper.

## 2   Related Work

Many works have been done to solve the empirical risk minimization problem under differential privacy. Generally, there are three types of algorithms proposed. Output perturbation algorithms perturb the model parameters based on sensitivity, for example, [1] analyzed the sensitivity of optimal solutions trained between neighboring databases; [5] tackled the case when full gradient descent is applied; and [16] and [7] analyzed the situation of applying stochastic gradient descent on permuting mini-batches. Objective perturbation algorithms perturb the training objective functions, and the privacy guarantee is subject to an exact solution of the ERM problem: [1] presented the first objective perturbation technique, and it is extended by [2]. Gradient perturbation algorithms perturb the (stochastic) gradients used for model updating by first-order optimization methods, and use a composition technique to quantify the overall privacy leak for multiple access of the data through gradient calculation. For example, [3] proposed "strong composition" theorem, then [4] proposed "moment accountant" method, which is also used in [6] and [17]. The Réyni differential privacy was introduced by [18], which can also be applied in gradient perturbation, especially after [15] proposed its amplification by subsampling results.

Alternating Direction Method of Multipliers (ADMM) is an old algorithm to solve optimization problems [19]. It has been extensively studied, and applied in many domains such as outlier recovery [20], image processing [21], and sensor detection [22]. In addition to its original version, many variations has been presented, such as [23, 24] and [12]. Several ADMM based differentially private algorithms have been presented, for example, [25] applied objective perturbation technique on the original ADMM problem, [26] and [27] applied output and objective perturbation technique, and [28] applied gradient perturbation technique on ADMM-based algorithms in distributed settings.

$L_1$ regularized ERM problem was first proposed for linear regression, that is least absolute shrinkage and selection operator (LASSO) [29]. Some variants of LASSO exists, such as [30] and [31]. It has been used for classification problems, and many algorithms for solving $L_1$ regularized generalized linear models were presented, such as [32], [33], and [34]. [35] and [36] has shown that $L_1$ regularized classification has good performance in feature selection. Limited to the assumption on the loss function, many differentially private ERM algorithms cannot be directly applied on $L_1$ regularized classification, with a few exceptions such as [4, 25], and [28].

# 3 Preliminaries

In this section we introduce relative background of this paper. We will start with definitions and lemmas in differential privacy and Rényi differential privacy, the $L_1$-regularized classification problem we aim to solve, and then the ADMM algorithm based on which we proposed our algorithms.

We assume a dataset $D = \{d_1, ..., d_n\} \sim \mathcal{D}^n$ is a set collected from $n$ individuals from an unknown population distribution $\mathcal{D}$, where $d_i = (s_i, l_i)$ for $i = 1, ..., n$ is a record of one individual, with $s_i$ being a vector of features of dimension $p$, and $l_i \in \{-1, +1\}$ being its label. Two datasets $D$ and $D'$ are considered *neighboring*, if $D'$ can be obtained by replacing one record with another one from $\mathcal{D}$, notated as $D \sim D'$. We use $x, y, z$ to denote model parameters, and $\| \cdot \|_1$ (resp. $\| \cdot \|_2$) as $L_1$ (resp. $L_2$) norm of a vector.

## 3.1 Differential Privacy

Differential privacy is so far the standard standard for protecting the privacy of sensitive datasets. Its formal definition is stated as:

**Definition 1** (($\epsilon, \delta$)-Differential Privacy (DP))**.** *[37] [38] Given privacy parameters $\epsilon \geq 0, 0 \leq \delta \leq 1$, a randomized mechanism (algorithm) $\mathcal{M}$ satisfies ($\epsilon, \delta$)-DP if for every event $S \subseteq range(M)$, and for every pair of neighboring datasets $D \sim D'$,*

$$\Pr[\mathcal{M}(D) \in S] \leq e^\epsilon \Pr[\mathcal{M}(D') \in S] + \delta \tag{2}$$

If $\delta = 0$, it is called *pure* differential privacy, and $\delta > 0$ is called *approximate* differential privacy.

With pure differential privacy, even the strongest attacker with arbitrary background information has limited ability to make inferences on the unknown record(s). With approximated differential privacy, this guarantee holds with a high chance, while failure of privacy preserving happens with probability at most $\delta$ (informally called "all-bets-are-off"). In practice, $\delta$ should be taken significantly small, such as $\Theta(n^{-2})$.

While approximate DP is a relaxation of pure DP, some other relaxations of pure DP also exists, such as zero-concentrated differential privacy (zCDP) [39] and Rényi Differential Privacy (RDP) [18]. These relaxations do not have such semantic meanings as approximate DP, but they are shown to stand between pure and approximate DP: they provide weaker protection than pure DP, but stronger protection than approximated DP, for any given $\delta > 0$. In this paper, we will focus on Rényi Differential Privacy.

## 3.2 Rényi Differential Privacy

Define $Z = \frac{\Pr[\mathcal{M}(D) \in S]}{\Pr[\mathcal{M}(D') \in S]}$ as the privacy loss random variable, instead of requiring it always inside range $[-\epsilon, \epsilon]$ as pure DP, Rényi differential privacy (RDP) constraints its expectation by Rényi divergence.

**Definition 2** (($\alpha, \epsilon$)-Rényi Differential Privacy (RDP))**.** *[18] Given a real number $\alpha \in (1, +\infty)$ and privacy parameter $\epsilon \geq 0$, a randomized mechanism (algorithm) $\mathcal{M}$ satisfies ($\alpha, \epsilon$)-RDP if for every pair of neighboring datasets $D \sim D'$, the Rényi $\alpha$-divergence between $\mathcal{M}(D)$ and $\mathcal{M}(D')$ satisfies*

$$D_\alpha[\mathcal{M}(D)\|\mathcal{M}(D')] \leq \epsilon \tag{3}$$

That is, the privacy parameter $\epsilon$ bounds the moment $\alpha$ of the Rényi divergence $D_\alpha$, which is defined as

**Definition 3** (Rényi Divergence)**.** *For probability distributions $\mathcal{M}(D)$ and $\mathcal{M}(D')$ over a set $\Omega$, and let $\alpha \in (1, +\infty)$. Then Rényi $\alpha$-divergence is*

$$D_\alpha(\mathcal{M}(D)\|\mathcal{M}(D')) := \frac{1}{\alpha - 1} \log \mathbb{E}_{x \sim \mathcal{M}(D')}\left[\left(\frac{P_{\mathcal{M}(D)}(x)}{P_{\mathcal{M}(D')}(x)}\right)^\alpha\right] \tag{4}$$

One method to achieve RDP is through the Gaussian mechanism: when a query $q(D)$ is taken over the dataset, the Gaussian mechanism adds a Gaussian noise $\gamma \sim \mathcal{N}(0, \sigma^2 \mathbf{I}_k)$, and release perturbed $q(D) + \gamma$.

**Lemma 1** (Gaussian Mechanism)**.** *[18] Let $q : \mathcal{D}^n \to \mathbb{R}^k$ be a vector-valued function over datasets. Let $\mathcal{M}$ be a mechanism releasing $q(D) + \gamma$ where $\gamma \sim \mathcal{N}(0, \sigma^2 \mathbf{I}_k)$, then for any $D \sim D'$ and any $\alpha \in (1, +\infty)$,*

$$D_\alpha(\mathcal{M}(D)\|\mathcal{M}(D')) \leq \alpha \Delta_2^2(q)/(2\sigma^2) \tag{5}$$

Gaussian mechanism relies on the $L_2$ sensitivity:

**Definition 4** ($L_2$ sensitivity). *Let $q : \mathcal{D}^n \to \mathbb{R}^k$ be a vector-valued function over datasets. The $L_2$ sensitivity of $q$, denoted as $\Delta_2(q)$, is defined as*

$$\Delta_2(q) = \sup_{D \sim D'} \|q(D) - q(D')\|_2 \tag{6}$$

Therefore, when scale the variance $\sigma^2 = \alpha \Delta_2^2(q)/(2\epsilon)$, then $\mathcal{M}$ satisfies $(\alpha, \epsilon)$-RDP.

