# Zeroth-Order Online Alternating Direction Method of Multipliers: Convergence Analysis and Applications

**Sijia Liu**
University of Michigan
IBM Research, Cambridge

**Jie Chen**
Northwestern Polytechnical
University, China

**Pin-Yu Chen**
IBM Research,
Yorktown Heights

**Alfred O. Hero**
University of Michigan

## Abstract

In this paper, we design and analyze a new zeroth-order online algorithm, namely, the zeroth-order online alternating direction method of multipliers (ZOO-ADMM), which enjoys dual advantages of being gradient-free operation and employing the ADMM to accommodate complex structured regularizers. Compared to the first-order gradient-based online algorithm, we show that ZOO-ADMM requires $\sqrt{m}$ times more iterations, leading to a convergence rate of $O(\sqrt{m}/\sqrt{T})$, where $m$ is the number of optimization variables, and $T$ is the number of iterations. To accelerate ZOO-ADMM, we propose two minibatch strategies: gradient sample averaging and observation averaging, resulting in an improved convergence rate of $O(\sqrt{1 + q^{-1}m}/\sqrt{T})$, where $q$ is the minibatch size. In addition to convergence analysis, we also demonstrate ZOO-ADMM to applications in signal processing, statistics, and machine learning.

## 1 Introduction

Online convex optimization (OCO) performs sequential inference in a data-driven adaptive fashion, and has found a wide range of applications (Hall and Willett, 2015; Hazan, 2016; Hosseini et al., 2016). In this paper, we focus on regularized convex optimization in the OCO setting, where a cumulative empirical loss is minimized together with a fixed regularization term. Regularized loss minimization is a common learning paradigm, which has been very effective in promotion of sparsity through $\ell_1$ or mixed $\ell_1/\ell_2$ regularization (Bach et al., 2012), low-rank matrix completion via nuclear norm regularization (Candès and Recht, 2009), graph signal recovery via graph Laplacian regularization (Chen and Liu, 2017), and constrained optimization by imposing indicator functions of constraint sets (Parikh and Boyd, 2014).

Several OCO algorithms have been proposed for regularized optimization, e.g., composite mirror descent, namely, proximal stochastic gradient descent (Duchi et al., 2010), regularized dual averaging (Xiao, 2010), and adaptive gradient descent (Duchi et al., 2011). However, the complexity of the aforementioned algorithms is dominated by the computation of the proximal operation with respect to the regularizers (Parikh and Boyd, 2014). An alternative is to use online alternating direction method of multipliers (O-ADMM) (Ouyang et al., 2013; Suzuki, 2013; Wang and Banerjee, 2013). Different from the algorithms in (Duchi et al., 2010, 2011; Xiao, 2010), the ADMM framework offers the possibility of splitting the optimization problem into a sequence of easily-solved subproblems. It was shown in (Ouyang et al., 2013; Suzuki, 2013; Wang and Banerjee, 2013) that the online variant of ADMM has convergence rate of $O(1/\sqrt{T})$ for convex loss functions and $O(\log T/T)$ for strongly convex loss functions, where $T$ is the number of iterations.

One limitation of existing O-ADMM algorithms is the need to compute and repeatedly evaluate the gradient of the loss function over the iterations. In many practical scenarios, an explicit expression for the gradient is difficult to obtain. For example, in bandit optimization (Agarwal et al., 2010), a player receives partial feedback in terms of loss function values revealed by her adversary, and making it impossible to compute the gradient of the full loss function. In adversarial black-box machine learning models, only the function values (e.g., prediction results) are provided (Chen et al., 2017). Moreover, in some high dimensional settings, acquiring the gradient information may be difficult, e.g., involving matrix inversion (Boyd and Vandenberghe, 2004). This motivates the development of gradient-free (zeroth-order) optimization algorithms.

arXiv:1710.07804v2 [stat.ML] 18 Feb 2018

Zeroth-order optimization approximates the full gradient via a randomized gradient estimate (Agarwal et al., 2010; Duchi et al., 2015; Ghadimi and Lan, 2013; Hajinezhad et al., 2017; Nesterov and Spokoiny, 2015; Shamir, 2017). For example, in (Agarwal et al., 2010; Shamir, 2017), zeroth-order algorithms were developed for bandit convex optimization with multi-point bandit feedback. In (Nesterov and Spokoiny, 2015), a zeroth-order gradient descent algorithm was proposed that has $O(m/\sqrt{T})$ convergence rate, where $m$ is the number of variables in the objective function. A similar convergence rate was found in (Ghadimi and Lan, 2013) for nonconvex optimization. This slowdown (proportional to the problem size $m$) in convergence rate was further improved to $O(\sqrt{m}/\sqrt{T})$ (Duchi et al., 2015), whose optimality was proved under the framework of mirror descent algorithms. A more recent relevant paper is (Gao et al., 2017), where a variant of the ADMM algorithm that uses gradient estimation was introduced. However, the ADMM algorithm presented in (Gao et al., 2017) was not customized for OCO. Furthermore, it only ensured that the linear equality constraints are satisfied in expectation; hence, a particular instance of the proposed solution could violate the constraints.

In this paper, we propose a zeroth-order online ADMM (called ZOO-ADMM) algorithm, and analyze its convergence rate under different settings, including stochastic optimization, learning with strongly convex loss functions, and minibatch strategies for convergence acceleration. We summarize our contributions as follows.

• We integrate the idea of zeroth-order optimization with online ADMM, leading to a new gradient-free OCO algorithm, ZOO-ADMM.

• We prove ZOO-ADMM yields a $O(\sqrt{m}/\sqrt{T})$ convergence rate for smooth+nonsmooth composite objective functions.

• We introduce a general hybrid minibatch strategy for acceleration of ZOO-ADMM, leading to an improved convergence rate $O(\sqrt{1+q^{-1}m}/\sqrt{T})$, where $q$ is the minibatch size.

• We illustrate the practical utility of ZOO-ADMM in machine leanring, signal processing and statistics.

## 2   ADMM: from First to Zeroth Order

In this paper, we consider the regularized loss minimization problem over a time horizon of length $T$

$$\begin{array}{ll}
\underset{\mathbf{x}\in\mathcal{X},\mathbf{y}\in\mathcal{Y}}{\text{minimize}} & \dfrac{1}{T}\sum_{t=1}^{T} f(\mathbf{x};\mathbf{w}_t)+\phi(\mathbf{y}) \\
\text{subject to} & \mathbf{Ax}+\mathbf{By}=\mathbf{c},
\end{array} \quad (1)$$

where $\mathbf{x}\in\mathbb{R}^m$ and $\mathbf{y}\in\mathbb{R}^d$ are optimization variables, $\mathcal{X}$ and $\mathcal{Y}$ are closed convex sets, $f(\cdot;\mathbf{w}_t)$ is a convex and smooth cost/loss function parameterized by $\mathbf{w}_t$ at time $t$, $\phi$ is a convex regularization function (possibly nonsmooth), and $\mathbf{A}\in\mathbb{R}^{l\times m}$, $\mathbf{B}\in\mathbb{R}^{l\times d}$, and $\mathbf{c}\in\mathbb{R}^l$ are appropriate coefficients associated with a system of $l$ linear constraints.

In problem (21), the use of time-varying cost functions $\{f(\mathbf{x};\mathbf{w}_t)\}_{t=1}^T$ captures possibly time-varying environmental uncertainties that may exist in the online setting (Hazan, 2016; Shalev-Shwartz, 2012). We can also write the online cost as $f_t(\mathbf{x})$ when it cannot be explicitly parameterized by $\mathbf{w}_t$. One interpretation of $\{f(\mathbf{x};\mathbf{w}_t)\}_{t=1}^T$ is the empirical approximation to the stochastic objective function $\mathbb{E}_{\mathbf{w}\sim P}[f(\mathbf{x};\mathbf{w})]$. Here $P$ is an empirical distribution with density $\sum_t \delta(\mathbf{w},\mathbf{w}_t)$, where $\{\mathbf{w}_t\}_{t=1}^T$ is a set of i.i.d. samples, and $\delta(\cdot,\mathbf{w}_t)$ is the Dirac delta function at $\mathbf{w}_t$. We also note that when $\mathcal{Y}=\mathcal{X}$, $l=m$, $\mathbf{A}=\mathbf{I}_m$, $\mathbf{B}=-\mathbf{I}_m$, $\mathbf{c}=\mathbf{0}_m$, the variable $\mathbf{y}$ and the linear constraint in (21) can be eliminated, leading to a standard OCO formulation. Here $\mathbf{I}_m$ denotes the $m\times m$ identity matrix, and $\mathbf{0}_m$ is the $m\times 1$ vector of all zeros[1].

### 2.1   Background on O-ADMM

O-ADMM (Ouyang et al., 2013; Suzuki, 2013; Wang and Banerjee, 2013) was originally proposed to extend batch-type ADMM methods to the OCO setting. For solving (21), a widely-used algorithm was developed by (Suzuki, 2013), which combines online proximal gradient descent and ADMM in the following form:

$$\begin{aligned}
\mathbf{x}_{t+1} = \underset{\mathbf{x}\in\mathcal{X}}{\arg\min}\Big\{ & \mathbf{g}_t^T\mathbf{x}-\boldsymbol{\lambda}_t^T(\mathbf{Ax}+\mathbf{By}_t-\mathbf{c}) \\
& +\frac{\rho}{2}\|\mathbf{Ax}+\mathbf{By}_t-\mathbf{c}\|_2^2+\frac{1}{2\eta_t}\|\mathbf{x}-\mathbf{x}_t\|_{\mathbf{G}_t}^2\Big\}, \quad (2)
\end{aligned}$$

$$\begin{aligned}
\mathbf{y}_{t+1} = \underset{\mathbf{y}\in\mathcal{Y}}{\arg\min}\Big\{ & \phi(\mathbf{y})-\boldsymbol{\lambda}_t^T(\mathbf{Ax}_{t+1}+\mathbf{By}-\mathbf{c}) \\
& +\frac{\rho}{2}\|\mathbf{Ax}_{t+1}+\mathbf{By}-\mathbf{c}\|_2^2\Big\}, \quad (3)
\end{aligned}$$

$$\boldsymbol{\lambda}_{t+1}=\boldsymbol{\lambda}_t-\rho(\mathbf{Ax}_{t+1}+\mathbf{By}_{t+1}-\mathbf{c}), \quad (4)$$

where $t$ is the iteration number (possibly the same as the time step), $\mathbf{g}_t$ is the gradient of the cost function $f(\mathbf{x};\mathbf{w}_t)$ at $\mathbf{x}_t$, namely, $\mathbf{g}_t=\nabla_\mathbf{x}f(\mathbf{x};\mathbf{w}_t)|_{\mathbf{x}=\mathbf{x}_t}$, $\boldsymbol{\lambda}_t$ is a Lagrange multiplier (also known as the dual variable), $\rho$ is a positive weight to penalize the augmented term associated with the equality constraint of (21), $\|\cdot\|_2$ denotes the $\ell_2$ norm, $\eta_t$ is a non-increasing sequence of positive step sizes, and $\|\mathbf{x}-\mathbf{x}_t\|_{\mathbf{G}_t}^2=(\mathbf{x}-\mathbf{x}_t)^T\mathbf{G}_t(\mathbf{x}-\mathbf{x}_t)$ is a Bregman divergence generated by the strongly convex function $(1/2)\mathbf{x}^T\mathbf{G}_t\mathbf{x}$ with a known symmetric positive definite coefficient matrix $\mathbf{G}_t$.

---

[1] In the sequel we will omit the dimension index $m$, which can be inferred from the context.

Similar to batch-type ADMM algorithms, the sub-problem in (23) is often easily solved via the proximal operator with respect to $\phi$ (Boyd et al., 2011). However, one limitation of O-ADMM is that it requires the gradient $\mathbf{g}_t$ in (2). We will develop the gradient-free (zeroth-order) O-ADMM algorithm below that relaxes this requirement.

## 2.2 Motivation of ZOO-ADMM

To avoid explicit gradient calculations in (2), we adopt a random gradient estimator to estimate the gradient of a smooth cost function (Duchi et al., 2015; Ghadimi and Lan, 2013; Nesterov and Spokoiny, 2015; Shamir, 2017). The gradient estimate of $f(\mathbf{w}; \mathbf{w}_t)$ is given by

$$\hat{\mathbf{g}}_t = \frac{f(\mathbf{x}_t + \beta_t \mathbf{z}_t; \mathbf{w}_t) - f(\mathbf{x}_t; \mathbf{w}_t)}{\beta_t} \mathbf{z}_t, \qquad (5)$$

where $\mathbf{z}_t \in \mathbb{R}^m$ is a random vector drawn independently at each iteration $t$ from a distribution $\mathbf{z} \sim \mu$ with $\mathbb{E}_\mu[\mathbf{z}\mathbf{z}^T] = \mathbf{I}$, and $\{\beta_t\}$ is a non-increasing sequence of small positive smoothing constants. Here for notational simplicity we replace $\{\}_{t=1}^T$ with $\{\}$. The rationale behind the estimator (5) is that $\hat{\mathbf{g}}_t$ becomes an unbiased estimator of $\mathbf{g}_t$ when the smoothing parameter $\beta_t$ approaches zero (Duchi et al., 2015).

