

# USING TAYLOR-APPROXIMATED GRADIENTS TO IMPROVE THE FRANK-WOLFE METHOD FOR EMPIRICAL RISK MINIMIZATION

ZIKAI XIONG\* AND ROBERT M. FREUND†

**Abstract.** The Frank-Wolfe method has become increasingly useful in statistical and machine learning applications, due to the structure-inducing properties of the iterates, and especially in settings where linear minimization over the feasible set is more computationally efficient than projection. In the setting of Empirical Risk Minimization – one of the fundamental optimization problems in statistical and machine learning – the computational effectiveness of Frank-Wolfe methods typically grows linearly in the number of data observations  $n$ . This is in stark contrast to the case for typical stochastic projection methods. In order to reduce this dependence on  $n$ , we look to second-order smoothness of typical smooth loss functions (least squares loss and logistic loss, for example) and we propose amending the Frank-Wolfe method with Taylor series-approximated gradients, including variants for both deterministic and stochastic settings. Compared with current state-of-the-art methods in the regime where the optimality tolerance  $\varepsilon$  is sufficiently small, our methods are able to simultaneously reduce the dependence on large  $n$  while obtaining optimal convergence rates of Frank-Wolfe methods, in both the convex and non-convex settings. We also propose a novel adaptive step-size approach for which we have computational guarantees. Last of all, we present computational experiments which show that our methods exhibit very significant speed-ups over existing methods on real-world datasets for both convex and non-convex binary classification problems.

**Key words.** Frank-Wolfe, linear minimization oracle, Empirical Risk Minimization, linear prediction, computational complexity, convex optimization

**AMS subject classifications.** 90C25, 90C26, 90C06, 90C60, 68Q25

**1. Introduction.** Our problem of interest is the following Empirical Risk Minimization (ERM) problem (also called the *finite-sum* problem) with constraints:

$$(1.1) \quad \text{ERM:} \quad \text{minimize}_{x \in \mathcal{C}} \quad F(x) := \frac{1}{n} \sum_{i=1}^n f_i(x) ,$$

where  $f_i(\cdot) : \mathbb{R}^p \rightarrow \mathbb{R}$  is the loss function of observation  $i$  for  $i = 1, \dots, n$ , and  $x$  is comprised of the (model) coefficients/parameters which are to be optimized over the feasibility set  $\mathcal{C} \subset \mathbb{R}^p$ . Problem (1.1) is a fundamental problem in statistical and machine learning, as it is the underlying optimization problem on which models are “trained” or “learned”, see [28, 10].

We are especially interested in and will focus in this study on those empirical risk minimization problems where the loss functions exhibit “linear prediction” [18], namely each loss function is of the form  $f_i(x) := l_i(w_i^\top x)$  where  $w_i$  is the  $i$ -th feature vector and  $l_i(\cdot) : \mathbb{R} \rightarrow \mathbb{R}$  is a univariate function; thus  $f_i$  is a scalar function of the linear function  $w_i^\top x$ . In this case (1.1) becomes

$$(1.2) \quad \text{ERM}_\ell : \quad \text{minimize}_{x \in \mathcal{C}} \quad F(x) := \frac{1}{n} \sum_{i=1}^n [f_i(x) = l_i(w_i^\top x)] ,$$

and we refer to (1.2) as “ERM<sub>ℓ</sub>” for “ERM with linear prediction”, because the  $i$ -th predicted value of interest is the linear function  $w_i^\top x$  of the model coefficients  $x$ . This

\*MIT Operations Research Center, 77 Massachusetts Avenue, Cambridge, MA 02139, USA (zikai@mit.edu). Research supported by AFOSR Grant No. FA9550-19-1-0240.

†MIT Sloan School of Management, 77 Massachusetts Avenue, Cambridge, MA 02139, USA (rfreund@mit.edu). Research supported by AFOSR Grant No. FA9550-19-1-0240.

structure arises in many canonical applications of statistical and machine learning problems, including support vector machines (SVMs), LASSO, logistic regression, and matrix completion, among others, see [20, 18] for other examples. Let  $W \in \mathbb{R}^{p \times n}$  be the matrix whose columns are comprised of the  $n$  feature vectors  $w_1, w_2, \dots, w_n$  in (1.2). We also define the function  $l(\cdot) : \mathbb{R}^n \rightarrow \mathbb{R}$  as  $l(\theta) := \frac{1}{n} \sum_{i=1}^n l_i(\theta_i)$ , and so the objective function of  $\text{ERM}_\ell$  can be written as  $F(x) = l(W^\top x)$  with  $(\theta_1, \dots, \theta_n) = \theta = W^\top x$ . Let  $w_i(\mathcal{C}) \subset \mathbb{R}$  denote the range of  $w_i^\top x$  over  $x \in \mathcal{C}$ . We make the following assumptions for  $\text{ERM}_\ell$  throughout this paper:

*Assumption 1.1.* The following conditions hold for the  $\text{ERM}_\ell$  problem (1.2):

- (i) the feasibility set  $\mathcal{C}$  is a compact convex set,
- (ii) for any  $g \in \mathbb{R}^p$ , the linear minimization problem  $\arg \min_{x \in \mathcal{C}} g^\top x$  can be easily solved,
- (iii) for all  $i \in [n]$  the univariate loss function  $l_i(\cdot)$  is twice-differentiable and its first-order derivative  $l'_i(\cdot)$  is  $L$ -Lipschitz continuous on the range of  $w_i^\top x$  over  $x \in \mathcal{C}$ , namely

$$(1.3) \quad |l'_i(\bar{\theta}) - l'_i(\hat{\theta})| \leq L|\bar{\theta} - \hat{\theta}| \quad \text{for any } \bar{\theta}, \hat{\theta} \in w_i(\mathcal{C}), \text{ and}$$

- (iv) for all  $i \in [n]$  the second-order derivative  $l''_i(\cdot)$  is  $\hat{L}$ -Lipschitz continuous on the range of  $w_i^\top x$  over  $x \in \mathcal{C}$ , namely

$$(1.4) \quad |l''_i(\bar{\theta}) - l''_i(\hat{\theta})| \leq \hat{L}|\bar{\theta} - \hat{\theta}| \quad \text{for any } \bar{\theta}, \hat{\theta} \in w_i(\mathcal{C}).$$

The first three conditions of Assumption 1.1 are canonical for Frank-Wolfe and other methods that rely on a Linear Minimization Oracle (LMO). However, condition (iv) of Assumption 1.1 concerns the smoothness of the second derivatives of the loss functions, and does not typically appear in the literature. In Section 1.1 we discuss and give some justification for this extra condition.

**1.1. Examples of oft-used loss functions with smooth second derivatives.** Here we discuss three oft-used loss functions that have smooth second derivatives and therefore satisfy condition (iv) of Assumption 1.1.

*Quadratic losses.* In very many practical instances of  $\text{ERM}_\ell$  the loss functions are of the form  $l_i(\theta_i) = \frac{1}{2}(y_i - \theta_i)^2$  for some observed value  $y_i$ , such as in LASSO instances [27], matrix completion with quadratic losses, and structured sparse matrix estimation with CUR factorization [19, 18]. We note that  $L = 1$  and  $\hat{L} = 0$  for quadratic losses.

*Logistic regression for binary classification.* In logistic regression we are given a training set  $\{(w_i, y_i)\}_{i=1}^n$  with  $w_i \in \mathbb{R}^p$  and  $y_i \in \{-1, 1\}$ , and we solve the maximum likelihood estimation problem by computing the optimal solution to the following logistic regression optimization problem:

$$(1.5) \quad \min_{x \in \mathcal{C}} F(x) := \frac{1}{n} \sum_{i=1}^n \ln(1 + \exp(-y_i w_i^\top x)).$$

In this problem,  $F(x) := \frac{1}{n} \sum_{i=1}^n l_i(w_i^\top x)$  and  $l_i(\theta_i) := \ln(1 + \exp(-y_i \theta_i))$ . For each  $i \in [n]$  and simplifying notation using  $v = \theta_i$ , we have  $l_i(v) = \ln(1 + \exp(-y_i v))$ ,  $l'_i(v) = y_i \frac{-1}{1 + \exp(y_i v)}$ ,  $l''_i(v) = y_i^2 \frac{\exp(y_i v)}{(1 + \exp(y_i v))^2}$ , and  $l'''_i(v) = y_i^3 \frac{\exp(y_i v) - \exp(2y_i v)}{(1 + \exp(y_i v))^3}$ . We note that  $|l''_i(v)|$  obtains its maximal value of  $\frac{1}{4}$  at  $v = 0$ , and hence  $L = \frac{1}{4}$ . Similarly,  $|l'''_i(v)|$  obtains its maximal value of  $\frac{1}{6\sqrt{3}}$  when  $\exp(y_i v) = 2 - \sqrt{3}$  and so  $\hat{L} = \frac{1}{6\sqrt{3}}$ . Summarizing,  $l'_i$  is  $\frac{1}{4}$ -Lipschitz continuous and  $l''_i$  is  $\frac{1}{6\sqrt{3}}$ -Lipschitz continuous for  $i \in [n]$ .