Gaussian mechanism makes the mechanism $\mathcal{M}$ satisfy $(\alpha, \epsilon)$-RDP for a series of $\alpha$, so we can use $\epsilon(\alpha)$ to denote the privacy $\epsilon$ under moment $\alpha$. In empirical risk minimization algorithms, it is common that the mechanism is taken over a randomized subsample of the dataset $B$, instead of the whole dataset $D$. Then, application Gaussian Mechanism on the subsample $B$ would satisfy $(\alpha, \epsilon(\alpha))$-RDP with respect to $B$. Due to the subsampling procedure, the mechanism would satisfy an amplified privacy with respect to the whole dataset $D$, as given by the following lemma:

**Lemma 2** (RDP for subsampled mechanism). *[15] For a randomized mechanism $\mathcal{M}$ and a dataset $D \sim \mathcal{D}^n$, define $\mathcal{M} \circ \text{SUBSAMPLE}$ as (1) subsample without replacement $m$ datapoints from the dataset (denote $q = m/n$ as sampling ratio); (2) apply $\mathcal{M}$ on the subsampled dataset as input, then if $\mathcal{M}$ satisfies $(\alpha, \epsilon(\alpha))$-RDP with respect to the subsample for all integers $\alpha > 2$, then the new randomized mechanism $\mathcal{M} \circ \text{SUBSAMPLE}$ satisfies $(\alpha, \epsilon'(\alpha))$-RDP with respect to $D$, where*

$$\begin{aligned}
\epsilon'(\alpha) \leq \frac{1}{\alpha - 1} \log \Big( 1 &+ q^2 \binom{\alpha}{2} \min\left\{ 4(e^{\epsilon(2)} - 1), 2e^{\epsilon(2)} \right\} \\
&+ \sum_{j=3}^{\alpha} q^j \binom{\alpha}{j} 2e^{(j-1)\epsilon(j)} \Big)
\end{aligned} \tag{7}$$

Similar as DP, RDP has below composition properties:

**Lemma 3** (RDP composition). *[18] For randomized mechanisms $\mathcal{M}_1$ and $\mathcal{M}_2$ applied on dataset $D$, if $\mathcal{M}_1$ satisfies $(\alpha, \epsilon_1)$-RDP and $M_2$ satisfies $(\alpha, \epsilon_2)$-RDP, then their composition $\mathcal{M}_1 \circ \mathcal{M}_2$ satisfies $(\alpha, \epsilon_1 + \epsilon_2)$-RDP.*

RDP is said to provide stronger protection than approximate DP, due to below conversion to $(\epsilon, \delta)$-DP:

**Proposition 1** (RDP to $(\epsilon, \delta)$-DP). *[18] If $\mathcal{M}$ satisfies $(\alpha, \epsilon)$-RDP, then it satisfies $(\epsilon(\delta), \delta)$-DP for $\epsilon(\delta) \geq \epsilon + \frac{\log(1/\delta)}{\alpha - 1}$.*

Therefore, when evaluating our proposed algorithms, to compare with other algorithms which satisfies $(\epsilon, \delta)$-DP, we keep track of $(\alpha, \epsilon)$ pairs which our algorithm satisfies for a series of $\alpha$ values, then convert each pair into a $(\epsilon(\delta), \delta)$ pair it satisfies by Proposition 1, for a pre-defined small $\delta$, and choose the smallest $\epsilon(\delta)$ as the $(\epsilon, \delta)$-DP it satisfies to compare with other algorithms.

### 3.3 Regularized Empirical Risk Minimization

Many problems in machine learning can be formulated as empirical risk minimization (ERM), which seek a solution $x^* \in \Theta$ that minimizes an empirical loss on the training data:

$$x^* = \arg\min_{x \in \Theta} F(x, D) := \arg\min_{x \in \Theta} \frac{1}{n} \sum_{i=1}^{n} \ell(x, d_i), \tag{8}$$

where $\Theta$ is a parameter space, $\ell$ is a *loss* function. To prevent overfitting, it is common to add a (data-independent) regularization term into the objective function, i.e. $\ell(x, d_i) = f(x, d_i) + R(x)$. For $L_1$ regularization, $R(x) = \lambda\|x\|_1$. For example, $L_1$ regularized logistic regression, one can fit the model by solving

$$x^* = \arg\min_{x \in \Theta} \frac{1}{n} \sum_{i=1}^{n} \log(1 + \exp(-l_i x^T s_i)) + \lambda\|x\|_1 \tag{9}$$

Recall that each datum $d_i = (s_i, l_i)$ as feature vector $s_i$ and label $l_i$. However, due to that many optimization algorithms assume the loss function to be doubly differentiable, it cannot be directly used on $L_1$ regularization problems. In this paper, we make the following assumptions on the loss function:

- **Convexity** Both the data-dependent function $f$ and regularization term $R$ are convex.
- **Differentiability** The non-regularized data-dependent function $f$ is continuously differentiable with respect to $x$.
- **Bounded gradient** There exists a constant $C > 0$ such that $\|\nabla f(x, d)\|_2 \leq C$ for all $x \in \Theta$ and $d \in \mathcal{D}$. Usually it is satisfied by preprocessing the data to ensure the feature $s_i$ of each data $d_i$ lies inside a ball of some radius $r$, or directly clip the $L_2$ norm of individual gradient by a threshold $C$.

4

### 3.4 Alternating Direction Method of Multipliers

The Alternating Direction Method of Multipliers (ADMM) algorithm was proposed decades ago, and has recently been widely used to solve optimization problems in machine learning [19]. Consider the optimization problem

$$
\begin{aligned}
\text{minimize} \quad & f(x) + h(z) \\
\text{subject to} \quad & Ax + Bz = c
\end{aligned}
\tag{10}
$$

where $f : \mathbb{R}^n \to \mathbb{R}$, $g : \mathbb{R}^m \to \mathbb{R}$, $A \in \mathbb{R}^{p \times n}$, $B \in \mathbb{R}^{p \times m}$, and $c \in \mathbb{R}^p$. ADMM forms the augmented Lagrangian of the problem:

$$
L_\rho(x, z, y) := f(x) + h(z) + y^T(Ax + Bz - c) + \frac{\rho}{2}\|Ax + Bz - c\|_2^2
\tag{11}
$$

where $x, z$ are called the *primal* variables, $y \in \mathbb{R}^p$ is called the *dual* variable, and $\rho > 0$ is a pre-selected *penalty* parameter.

ADMM algorithm solves the optimization problem by alternating the iterations below

$$
x\text{-minimization step: } x^{k+1} \leftarrow \arg\min_x L_\rho(x, z^k, y^k)
\tag{12}
$$

$$
z\text{-minimization step: } z^{k+1} \leftarrow \arg\min_z L_\rho(x^{k+1}, z, y^k)
\tag{13}
$$

$$
\text{dual variable update: } y^{k+1} \leftarrow y^k + \rho(Ax^{k+1} + Bz^{k+1} - c)
\tag{14}
$$

Therefore, $x$ and $z$ are updated in an alternating fashion, and separating minimization over $x$ and $z$ into two steps can make the otherwise hard-to-solve optimization problem solvable in a sequential manner.

### 3.5 Stochastic ADMM

One variant of ADMM, stochastic ADMM (sADMM), was proposed by [12] and tested on $L_1$ regularized linear regression (LASSO). This variant was proposed based on the observation that, for ADMM problems, usually one of $f(x)$ and $h(z)$ is data-dependent, and it is both expensive and unnecessary to exactly solve its minimization step for each iteration. To be specific, let $f$ be data-dependent, and $h$ be data-independent, then the optimization problem becomes $f(x, D) + h(z)$, and sADMM approximate $L_\rho$ by *approximated* augmented Lagrangian $\hat{L}_\rho$, defined at iteration $k$ as

$$
\begin{aligned}
\hat{L}_\rho(x, z, y) := & f(x^k) + \langle \nabla f(x^k, B_k), x \rangle + \frac{\|x - x^k\|_2^2}{2\eta^k} \\
& + h(z) + y^T(Ax + Bz - c) + \frac{\rho}{2}\|Ax + Bz - c\|_2^2
\end{aligned}
\tag{15}
$$

where $B_k$ is a portion of the data accessed at iteration $k$, and $\eta^k$ is the learning rate at iteration $k$. After this approximation of $L_\rho$ by $\hat{L}_\rho$, one can derive an exact solution for each $x$-minimization step in (12), instead of solving a computationally expensive ERM problem.