After replacing $\mathbf{g}_t$ with $\hat{\mathbf{g}}_t$ in (5), the resulting algorithm (2)-(24) can be implemented without explicit gradient computation. This extension is called zeroth-order O-ADMM (ZOO-ADMM) that involves a modification of step (2) :

$$\mathbf{x}_{t+1} = \operatorname*{arg\,min}_{\mathbf{x} \in \mathcal{X}} \left\{ \hat{\mathbf{g}}_t^T \mathbf{x} - \boldsymbol{\lambda}_t^T (\mathbf{A}\mathbf{x} + \mathbf{B}\mathbf{y}_t - \mathbf{c}) \right.$$
$$\left. + \frac{\rho}{2} \|\mathbf{A}\mathbf{x} + \mathbf{B}\mathbf{y}_t - \mathbf{c}\|_2^2 + \frac{1}{2\eta_t} \|\mathbf{x} - \mathbf{x}_t\|_{\mathbf{G}_t}^2 \right\}. \quad (6)$$

In (22), we can specify the matrix $\mathbf{G}_t$ in such a way as to cancel the term $\|\mathbf{A}\mathbf{x}\|_2^2$. This technique has been used in the linearized ADMM algorithms (Parikh and Boyd, 2014; Zhang et al., 2011) to avoid matrix inversions. Defining $\mathbf{G}_t = \alpha \mathbf{I} - \rho \eta_t \mathbf{A}^T \mathbf{A}$, the update rule (22) simplifies to a projection operator

$$\mathbf{x}_{t+1} = \operatorname*{arg\,min}_{\mathbf{x} \in \mathcal{X}} \left\{ \|\mathbf{x} - \boldsymbol{\omega}\|_2^2 \right\} \text{ with} \qquad (7)$$
$$\boldsymbol{\omega} := \left[ \frac{\eta_t}{\alpha} \left( -\hat{\mathbf{g}}_t + \mathbf{A}^T (\boldsymbol{\lambda}_t - \rho(\mathbf{A}\mathbf{x}_t + \mathbf{B}\mathbf{y}_t - \mathbf{c})) \right) + \mathbf{x}_t \right],$$

where $\alpha > 0$ is a parameter selected to ensure $\mathbf{G}_t \succeq \mathbf{I}$. Here $\mathbf{X} \succeq \mathbf{Y}$ signifies that $\mathbf{X} - \mathbf{Y}$ is positive semidefinite.

To evaluate the convergence behavior of ZOO-ADMM, we will derive its expected average regret (Hazan,

2016)

$$\overline{\text{Regret}}_T(\mathbf{x}_t, \mathbf{y}_t, \mathbf{x}^*, \mathbf{y}^*) := \mathbb{E}\left[ \frac{1}{T} \sum_{t=1}^T (f(\mathbf{x}_t; \mathbf{w}_t) + \phi(\mathbf{y}_t)) \right.$$
$$\left. - \frac{1}{T} \sum_{t=1}^T (f(\mathbf{x}^*; \mathbf{w}_t) + \phi(\mathbf{y}^*)) \right], \qquad (8)$$

where $(\mathbf{x}^*, \mathbf{y}^*)$ denotes the best batch offline solution.

# 3 Algorithm and Convergence Analysis of ZOO-ADMM

In this section, we begin by stating assumptions used in our analysis. We then formally define the ZOO-ADMM algorithm and derive its convergence rate.

We assume the following conditions in our analysis.
- *Assumption A:* In problem (21), $\mathcal{X}$ and $\mathcal{Y}$ are bounded with finite diameter $R$, and at least one of $\mathbf{A}$ and $\mathbf{B}$ in $\mathbf{A}\mathbf{x} + \mathbf{B}\mathbf{y} = \mathbf{c}$ is invertible.

- *Assumption B:* $f(\cdot; \mathbf{w}_t)$ is convex and Lipschitz continuous with $\sqrt{\mathbb{E}[\|\nabla_\mathbf{x} f(\mathbf{x}; \mathbf{w}_t)\|_2^2]} \leq L_1$ for all $t$ and $\mathbf{x} \in \mathcal{X}$.

- Assumption C: $f(\cdot; \mathbf{w}_t)$ is $L_g(\mathbf{w}_t)$-smooth with $L_g = \sqrt{\mathbb{E}[(L_g(\mathbf{w}_t)^2)]}$.

- *Assumption D:* $\phi$ is convex and $L_2$-Lipschitz continuous with $\|\partial\phi(\mathbf{y})\|_2 \leq L_2$ for all $\mathbf{y} \in \mathcal{Y}$, where $\partial\phi(\mathbf{y})$ denotes the subgradient of $\phi$.

- *Assumption E:* In (5), given $\mathbf{z} \sim \mu$, the quantity $M(\mu) := \sqrt{\mathbb{E}[\|\mathbf{z}\|_2^6]}$ is finite, and there is a function $s : \mathbb{N} \to \mathbb{R}_+$ satisfying $\mathbb{E}[\|\langle \mathbf{a}, \mathbf{z}\rangle \mathbf{z}\|_2^2] \leq s(m)\|\mathbf{a}\|_2^2$ for all $\mathbf{a} \in \mathbb{R}^m$, where $\langle \cdot, \cdot \rangle$ denotes the inner product of two vectors.

We remark that Assumptions A-D are standard for stochastic gradient-based and ADMM-type methods (Boyd et al., 2011; Hazan, 2016; Shalev-Shwartz, 2012; Suzuki, 2013). We elaborate on the rationale behind them in Sec. 8.1. Assumption E places moment constraints on the distribution $\mu$ that will allow us to derive the necessary concentration bounds for our convergence analysis. If $\mu$ is uniform on the surface of the Euclidean-ball of radius $\sqrt{m}$, we have $M(\mu) = m^{1.5}$ and $s(m) = m$. And if $\mu = \mathcal{N}(\mathbf{0}, \mathbf{I}_{m \times m})$, we have $M(\mu) \approx m^{1.5}$ and $s(m) \approx m$ (Duchi et al., 2015). For ease of representation, we restrict our attention to the case that $s(m) = m$ in the rest of the paper. It is also worth mentioning that the convex and strongly convex conditions of $f(\cdot; \mathbf{w}_t)$ can be described as

$$f(\mathbf{x}; \mathbf{w}_t) \geq f(\tilde{\mathbf{x}}; \mathbf{w}_t) + (\mathbf{x} - \tilde{\mathbf{x}})^T \nabla_\mathbf{x} f(\tilde{\mathbf{x}}; \mathbf{w}_t)$$
$$+ \frac{\sigma}{2} \|\mathbf{x} - \tilde{\mathbf{x}}\|^2, \ \forall \mathbf{x}, \tilde{\mathbf{x}}, \qquad (9)$$

where $\sigma \geq 0$ is a parameter controlling convexity. If

$\sigma > 0$, then $f(\cdot; \mathbf{w}_t)$ is strongly convex with parameter $\sigma$. Otherwise ($\sigma = 0$), (9) implies convexity of $f(\cdot; \mathbf{w}_t)$.

The ZOO-ADMM iterations are given as Algorithm 1. Compared to O-ADMM in (Suzuki, 2013), we only require querying two function values for the generation of gradient estimate at step 3. Also different from (Gao et al., 2017), steps 7-11 of Algorithm 1 imply that the equality constraint of problem (21) is always satisfied at $\{\mathbf{x}_t, \mathbf{y}'_t\}$ or $\{\mathbf{x}'_t, \mathbf{y}_t\}$. The average regret of ZOO-ADMM is bounded in Theorem 1.

**Theorem 1** *Suppose $\mathbf{B}$ is invertible in problem (21). For $\{\mathbf{x}_t, \mathbf{y}'_t\}$ generated by ZOO-ADMM, the expected average regret is bounded as*

$$\overline{\text{Regret}}_T(\mathbf{x}_t, \mathbf{y}'_t, \mathbf{x}^*, \mathbf{y}^*)$$

$$\leq \frac{1}{T}\sum_{t=2}^{T}\max\left\{\frac{\alpha}{2\eta_t} - \frac{\alpha}{2\eta_{t-1}} - \frac{\sigma}{2}, 0\right\}R^2 + \frac{mL_1^2}{T}\sum_{t=1}^{T}\eta_t$$

$$+ \frac{M(\mu)^2 L_g^2}{4T}\sum_{t=1}^{T}\eta_t\beta_t^2 + \frac{K}{T}, \tag{10}$$

*where $\alpha$ is introduced in (7), $R$, $L_1$, $L_g$, $s(m)$ and $M(\mu)$ are defined in Assumptions A-E, and $K$ denotes a constant term that depends on $\alpha$, $R$, $\eta_1$, $\mathbf{A}$, $\mathbf{B}$, $\boldsymbol{\lambda}$, $\rho$ and $L_2$. Suppose $\mathbf{A}$ is invertible in problem (21). For $\{\mathbf{x}'_t, \mathbf{y}_t\}$, the regret $\overline{\text{Regret}}_T(\mathbf{x}'_t, \mathbf{y}_t, \mathbf{x}^*, \mathbf{y}^*)$ obeys the same bounds as (47).*

**Proof:** See Sec. 8.2. ∎

In Theorem 1, if the step size $\eta_t$ and the smoothing parameter $\beta_t$ are chosen as

$$\eta_t = \frac{C_1}{m\sqrt{t}}, \ \beta_t = \frac{C_2}{M(\mu)t} \tag{11}$$

for some constant $C_1 > 0$ and $C_2 > 0$, then the regret bound (47) simplifies to

$$\overline{\text{Regret}}_T(\mathbf{x}_t, \mathbf{y}'_t, \mathbf{x}^*, \mathbf{y}^*) \leq \frac{\alpha R^2}{2C_1}\frac{\sqrt{m}}{\sqrt{T}}$$

$$+ 2C_1 L_1^2 \frac{\sqrt{m}}{\sqrt{T}} + \frac{5C_1 C_2^2 L_g^2}{12}\frac{1}{T} + \frac{K}{T}. \tag{12}$$

The above simplification is derived in Sec. 8.3.

It is clear from (12) that ZOO-ADMM converges at least as fast as $O(\sqrt{m}/\sqrt{T})$, which is similar to the convergence rate of O-ADMM found by (Suzuki, 2013) but involves an additional factor $\sqrt{m}$. Such a dimension-dependent effect on the convergence rate has also been reported for other zeroth-order optimization algorithms (Duchi et al., 2015; Ghadimi and Lan, 2013; Shamir, 2017), leading to the same convergence rate as ours. In (12), even if we set $C_2 = 0$ (namely, $\beta_t = 0$) for an unbiased gradient estimate (5), the dimension-dependent factor $\sqrt{m}$ is not eliminated.

---

**Algorithm 1** ZOO-ADMM for solving problem (21)

1: Input: $\mathbf{x}_1 \in \mathcal{X}$, $\mathbf{y}_1 \in \mathcal{Y}$, $\boldsymbol{\lambda}_1 = \mathbf{0}$, $\rho > 0$, step sizes $\{\eta_t\}$, smoothing constants $\{\beta_t\}$, distribution $\mu$, and $\alpha \geq \rho\eta_t\lambda_{\max}(\mathbf{A}^T\mathbf{A}) + 1$ so that $\mathbf{G}_t \succeq \mathbf{I}$, where $\lambda_{\max}(\cdot)$ denotes the maximum eigenvalue of a symmetric matrix
2: **for** $t = 1, 2, \ldots, T$ **do**
3:     sample $\mathbf{z}_t \sim \mu$ to generate $\hat{\mathbf{g}}_t$ using (5)
4:     update $\mathbf{x}_{t+1}$ via (7) under $\hat{\mathbf{g}}_t$ and $(\mathbf{x}_t, \mathbf{y}_t, \boldsymbol{\lambda}_t)$
5:     update $\mathbf{y}_{t+1}$ via (23) under $(\mathbf{x}_{t+1}, \boldsymbol{\lambda}_t)$
6:     update $\boldsymbol{\lambda}_{t+1}$ via (24) under $(\mathbf{x}_{t+1}, \mathbf{y}_{t+1}, \boldsymbol{\lambda}_t)$
7:     **if** $\mathbf{B}$ is invertible **then**
8:         compute $\mathbf{y}'_{t+1} := \mathbf{B}^{-1}(\mathbf{c} - \mathbf{A}\mathbf{x}_{t+1})$
9:     **else**
10:        compute $\mathbf{x}'_{t+1} := \mathbf{A}^{-1}(\mathbf{c} - \mathbf{B}\mathbf{y}_{t+1})$
11:     **end if**
12: **end for**
13: output: $\{\mathbf{x}_t, \mathbf{y}'_t\}$ or $\{\mathbf{x}'_t, \mathbf{y}_t\}$, running average $(\bar{\mathbf{x}}_T, \bar{\mathbf{y}}'_T)$ or $(\bar{\mathbf{x}}'_T, \bar{\mathbf{y}}_T)$, where $\bar{\mathbf{x}}_T = \frac{1}{T}\sum_{k=1}^{T}\mathbf{x}_k$.

---

That is because the second moment of the gradient estimate also depends on the number of optimization variables. In the next section, we will propose two minibatch strategies that can be used to reduce the variance of the gradient estimate and to improve the convergence speed of ZOO-ADMM.

## 4 Convergence for Special Cases

In this section, we specialize ZOO-ADMM to three cases: a) stochastic optimization, b) strongly convex cost function in (21), and c) the use of minibatch strategies for evaluation of gradient estimates. Without loss of generality, we restrict analysis to the case that $\mathbf{B}$ is invertible in (21).

The stochastic optimization problem is a special case of the OCO problem (21). If the objective function becomes $F(\mathbf{x}, \mathbf{y}) := \mathbb{E}_{\mathbf{w}}[f(\mathbf{x}; \mathbf{w})] + \phi(\mathbf{y})$ then we can link the regret with the optimization error at the running average $\bar{\mathbf{x}}_T$ and $\bar{\mathbf{y}}_T$ under the condition that $F$ is convex. We state our results as Corollary 1.