*Binary linear classification with non-convex losses* [20]. We are given a training set  $\{(w_i, y_i)\}_{i=1}^n$  with  $w_i \in \mathbb{R}^p$  and  $y_i \in \{0, 1\}$ , and the task is to learn  $\mathbb{P}[Y = 1 | w = w_i] = \sigma(w_i^\top x)$  for a given threshold function  $\sigma : \mathbb{R} \rightarrow [0, 1]$ . This leads to the optimization problem:

$$(1.6) \quad \min_{x \in \mathcal{C}} \frac{1}{n} \sum_{i=1}^n (y_i - \sigma(w_i^\top x))^2,$$

where  $F(x) := \frac{1}{n} \sum_{i=1}^n l_i(w_i^\top x)$ ,  $l_i(\cdot) := (y_i - \sigma(\cdot))^2$ , and  $\sigma : \mathbb{R} \rightarrow [0, 1]$  is the given threshold function, such as the sigmoid function  $\sigma(v) = (1 + \exp(-v))^{-1}$ . We note in general that (1.6) is a non-convex optimization problem. It is straightforward to derive  $l'_i(v) = -2\sigma'(v)(y_i - \sigma(v))$ ,  $l''_i(v) = 2(\sigma'(v))^2 - 2\sigma''(v)(y_i - \sigma(v))$ , and  $l'''_i(v) = 6\sigma'(v)\sigma''(v) - 2\sigma'''(v)(y_i - \sigma(v))$ . Suppose that there exists a constant  $L_\sigma$  such that  $|\sigma'(v)|$ ,  $|\sigma''(v)|$ , and  $|\sigma'''(v)| \leq L_\sigma$  for any  $v \in \mathbb{R}$ . Then it follows that  $l'_i$  is  $L := (2L_\sigma^2 + 2L_\sigma)$ -Lipschitz continuous and  $l''_i$  is  $\hat{L} := (6L_\sigma^2 + 2L_\sigma)$ -Lipschitz continuous.

In the often-encountered case when  $\sigma$  is the sigmoid function, we have  $\sigma'(x) = \frac{\exp(x)}{(1+\exp(x))^2}$ ,  $\sigma''(x) = \frac{\exp(x)(\exp(x)-1)}{(1+\exp(x))^3}$ , and  $\sigma'''(x) = \frac{\exp(x)(1-4\exp(x)+\exp(2x))}{(1+\exp(x))^4}$ , and thus  $|\sigma'(x)| \leq 1/4$ ,  $|\sigma''(x)| \leq 1/6\sqrt{3}$  and  $|\sigma'''(x)| \leq 1/24$ . Because  $y_i \in \{0, 1\}$  and  $\sigma(x) \in [0, 1]$ ,  $|y_i - \sigma(x)| \leq 1$ , it holds that  $|l'_i(x)| \leq \frac{1}{8} + \frac{1}{3\sqrt{3}}$  and  $|l''_i(x)| \leq \frac{1}{4\sqrt{3}} + \frac{1}{12}$ . Therefore for each  $i \in [n]$  we have  $l'_i$  is  $(\frac{1}{8} + \frac{1}{3\sqrt{3}})$ -Lipschitz continuous and  $l''_i$  is  $(\frac{1}{4\sqrt{3}} + \frac{1}{12})$ -Lipschitz continuous.

**1.2. Motivation and literature review.** With an ever-increasing number  $n$  of data observations in statistical and machine learning applications in ERM and  $\text{ERM}_\ell$ , there is growing importance in solving models where the number of observations  $n$  is very large (compared, say, to the dimension  $p$  of the features). The presence of large  $n$  does not pose a particular problem (at least in theory) for a typical version of the stochastic gradient descent method (SGD), which forms a gradient estimate at each iteration by randomly sampling a mini-batch of observations and then projects back onto the feasibility set  $\mathcal{C}$  after taking a gradient-estimate-based step. The computational complexity of SGD is not adversely affected by large values of  $n$ .