For $L_1$ regularized ERM, let $h(z)$ be the regularization term $R(z) = \lambda\|z\|_1$, the constraint $Ax + Bz = c$ reduces to $x = z$, then by taking derivative of $\hat{L}_\rho(x, z^k, y^k)$ and set to zero, one get

$$
x^{k+1} \leftarrow \frac{1}{\rho + 1/\eta^k}\left(-\nabla f(x, B_k) - y^k + \rho z^k + x^k/\eta^k\right)
\tag{16}
$$

as the exact solution to minimize $\hat{L}_\rho(x, z^k, y^k)$, and

$$
y^{k+1} \leftarrow y^k + \rho(x - z)
\tag{17}
$$

to update the dual variable $y$.

## 4 Algorithm

In this section we propose the main algorithms. We propose two sADMM based $L_1$ regularized classification algorithms, both satisfies Rényi differential privacy. One achieves privacy by gradient perturbation relying on randomized subsampling; the other is through model perturbation after each epoch relying on sensitivity calculation. Both algorithms assume a centralized computing: all training data were collected in a center, which performs the computation locally. This is because we assume the data is small-to-median sized, where $L_1$ regularization are usually applied on.

### 4.1 Rényi differentially private subsampling algorithm

Our subsampling private sADMM algorithm (ssADMM) is presented in Algorithm 1. This algorithm is inspired by the gradient perturbation technique proposed in [4], on differentially private stochastic gradient descent (DP-SGD).

Similar as DP-SGD, our ssADMM algorithm perturbs the mini-batch gradient by Gaussian noise right after gradient evaluation in line 6. However, Algorithm 1 differs from DP-SGD for the following aspects: (i) By utilizing ADMM, we are able separate gradient descent and $L_1$ regularization into two steps, so that pure gradient can be computed and perturbed in $x$-minimization step; for DP-SGD, proximal gradient has to be used to handle $L_1$ regularization; (ii) while DP-SGD suggest using constant learning rate, we proved that using decreasing step size in Algorithm 1 help accelerate convergence, as in Theorem 2 and numerical experiments; (iii) authors of DP-SGD proposed the moment accountant (MA) method to analyze the privacy loss, and convert to $(\epsilon, \delta)$-DP; we use the most recent RDP for subsampling mechanism, which is a more advanced technique to analyze privacy loss, and also easier to implement.

---

**Algorithm 1** RDP subsampling sADMM $L_1$ regularized ERM algorithm (ssADMM)

---

1: **Input**: Dataset $D = \{d_1, ..., d_n\}$. Penalty parameter $\rho$, mini-batch size $m$, total iterations $T$.
2: **Initialize**: primal variables $x^0, z^0$, dual variable $y^0$.
3: **for** iteration $k = 0, 1, ..., T-1$ **do**
4:    Sample mini-batch $B_k$ from $D$ of size $m$.
5:    $g_k \leftarrow \frac{1}{m} \sum_{d_i \in B_k} \nabla f(x^k, d_i)$            ▷ compute gradient
6:    $\tilde{g}_k \leftarrow g_k + \gamma$ where $\gamma \sim N(0, \sigma^2 \mathbf{I}_p)$       ▷ perturb gradient by Gaussian noise
7:    Compute $x^{k+1}$ by (16) using $\tilde{g}_k$          ▷ primal variable $x$
8:    Compute $z^{k+1}$ by (18)             ▷ primal variable $z$
9:    Compute $y^{k+1}$ by (17)             ▷ dual variable $y$
10: **end for**
11: **Output**: $x^T$

---

Since the regularization is data-independent, it does not cause any privacy leak. Therefore, any (non-) smooth regularizers are applicable for Algorithm 1, with the same privacy guarantee. Since in this paper we use $L_1$ regularization as an example, for the $z$-minimization step, we utilize soft-thresholding technique from [19] to acquire the solution to minimize $L_\rho(x^{k+1}, z, y^k)$:

$$z^{k+1} \leftarrow \mathcal{S}_{\frac{\lambda}{\rho}}(x^{k+1} + y^k/\rho) \tag{18}$$

where soft-thresholding operator is defined as

$$\mathcal{S}_t(x)_i = \begin{cases} x_i - t & \text{if } x_i > t \\ x_i + t & \text{if } x_i < -t \\ 0 & \text{otherwise} \end{cases} \tag{19}$$

Similar technique has been used in [12] and [25].

Another ADMM based algorithm proposed in [28] (DP-ADMM) also used gradient perturbation technique. Our method differed from theirs for the following aspects: (i) DP-ADMM is used for distributed learning, so that the training objective is assigned into multiple parties each holding a portion of the data, instead in ssADMM it is the data dependent loss and regularization that are separated; (ii) in DP-ADMM, each party is perturbing full gradient and transmit to the center, so that there is no privacy amplification effect, therefore although both algorithms solve optimization approximately, their privacy loss is higher than ours at each step. Our methods differ from the ADMM-objP method (DPLL in [25]) for the following aspect: (i) ADMM-objP perturb the training objective at each iteration, and use full gradient descent multiple times to acquire exact solution at each iteration, which is not as efficient as ours, since our method only access a portion of data once at each step; (ii) ADMM-objP guarantees privacy only if exact solution is acquired at each step, therefore the privacy guarantee is only theoretically true. The privacy guarantee of ssADMM is given by Theorem 1.

**Theorem 1.** *Algorithm 1 is $(\alpha, \epsilon)$-RDP.*

*Proof.* We first show the $L_2$ sensitivity of batch gradient $g_k$. Assume neighboring mini-batches $B_i$ and $B_i'$ differ by one record $d_s \in B$ and $d_s \in B'$, by Definition 4,

$$
\begin{aligned}
\Delta_2^k(g) =& \Delta_2 \left[ \frac{1}{m} \sum_{d_i \in B_k} \nabla f(x^k, d_i) \right] \\
=& \sup_{B_k \sim B_k'} \left\| \frac{1}{m} \sum_{d_i \in B_k} \nabla f(x^k, d_i) - \frac{1}{m} \sum_{d_i \in B_k'} \nabla f(x^k, d_i) \right\|_2 \\
=& \frac{1}{m} \sup \| \nabla f(x^k, d_s) - \nabla f(x^k, d_s') \|_2 \leq \frac{2C}{m}
\end{aligned}
\tag{20}
$$

Let $\epsilon_k(\alpha) = \alpha(\Delta_2^k(g))^2/2\sigma^2$. So each iteration is $(\alpha, \epsilon_k(\alpha))$-RDP by Lemma 1, with respect to the batch $B_k$. Since $B_k$ is a randomized subsample of $D$, by Lemma 2, we can calculate $\epsilon_k'(\alpha)$ so that each iteration is $(\alpha, \epsilon_k'(\alpha))$-RDP with respect to $D$. Since the algorithm has run $T$ iterations, let $\epsilon = \sum_{k=0}^{T-1} \epsilon_k'(\alpha)$, by Lemma 4, Algorithm 1 is $(\alpha, \epsilon)$-RDP. $\square$

**Theorem 2.** *If we choose $\eta^k = O(1/\sqrt{k})$, and train for $t$ iterations, then Algorithm 1 has the expected convergence rate of $O(1/\sqrt{t})$.*

*Proof.* See proof in appendix. $\square$

### 4.2 Rényi differentially private model perturbation algorithm

Our model perturbation private sADMM algorithm (mpADMM) is presented in Algorithm 2. Different from perturbing the gradients, this algorithm use the unperturbed gradients to do model calculation for a whole step, and keep track of the $L_2$ sensitivity of all data-dependent model vectors. After each epoch, Gaussian noises are injected into model vectors $x, y, z$, and total privacy $\epsilon$ is updated, according to sensitivity and $\sigma^2$. Due to it is difficult to calculate the sensitivity over multiple epochs, we perform output perturbation after each epoch. Therefore, this algorithm can be considered as multiple-time output perturbation algorithm.