**Corollary 1** *Consider the stochastic optimization problem with the objective function $F(\mathbf{x}, \mathbf{y}) := \mathbb{E}_{\mathbf{w}}[f(\mathbf{x}; \mathbf{w})] + \phi(\mathbf{y})$, and set $\eta_t$ and $\beta_t$ using (11). For $\{\bar{\mathbf{x}}_t, \bar{\mathbf{y}}'_t\}$ generated by ZOO-ADMM, the optimization error $\mathbb{E}\left[F(\bar{\mathbf{x}}_T, \bar{\mathbf{y}}'_T) - F(\mathbf{x}^*, \mathbf{y}^*)\right]$ obeys the same bound as (12).*

**Proof:** See Sec. 8.4. ∎

We recall from (9) that $\sigma$ controls the convexity of $f_t$, where $\sigma > 0$ if $f_t$ is strongly convex. In Corollary 2, we show that $\sigma$ affects the average regret of ZOO-ADMM.

**Corollary 2** *Suppose $f(\cdot; \mathbf{w}_t)$ is strongly convex, and*

Sijia Liu, Jie Chen, Pin-Yu Chen, Alfred O. Hero

*the step size $\eta_t$ and the smoothing parameter $\beta_t$ are chosen as $\eta_t = \frac{\alpha}{\sigma t}$ and $\beta_t = \frac{C_2}{M(\mu)t}$ for $C_2 > 0$. Given $\{\mathbf{x}_t, \mathbf{y}_t'\}$ generated by ZOO-ADMM, the expected average regret can be bounded as*

$$\overline{\mathrm{Regret}}_T(\mathbf{x}_t, \mathbf{y}_t', \mathbf{x}^*, \mathbf{y}^*) \leq \frac{\alpha L_1^2}{\sigma} \frac{m \log T}{T}$$
$$+ \frac{3\alpha C_2^2 L_g^2}{8\sigma} \frac{1}{T} + \frac{K}{T}. \quad (13)$$

**Proof:** See Sec. 8.5. ∎

Corollary 2 implies that when the cost function is strongly convex, the regret bound of ZOO-ADMM could achieve $O(m/T)$ up to a logarithmic factor $\log T$. Compared to the regret bound $O(\sqrt{m}/\sqrt{T})$ in the general case (12), the condition of strong convexity improves the regret bound in terms of the number of iterations $T$, but the dimension-dependent factor now becomes linear in the dimension $m$ due to the effect of the second moment of gradient estimate.

The use of a gradient estimator makes the convergence rate of ZOO-ADMM dependent on the dimension $m$, i.e., the number of optimization variables. Thus, it is important to study the impact of minibatch strategies on the acceleration of the convergence speed (Cotter et al., 2011; Duchi et al., 2015; Li et al., 2014; Suzuki, 2013). Here we present two minibatch strategies: gradient sample averaging and observation averaging. In the first strategy, instead of using a single sample as in (5), the average of $q$ sub-samples $\{\mathbf{z}_{t,i}\}_{i=1}^q$ are used for gradient estimation

$$\hat{\mathbf{g}}_t = \frac{1}{q} \sum_{i=1}^q \frac{f(\mathbf{x}_t + \beta_t \mathbf{z}_{t,i}; \mathbf{w}_t) - f(\mathbf{x}_t; \mathbf{w}_t)}{\beta_t} \mathbf{z}_{t,i}, \quad (14)$$

where $q$ is called the batch size. The use of (14) is analogous to the use of an average gradient in incremental gradient (Blatt et al., 2007) and stochastic gradient (Roux et al., 2012). In the second strategy, we use a subset of observations $\{\mathbf{w}_{t,i}\}_{i=1}^q$ to reduce the gradient variance,

$$\hat{\mathbf{g}}_t = \frac{1}{q} \sum_{i=1}^q \frac{f(\mathbf{x}_t + \beta_t \mathbf{z}_t; \mathbf{w}_{t,i}) - f(\mathbf{x}_t; \mathbf{w}_{t,i})}{\beta_t} \mathbf{z}_t. \quad (15)$$

We note that in the online setting, the subset of observations $\{\mathbf{w}_{t,i}\}_{i=1}^q$ can be obtained via a sliding time window of length $q$, namely, $\mathbf{w}_{i,t} = \mathbf{w}_{t-i+1}$ for $i = 1, 2, \ldots, q$.

Combination of (14) and (15) yields a hybrid strategy

$$\hat{\mathbf{g}}_t = \frac{1}{q_1 q_2} \sum_{j=1}^{q_1} \sum_{i=1}^{q_2} \frac{f(\mathbf{x}_t + \beta_t \mathbf{z}_{t,j}; \mathbf{w}_{t,i}) - f(\mathbf{x}_t; \mathbf{w}_{t,i})}{\beta_t} \mathbf{z}_{t,j}.$$
$$(16)$$

In Corollary 3, we demonstrate the convergence behavior of the general hybrid ZOO-ADMM.

**Corollary 3** *Consider the hybrid minibatch strategy (51) in ZOO-ADMM, and set $\eta_t = \frac{C_1}{\sqrt{1+\frac{m}{q_1 q_2}}\sqrt{t}}$ and $\beta_t = \frac{C_2}{M(\mu)t}$. The expected average regret is bounded as*

$$\overline{\mathrm{Regret}}_T(\mathbf{x}_t, \mathbf{y}_t', \mathbf{x}^*, \mathbf{y}^*) \leq \frac{\alpha R^2}{2C_1} \frac{\sqrt{1+\frac{s(m)}{q_1 q_2}}}{\sqrt{T}}$$
$$+ 2C_1 L_1^2 \frac{\sqrt{1+\frac{s(m)}{q_1 q_2}}}{\sqrt{T}} + \frac{5C_1 C_2^2 L_g^2}{6} \frac{1}{T} + \frac{K}{T}, \quad (17)$$

*where $q_1$ and $q_2$ are number of sub-samples $\{\mathbf{z}_{t,i}\}$ and $\{\mathbf{w}_{t,i}\}$, respectively.*

**Proof:** See Sec. 8.6. ∎

It is clear from Corollary 3 that the use of minibatch strategies can alleviate the dimension dependency, leading to the regret bound $O(\sqrt{1+m/(q_1 q_2)}/\sqrt{T})$. The regret bound in (17) also implies that the convergence behavior of ZOO-ADMM is similar using either gradient sample averaging minibatch (14) or observation averaging minibatch (15). If $q_1 = 1$ and $q_2 = 1$, the regret bound (17) reduces to $O(\sqrt{m}/\sqrt{T})$, which is the general case in (12). If $q_1 q_2 = O(m)$, we obtain the regret error $O(1/\sqrt{T})$ as in the case where an explicit expression for the gradient is used in the OCO algorithms.

## 5 Applications of ZOO-ADMM

In this section, we demonstrate several applications of ZOO-ADMM in signal processing, statistics and machine learning.

### 5.1 Black-box optimization

In some OCO problems, explicit gradient calculation is impossible due to the lack of a mathematical expression for the loss function. For example, commercial recommender systems try to build a representation of a customer's buying preference function based on a discrete number of queries or purchasing history, and the system never has access to the gradient of the user's preference function over their product line, which may even be unknown to the user. Gradient-free methods are therefore necessary. A specific example is the Yahoo! music recommendation system (Dror et al., 2012), which will be further discussed in the Sec. 6. In these examples, one can consider each user as a black-box model that provides feedback on the value of an objective function, e.g., relative preferences over all products, based on an online evaluation of the objective function at discrete points on its domain. Such a system can benefit from ZOO-ADMM.

## 5.2 Sensor selection

Sensor selection for parameter estimation is a fundamental problem in smart grids, communication systems, and wireless sensor networks (Hero and Cochran, 2011; Liu et al., 2016). The goal is to seek the optimal tradeoff between sensor activations and the estimation accuracy. The sensor selection problem is also closely related to leader selection (Lin et al., 2014) and experimental design (Boyd and Vandenberghe, 2004).

For sensor selection, we often solve a (relaxed) convex program of the form (Joshi and Boyd, 2009)

$$\underset{\mathbf{x}}{\text{minimize}} \quad \frac{1}{T} \sum_{t=1}^{T} \left[ -\text{logdet}\left( \sum_{i=1}^{m} x_i \mathbf{a}_{i,t} \mathbf{a}_{i,t}^T \right) \right] \quad (18)$$
$$\text{subject to} \quad \mathbf{1}^T \mathbf{x} = m_0, \ \mathbf{0} \le \mathbf{x} \le \mathbf{1},$$

where $\mathbf{x} \in \mathbb{R}^m$ is the optimization variable, $m$ is the number of sensors, $\mathbf{a}_{i,t} \in \mathbb{R}^n$ is the observation coefficient of sensor $i$ at time $t$, and $m_0$ is the number of selected sensors. The objective function of (18) can be interpreted as the log determinant of error covariance associated with the maximum likelihood estimator for parameter estimation (Rao, 1973). The constraint $\mathbf{0} \le \mathbf{x} \le \mathbf{1}$ is a relaxed convex hull of the Boolean constraint $\mathbf{x} \in \{0,1\}^m$, which encodes whether or not a sensor is selected.

Conventional methods such as projected gradient (first-order) and interior-point (second-order) algorithms can be used to solve problem (18). However, both of them involve calculation of inverse matrices necessary to evaluate the gradient of the cost function. By contrast, we can rewrite (18) in a form amenable to ZOO-ADMM that avoids matrix inversion,

$$\underset{\mathbf{x},\mathbf{y}}{\text{minimize}} \quad \frac{1}{T} \sum_{t=1}^{T} f(\mathbf{x}; \mathbf{w}_t) + \mathcal{I}_1(\mathbf{x}) + \mathcal{I}_2(\mathbf{y}) \quad (19)$$
$$\text{subject to} \quad \mathbf{x} - \mathbf{y} = \mathbf{0},$$

where $\mathbf{y} \in \mathbb{R}^m$ is an auxiliary variable, $f(\mathbf{x}; \mathbf{w}_t) = -\text{logdet}(\sum_{i=1}^{m} x_i \mathbf{a}_{i,t} \mathbf{a}_{i,t}^T)$ with $\mathbf{w}_t = \{\mathbf{a}_{i,t}\}_{i=1}^{m}$, and $\{\mathcal{I}_i\}$ are indicator functions

$$\mathcal{I}_1(\mathbf{x}) = \left\{ \begin{array}{ll} 0 & \mathbf{0} \le \mathbf{x} \le \mathbf{1} \\ \infty & \text{otherwise,} \end{array} \right. \mathcal{I}_2(\mathbf{y}) = \left\{ \begin{array}{ll} 0 & \mathbf{1}^T \mathbf{y} = m_0 \\ \infty & \text{otherwise.} \end{array} \right.$$

We specify the ZOO-ADMM algorithm for solving (59) in Sec. 8.7.

## 5.3 Sparse Cox regression

In survival analysis, Cox regression (also known as proportional hazards regression) is a method to investigate effects of variables of interest upon the amount of time that elapses before a specified event occurs, e.g., relating gene expression profiles to survival time (time

to cancer recurrence or death) (Sohn et al., 2009). Let $\{\mathbf{a}_i \in \mathbb{R}^m, \delta_i \in \{0,1\}, t_i \in \mathbb{R}_+\}_{i=1}^{n}$ be $n$ triples of $m$ covariates, where $\mathbf{a}_i$ is a vector of covariates or factors for subject $i$, $\delta_i$ is a censoring indicator variable taking 1 if an event (e.g., death) is observed and 0 otherwise, and $t_i$ denotes the censoring time.

This sparse regression problem can be formulated as the solution to an $\ell_1$ penalized optimization problem (Park and Hastie, 2007; Sohn et al., 2009), which yields

$$\underset{\mathbf{x}}{\text{minimize}} \quad \frac{1}{n} \sum_{i=1}^{n} \delta_i \left\{ -\mathbf{a}_i^T \mathbf{x} + \log\left( \sum_{j \in \mathcal{R}_i} e^{\mathbf{a}_j^T \mathbf{x}} \right) \right\}$$
$$+ \gamma \|\mathbf{x}\|_1 \quad (20)$$

where $\mathbf{x} \in \mathbb{R}^m$ is the vector of covariates coefficients to be designed, $\mathcal{R}_i$ is the set of subjects at risk at time $t_i$, namely, $\mathcal{R}_i = \{j : t_j \ge t_i\}$, and $\gamma > 0$ is a regularization parameter. In the objective function of (20), the first term corresponds to the (negative) log partial likelihood for the Cox proportional hazards model (Cox, 1972), and the second term encourages sparsity of the covariate coefficients.

By introducing a new variable $\mathbf{y} \in \mathbb{R}^m$ together with the constraint $\mathbf{x} - \mathbf{y} = \mathbf{0}$, problem (20) can be cast as the canonical form (21) amenable to the ZOO-ADMM algorithm. This helps us to avoid the gradient calculation for the involved objective function in Cox regression. We specify the ZOO-ADMM algorithm for solving (20) in Sec. 8.8.