However, in many applications of ERM and  $\text{ERM}_\ell$  it is more attractive to use the Frank-Wolfe method or a version thereof, due to the structure-inducing properties of the iterates (sparsity, low-rank) especially when combined with in-face steps, see for example [7]. The Frank-Wolfe method is also attractive in settings where a Linear Minimization Oracle (LMO) on the feasibility set  $\mathcal{C}$  is more computationally efficient than projection onto  $\mathcal{C}$ . Here an LMO is defined to be a subroutine that solves and returns

$$\arg \min_{x \in \mathcal{C}} \langle g, x \rangle$$

for any given objective function vector  $g \in \mathbb{R}^p$ . For example, when  $\mathcal{C}$  is the nuclear norm ball, the LMO only requires computing the leading singular vector of a matrix, as compared to a projection oracle which requires a full SVD of a matrix, see [13] for example. As a result, the Frank-Wolfe method, which replaces gradient steps and projection by an LMO instead, has become increasingly useful in statistical and machine learning applications, for both convex and non-convex loss functions. However, when the number of observations  $n$  is large, the Frank-Wolfe method has lacked “good” efficiency as compared to SGD, as the number of operations of Frank-Wolfe typically

grows as  $O(n/\varepsilon)$ . The primary motivation of the research in this paper therefore is to address the following question:

*Are there efficient stochastic (or deterministic) Frank-Wolfe methods that eliminate or reduce the dependence on the number of observations  $n$ , in theory and in practice?*

In the last decade many researchers have studied stochastic versions of the Frank-Wolfe method. Before discussing these related works, let us introduce some useful measurements for Frank-Wolfe methods. We use the following Frank-Wolfe gap  $\mathcal{G}(x)$  as a measure of non-stationarity:

$$(1.7) \quad \mathcal{G}(x) := \max_{s \in \mathcal{C}} \langle x - s, \nabla F(x) \rangle .$$

The Frank-Wolfe gap is always nonnegative and is 0 if and only if  $x$  is a stationary point of (1.2). For problems with convex losses  $F$ , supposing that  $x^*$  is an optimal solution, we say  $x$  is an  $\varepsilon$ -optimal solution if and only if the optimality gap  $F(x) - F(x^*)$  is no larger than  $\varepsilon$ . Since the Frank-Wolfe gap is an upper bound on the optimality gap when  $F$  is convex,  $\mathcal{G}(x^k)$  is a conservative measure of the optimality gap at  $x^k$ , see [13] for example. The Frank-Wolfe gap does not necessarily bound the optimality gap in the non-convex setting; nevertheless  $\mathcal{G}(x)$  is always nonnegative and  $\mathcal{G}(x) = 0$  if and only if  $x$  is a stationary point of (1.2), and for this reason  $\mathcal{G}(x)$  is often used as a measure of non-stationarity at  $x$ , see [15, 24, 29, 22]. For an instance of (1.2) with non-convex objective function  $F$ , we therefore say that a point  $x \in \mathcal{C}$  is an  $\varepsilon$ -stationarity point of (1.2) if  $\mathcal{G}(x) \leq \varepsilon$ .

Generally speaking, for deterministic versions of the Frank-Wolfe method the lower bound LMO complexity to obtain an  $\varepsilon$ -optimal solution is  $O(1/\varepsilon)$ , see [13, 16]. And in the case when  $F$  is non-convex, the LMO complexity is at least  $O(1/\varepsilon^2)$ , see [26]. These deterministic methods require access to exact full-batch gradients per iteration, whose complexity scales with the number of observations  $n$ . When  $n$  is very large, computing exact gradients might be no longer efficient. We refer the reader to [6, 13, 17, 15] among other recent deterministic versions of the Frank-Wolfe method in this context.

Stochastic Frank-Wolfe methods avoid computing the exact gradient in each iteration by doing stochastic gradient estimation aimed at reducing the dependence on  $n$  in the overall computational complexity. Most existing stochastic Frank-Wolfe methods typically replace exact gradients by the average of sampled gradients as recorded in the current or previous iterates, so that they can reduce (or eliminate) the dependence on  $n$  in the overall complexity. Unfortunately the approaches that eliminate the dependence on  $n$  do so at the expense of a higher overall LMO complexity which tends to grow much faster than  $O(1/\varepsilon)$ . Representative works in this vein include [11, 17, 21]. A second approach is to estimate gradients by the variance reduction techniques which emerged from other stochastic first-order methods, such as SVRG [14] and SPIDER [4]. These methods still periodically require access to the exact full gradient, and they also compute averages from a gradually increasing number of gradient samples, see [11, 24, 29, 26, 9]. Also unfortunately, some of these methods perform worse in practice than the methods with even weaker theoretical guarantees [11]. A third approach is based on the primal-dual structure of  $\text{ERM}_\ell$  when  $f$  is convex. However, depending on the specific assumptions and problem set-up, the overall dependence of these methods on  $\varepsilon$  and  $n$  is still  $O(n/\varepsilon)$ , which is the same as the