---

**Algorithm 2** RDP model perturbation sADMM $L_1$ regularized ERM algorithm (mpADMM)

---

1: **Input**: Dataset $D = \{d_1, ..., d_n\}$. Penalty parameter $\rho$, total epochs $T$.
2: **Initialize**: primal variables $x^0, z^0$, dual variable $y^0$.
3: **for** epoch $k = 0, 1, ..., T-1$ **do**
4:     $g_k \leftarrow \frac{1}{n} \sum_{d_i \in D} \nabla f(x^k, d_i)$                                     ▷ compute gradient
5:     Compute $x^{k+1}$ by (16)                                                                                ▷ primal variable $x$
6:     Compute $z^{k+1}$ by (18)                                                                                ▷ primal variable $z$
7:     Compute $y^{k+1}$ by (17)                                                                                ▷ dual variable $y$
8:     Sample $\gamma_1, \gamma_2, \gamma_3 \sim N(0, \sigma^2 \mathbf{I}_p)$
9:     $x^{k+1} \leftarrow x^{k+1} + \gamma_1, y^{k+1} = y^{k+1} + \gamma_2, z^{k+1} = z^{k+1} + \gamma_3$       ▷ perturb the model
10: **end for**
11: **Output**: $x^T$

---

To calculate the sensitivity, since unperturbed batch gradient is used here, after one epoch, all primal and dual variables are data-dependent. Assume neighboring datasets $D$ and $D'$ differ at position $s$: $d_s \in D$ and $d_s' \in D'$. We define $\delta_x := x - (x')$ where $x$ and $(x')$ are primal variables evaluated on $D$ and $D'$, respectively, after one epoch. Also, define $\delta_z^k$ and $\delta_y^k$ similarly. Then, after epoch $k$,

$$
\begin{aligned}
\delta_x^{k+1} =& x^{k+1} - (x')^{k+1} \\
=& \frac{1}{\rho + 1/\eta^k} \left( -\frac{1}{n} \sum_{d_i \in D} \nabla f(x^k, d_i) - y^k + \rho z^k + x^k/\eta^k \right) - \\
& \frac{1}{\rho + 1/\eta^k} \left( -\frac{1}{n} \sum_{d_i \in D'} \nabla f(x^k, d_i) - y^k + \rho z^k + x^k/\eta^k \right) \\
=& (\nabla f(x^k, d_s') - \nabla f(x^k, d_s))/n(1 + \eta^{k+1}\rho)
\end{aligned}
\tag{21}
$$

7

Consider when the soft-thresholding operator $\mathcal{S}_t$ (19) applied on two vectors $w$ and $w'$, and compare $\mathcal{S}_t(w) - \mathcal{S}_t(w')$ with $w - w'$ element-wise:

- If $w_i$ and $w'_i$ are of different signs, applying $\mathcal{S}$ on $w_i$ and $w'_i$ would bring them closer, therefore $|\mathcal{S}_t(w_i) - \mathcal{S}_t(w'_i)| < |w_i - w'_i|$;
- If $w_i$ and $w'_i$ are of the same sign, without loss of generality, let $|w_i| \leq |w'_i|$. One can easily observe that
  - If $t \leq |w_i| \leq |w'_i|$, then $|\mathcal{S}_t(w_i) - \mathcal{S}_t(w'_i)| = |(|w_i| - t) - (|w'_i| - t)| = |w_i - w'_i|$;
  - If $|w_i| < t < |w'_i|$, then $|\mathcal{S}_t(w_i) - \mathcal{S}_t(w'_i)| = |0 - (|w'_i| - t)| < |w_i - w'_i|$ since $t < |w'_i|$;
  - If $|w_i| \leq |w'_i| \leq t$, then $|\mathcal{S}_t(w_i) - \mathcal{S}_t(w'_i)| = 0 \leq |w_i - w'_i|$;

For vectors $u, v$, we can use $u \preccurlyeq v$ to represent $|u_i| < |v_i|$ and $u_i, v_i$ have the same sign, for each index $i$. Obviously $u \preccurlyeq v$ indicates $\|u\|_2 \leq \|v\|_2$. In either case above, we have $|\mathcal{S}_t(w_i) - \mathcal{S}_t(w'_i)| \leq |w_i - w'_i|$, and sign preserves (or becomes zero), so $\mathcal{S}_t(w) - \mathcal{S}_t(w') \preccurlyeq w - w'$ for any threshold $t$. Therefore,

$$\begin{aligned}
\delta_z^{k+1} &= z^{k+1} - (z')^{k+1} \\
&= \mathcal{S}_{\frac{\lambda}{\rho}}(x^{k+1} + y^k/\rho) - \mathcal{S}_{\frac{\lambda}{\rho}}((x')^{k+1} + y^k/\rho) \\
&\preccurlyeq x^{k+1} + y^k/\rho - ((x')^{k+1} + y^k/\rho) = \delta_x^{k+1}
\end{aligned} \tag{22}$$

and

$$\begin{aligned}
\delta_y^{k+1} &= y^{k+1} - (y')^{k+1} \\
&= y^k + \rho(x^{k+1} - z^{k+1}) - \left(y^k + \rho((x')^{k+1} - (z')^{k+1})\right) \\
&= \rho(\delta_x^{k+1} - \delta_z^{k+1}) \preccurlyeq \rho \delta_x^{k+1}
\end{aligned} \tag{23}$$

The last $\preccurlyeq$ holds because $\delta_z^{k+1} \preccurlyeq \delta_x^{k+1}$, the subtraction by $\delta_z^{k+1}$ only pushes each element of $\delta_x^{k+1}$ towards zero. So we have below conclusions for sensitivities of $x, z, y$ after epoch $k$:

$$\Delta_2^{k+1}(x) = \|\delta_x^{k+1}\|_2 \leq \frac{2C}{n(1 + \eta^{k+1}\rho)} \tag{24}$$

$$\Delta_2^{k+1}(z) = \|\delta_z^{k+1}\|_2 \leq \|\delta_x^{k+1}\|_2 \leq \frac{2C}{n(1 + \eta^{k+1}\rho)} \tag{25}$$

$$\Delta_2^{k+1}(y) = \|\delta_y^{k+1}\|_2 \leq \rho\|\delta_x^{k+1}\|_2 \leq \frac{2\rho C}{n(1 + \eta^{k+1}\rho)} \tag{26}$$

**Theorem 3.** *Algorithm 2 is $(\alpha, \epsilon)$-RDP.*

*Proof.* Let $\epsilon_{k+1,w}(\alpha) = \alpha(\Delta_2^{k+1}(w))^2/2\sigma^2$ for $w \in \{x, z, y\}$. By Lemma 1, each epoch is $(\alpha, \sum_{w \in \{x,z,y\}} \epsilon_{k+1,w}(\alpha))$-RDP, with respect to $D$. Since the algorithm has run $T$ epochs, by Lemma 4, let $\epsilon = \sum_{k=1}^{T} \sum_{w \in \{x,z,y\}} \epsilon_{k,w}(\alpha))$, then Algorithm 2 is $(\alpha, \epsilon)$-RDP. $\qquad\square$

## 5 Experimental Results

In this section we will present our experimental results on both real and simulated datasets. We will first show performance of classification on two real datasets, then show performance of both classification and feature selection on a synthetic dataset.

### 5.1 ERM models

We perform our experiments on $L_1$ regularized logistic regression and huberized SVM. The objective function of logistic regression is in (9). For huberized SVM, the objection function is

$$F(x, D) := \frac{1}{n}\sum_{i=1}^{n} \ell_{\text{huber}}(l_i x^T s_i) + \lambda\|x\|_1 \tag{27}$$

where

$$\ell_{\text{huber}}(z) := \begin{cases} 0 & \text{if } z > 1 + h \\ \frac{1}{4h}(1 + h - z)^2 & \text{if } |1 - z| \leq h \\ 1 - z & \text{otherwise} \end{cases} \tag{28}$$

is the huberized hinge loss (we set $h = 0.5$ in all experiments).