## 6 Experiments

In this section, we demonstrate the effectiveness of ZOO-ADMM, and validate its convergence behavior for the applications introduced in Sec. 5. In Algorithm 1, we set $\mathbf{x}_1 = \mathbf{0}$, $\mathbf{y}_1 = \mathbf{0}$, $\boldsymbol{\lambda}_1 = \mathbf{0}$, $\rho = 10$, $\eta_t = 1/\sqrt{mt}$, $\beta_t = 1/(m^{1.5}t)$, $\alpha = \rho \eta_t \lambda_{\max}(\mathbf{A}^T \mathbf{A}) + 1$, and the distribution $\mu$ is chosen to be uniform on the surface of the Euclidean-ball of radius $\sqrt{m}$. Unless specified otherwise, we use the gradient sample averaging minibatch of size 30 in ZOO-ADMM. Through this section, we compare ZOO-ADMM with the conventional O-ADMM algorithm in (Suzuki, 2013) under the same parameter settings. Our experiments are performed on a synthetic dataset for sensor selection, and on real datasets for black-box optimization and Cox regression. Experiments were conducted by Matlab R2016 on a machine with 3.20 GHz CPU and 8 GB RAM.

**Black-box optimization:** We consider prediction of users' ratings in the Yahoo! music system (Dror et al., 2012). Our dataset, provided by (Lian et al., 2016), include $n' = 131072$ true music ratings $\mathbf{r} \in \mathbb{R}^{n'}$, and the predicted ratings of $m = 237$ individual models
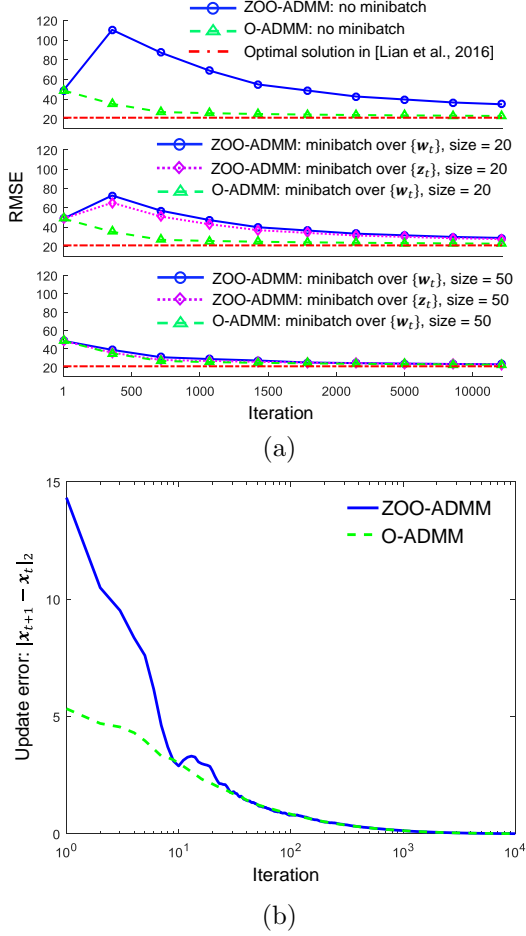
Sijia Liu, Jie Chen, Pin-Yu Chen, Alfred O. Hero

(a)



(b)

**Figure 1:** Convergence of ZOO-ADMM: a) RMSE under different minibatch strategies, b) update error with minibatch size equal to 50.



(a)



(b)

**Figure 2:** ZOO-ADMM for sensor selection: a) MSE versus number of selected sensors $m_0$, b) computation time versus number of optimization variables.

created from the NTU KDD-Cup team (Chen et al., 2011). Let $\mathbf{C} \in \mathbb{R}^{n \times m}$ represent a matrix of each models' predicted ratings on Yahoo! music data sample. We split the dataset $(\mathbf{C}, \mathbf{r})$ into two equal parts, leading to the training dataset $(\mathbf{C}_1 \in \mathbb{R}^{n \times m}, \mathbf{r}_1 \in \mathbb{R}^n)$ and the test dataset $(\mathbf{C}_2 \in \mathbb{R}^{n \times m}, \mathbf{r}_2 \in \mathbb{R}^n)$, where $n = n'/2$.

Our goal is to find the optimal coefficients $\mathbf{x}$ to blend $m$ individual models such that the mean squared error $f(\mathbf{x}) := \frac{1}{n} \sum_{i=1}^{n} f(\mathbf{x}; \mathbf{w}_i) = \frac{1}{n} \sum_{i=1}^{n} ([\mathbf{C}_1]_i^T \mathbf{x} - [\mathbf{r}_1]_i)^2$ is minimized, where $\mathbf{w}_i = ([\mathbf{C}_1]_i, [\mathbf{r}_1]_i)$, $[\mathbf{C}_1]_i$ is the $i$th row vector of $\mathbf{C}_1$, and $[\mathbf{r}_1]_i$ is the $i$th entry of $\mathbf{r}_1$. Since $(\mathbf{C}, \mathbf{r})$ includes predicted ratings on Yahoo! Music data using NTU KDD-Cup team's models, it is private information known only to other users. Therefore, the information $(\mathbf{C}, \mathbf{r})$ cannot be accessed directly (Lian et al., 2016), and explicit gradient calculation for $f$ is not possible. We thus treat the loss function as a black box, where it is evaluated at individual points $\mathbf{x}$ in its domain but not over any open region of its domain.

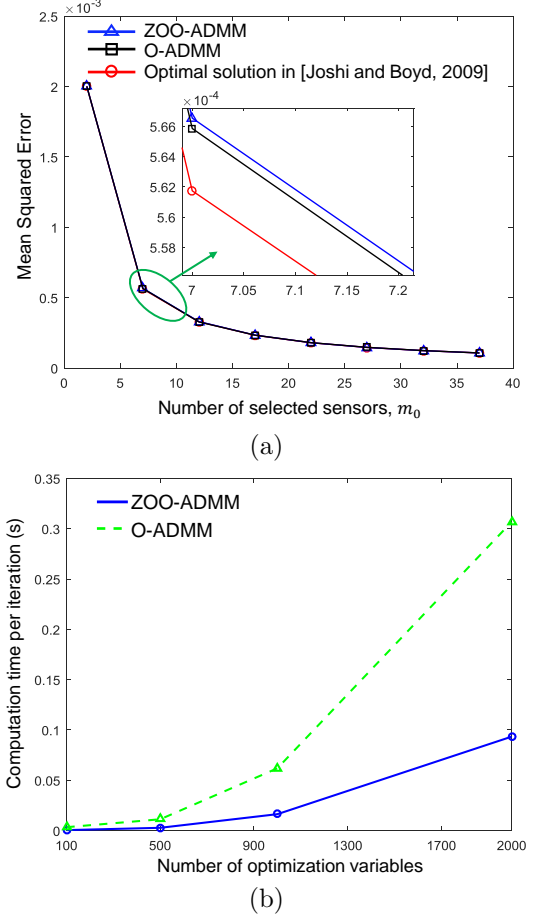As discussed in Sec. 5.1, we can apply ZOO-ADMM

to solve the proposed linear blending problem, and the prediction accuracy can be measured by the root mean squared error (RMSE) of the test data RMSE = $\sqrt{\|\mathbf{r}_2 - \mathbf{C}_2 \mathbf{x}\|_2^2 / n}$, where an update of $\mathbf{x}$ is obtained at each iteration.

In Fig. 1, we compare the performance of ZO-ADMM with O-ADMM and the optimal solution provided by (Lian et al., 2016). In Fig. 1-(a), we present RMSE as a function of iteration number under different minibatch schemes. As we can see, both gradient sample averaging (over $\{\mathbf{z}_t\}$) and observation averaging (over $\{\mathbf{w}_t\}$) significantly accelerate the convergence speed of ZOO-ADMM. In particular, when the minibatch size $q$ is large enough (50 in our example), the dimension-dependent slowdown factor of ZOO-ADMM can be mitigated. We also observe that ZOO-ADMM reaches the best RMSE in (Lian et al., 2016) after 10000 iterations. In Fig. 1-(b), we show the convergence error $\|\mathbf{x}_{t+1} - \mathbf{x}_t\|_2$ versus iteration number using gradient sample averaging minibatch of size 50. Compared to O-ADMM, ZOO-ADMM has a larger performance gap in its first few iterations, but it thereafter con-
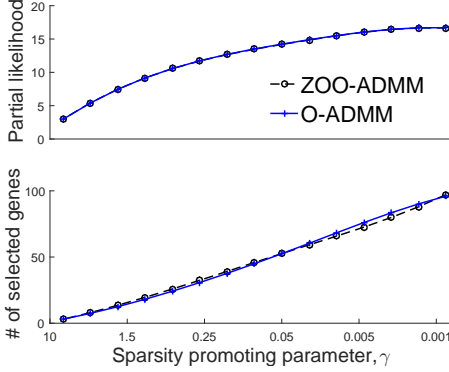
**Figure 3:** Partial likelihood and number of selected genes versus sparsity promoting parameter $\gamma$.

verges quickly resulting in comparable performance to O-ADMM.

**Sensor selection:** We consider an example of estimating a spatial random field based on measurements of the field at a discrete set of sensor locations. Assume that $m = 100$ sensors are randomly deployed over a square region to monitor a vector of field intensities (e.g., temperature values). The objective is to estimate the field intensity at $n = 5$ locations over a time period of $T = 1000$ secs. In (18), the observation vectors $\{\mathbf{a}_{i,t}\}$ are chosen randomly, and independently, from a distribution $\mathcal{N}(\mu_i \mathbf{1}_n, \mathbf{I}_n)$. Here $\mu_i$ is generated by an exponential model (Liu et al., 2016), $\mu_i = 5e^{\sum_{j=1}^{n} \|\hat{\mathbf{s}}_j - \tilde{\mathbf{s}}_i\|_2/n}$, where $\hat{\mathbf{s}}_j$ is the $j$-th spatial location at which the field intensity is to be estimated and $\tilde{\mathbf{s}}_i$ is the spatial location of the $i$ sensor.

In Fig. 2, we present the performance of ZOO-ADMM for sensor selection. In Fig. 2-(a), we show the mean squared error (MSE) averaged over 50 random trials for different number of selected sensors $m_0$ in (18). We compare our approach with O-ADMM and the method in (Joshi and Boyd, 2009). The figure shows that ZOO-ADMM yields almost the same MSE as O-ADMM. The method in (Joshi and Boyd, 2009) yields slightly better estimation performance, since it uses the second-order optimization method for sensor selection. In Fig. 2-(b), we present the computation time of ZOO-ADMM versus the number of optimization variables $m$. The figure shows that ZOO-ADMM becomes much more computationally efficient as $m$ increases since no matrix inversion is required.

**Sparse Cox regression:** We next employ ZOO-ADMM to solve problem (20) for building a sparse predictor of patient survival using the Kidney renal clear cell carcinoma dataset[2]. The aforementioned dataset includes clinical data (survival time and censoring information) and gene expression data for 606 patients

_____

[2] Available at `http://gdac.broadinstitute.org/`

**Table 1:** Percentage of common genes found using ZOO-ADMM and Cox scores (Witten and Tibshirani, 2010).

|  | $\gamma = 1.5$ | $\gamma = 0.05$ | $\gamma = 0.001$ |
|---|---|---|---|
| # selected genes | 19 | 56 | 93 |
| Overlapping (%) | 80.1% | 87.5% | 92.3% |

(534 with tumor and 72 without tumor). Our goal is to seek the best subset of genes (in terms of optimal sparse covariate coefficients) that make the most significant impact on the survival time.

In Fig. 3, we show the partial likelihood and number of selected genes as functions of the regularization parameter $\gamma$. The figure shows that ZOO-ADMM nearly attains the accuracy of O-ADMM. Furthermore, the likelihood increases as the number of selected genes increases. There is thus a tradeoff between the (negative) log partial likelihood and the sparsity of covariate coefficients in problem (20). To test the significance of our selected genes, we compare our approach with the significance analysis based on univariate Cox scores used in (Witten and Tibshirani, 2010). The percentage of overlap between the genes identified by each method is shown in Table 1 under different values of $\gamma$. Despite its use of a zeroth order approximation to the gradient, the ZOO-ADMM selects at least 80% of the genes selected by the gradient-based Cox scores of (Witten and Tibshirani, 2010).

## 7 Conclusion

In this paper, we proposed and analyzed a gradient-free (zeroth-order) online optimization algorithm, ZOO-ADMM. We showed that the regret bound of ZOO-ADMM suffers an additional dimension-dependent factor in convergence rate over gradient-based online variants of ADMM, leading to $O(\sqrt{m}/\sqrt{T})$ convergence rate, where $m$ is the number of optimization variables. To alleviate the dimension dependence, we presented two minibatch strategies that yield an improved convergence rate of $O(\sqrt{1 + q^{-1}m}/\sqrt{T})$, where $q$ is the minibatch size. We illustrated the effectiveness of ZOO-ADMM via multiple applications using both synthetic and real-world datasets. In the future, we would like to relax the assumptions on smoothness and convexity of the cost function in ZOO-ADMM.