## 5.2 Baselines

Many differentially private ERM algorithms cannot be applied to $L_1$ regularized classification, such as ObjPert [1], [2], OutPert [40], PVP and DVP [5], PSGD [16], and RSGD [7]. Therefore, we compare our proposed algorithms with these baselines: DP-SGD [4], DP-ADMM [25], ADMM-objP [28], and Non-Private approach.

DP-SGD performs stochastic gradient descent with Gaussian perturbation. (Although their paper proposed moment accountant approach to analyze the privacy leak, we use Lemma 2 to analyze as we do on ssADMM, since it gives tighter bound on $\epsilon$.) For DP-SGD, when the algorithm requires taking gradient on $f(x^k, B_k) + \lambda\|x^k\|_1$, we use the proximal gradient technique

$$x^{k+1} \leftarrow \mathcal{S}_{\lambda\eta^k}[x^k - \eta^k \nabla f(x^k, B_k)] \tag{29}$$

to update $x^{k+1}$, as suggested in [41] and [42]. DP-ADMM is a distributed learning version of ADMM, where each party transfers perturbed primal variables to the center, and the center draw a consensus of the parties then transfer primal and dual variable back to each party. ADMM-objP is an ADMM version of the objective perturbation algorithm. At each iteration, the trainer optimize a perturbed unregulated objective function, therefore although the algorithm satisfies pure $\epsilon$-DP, in practice it is not really differentially private due to the objective function can only be approximately solved. According to their paper, we apply gradient descent enough times and assume the optimization problem is exactly solved at each iteration.

The DP-SVRG algorithm presented in [6] can also be applied on non-smooth regularizers, but we have implemented and found that, due to the extra privacy budget required to spent on perturbing the full gradient, with the high privacy range ($\epsilon \leq 1$), if we choose a large $\sigma^2$, the perturbed full gradient cannot help as a control variant to fasten the training, but actually slows down the minimization of empirical loss; if we choose a small $\sigma^2$, the privacy budget accumulates too fast and exceed our range in a few iterations. Therefore we have dropped this algorithm in our comparisons.

## 5.3 Datasets and Pre-proessing

Two real datasets on human subjects were used in our study: (i) the Adult dataset [43] was generated from 1994 US Census, with $n = 48,842$, $p = 124$, and the frequency of the majority label is 0.761; (ii) the IPUMS-BR dataset [44] was extracted from IPUMS data, with $n = 38,000$, $p = 53$, and the frequency of the majority label is 0.507.

To test the performance on feature selection, we created a synthetic dataset with many irrelevant features, using similar strategy as in [25]. To be specific, we generate a 100-dimension data $s_i \sim \mathcal{N}(0_{100}, \Sigma)$ where $\Sigma_{i,j} = 0.5^{|i-j|}$. Let $x$ be the true model, defined as $x_{1:10} = (0.5, 1, 1.5, 2, 2.5, 3, 3.5, 4, 4.5, 5)$, $x_{11:20} = -x_{1:10}$, and $x_{21:100} = (0, ..., 0)$. For the label of each row $l_i$, we sample the Bernoulli distribution with $P(l_i = 1) = 1/(1 + \exp(-x^T s_i + \iota))$, where $\iota \sim \mathcal{N}(0, 1)$ is a random noise. Therefore, to predict $l_i$, $s_i$ contains 20 relevant features and 80 irrelevant features. We generate 40,000 samples to constitute one dataset, the frequency of the majority label is 0.500. We only perform logistic regression on simulated data, since it is usually used for attribute selection.

We did 10-fold cross validation on each experiment for each algorithm, and due to randomness from noisy injection, we repeat each fold 10 times and report average classification accuracy and objective value on testing data. For the simulated data, we generated 10 datasets using the simulation strategy, and report the average performance.

An intercept is added into each dataset. All numerical attributes are re-scaled into [0, 1] by Min-Max scalar. For the algorithms requiring feature vector to have bounded $L_2$ norm, we normalize to make $\|x_i\| \leq 1$ for $i = 1, ..., n$.

## 5.4 Parameter setting

We keep $\delta = 10^{-8}$ for all experiments. For those algorithms satisfying RDP, we choose the best conversion to $(\epsilon, \delta)$-DP. In non-private settings, model users usually train a series models with different candidates of regularization coefficient $\lambda$, and select the one with highest testing performance. However, this process is data-dependent, therefore in private settings we cannot take a "best performing" coefficient for granted. Instead, we performed two group of experiments by two frequently using coefficients: low regularization with $\lambda = 0.0001$ and high regularization with $\lambda = 0.001$.

For ssADMM and DP-SGD, we set mini-batch size $m = \sqrt{n}$. We choose $\eta^k = \eta^0/h$ where $h$ is the current expected epoch (we consider every $n/m$ iterations as one expected epoch), since we find this schedule has the best performance for both algorithms, compare to a constant learning rate, or a decreasing one at a rate of $O(1/\sqrt{k})$. After tuning on the simulated data, we set penalty term $\rho = 0.25$ for ssADMM and $\rho = 0.5$ for mpADMM. For mpADMM, we use a constant learning rate $\eta$. For DP-ADMM, we assume there are 2 parties, each holding half of the data. (If there is only one party, DP-ADMM will reduce to DP-SGD with sampling ratio=1.) For ADMM-objP, at each iteration we optimize the perturbed objective function by full gradient descent running 20 epochs. Other parameters for DP-ADMM and ADMM-objP are set according to their paper.

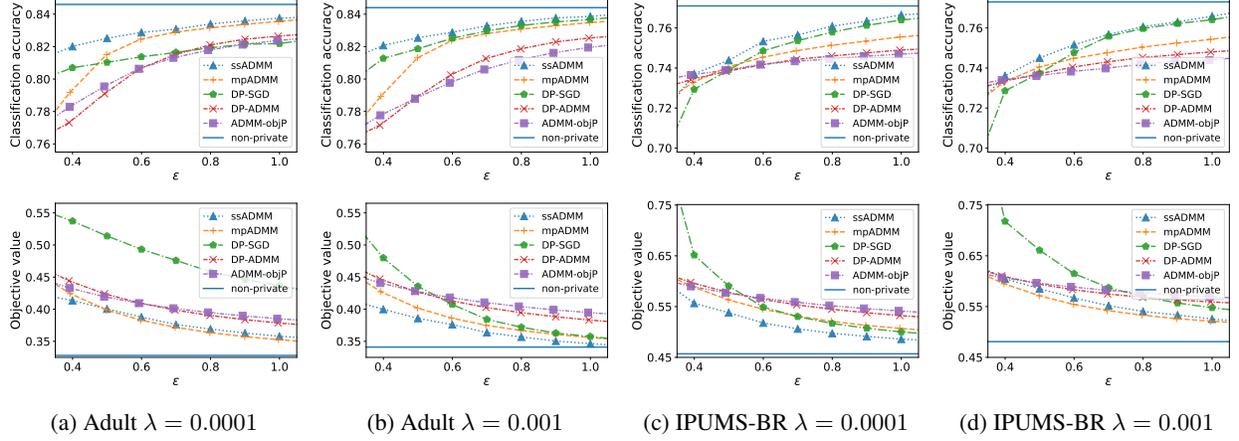## 5.5 Classification Performance on Real Data