## 8 Supplementary Material

### 8.1 Assumptions and Key Notations

Recall that we consider the regularized loss minimization problem over a time horizon of length $T$,

$$\underset{\mathbf{x}\in\mathcal{X},\mathbf{y}\in\mathcal{Y}}{\text{minimize}} \quad \frac{1}{T}\sum_{t=1}^{T} f_t(\mathbf{x};\mathbf{w}_t) + \phi(\mathbf{y}) \tag{21}$$
$$\text{subject to} \quad \mathbf{A}\mathbf{x} + \mathbf{B}\mathbf{y} = \mathbf{c}.$$

ZOO-ADMM is given by

$$\mathbf{x}_{t+1} = \underset{\mathbf{x}\in\mathcal{X}}{\arg\min} \left\{ \hat{\mathbf{g}}_t^T \mathbf{x} - \boldsymbol{\lambda}_t^T(\mathbf{A}\mathbf{x} + \mathbf{B}\mathbf{y}_t - \mathbf{c}) \right.$$
$$\left. + \frac{\rho}{2}\|\mathbf{A}\mathbf{x} + \mathbf{B}\mathbf{y}_t - \mathbf{c}\|_2^2 + \frac{1}{2\eta_t}\|\mathbf{x} - \mathbf{x}_t\|_{\mathbf{G}_t}^2 \right\}, \tag{22}$$

$$\mathbf{y}_{t+1} = \underset{\mathbf{y}\in\mathcal{Y}}{\arg\min} \left\{ \phi(\mathbf{y}) - \boldsymbol{\lambda}_t^T(\mathbf{A}\mathbf{x}_{t+1} + \mathbf{B}\mathbf{y} - \mathbf{c}) \right.$$
$$\left. + \frac{\rho}{2}\|\mathbf{A}\mathbf{x}_{t+1} + \mathbf{B}\mathbf{y} - \mathbf{c}\|_2^2 \right\}, \tag{23}$$

$$\boldsymbol{\lambda}_{t+1} = \boldsymbol{\lambda}_t - \rho(\mathbf{A}\mathbf{x}_{t+1} + \mathbf{B}\mathbf{y}_{t+1} - \mathbf{c}), \tag{24}$$

where $\mathbf{G}_t = \alpha\mathbf{I} - \rho\eta_t\mathbf{A}^T\mathbf{A}$.

We first elaborate on our assumptions.

- Assumption A implies that $\|\mathbf{x} - \mathbf{x}'\|_2 \leq R$ and $\|\mathbf{y} - \mathbf{y}'\|_2 \leq R$ for all $\mathbf{x}, \mathbf{x}' \in \mathcal{X}$ and for all $\mathbf{y}, \mathbf{y}' \in \mathcal{Y}$.

- Based on Jensen's inequality, Assumptions B implies that $\|\mathbb{E}[\nabla_{\mathbf{x}}f(\mathbf{x};\mathbf{w}_t)]\|_2 \leq L_1$.

- Assumption C implies a Lipschitz condition over the gradient $\nabla_{\mathbf{x}}f(\mathbf{x};\mathbf{w}_t)$ with constant $L_g(\mathbf{w}_t)$ (Bubeck et al., 2015; Hazan, 2016). Also based on Jensen's inequality, we have $|\mathbb{E}[L_g(\mathbf{w}_t)]| \leq L_g$.

We next introduce key notations used in our analysis. Given the primal-dual variables $\mathbf{x}$, $\mathbf{y}$ and $\boldsymbol{\lambda}$ of problem (21), we define $\mathbf{v} := [\mathbf{x}^T, \mathbf{y}^T, \boldsymbol{\lambda}^T]$, and a primal-dual mapping $H$

$$H(\mathbf{v}) := \mathbf{C}\mathbf{v} - \begin{bmatrix} 0 \\ 0 \\ \mathbf{c} \end{bmatrix}, \quad \mathbf{C} := \begin{bmatrix} 0 & 0 & -\mathbf{A}^T \\ 0 & 0 & -\mathbf{B}^T \\ \mathbf{A} & \mathbf{B} & 0 \end{bmatrix}, \tag{25}$$

where $\mathbf{C}$ is skew symmetric, namely, $\mathbf{C}^T = -\mathbf{C}$. An important property of the affine mapping $H$ is that $\langle \mathbf{v}_1 - \mathbf{v}_2, H(\mathbf{v}_1) - H(\mathbf{v}_2) \rangle = 0$ for every $\mathbf{v}_1$ and $\mathbf{v}_2$. Supposing the sequence $\{\mathbf{v}_t\}$ is generated by an algorithm, we introduce the auxiliary sequence

$$\tilde{\mathbf{v}}_t := [\mathbf{x}_t^T, \mathbf{y}_t^T, \tilde{\boldsymbol{\lambda}}_t^T]^T, \tag{26}$$

where $\tilde{\boldsymbol{\lambda}}_t := \boldsymbol{\lambda}_t - \rho(\mathbf{A}\mathbf{x}_{t+1} + \mathbf{B}\mathbf{y}_t - \mathbf{c})$.

### 8.2 Proof of Theorem 1

Since the sequences $\{\mathbf{x}_t\}$, $\{\mathbf{y}_t\}$ and $\{\boldsymbol{\lambda}_t\}$ produced from (22)-(24) have the same structure as the ADMM/O-ADMM steps, the property of ADMM given by Theorem 4 of (Suzuki, 2013) is directly applicable to our case, yielding

$$\sum_{t=1}^{T}(f_t(\mathbf{x}_t) + \phi(\mathbf{y}_t)) - \sum_{t=1}^{T}(f_t(\mathbf{x}) + \phi(\mathbf{y}))$$
$$+ \sum_{t=1}^{T}(\tilde{\mathbf{v}}_t - \mathbf{v})^T H(\tilde{\mathbf{v}}_t) \leq \frac{\|\mathbf{x}_1 - \mathbf{x}\|_{\mathbf{G}_1}^2}{2\eta_1}$$
$$+ \sum_{t=2}^{T}\left(\frac{\|\mathbf{x}_t - \mathbf{x}\|_{\mathbf{G}_t}^2}{2\eta_t} - \frac{\|\mathbf{x}_t - \mathbf{x}\|_{\mathbf{G}_{t-1}}^2}{2\eta_{t-1}}\right)$$
$$+ \langle \boldsymbol{\lambda}, \mathbf{A}(\mathbf{x}_{T+1} - \mathbf{x}_1) \rangle + \frac{\rho}{2}\|\mathbf{y}_1 - \mathbf{y}\|_{\mathbf{B}^T\mathbf{B}}^2 + \frac{\|\boldsymbol{\lambda}_1 - \boldsymbol{\lambda}\|_2^2}{2\rho}$$
$$- \frac{\|\boldsymbol{\lambda}_{T+1} - \boldsymbol{\lambda}\|_2^2}{2\rho} + \langle \mathbf{B}(\mathbf{y} - \mathbf{y}_{T+1}), \boldsymbol{\lambda}_{T+1} - \boldsymbol{\lambda} \rangle$$
$$- \langle \mathbf{B}(\mathbf{y} - \mathbf{y}_1), \boldsymbol{\lambda}_1 - \boldsymbol{\lambda} \rangle - \sum_{t=1}^{T}\frac{\|\boldsymbol{\lambda}_t - \boldsymbol{\lambda}_{t+1}\|_2^2}{2\rho}$$
$$- \sum_{t=1}^{T}\frac{\sigma}{2}\|\mathbf{x}_t - \mathbf{x}\|_2^2 + \sum_{t=1}^{T}\frac{\eta_t}{2}\|\hat{\mathbf{g}}_t\|_{\mathbf{G}_t^{-1}}^2. \tag{27}$$

Here for notational simplicity we have used, and henceforth will continue to use, $f_t(\mathbf{x}_t)$ instead of $f(\mathbf{x}_t;\mathbf{w}_t)$.

In (27), based on $\mathbf{G}_t = \alpha\mathbf{I} - \rho\eta_t\mathbf{A}^T\mathbf{A}$, we have

$$\frac{\|\mathbf{x}_t - \mathbf{x}\|_{\mathbf{G}_t}^2}{2\eta_t} - \frac{\|\mathbf{x}_t - \mathbf{x}\|_{\mathbf{G}_{t-1}}^2}{2\eta_{t-1}}$$
$$= \left(\frac{\alpha}{2\eta_t} - \frac{\alpha}{2\eta_{t-1}}\right)\|\mathbf{x}_t - \mathbf{x}\|_2^2,$$

which yields

$$\sum_{t=2}^{T}\left(\frac{\|\mathbf{x}_t - \mathbf{x}\|_{\mathbf{G}_t}^2}{2\eta_t} - \frac{\|\mathbf{x}_t - \mathbf{x}\|_{\mathbf{G}_{t-1}}^2}{2\eta_{t-1}}\right)$$
$$- \sum_{t=1}^{T}\frac{\sigma}{2}\|\mathbf{x}_t - \mathbf{x}\|_2^2 \leq$$
$$\sum_{t=2}^{T}\max\{\frac{\alpha}{2\eta_t} - \frac{\alpha}{2\eta_{t-1}} - \frac{\sigma}{2}, 0\}R^2. \tag{28}$$

We also note that the terms $\frac{1}{2\eta_1}\|\mathbf{x}_1 - \mathbf{x}\|_{\mathbf{G}_1}^2$, $\langle \boldsymbol{\lambda}, \mathbf{A}(\mathbf{x}_{T+1} - \mathbf{x}_1) \rangle$, $\frac{\rho}{2}\|\mathbf{y}_1 - \mathbf{y}\|_{\mathbf{B}^T\mathbf{B}}$, $\frac{1}{2\rho}(\|\boldsymbol{\lambda}_1 - \boldsymbol{\lambda}\|_2^2 - \|\boldsymbol{\lambda}_{T+1} - \boldsymbol{\lambda}\|_2^2)$, $\langle \mathbf{B}(\mathbf{y} - \mathbf{y}_{T+1}), \boldsymbol{\lambda}_{T+1} - \boldsymbol{\lambda} \rangle$, and $\langle \mathbf{B}(\mathbf{y} - \mathbf{y}_1), \boldsymbol{\lambda}_1 - \boldsymbol{\lambda} \rangle$ are *independent* of time $t$. In particular,

we have

$$\|\mathbf{x}_1 - \mathbf{x}\|_{\mathbf{G}_1}^2 \le \alpha R^2,$$
$$\langle \boldsymbol{\lambda}, \mathbf{A}(\mathbf{x}_{T+1} - \mathbf{x}_1)\rangle \le R\|\boldsymbol{\lambda}\|_2\|\mathbf{A}\|_F,$$
$$(\|\boldsymbol{\lambda}_1 - \boldsymbol{\lambda}\|_2^2 - \|\boldsymbol{\lambda}_{T+1} - \boldsymbol{\lambda}\|_2^2) \le \|\boldsymbol{\lambda}\|_2^2,$$
$$\langle \mathbf{B}(\mathbf{y} - \mathbf{y}_1), \boldsymbol{\lambda} - \boldsymbol{\lambda}_1\rangle \le R\|\mathbf{B}\|_F\|\boldsymbol{\lambda}\|_2, \qquad (29)$$

where $\|\cdot\|_F$ denotes the Frobenius norm of a matrix, and we have used the facts that $\mathbf{G}_t \preceq \alpha\mathbf{I}$ and $\boldsymbol{\lambda}_1 = \mathbf{0}$.

Based on the optimality condition of $\mathbf{y}_{t+1}$ in (23), we have $\langle \partial\phi(\mathbf{y}_{t+1}) - \mathbf{B}^T\boldsymbol{\lambda}_t + \rho\mathbf{B}^T(\mathbf{A}\mathbf{x}_{t+1} + \mathbf{B}\mathbf{y}_{t+1} - \mathbf{c}), \mathbf{y} - \mathbf{y}_{t+1}\rangle \ge 0$ , $\forall \mathbf{y} \in \mathcal{Y}$, which is equivalent to $\langle \partial\phi(\mathbf{y}_{t+1}) - \mathbf{B}^T\boldsymbol{\lambda}_{t+1}, \mathbf{y} - \mathbf{y}_{t+1}\rangle \ge 0$. And thus, we obtain

$$\langle \boldsymbol{\lambda}_{t+1}, \mathbf{B}(\mathbf{y} - \mathbf{y}_{t+1})\rangle - \langle \boldsymbol{\lambda}, \mathbf{B}(\mathbf{y} - \mathbf{y}_{t+1})\rangle$$
$$\le \langle \partial\phi(\mathbf{y}_{t+1}), \mathbf{y} - \mathbf{y}_{t+1}\rangle - \langle \boldsymbol{\lambda}, \mathbf{B}(\mathbf{y} - \mathbf{y}_{t+1})\rangle,$$

which yields

$$\langle \mathbf{B}(\mathbf{y} - \mathbf{y}_{t+1}), \boldsymbol{\lambda}_{t+1} - \boldsymbol{\lambda}\rangle$$
$$\le \langle \mathbf{y} - \mathbf{y}_{t+1}, \partial\phi(\mathbf{y}_{t+1}) - \mathbf{B}^T\boldsymbol{\lambda}\rangle$$
$$\le R(L_2 + \|\mathbf{B}^T\boldsymbol{\lambda}\|_2), \qquad (30)$$

where we have used the fact that $\|\partial\phi(\mathbf{y}_{t+1})\|_2 \le L_2$.

Substituting (28)-(30) into (27), we then obtain

$$\frac{1}{T}\sum_{t=1}^{T}(f_t(\mathbf{x}_t) + \phi(\mathbf{y}_t)) - \frac{1}{T}\sum_{t=1}^{T}(f_t(\mathbf{x}) + \phi(\mathbf{y}))$$
$$+ \frac{1}{T}\sum_{t=1}^{T}(\tilde{\mathbf{v}}_t - \mathbf{v})^T H(\tilde{\mathbf{v}}_t) + \frac{1}{T}\sum_{t=1}^{T}\frac{\|\boldsymbol{\lambda}_{t+1} - \boldsymbol{\lambda}_t\|_2^2}{2\rho}$$
$$\le \frac{1}{T}\sum_{t=2}^{T}\max\{\frac{\alpha}{2\eta_t} - \frac{\alpha}{2\eta_{t-1}} - \frac{\sigma}{2}, 0\}R^2$$
$$+ \frac{1}{T}\sum_{t=1}^{T}\frac{\eta_t}{2}\|\hat{\mathbf{g}}_t\|^2 + \frac{K}{T}, \qquad (31)$$

where $K$ is a constant term related to $\alpha$, $R$, $\eta_1$, $\mathbf{A}$, $\mathbf{B}$, $\boldsymbol{\lambda}$, $\rho$ and $L_2$, $K = \frac{\alpha R^2}{2\eta_1} + R\|\boldsymbol{\lambda}\|_2\|\mathbf{A}\|_F + \frac{1}{2\rho}\|\boldsymbol{\lambda}\|_2^2 + R\|\mathbf{B}\|_F\|\boldsymbol{\lambda}\|_2 + R(L_2 + \|\mathbf{B}^T\boldsymbol{\lambda}\|_2)$, and we have used the fact that $\|\hat{\mathbf{g}}_t\|_{\mathbf{G}_t^{-1}}^2 \le \|\hat{\mathbf{g}}_t\|_2^2$ (due to $\mathbf{G}_t^{-1} \preceq \mathbf{I}$).

Based on (31) we continue to prove Theorem 1. When $\mathbf{B}$ is invertible and $\mathbf{y}_t' = \mathbf{B}^{-1}(\mathbf{c} - \mathbf{A}\mathbf{x}_t)$, we obtain

$$\mathbf{B}(\mathbf{y}_t' - \mathbf{y}_t) = \frac{1}{\rho}(\boldsymbol{\lambda}_t - \boldsymbol{\lambda}_{t-1}). \qquad (32)$$

Based on the convexity of $f$ and $\phi$, we obtain

$$f_t(\mathbf{x}_t) + \phi(\mathbf{y}_t') \le f_t(\mathbf{x}_t) + \phi(\mathbf{y}_t) + \langle \partial\phi(\mathbf{y}_t'), \mathbf{y}_t' - \mathbf{y}_t\rangle$$
$$= f_t(\mathbf{x}_t) + \phi(\mathbf{y}_t) + \frac{1}{\rho}\langle (\mathbf{B}^{-1})^T\partial\phi(\mathbf{y}_t'), \boldsymbol{\lambda}_t - \boldsymbol{\lambda}_{t-1}\rangle, \qquad (33)$$

where the last equality holds due to (32).

Let $(\mathbf{x}^*, \mathbf{y}^*)$ be the optimal solution (implying $\mathbf{A}\mathbf{x}^* + \mathbf{B}\mathbf{y}^* - \mathbf{c} = \mathbf{0}$). For any dual variable $\boldsymbol{\lambda}^*$ and $\tilde{\mathbf{v}}_t = [\mathbf{x}_t^T, \mathbf{y}_t^T, \tilde{\boldsymbol{\lambda}}_t^T]^T$, we have

$$(\tilde{\mathbf{v}}_t - \mathbf{v}^*)^T H(\tilde{\mathbf{v}}_t) = H(\mathbf{v}^*)^T(\tilde{\mathbf{v}}_t - \mathbf{v}^*)$$
$$= \begin{bmatrix} -\mathbf{A}^T\boldsymbol{\lambda}^* \\ -\mathbf{B}^T\boldsymbol{\lambda}^* \\ \mathbf{A}\mathbf{x}^* + \mathbf{B}\mathbf{y}^* - \mathbf{c} \end{bmatrix}^T \begin{bmatrix} \mathbf{x}_t - \mathbf{x}^* \\ \mathbf{y}_t - \mathbf{y}^* \\ \tilde{\boldsymbol{\lambda}}_t - \boldsymbol{\lambda}^* \end{bmatrix}$$
$$= \langle \boldsymbol{\lambda}^*, \mathbf{c} - \mathbf{A}\mathbf{x}_t - \mathbf{B}\mathbf{y}_t\rangle = \frac{1}{\rho}\langle \boldsymbol{\lambda}^*, \boldsymbol{\lambda}_t - \boldsymbol{\lambda}_{t-1}\rangle \qquad (34)$$

where $\mathbf{v}^* := [(\mathbf{x}^*)^T, (\mathbf{y}^*)^T, (\boldsymbol{\lambda}^*)^T]^T$, and the affine mapping $H(\cdot)$ is given by (25).

Setting $\boldsymbol{\lambda}^* = (\mathbf{B}^{-1})^T\partial\phi(\mathbf{y}_t')$, based on (33) and (34) we have

$$f_t(\mathbf{x}_t) + \phi(\mathbf{y}_t') - (f_t(\mathbf{x}^*) + \phi(\mathbf{y}^*))$$
$$\le f_t(\mathbf{x}_t) + \phi(\mathbf{y}_t) + (\tilde{\mathbf{v}}_t - \mathbf{v}^*)^T H(\tilde{\mathbf{v}}_t)$$
$$- (f_t(\mathbf{x}^*) + \phi(\mathbf{y}^*)). \qquad (35)$$

Combining (31) and (35) yields

$$\frac{1}{T}\sum_{t=1}^{T}(f_t(\mathbf{x}_t) + \phi(\mathbf{y}_t')) - \frac{1}{T}\sum_{t=1}^{T}(f_t(\mathbf{x}^*) + \phi(\mathbf{y}^*))$$
$$+ \frac{1}{T}\sum_{t=1}^{T}\frac{\|\boldsymbol{\lambda}_{t+1} - \boldsymbol{\lambda}_t\|_2^2}{2\rho} \le \frac{1}{T}\sum_{t=1}^{T}(f_t(\mathbf{x}_t) + \phi(\mathbf{y}_t))$$
$$- \frac{1}{T}\sum_{t=1}^{T}(f_t(\mathbf{x}^*) + \phi(\mathbf{y}^*)) + \frac{1}{T}\sum_{t=1}^{T}(\tilde{\mathbf{v}}_t - \mathbf{v}^*)^T H(\tilde{\mathbf{v}}_t)$$
$$+ \frac{1}{T}\sum_{t=1}^{T}\frac{\|\boldsymbol{\lambda}_{t+1} - \boldsymbol{\lambda}_t\|_2^2}{2\rho}$$
$$\le \frac{1}{T}\sum_{t=2}^{T}\max\{\frac{\alpha}{2\eta_t} - \frac{\alpha}{2\eta_{t-1}} - \frac{\sigma}{2}, 0\}R^2$$
$$+ \frac{1}{T}\sum_{t=1}^{T}\frac{\eta_t}{2}\|\hat{\mathbf{g}}_t\|_2^2 + \frac{K}{T}. \qquad (36)$$

Since $\boldsymbol{\lambda}_{t+1} - \boldsymbol{\lambda}_t = \rho(\mathbf{A}\mathbf{x}_{t+1} + \mathbf{B}\mathbf{y}_{t+1} - \mathbf{c})$, from (36) we have

$$\frac{1}{T}\sum_{t=1}^{T}(f_t(\mathbf{x}_t) + \phi(\mathbf{y}_t')) - \frac{1}{T}\sum_{t=1}^{T}(f_t(\mathbf{x}^*) + \phi(\mathbf{y}^*))$$
$$+ \frac{\rho}{2T}\sum_{t=1}^{T}\|\mathbf{A}\mathbf{x}_{t+1} + \mathbf{B}\mathbf{y}_{t+1} - \mathbf{c}\|_2^2$$
$$\le \frac{1}{T}\sum_{t=2}^{T}\max\{\frac{\alpha}{2\eta_t} - \frac{\alpha}{2\eta_{t-1}} - \frac{\sigma}{2}, 0\}R^2$$
$$+ \frac{1}{T}\sum_{t=1}^{T}\frac{\eta_t}{2}\|\hat{\mathbf{g}}_t\|_2^2 + \frac{K}{T}. \qquad (37)$$

Taking expectations for both sides of (37) with respect to its randomness, we have

$$\mathbb{E}\left[\frac{1}{T}\sum_{t=1}^{T}(f_t(\mathbf{x}_t)+\phi(\mathbf{y}_t'))-\frac{1}{T}\sum_{t=1}^{T}(f_t(\mathbf{x}^*)+\phi(\mathbf{y}^*))\right]$$

$$+\mathbb{E}\left[\frac{\rho}{2T}\sum_{t=1}^{T}\|\mathbf{A}\mathbf{x}_{t+1}+\mathbf{B}\mathbf{y}_{t+1}-\mathbf{c}\|_2^2\right]$$

$$\leq\frac{1}{T}\sum_{t=2}^{T}\max\{\frac{\alpha}{2\eta_t}-\frac{\alpha}{2\eta_{t-1}}-\frac{\sigma}{2},0\}R^2$$

$$+\frac{1}{T}\sum_{t=1}^{T}\frac{\eta_t}{2}\mathbb{E}[\|\hat{\mathbf{g}}_t\|_2^2]+\frac{K}{T}. \tag{38}$$

Based on (Duchi et al., 2015, Lemma 1), the second-order statistics of the gradient estimate $\hat{\mathbf{g}}_t$ is given by

$$\mathbb{E}_{\mathbf{z}_t}[\hat{\mathbf{g}}_t]=\mathbf{g}_t+\beta_t L_g(\mathbf{w}_t)\nu(\mathbf{x}_t,\beta_t), \tag{39}$$

$$\mathbb{E}_{\mathbf{z}_t}[\|\hat{\mathbf{g}}_t\|_2^2]\leq 2s(m)\|\mathbf{g}_t\|_2^2+\frac{1}{2}\beta_t^2 L_g(\mathbf{w}_t)^2 M(\mu)^2, \tag{40}$$

where $\mathbf{g}_t = \nabla_{\mathbf{x}}f(\mathbf{x};\mathbf{w}_t)|_{\mathbf{x}=\mathbf{x}_t}$, $\|\nu(\mathbf{x}_t,\beta_t)\|_2 \leq \frac{1}{2}\mathbb{E}_{\mathbf{z}}[\|\mathbf{z}\|_2^3]$, $L_g(\mathbf{w}_t)$ is defined in Assumption C, and $s(m)$ and $M(\mu)$ are introduced in Assumption E. According to (40), we have

$$\mathbb{E}[\|\hat{\mathbf{g}}_t\|_2^2]=\mathbb{E}\left[\mathbb{E}_{\mathbf{z}}[\|\hat{\mathbf{g}}_t\|_2^2]\right]$$

$$\leq\mathbb{E}\left[2s(m)\|\mathbf{g}_t\|_2^2+\frac{1}{2}\beta_t^2 L_{g,t}^2 M(\mu)^2\right]$$

$$\leq 2s(m)L_1^2+\frac{1}{2}\beta_t^2 L_g^2 M(\mu)^2, \tag{41}$$

where for ease of notation, we have replaced $L_g(\mathbf{w}_t)$ with $L_{g,t}$, and the last inequality holds due to Assumptions B and C.

Substituting (41) into (38), the expected average regret can be bounded as

$$\overline{\text{Regret}}_T(\mathbf{x}_t,\mathbf{y}_t',\mathbf{x}^*,\mathbf{y}^*)$$

$$\leq\frac{1}{T}\sum_{t=2}^{T}\max\{\frac{\alpha}{2\eta_t}-\frac{\alpha}{2\eta_{t-1}}-\frac{\sigma}{2},0\}R^2+\frac{s(m)L_1^2}{T}\sum_{t=1}^{T}\eta_t$$

$$+\frac{M(\mu)^2 L_g^2}{4T}\sum_{t=1}^{T}\eta_t\beta_t^2+\frac{K}{T}. \tag{42}$$

On the other hand, when $\mathbf{A}$ is invertible and $\mathbf{x}_t' = \mathbf{A}^{-1}(\mathbf{c}-\mathbf{B}\mathbf{y}_t)$, we obtain

$$\mathbf{A}(\mathbf{x}_t'-\mathbf{x}_t)=\frac{1}{\rho}(\boldsymbol{\lambda}_t-\boldsymbol{\lambda}_{t-1}).$$

Based on the convexity of $f$ and $\phi$, we obtain

$$f_t(\mathbf{x}_t')+\phi(\mathbf{y}_t)$$

$$\leq f_t(\mathbf{x}_t)+\phi(\mathbf{y}_t)+\langle\nabla f_t(\mathbf{x}_t'),\mathbf{x}_t'-\mathbf{x}_t\rangle$$

$$=f_t(\mathbf{x}_t)+\phi(\mathbf{y}_t)+\frac{1}{\rho}\langle(\mathbf{A}^{-1})^T\nabla f_t(\mathbf{x}_t'),\boldsymbol{\lambda}_t-\boldsymbol{\lambda}_{t-1}\rangle. \tag{43}$$

Setting $\boldsymbol{\lambda}^* = (\mathbf{A}^{-1})^T\nabla f_t(\mathbf{x}_t')$, based on (43) and (34) we have

$$f_t(\mathbf{x}_t')+\phi(\mathbf{y}_t)-(f_t(\mathbf{x}^*)+\phi(\mathbf{y}^*))\leq f_t(\mathbf{x}_t)$$

$$+\phi(\mathbf{y}_t)+(\tilde{\mathbf{v}}_t-\mathbf{v}^*)^T H(\tilde{\mathbf{v}}_t)-(f_t(\mathbf{x}^*)+\phi(\mathbf{y}^*)). \tag{44}$$

Since the right hand side (RHS) of (44) and RHS of (35) are same, we can then mimic the aforementioned procedure to prove that the regret $\overline{\text{Regret}}_T(\mathbf{x}_t',\mathbf{y}_t,\mathbf{x}^*,\mathbf{y}^*)$ obeys the same bounds as (42).