Figure 1: Logistic regression result by $\epsilon$ (Top: Classification accuracy; Bottom: Objective value)
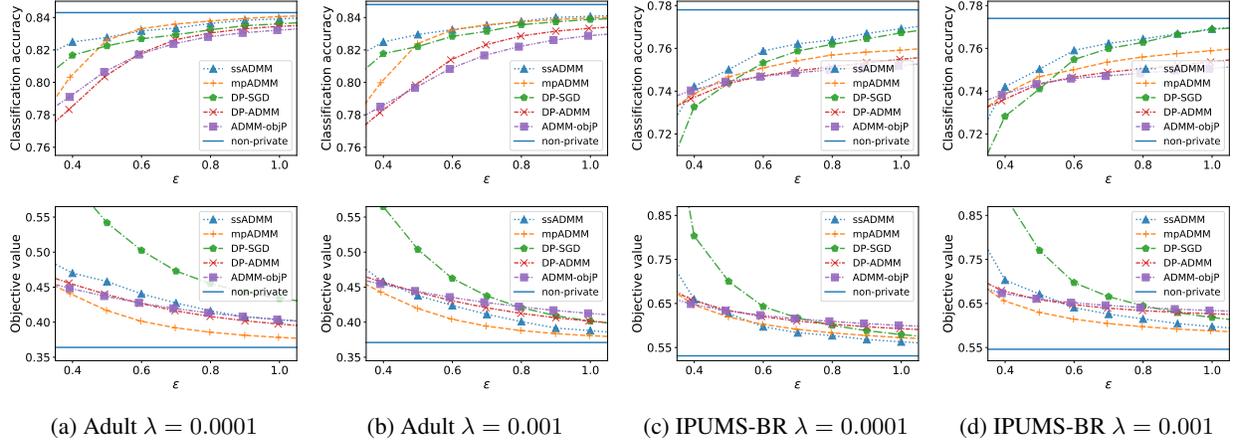
(a) Adult $\lambda = 0.0001$    (b) Adult $\lambda = 0.001$    (c) IPUMS-BR $\lambda = 0.0001$    (d) IPUMS-BR $\lambda = 0.001$



Figure 2: Huberized SVM result by $\epsilon$ (Top: Classification accuracy; Bottom: Objective value)

(a) Adult $\lambda = 0.0001$    (b) Adult $\lambda = 0.001$    (c) IPUMS-BR $\lambda = 0.0001$    (d) IPUMS-BR $\lambda = 0.001$

Figure 1 and Figure 2 plots the testing data accuracy (top) and objective values (bottom) of the algorithms trading off with privacy parameter $\epsilon$, for $L_1$ regularized logistic regression and huberized SVM, respectively. We can see for classification accuracy, ssADMM outperforms other algorithms in most cases. This is in accordance with the experiment in [12] that sADMM outperforms proximal gradient in non-private setting. [45] also show that ADMM based algorithms are more robust to noisy data with outliers. Although DP-SGD has better classification accuracy than mpADMM in some cases, its objective value is usually higher. DP-ADMM and ADMM-objP can achieve high utility when $\epsilon$ gets high, but in our testing range of $\epsilon$, they cannot perform as good as other algorithms. mpADMM performs better in adult dataset than in IPUMS-BR dataset, probably because Adult dataset is more sparse compare to IPUMS-BR, due to it is binary transferred through one-hot encoding. And that model perturbation are more robust to data with irrelevant attributes is in accordance with our observations on the simulated data.

## 5.6 Performance on Simulated Data

To measure the attribute selection performance, we test how many relevant attributes are selected by each algorithm for $L_1$ regularized logistic regression. Since the dataset is standardized, we can use the magnitude of the coefficient to rank the attributes, due to that noisy perturbation might cause the coefficients of irrelevant attributes slightly differ from zero.

We define a criterion $\xi_k$ to measure the coverage of relevant attributes if top $k$ attributes suggested by the algorithm were selected. For example, since we know there are 20 relevant attributes in the simulated data, if we select $k = 30$ attributes by magnitude of coefficient, 16 of them are the true relevant ones (i.e. among $x_1, ..., x_{20}$), then $\xi_{30} = 16/20 = 0.8$.
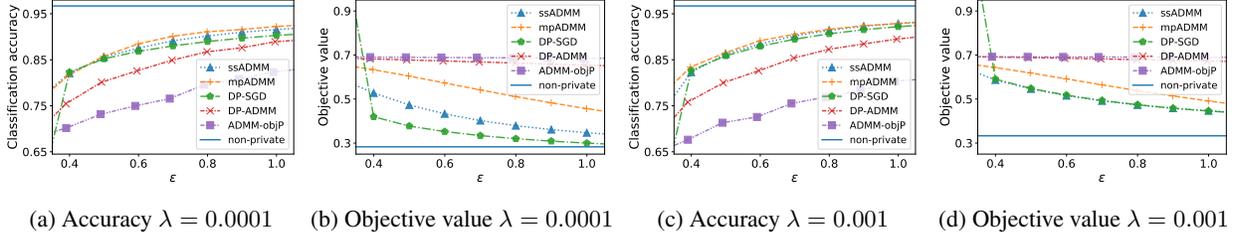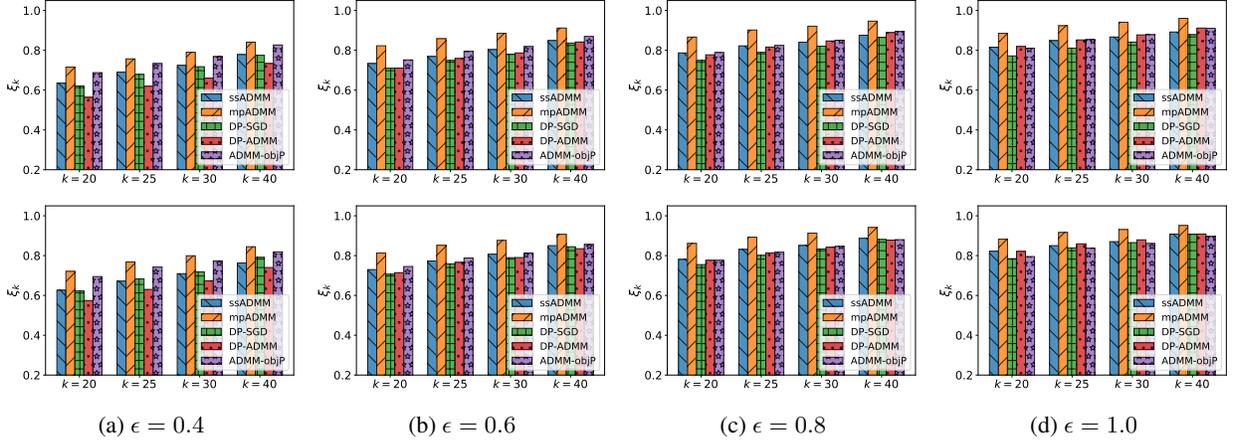
Figure 3: Classification performance on simulated data



Figure 4: Attribute selection performance on simulated data (Top: $\lambda = 0.0001$; Bottom: $\lambda = 0.001$)

This make sense because in real case, the number of attributes we choose to select from an attribute ranker depends on the budget we can spend to collect data. We test all algorithms for $k = 20, 25, 30$, and $40$.

Figure 3 shows the classification performance of each algorithm on the simulated data. For non-private performance, we assume the true model is known. We can see that ssADMM, mpADMM, and DP-SGD have similar performance in classification. Figure 4 shows the performance of attribute selection. Although classification accuracy are close, we can see that mpADMM can detect more relevant attributes, especially in the lower $\epsilon$ range. ADMM-objP, which was originally proposed for feature selection, can outperform ssADMM and DP-SGD for feature selection in low $\epsilon$ while its classification accuracy is behind ssADMM and DP-SGD. However, ADMM-objP usually require much more epochs in training compare to the other algorithms. Therefore, if we know the data is sparse and the major goal is focused on attribute selection, mpADMM is more preferable.

## 6 Conclusions

We present two privatizations of stochastic ADMM under Rényi differential privacy. One algorithm combines gradient perturbation technique with privacy amplification result to reduce the total privacy loss throughout the execution. The other algorithm uses the output perturbation (with numerical computation of sensitivity) to privately release the solution at the end of each training epoch. These algorithms can be used to solve optimization problems with complex structural regularization that induces sparsity.

## References

[1] Kamalika Chaudhuri, Claire Monteleoni, and Anand D Sarwate. Differentially private empirical risk minimization. *Journal of Machine Learning Research*, 12(Mar):1069–1109, 2011.

[2] Daniel Kifer, Adam Smith, and Abhradeep Thakurta. Private convex empirical risk minimization and high-dimensional regression. In *Conference on Learning Theory*, pages 25–1, 2012.

[3] Raef Bassily, Adam Smith, and Abhradeep Thakurta. Private empirical risk minimization: Efficient algorithms and tight error bounds. In *Foundations of Computer Science (FOCS), 2014 IEEE 55th Annual Symposium on*, pages 464–473. IEEE, 2014.

[4] Martin Abadi, Andy Chu, Ian Goodfellow, H Brendan McMahan, Ilya Mironov, Kunal Talwar, and Li Zhang. Deep learning with differential privacy. In *Proceedings of the 2016 ACM SIGSAC Conference on Computer and Communications Security*, pages 308–318. ACM, 2016.

[5] Jiaqi Zhang, Kai Zheng, Wenlong Mou, and Liwei Wang. Efficient private erm for smooth objectives. *arXiv preprint arXiv:1703.09947*, 2017.

[6] Di Wang, Minwei Ye, and Jinhui Xu. Differentially private empirical risk minimization revisited: Faster and more general. In *Advances in Neural Information Processing Systems*, pages 2722–2731, 2017.

[7] Chen Chen, Jaewoo Lee, and Dan Kifer. Renyi differentially private erm for smooth objectives. In *The 22nd International Conference on Artificial Intelligence and Statistics*, pages 2037–2046, 2019.

[8] Ryan J Tibshirani, Jonathan Taylor, et al. The solution path of the generalized lasso. *The Annals of Statistics*, 39(3):1335–1371, 2011.

[9] Xiao Zhang, Lingxiao Wang, Yaodong Yu, and Quanquan Gu. A primal-dual analysis of global optimality in nonconvex low-rank matrix recovery. In *International conference on machine learning*, pages 5857–5866, 2018.

[10] Guangcan Liu, Qingshan Liu, and Ping Li. Blessing of dimensionality: Recovering mixture data via dictionary pursuit. *IEEE transactions on pattern analysis and machine intelligence*, 39(1):47–60, 2016.

[11] Daniel Gabay and Bertrand Mercier. A dual algorithm for the solution of nonlinear variational problems via finite element approximation. *Computers & Mathematics with Applications*, 2(1):17–40, 1976.

[12] Hua Ouyang, Niao He, Long Tran, and Alexander Gray. Stochastic alternating direction method of multipliers. In *International Conference on Machine Learning*, pages 80–88, 2013.

[13] Jianqing Fan and Runze Li. Variable selection via nonconcave penalized likelihood and its oracle properties. *Journal of the American statistical Association*, 96(456):1348–1360, 2001.

[14] Cun-Hui Zhang et al. Nearly unbiased variable selection under minimax concave penalty. *The Annals of statistics*, 38(2):894–942, 2010.

[15] Yu-Xiang Wang, Borja Balle, and Shiva Kasiviswanathan. Subsampled r\'enyi differential privacy and analytical moments accountant. *arXiv preprint arXiv:1808.00087*, 2018.

[16] Xi Wu, Fengan Li, Arun Kumar, Kamalika Chaudhuri, Somesh Jha, and Jeffrey Naughton. Bolt-on differential privacy for scalable stochastic gradient descent-based analytics. In *Proceedings of the 2017 ACM International Conference on Management of Data*, pages 1307–1322. ACM, 2017.

[17] Antti Koskela and Antti Honkela. Learning rate adaptation for differentially private stochastic gradient descent. *arXiv preprint arXiv:1809.03832*, 2018.

[18] Ilya Mironov. Rényi differential privacy. In *2017 IEEE 30th Computer Security Foundations Symposium (CSF)*, pages 263–275. IEEE, 2017.

[19] Stephen Boyd, Neal Parikh, Eric Chu, Borja Peleato, Jonathan Eckstein, et al. Distributed optimization and statistical learning via the alternating direction method of multipliers. *Foundations and Trends® in Machine learning*, 3(1):1–122, 2011.

[20] Huachun Tan, Jianshuai Feng, Guangdong Feng, Wuhong Wang, and Yu-Jin Zhang. Traffic volume data outlier recovery via tensor model. *Mathematical Problems in Engineering*, 2013, 2013.

[21] Stanley H Chan, Xiran Wang, and Omar A Elgendy. Plug-and-play admm for image restoration: Fixed-point convergence and applications. *IEEE Transactions on Computational Imaging*, 3(1):84–98, 2016.

[22] Neil K Dhingra, Mihailo R Jovanović, and Zhi-Quan Luo. An admm algorithm for optimal sensor and actuator selection. In *53rd IEEE Conference on Decision and Control*, pages 4039–4044. IEEE, 2014.

[23] Ernie Esser. Applications of lagrangian-based alternating direction methods and connections to split bregman. *CAM report*, 9:31, 2009.

[24] Junfeng Yang and Xiaoming Yuan. Linearized augmented lagrangian and alternating direction methods for nuclear norm minimization. *Mathematics of computation*, 82(281):301–329, 2013.

[25] Puyu Wang and Hai Zhang. Differential privacy for sparse classification learning. *arXiv preprint arXiv:1908.00780*, 2019.

[26] Tao Zhang and Quanyan Zhu. Dynamic differential privacy for admm-based distributed classification learning. *IEEE Transactions on Information Forensics and Security*, 12(1):172–187, 2017.

[27] Xueru Zhang, Mohammad Mahdi Khalili, and Mingyan Liu. Improving the privacy and accuracy of admm-based distributed algorithms. *arXiv preprint arXiv:1806.02246*, 2018.

[28] Zonghao Huang, Rui Hu, Yuanxiong Guo, Eric Chan-Tin, and Yanmin Gong. Dp-admm: Admm-based distributed learning with differential privacy. *IEEE Transactions on Information Forensics and Security*, 2019.

[29] Robert Tibshirani. Regression shrinkage and selection via the lasso. *Journal of the Royal Statistical Society: Series B (Methodological)*, 58(1):267–288, 1996.

[30] Hui Zou. The adaptive lasso and its oracle properties. *Journal of the American statistical association*, 101(476):1418–1429, 2006.

[31] Trevor Park and George Casella. The bayesian lasso. *Journal of the American Statistical Association*, 103(482):681–686, 2008.

[32] Su-In Lee, Honglak Lee, Pieter Abbeel, and Andrew Y Ng. Efficient l˜ 1 regularized logistic regression. In *AAAI*, volume 6, pages 401–408, 2006.

[33] Mee Young Park and Trevor Hastie. L1-regularization path algorithm for generalized linear models. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, 69(4):659–677, 2007.

[34] Yatao An Bian, Xiong Li, Yuncai Liu, and Ming-Hsuan Yang. Parallel coordinate descent newton method for efficient $l\_1$-regularized loss minimization. *IEEE transactions on neural networks and learning systems*, 2019.

[35] Andrew Y Ng. Feature selection, l 1 vs. l 2 regularization, and rotational invariance. In *Proceedings of the twenty-first international conference on Machine learning*, page 78. ACM, 2004.

[36] Joshua Goodman. Exponential priors for maximum entropy models. In *Proceedings of the Human Language Technology Conference of the North American Chapter of the Association for Computational Linguistics: HLT-NAACL 2004*, pages 305–312, 2004.

[37] Cynthia Dwork, Frank McSherry, Kobbi Nissim, and Adam Smith. Calibrating noise to sensitivity in private data analysis. In *Theory of cryptography conference*, pages 265–284. Springer, 2006.

[38] Cynthia Dwork, Krishnaram Kenthapadi, Frank McSherry, Ilya Mironov, and Moni Naor. Our data, ourselves: Privacy via distributed noise generation. In *Annual International Conference on the Theory and Applications of Cryptographic Techniques*, pages 486–503. Springer, 2006.

[39] Mark Bun and Thomas Steinke. Concentrated differential privacy: Simplifications, extensions, and lower bounds. In *Theory of Cryptography Conference*, pages 635–658. Springer, 2016.

[40] Tao Zhang and Quanyan Zhu. Dynamic differential privacy for admm-based distributed classification learning. *IEEE Transactions on Information Forensics and Security*, 12(1):172–187, 2016.