## 8.3 Simplification of Regret Bound

Consider terms in right hand side (RHS) of (42) together with $\eta_t = \frac{C_1}{\sqrt{s(m)}\sqrt{t}}$ and $\beta_t = \frac{C_2}{M(\mu)t}$, we have

$$\frac{1}{T}\sum_{t=2}^{T}\max\{\frac{\alpha}{2\eta_t}-\frac{\alpha}{2\eta_{t-1}}-\frac{\sigma}{2},0\}R^2$$

$$\leq\frac{1}{T}\sum_{t=2}^{T}(\frac{\alpha}{2\eta_t}-\frac{\alpha}{2\eta_{t-1}})R^2\leq\frac{1}{\sqrt{T}}\frac{\alpha R^2\sqrt{s(m)}}{2C_1},$$

$$\frac{s(m)L_1^2}{T}\sum_{t=1}^{T}\eta_t\leq\frac{2C_1\sqrt{s(m)}L_1^2}{\sqrt{T}},$$

$$\frac{M(\mu)^2 L_g^2}{4T}\sum_{t=1}^{T}\eta_t\beta_t^2=\frac{C_1 C_2^2 L_g^2}{4\sqrt{s(m)}T}\sum_{t=1}^{T}\frac{1}{t^{5/2}}$$

$$\leq\frac{5C_1 C_2^2 L_g^2}{12T}, \tag{45}$$

where we have used the facts that $\sum_{t=1}^{T}\frac{1}{\sqrt{t}}\leq 2\sqrt{T}$,

$$\sum_{t=1}^{T}(1/t^a)=1+\sum_{t=2}^{T}(1/t^a)$$

$$\leq 1+\int_{1}^{\infty}(1/t^a)=a/(a-1),\ \forall a>1, \tag{46}$$

and we recall that $s(m) = m \geq 1$. Substituting (45) into RHS of (42), we conclude that the expected average regret $\overline{\text{Regret}}_T(\mathbf{x}_t,\mathbf{y}_t',\mathbf{x}^*,\mathbf{y}^*)$ is upper bounded by

$$\frac{1}{\sqrt{T}}\frac{\alpha R^2\sqrt{s(m)}}{2C_1}+\frac{2C_1\sqrt{s(m)}L_1^2}{\sqrt{T}}+\frac{5C_1 C_2^2 L_g^2}{12T}+\frac{K}{T}. \tag{47}$$

## 8.4 Proof of Corollary 1

Given i.i.d. samples $\{\mathbf{w}_t\}$ drawn from the probability distribution $P$, from Theorem 1 we have

$$
\begin{aligned}
&\mathbb{E}\left[\frac{1}{T}\sum_{t=1}^{T}\left(f(\mathbf{x}_t;\mathbf{w}_t)+\phi(\mathbf{y}_t')\right)\right.\\
&\left.-\frac{1}{T}\sum_{t=1}^{T}\left(f(\mathbf{x}^*;\mathbf{w}_t)+\phi(\mathbf{y}^*)\right)\right]\\
&\leq\frac{1}{\sqrt{T}}\frac{\alpha R^2\sqrt{s(m)}}{2C_1}+\frac{2C_1\sqrt{s(m)}L_1^2}{\sqrt{T}}\\
&\quad+\frac{5C_1C_2^2L_g^2}{12}\frac{1}{T}+\frac{K}{T}.
\end{aligned}\tag{48}
$$

Based on $F(\mathbf{x},\mathbf{y})=\mathbb{E}_{\mathbf{w}}[f(\mathbf{x};\mathbf{w})]+\phi(\mathbf{y})$, from (48) we have

$$
\begin{aligned}
&\mathbb{E}\left[F(\bar{\mathbf{x}}_t,\bar{\mathbf{y}}_t)-F(\mathbf{x}^*,\mathbf{y}^*)\right]\\
&\leq\mathbb{E}\left[\frac{1}{T}\sum_{t=1}^{T}F(\mathbf{x}_t,\mathbf{y}_t)-F(\mathbf{x}^*,\mathbf{y}^*)\right]\\
&=\mathbb{E}_{\mathbf{z}_{1:T}}\left[\mathbb{E}_{\mathbf{w}_{1:T}}\left[\frac{1}{T}\sum_{t=1}^{T}\left(f(\mathbf{x}_t;\mathbf{w}_t)+\phi(\mathbf{y}_t')\right)\right.\right.\\
&\left.\left.\quad-\frac{1}{T}\sum_{t=1}^{T}\left(f(\mathbf{x}^*;\mathbf{w}_t)+\phi(\mathbf{y}^*)\right)\right]\right]\\
&\leq\frac{1}{\sqrt{T}}\frac{\alpha R^2\sqrt{s(m)}}{2C_1}+\frac{2C_1\sqrt{s(m)}L_1^2}{\sqrt{T}}\\
&\quad+\frac{5C_1C_2^2L_g^2}{12}\frac{1}{T}+\frac{K}{T},
\end{aligned}\tag{49}
$$

where the first inequality holds due to the convexity of $F$, and the second equality holds since $\mathbf{x}_t$ and $\mathbf{y}_t$ are implicit functions of i.i.d. random variables $\{\mathbf{w}_k\}_{k=1}^{t-1}$ and $\{\mathbf{z}_k\}_{k=1}^{t-1}$, and $\{\mathbf{w}_t\}$ and $\{\mathbf{z}_t\}$ are independent of each other.

## 8.5 Proof of Corollary 2

Substituting $\eta_t=\frac{\alpha}{\sigma t}$ and $\beta_t=\frac{C_2}{M(\mu)t}$ into RHS of (42), we have

$$
\begin{aligned}
&\frac{1}{T}\sum_{t=2}^{T}\max\{\frac{\alpha}{2\eta_t}-\frac{\alpha}{2\eta_{t-1}}-\frac{\sigma}{2},0\}R^2=0,\\
&\frac{s(m)L_1^2}{T}\sum_{t=1}^{T}\eta_t\leq\frac{\alpha s(m)L_1^2\log T}{\sigma T},\\
&\frac{M(\mu)^2L_g^2}{4T}\sum_{t=1}^{T}\eta_t\beta_t^2=\frac{\alpha C_2^2L_g^2}{4\sigma T}\sum_{t=1}^{T}\frac{1}{t^3}\leq\frac{3\alpha C_2^2L_g^2}{8\sigma T},
\end{aligned}\tag{50}
$$

where we have used the facts that $\sum_{t=1}^{T}\frac{1}{t}\leq 1+\log T$ and (46). Based on (50) and (47), we complete the proof.

## 8.6 Proof of Corollary 3

We consider the hybrid minibatch strategy

$$
\hat{\mathbf{g}}_t=\frac{1}{q_1q_2}\sum_{j=1}^{q_1}\sum_{i=1}^{q_2}\frac{f(\mathbf{x}_t+\beta_t\mathbf{z}_{t,j};\mathbf{w}_{t,i})-f(\mathbf{x}_t;\mathbf{w}_{t,i})}{\beta_t}\mathbf{z}_{t,j}
\tag{51}
$$

with $\hat{\mathbf{g}}_{t,ij}:=\frac{f(\mathbf{x}_t+\beta_t\mathbf{z}_{t,j};\mathbf{w}_{t,i})-f(\mathbf{x}_t;\mathbf{w}_{t,i})}{\beta_t}\mathbf{z}_{t,j}$. Based on (39) and i.i.d. samples $\{\mathbf{w}_{t,i}\}$ and $\{\mathbf{z}_{t,j}\}$, we have

$$
\bar{\mathbf{g}}_t:=\mathbb{E}[\hat{\mathbf{g}}_{t,ij}]=\mathbb{E}[\mathbf{g}_t]+\beta_t\mathbb{E}[L_{g,t}\nu(\mathbf{x}_t,\beta_t)],\ \forall i,j.\tag{52}
$$

where for ease of notation we have replaced $L_g(\mathbf{w}_t)$ with $L_{g,t}$, $\|\nu(\mathbf{x}_t,\beta_t)\|_2\leq\frac{1}{2}\mathbb{E}[\|\mathbf{z}\|_2^3]\leq M(\mu)$ due to Assumption E. From (51), we obtain

$$
\begin{aligned}
\mathbb{E}[\|\hat{\mathbf{g}}_t\|_2^2]&=\mathbb{E}\left[\left\|\frac{1}{q_1q_2}\sum_{i=1}^{q_1}\sum_{j=1}^{q_2}(\hat{\mathbf{g}}_{t,ij}-\bar{\mathbf{g}}_t)+\bar{\mathbf{g}}_t\right\|_2^2\right]\\
&=\|\bar{\mathbf{g}}_t\|_2^2+\mathbb{E}\left[\left\|\frac{1}{q_1q_2}\sum_{i=1}^{q_1}\sum_{j=1}^{q_2}(\hat{\mathbf{g}}_{t,ij}-\bar{\mathbf{g}}_t)\right\|_2^2\right]\\
&=\|\bar{\mathbf{g}}_t\|_2^2+\frac{1}{q_1q_2}\mathbb{E}[\|\hat{\mathbf{g}}_{t,11}-\bar{\mathbf{g}}_t\|_2^2]=\|\bar{\mathbf{g}}_t\|^2\\
&\quad+\frac{1}{q_1q_2}\mathbb{E}[\|\hat{\mathbf{g}}_{t,11}\|^2]-\frac{1}{q_1q_2}\|\bar{\mathbf{g}}_t\|^2,
\end{aligned}\tag{53}
$$

where we have used the fact that $\mathbb{E}[\hat{\mathbf{g}}_{t,ij}]=\mathbb{E}[\hat{\mathbf{g}}_{t,11}]$ for any $i$ and $j$.

The definition of $\bar{\mathbf{g}}_t$ in (52) yields

$$
\begin{aligned}
\|\bar{\mathbf{g}}_t\|^2&\leq 2\|\mathbb{E}[\mathbf{g}_t]\|_2^2+2\|\beta_t\mathbb{E}[L_{g,t}\nu(\mathbf{x}_t,\beta_t)]\|_2^2\\
&\leq 2\mathbb{E}[\|\mathbf{g}_t\|_2^2]+2\beta_t^2\mathbb{E}[L_{g,t}^2]\mathbb{E}[\|\nu(\mathbf{x}_t,\beta_t)\|_2^2]\\
&\leq 2\mathbb{E}[\|\mathbf{g}_t\|_2^2]+\frac{1}{2}\beta_t^2L_g^2M(\mu)^2,
\end{aligned}\tag{54}
$$

where the first inequality holds due to Cauchy-Schwarz inequality, and the second inequality holds due to Jensen's inequality. From (40), we obtain

$$
\mathbb{E}[\|\hat{\mathbf{g}}_{t,11}\|^2]\leq 2s(m)\mathbb{E}[\|\mathbf{g}_t\|_2^2]+\frac{1}{2}\beta_t^2L_g^2M(\mu)^2.\tag{55}
$$

Substituting (54) and (55) into (53), we obtain

$$
\begin{aligned}
\mathbb{E}[\|\hat{\mathbf{g}}_t\|_2^2]&\leq\|\bar{\mathbf{g}}_t\|_2^2+\frac{1}{q_1q_2}\mathbb{E}[\|\hat{\mathbf{g}}_{t,11}\|_2^2]\\
&\leq 2(1+\frac{s(m)}{q_1q_2})\mathbb{E}[\|\mathbf{g}_t\|_2^2]+\frac{q_1q_2+1}{2q_1q_2}\beta_t^2L_g^2M(\mu)^2.
\end{aligned}\tag{56}
$$

Similar to proof of Theorem 1, substituting (56) into

(38), we obtain

$$\overline{\text{Regret}}_T(\mathbf{x}_t, \mathbf{y}'_t, \mathbf{x}^*, \mathbf{y}^*)$$

$$\leq \frac{1}{T}\sum_{t=2}^{T}\max\{\frac{\alpha}{2\eta_t} - \frac{\alpha}{2\eta_{t-1}} - \frac{\sigma}{2}, 0\}R^2$$

$$+ \frac{1}{T}\sum_{t=1}^{T}\frac{\eta_t}{2}\mathbb{E}[\|\hat{\mathbf{g}}_t\|_2^2] + \frac{K}{T}$$

$$\leq \frac{1}{T}\sum_{t=2}^{T}\max\{\frac{\alpha}{2\eta_t} - \frac{\alpha}{2\eta_{t-1}} - \frac{\sigma}{2}, 0\}R^2$$

$$+ \frac{(q_1 q_2 + s(m))L_1^2}{q_1 q_2 T}\sum_{t=1}^{T}\eta_t$$

$$+ \frac{(q_1 q_2 + 1)L_g^2 M(\mu)^2}{4 q_1 q_2 T}\sum_{t=1}^{T}\eta_t \beta_t^2 + \frac{K}{T}. \qquad (57)$$

Substituting $\eta_t = \frac{C_1}{\sqrt{1 + \frac{s(m)}{q_1 q_2}}\sqrt{t}}$ and $\beta_t = \frac{C_2}{M(\mu)t}$ into (57), we obtain

$$\overline{\text{Regret}}_T(\mathbf{x}_t, \mathbf{y}'_t, \mathbf{x}^*, \mathbf{y}^*)$$

$$\leq \frac{\alpha R^2}{2 C_1}\frac{\sqrt{1 + \frac{s(m)}{q_1 q_2}}}{\sqrt{T}} + 2 C_1 L_1^2 \frac{\sqrt{1 + \frac{s(m)}{q_1 q_2}}}{\sqrt{T}}$$

$$+ \frac{5 C_1 C_2^2 L_g^2}{12 T}\frac{q_1 q_2 + 1}{q_1 q_2 \sqrt{1 + \frac{s(m)}{q_1 q_2}}} + \frac{K}{T}$$

$$\leq \frac{\alpha R^2}{2 C_1}\frac{\sqrt{1 + \frac{s(m)}{q_1 q_2}}}{\sqrt{T}} + 2 C_1 L_1^2 \frac{\sqrt{1 + \frac{s(m)}{q_1 q_2}}}{\sqrt{T}}$$

$$+ \frac{5 C_1 C_2^2 L_g^2}{6}\frac{1}{T} + \frac{K}{T}, \qquad (58)$$

which then completes the proof.

### 8.7 ZOO-ADMM for Sensor Selection

We recall that the sensor selection problem can be cast as

$$\underset{\mathbf{x}, \mathbf{y}}{\text{minimize}} \quad \frac{1}{T}\sum_{t=1}^{T}f(\mathbf{x}; \mathbf{w}_t) + \mathcal{I}_1(\mathbf{x}) + \mathcal{I}_2(\mathbf{y})$$
$$\text{subject to} \quad \mathbf{x} - \mathbf{y} = \mathbf{0}, \qquad (59)$$

where $\mathbf{y} \in \mathbb{R}^m$ is an auxiliary variable, $f(\mathbf{x}; \mathbf{w}_t) = -\text{logdet}(\sum_{i=1}^{m}x_i \mathbf{a}_{i,t}\mathbf{a}_{i,t}^T)$ with $\mathbf{w}_t = \{\mathbf{a}_{i,t}\}_{i=1}^{m}$, and $\{\mathcal{I}_i\}$ are indicator functions

$$\mathcal{I}_1(\mathbf{x}) = \begin{cases} 0 & \mathbf{0} \leq \mathbf{x} \leq \mathbf{1} \\ \infty & \text{otherwise,} \end{cases} \quad \mathcal{I}_2(\mathbf{y}) = \begin{cases} 0 & \mathbf{1}^T \mathbf{y} = m_0 \\ \infty & \text{otherwise.} \end{cases}$$

Based on (59), two key steps of ZOO-ADMM (22)-(23)

are given by

$$\mathbf{x}_{t+1} = \arg\min_{\mathbf{0} \leq \mathbf{x} \leq \mathbf{1}}\left\{\|\mathbf{x} - \mathbf{d}_t\|_2^2\right\}, \qquad (60)$$

$$\mathbf{y}_{t+1} = \arg\min_{\mathbf{1}^T \mathbf{y} = m_0}\left\{\|\mathbf{y} - (\mathbf{x}_{t+1} - (1/\rho)\boldsymbol{\lambda}_t)\|_2^2\right\}, \quad (61)$$

where $\hat{\mathbf{g}}_t$ is the gradient estimate, and $\mathbf{d}_t := \frac{\eta_t}{\alpha}(-\hat{\mathbf{g}}_t + \boldsymbol{\lambda}_t - \rho\mathbf{x}_t + \rho\mathbf{y}_t) + \mathbf{x}_t$. Sub-problems (60) and (61) yield closed-form solutions as below (Parikh and Boyd, 2014)

$$[\mathbf{x}_{t+1}]_i = \begin{cases} 0 & [\mathbf{d}_t]_i < 0 \\ [\mathbf{d}_t]_i & [\mathbf{d}_t]_i \in [0,1] \\ 1 & [\mathbf{d}_t]_i > 1, \end{cases} \quad \text{and} \qquad (62)$$

$$\mathbf{y}_{t+1} = \mathbf{x}_{t+1} - \frac{1}{\rho}\boldsymbol{\lambda}_t + \frac{m_0 - \mathbf{1}^T(\mathbf{x}_{t+1} - \boldsymbol{\lambda}_t/\rho)}{m}\mathbf{1}_m, \qquad (63)$$

where $[\mathbf{x}]_i$ denote the $i$th entry of $\mathbf{x}$.

### 8.8 ZOO-ADMM for Sparse Cox Regression

This sparse regression problem can formulated as

$$\underset{\mathbf{x}, \mathbf{y}}{\text{minimize}} \quad \frac{1}{n}\sum_{i=1}^{n}f(\mathbf{x}; \mathbf{w}_i) + \gamma\|\mathbf{y}\|_1$$
$$\text{subject to} \quad \mathbf{x} - \mathbf{y} = \mathbf{0}, \qquad (64)$$

where $f(\mathbf{x}; \mathbf{w}_i) = \delta_i\left\{-\mathbf{a}_i^T\mathbf{x} + \log\left(\sum_{j \in \mathcal{R}_i}e^{\mathbf{a}_j^T\mathbf{x}}\right)\right\}$ with $\mathbf{w}_i = \mathbf{a}_i$. By using the ZOO-ADMM algorithm, we can avoid the gradient calculation for the involved objective function in Cox regression. The two key steps of ZOO-ADMM (22)-(23) at iteration $i$ become

$$\mathbf{x}_{i+1} = \frac{\eta_t}{\alpha}(-\hat{\mathbf{g}}_i + \boldsymbol{\lambda}_i - \rho\mathbf{x}_i + \rho\mathbf{y}_i) + \mathbf{x}_i, \qquad (65)$$

$$\mathbf{y}_{i+1} = \arg\min_{\mathbf{y}}\left\{\|\mathbf{y}\|_1 + \frac{\rho}{2\gamma}\|\mathbf{y} - \mathbf{d}_i\|_2^2\right\}, \qquad (66)$$

where $\hat{\mathbf{g}}_i$ is the gradient estimate, $\mathbf{d}_i = (\mathbf{x}_{i+1} - (1/\rho)\boldsymbol{\lambda}_i)$, and the solution of sub-problem (66) is given by the soft-thresholding operator at the point $\mathbf{d}_i$ with parameter $\rho/\gamma$ (Parikh and Boyd, 2014, Sec. 6)

$$[\mathbf{y}_{i+1}]_k = \begin{cases} (1 - \frac{\gamma}{\rho|[\mathbf{d}_i]_k|})[\mathbf{d}_i]_k & [\mathbf{d}_i]_k > \frac{\gamma}{\rho} \\ 0 & [\mathbf{d}_i]_k \leq \frac{\gamma}{\rho}, \end{cases}$$

for $k = 1, 2, \ldots, m$.

## References

A. Agarwal, O. Dekel, and L. Xiao. Optimal algorithms for online convex optimization with multipoint bandit feedback. In *COLT*, pages 28–40, 2010.

F. Bach, R. Jenatton, J. Mairal, and G. Obozinski. Optimization with sparsity-inducing penalties. *Foundations and Trends® in Machine Learning*, 4 (1):1–106, 2012.

D. Blatt, A. O. Hero, and H. Gauchman. A convergent incremental gradient method with a constant step size. *SIAM Journal on Optimization*, 18(1):29–51, 2007.

S. Boyd and L. Vandenberghe. *Convex optimization.* Cambridge university press, 2004.

S. Boyd, N. Parikh, E. Chu, B. Peleato, and J. Eckstein. Distributed optimization and statistical learning via the alternating direction method of multipliers. *Foundations and Trends® in Machine Learning*, 3(1):1–122, 2011.

S. Bubeck et al. Convex optimization: Algorithms and complexity. *Foundations and Trends® in Machine Learning*, 8(3-4):231–357, 2015.

E. J. Candès and B. Recht. Exact matrix completion via convex optimization. *Foundations of Computational mathematics*, 9(6):717, 2009.

P.-L. Chen, C.-T. Tsai, Y.-N. Chen, K.-C. Chou, C.-L. Li, C.-H. Tsai, K.-W. Wu, Y.-C. Chou, C.-Y. Li, W.-S. Lin, et al. A linear ensemble of individual and blended models for music rating prediction. In *Proceedings of the 2011 International Conference on KDD Cup*, pages 21–60. JMLR. org, 2011.

P.-Y. Chen and S. Liu. Bias-variance tradeoff of graph laplacian regularizer. *IEEE Signal Processing Letters*, 2017.

P.-Y. Chen, H. Zhang, Y. Sharma, J. Yi, and C.-J. Hsieh. Zoo: Zeroth order optimization based black-box attacks to deep neural networks without training substitute models. *arXiv preprint arXiv:1708.03999*, 2017.

A. Cotter, O. Shamir, N. Srebro, and K. Sridharan. Better mini-batch algorithms via accelerated gradient methods. In *Advances in neural information processing systems*, pages 1647–1655, 2011.

D. R. Cox. Regression models and life-tables. *Journal of the Royal Statistical Society. Series B (Methodological)*, 34(2):187–220, 1972.

G. Dror, N. Koenigstein, Y. Koren, and M. Weimer. The yahoo! music dataset and kdd-cup11. In *Proceedings of KDD Cup 2011*, pages 3–18, 2012.

J. Duchi, S. Shalev-Shwartz, Y. Singer, and A. Tewari. Composite objective mirror descent. In *COLT*, pages 14–26, 2010.

J. Duchi, E. Hazan, and Y. Singer. Adaptive subgradient methods for online learning and stochastic optimization. *Journal of Machine Learning Research*, 12(Jul):2121–2159, 2011.

J. C. Duchi, M. I. Jordan, M. J. Wainwright, and A. Wibisono. Optimal rates for zero-order convex optimization: The power of two function evaluations. *IEEE Transactions on Information Theory*, 61(5):2788–2806, 2015.

X. Gao, B. Jiang, and S. Zhang. On the information-adaptive variants of the admm: An iteration complexity perspective. *Journal of Scientific Computing*, Dec 2017. ISSN 1573-7691.

S. Ghadimi and G. Lan. Stochastic first-and zeroth-order methods for nonconvex stochastic programming. *SIAM Journal on Optimization*, 23(4):2341–2368, 2013.

D. Hajinezhad, M. Hong, and A. Garcia. Zenith: A zeroth-order distributed algorithm for multi-agent nonconvex optimization. 2017.

E. C. Hall and R. M. Willett. Online convex optimization in dynamic environments. *IEEE Journal of Selected Topics in Signal Processing*, 9(4):647–662, June 2015. ISSN 1932-4553.

E. Hazan. Introduction to online convex optimization. *Foundations and Trends® in Optimization*, 2(3-4): 157–325, 2016.

A. O. Hero and D. Cochran. Sensor management: Past, present, and future. *IEEE Sensors Journal*, 11(12):3064–3075, 2011.

S. Hosseini, A. Chapman, and M. Mesbahi. Online distributed convex optimization on dynamic networks. *IEEE Transactions on Automatic Control*, 61(11): 3545–3550, 2016.

S. Joshi and S. Boyd. Sensor selection via convex optimization. *IEEE Transactions on Signal Processing*, 57(2):451–462, 2009.

M. Li, T. Zhang, Y. Chen, and A. J. Smola. Efficient mini-batch training for stochastic optimization. In *Proceedings of the 20th ACM SIGKDD international conference on Knowledge discovery and data mining*, pages 661–670. ACM, 2014.

X. Lian, H. Zhang, C.-J. Hsieh, Y. Huang, and J. Liu. A comprehensive linear speedup analysis for asynchronous stochastic parallel optimization from zeroth-order to first-order. In *Advances in Neural Information Processing Systems*, pages 3054–3062, 2016.

F. Lin, M. Fardad, and M. R. Jovanovic. Algorithms for leader selection in stochastically forced consensus networks. *IEEE Transactions on Automatic Control*, 59(7):1789–1802, 2014.

S. Liu, S. P. Chepuri, M. Fardad, E. Maşazade, G. Leus, and P. K. Varshney. Sensor selection for estimation with correlated measurement noise. *IEEE Transactions on Signal Processing*, 64(13): 3509–3522, 2016.

Y. Nesterov and V. Spokoiny. Random gradient-free minimization of convex functions. *Foundations of Computational Mathematics*, 2(17):527–566, 2015.

H. Ouyang, N. He, L. Tran, and A. Gray. Stochastic alternating direction method of multipliers. In *International Conference on Machine Learning*, pages 80–88, 2013.

N. Parikh and S. Boyd. Proximal algorithms. *Foundations and Trends® in Optimization*, 1(3):127–239, 2014.

M. Y. Park and T. Hastie. L1-regularization path algorithm for generalized linear models. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, 69(4):659–677, 2007.

C. R. Rao. *Linear statistical inference and its applications*, volume 2. Wiley New York, 1973.

N. L. Roux, M. Schmidt, and F. R. Bach. A stochastic gradient method with an exponential convergence rate for finite training sets. In *Advances in Neural Information Processing Systems*, pages 2663–2671, 2012.

S. Shalev-Shwartz. Online learning and online convex optimization. *Foundations and Trends® in Machine Learning*, 4(2):107–194, 2012.

O. Shamir. An optimal algorithm for bandit and zero-order convex optimization with two-point feedback. *Journal of Machine Learning Research*, 18(52):1–11, 2017.

I. Sohn, J. Kim, S.-H. Jung, and C. Park. Gradient lasso for cox proportional hazards model. *Bioinformatics*, 25(14):1775–1781, 2009.

T. Suzuki. Dual averaging and proximal gradient descent for online alternating direction multiplier method. In *International Conference on Machine Learning*, pages 392–400, 2013.

H. Wang and A. Banerjee. Online alternating direction method (longer version). *arXiv preprint arXiv:1306.3721*, 2013.

D. M. Witten and R. Tibshirani. Survival analysis with high-dimensional covariates. *Statistical methods in medical research*, 19(1):29–51, 2010.

L. Xiao. Dual averaging methods for regularized stochastic learning and online optimization. *Journal of Machine Learning Research*, 11(Oct.):2543–2596, 2010.

X. Zhang, M. Burger, and S. Osher. A unified primal-dual algorithm framework based on bregman iteration. *Journal of Scientific Computing*, 46(1):20–46, 2011.