[41] John Duchi and Yoram Singer. Efficient online and batch learning using forward backward splitting. *Journal of Machine Learning Research*, 10(Dec):2899–2934, 2009.

[42] Patrick L Combettes and Jean-Christophe Pesquet. Proximal splitting methods in signal processing. In *Fixed-point algorithms for inverse problems in science and engineering*, pages 185–212. Springer, 2011.

[43] Chih-Chung Chang and Chih-Jen Lin. Libsvm: A library for support vector machines. *ACM transactions on intelligent systems and technology (TIST)*, 2(3):27, 2011.

[44] Steven Ruggles, Katie Genadek, Ronald Goeken, Josiah Grover, and Matthew Sobek. Integrated public use microdata series: Version 6.0 [dataset]. *Minneapolis: University of Minnesota*, 23:56, 2015.

[45] Samaneh Azadi, Jiashi Feng, Stefanie Jegelka, and Trevor Darrell. Auxiliary image regularization for deep cnns with noisy labels. *arXiv preprint arXiv:1511.07069*, 2015.

# A Proof of Theorem 2

The proof is done by applying similar technique for Theorem 1 in [12], plus considering the Gaussian noise term added. Define

$$u := \begin{pmatrix} x \\ z \end{pmatrix}, \overline{u}^k := \begin{pmatrix} \frac{1}{k} \sum_{i=1}^{k-1} x^i \\ \frac{1}{k} \sum_{i=1}^{k-1} z^i \end{pmatrix}, \theta(u) := f(x) + h(z),$$

and define

$$w := \begin{pmatrix} x \\ z \\ y \end{pmatrix}, \overline{w}^k := \begin{pmatrix} \frac{1}{k}\sum_{i=1}^{k-1} x^i \\ \frac{1}{k}\sum_{i=1}^{k-1} z^i \\ \frac{1}{k}\sum_{i=1}^{k-1} y^i \end{pmatrix}, F(w) := \begin{pmatrix} -y \\ y \\ x - z \end{pmatrix}$$

Denote $u^* := \begin{pmatrix} x^* \\ z^* \end{pmatrix}$ as the optimal solution, and $\delta_{k+1} := \nabla f(x^k, B_k) - \nabla f(x^k, D)$, $d_{\mathcal{X}} := \sup_{x_a, x_b \in \mathcal{X}} \|x_a - x_b\|$, $d_{y^*} := \|y^0 - y^*\|$.

Therefore, consider the expectation of $\theta(\overline{u}^t) - \theta(u^*)$ after $t$ iterations,

$$\mathbb{E}\left[\theta(\overline{u}^t) - \theta(u^*) + (\overline{w}^t - w^*)^T F(\overline{w}^t)\right]$$

$$= \mathbb{E}\left[\theta(\overline{u}^t) - \theta(u^*) + (\overline{x}^t - x^*)^T(-\overline{y}^t) + (\overline{z}^t - z^*)^T(\overline{y}^t)\right.$$

$$\left. + (\overline{y} - y)^T(\overline{x}^t - \overline{z}^t)\right]$$

$$\leq \mathbb{E}\left[\frac{1}{t}\sum_{k=0}^{t-1}\left[\frac{\eta^k}{2}\|\nabla f(x^k, B_k) + \gamma^k\|^2 + \frac{1}{2\eta^k}(\|x^k - x^*\|^2 - \|x^{k+1} - x^*\|^2)\right.\right.$$

$$\left.\left. + \langle \delta_{k+1}, x^* - x^k \rangle\right] + \frac{1}{t}\left(\frac{\rho}{2}\|x^* - z^0\|^2 + \frac{1}{2\rho}\|y - y^0\|^2\right)\right]$$

$$\leq \mathbb{E}\left[\frac{1}{t}\sum_{k=0}^{t-1}\left[\frac{\eta^k(C^2 + p\sigma^2)}{2} + \langle \delta_{k+1}, x^* - x^k \rangle\right]\right.$$

$$\left. + \frac{1}{t}\left(\frac{d_{\mathcal{X}}^2}{2\eta^{t-1}} + \frac{\rho}{2}d_{y^*}^2 + \frac{1}{2\rho}\|y - y^0\|^2\right)\right]$$

$$= \mathbb{E}\left[\frac{1}{t}\sum_{k=0}^{t-1}\left[\frac{\eta^k(C^2 + p\sigma^2)}{2}\right] + \frac{1}{t}\left(\frac{d_{\mathcal{X}}^2}{2\eta^{t-1}} + \frac{\rho}{2}d_{y^*}^2 + \frac{1}{2\rho}\|y - y^0\|^2\right)\right] \qquad (30)$$

while the first inequality holds by applying an expected version of Lemma 2 in [12], note that since noisy perturbation $\gamma \sim \mathcal{N}(0, \sigma^2 \mathbf{I}_p)$, $\mathbb{E}[\nabla f(x^k, B_k) + \gamma] = \nabla f(x^k, B_k)$, and $\mathbb{E}[\|\nabla f(x^k, B_k) + \gamma^k\|^2] \leq \mathbb{E}[\|\nabla f(x^k, B_k)\|^2] + \mathbb{E}[\|\gamma\|^2] + 2\mathbb{E}[\|\nabla f(x^k, B_k)\|]\mathbb{E}[\gamma] \leq C^2 + p\sigma^2$. The last equality holds because we assume $x^k$ is independent of $B_k$ (which was used to calculate $x^{k+1}$) is independent of $x^k$, hence $\mathbb{E}_{B_k|B_{[0:k-1]}}\langle \delta_{k+1}, x^* - x^k \rangle = 0$.

The above holds for all dual variable $y$, hence it holds for $y$ in a ball $\mathcal{B}_0 = \{y : \|y\|_2 \leq \beta\}$. According to (33) in [12],

$$\max_{y \in \mathcal{B}_0}\{\theta(\overline{u}^t) - \theta(u^*) + (\overline{w}^t - w^*)^T F(\overline{w}^t)\} = \theta(\overline{u}^t) - \theta(u^*) + \beta\|\overline{x}_t - \overline{z}_t\| \qquad (31)$$

Therefore, continue on (30), we can have

$$\mathbb{E}\left[\theta(\overline{u}^t) - \theta(u^*) + \beta\|\overline{x}_t - \overline{z}_t\|\right]$$

$$\leq \mathbb{E}\left[\frac{1}{t}\sum_{k=0}^{t-1}\left[\frac{\eta^k(C^2 + p\sigma^2)}{2}\right] + \frac{1}{t}\left(\frac{d_{\mathcal{X}}^2}{2\eta^{t-1}} + \frac{\rho}{2}d_{y^*}^2 + \frac{1}{2\rho}\|y - y^0\|^2\right)\right]$$

$$\leq \mathbb{E}\left[\frac{1}{t}\sum_{k=0}^{t-1}\left[\frac{\eta^k(C^2 + p\sigma^2)}{2}\right] + \frac{1}{t}\left(\frac{d_{\mathcal{X}}^2}{2\eta^{t-1}} + \frac{\rho}{2}d_{y^*}^2\right)\right] \qquad (32)$$

$$+ \mathbb{E}\left[\max_{y \in \mathcal{B}_0}\{\frac{1}{2\rho t}\|y - y_0\|^2\right]$$

$$\leq \frac{1}{t}\left(\frac{C^2 + p\sigma^2}{2}\sum_{k=1}^{t}\eta^k + \frac{d_{\mathcal{X}}^2}{2\eta^{t-1}}\right) + \frac{\rho d_{y^*}^2}{2t} + \frac{\beta^2}{2\rho t}$$

So if we choose $\eta^k = \frac{d_{\mathcal{X}}}{\sqrt{2(C^2 + p\sigma^2)k}} = O(1/\sqrt{k})$, $\mathbb{E}\left[\theta(\overline{u}^t) - \theta(u^*) + \beta\|\overline{x}_t - \overline{z}_t\|\right] \leq \frac{d_{\mathcal{X}}\sqrt{2(C^2 + p\sigma^2)}}{\sqrt{t}} + \frac{\rho d_{y^*}^2}{2t} + \frac{\beta^2}{2\rho t} = O(1/\sqrt{t})$.