traditional (deterministic) Frank-Wolfe method, see [18, 22]. In addition to the above approaches, some methods avoid directly computing gradients by estimating them via zeroth-order information (see, e.g., [8, 23, 12, 1, 25]). The theoretical complexities of these methods are usually no better than the methods using gradients. Recently, a new line of research has emerged that seeks to estimate gradients with the aid of second-order information, see for example [26, 30, 9, 30]. However, the cost per iteration in these methods may still grow fast, or they still periodically require access to the exact gradient. In summary, there are no existing stochastic Frank-Wolfe methods (that we are aware of) whose overall complexity does not scale with  $n$  while maintaining a dependence on  $\varepsilon$  of only  $O(1/\varepsilon)$  (for convex losses) or  $O(1/\varepsilon^2)$  (for non-convex losses).

**1.3. Contributions.** In the setting of  $\text{ERM}_\ell$  the contributions of this paper are as follows:

- We develop a new family of Frank-Wolfe methods, which we call TUFW for “Taylor-point Updating Frank-Wolfe,” that replaces the exact gradient computation by a sum of (second order) Taylor-approximated gradients around some current and previous iterates, which we call the *Taylor points*.
- Different versions of TUFW are developed based on different rules for constructing the batches of observations for Taylor-point updating at each iteration. For both convex and non-convex losses, we propose both stochastic and deterministic rules. All of our rules exhibit a *decreasing* number of gradient calls over the course of the intended iterations, while retaining the optimal LMO complexity of  $O(1/\varepsilon)$  (for convex losses) and  $O(1/\varepsilon^2)$  (for non-convex losses), together with other flops ( $O(n/\sqrt{\varepsilon})$  for convex losses and  $O(n/\varepsilon^{3/2})$  for non-convex losses). In the regime when  $\varepsilon$  is sufficiently small these other flops are minor compared with the LMO complexity. Our stochastic TUFW also avoids periodically computing exact gradients or gradually increasing the sample size.
- We present computational experiments which show that our methods exhibit very significant speed-ups over existing methods on real-world datasets for convex and non-convex binary classification problems.
- Our TUFW method easily extends to the more general ERM problem, with corresponding versions of our theoretical guarantees.
- We also propose a novel adaptive step-size method which has similar theoretical guarantees, and for which our computational experiments indicate superior performance in practice.

**1.4. Outline.** In [section 2](#) we present our new type of Frank-Wolfe method, which we call TUFW for “Taylor-point Updating Frank-Wolfe”, along with some elementary properties of the method. In [section 3](#) we present computational guarantees and overall complexity of TUFW for convex losses, and in [section 4](#) we present computational guarantees and overall complexity for non-convex losses. We present computational experiments for both convex and non-convex instances of  $\text{ERM}_\ell$  in [section 6](#).

**1.5. Notation.** Let  $\mathbb{N}$  denote the natural numbers. For any  $n \in \mathbb{N}$ , we use  $[n]$  to denote the set  $\{1, 2, \dots, n\}$ . For any real number  $v$ , we use  $\lfloor v \rfloor$  to denote the largest integer number less than or equal to  $v$ . Depending on the context, we will use either  $\langle a, b \rangle$  or  $a^\top b$  to denote the inner product between  $a$  and  $b$  in Euclidean space. For a vector  $x$  in  $\mathbb{R}^p$ , the norm  $\|x\|_q$  is given by  $(\sum_{j=1}^p x_j^q)^{1/q}$  for  $q \in [1, \infty)$  and

215  $\|x\|_\infty := \max_j |x_j|$ . The dual of norm  $\|\cdot\|$  is defined as  $\|z\|_* := \sup\{z^\top x : \|x\| \leq 1\}$ .  
 216 For a matrix  $A \in \mathbb{R}^{p \times p}$ , the matrix operator norm  $\|A\|$  induced by the norm  $\|\cdot\|$  is  
 217 defined as  $\sup_{x: \|x\| \leq 1} \|Ax\|_*$ . For any symmetric positive definite matrix  $A \in \mathbb{S}^{p \times p}$ , we  
 218 define its induced norm on  $\mathbb{R}^p$  by  $\|z\|_A := \sqrt{z^\top A z}$ . We define  $\sum_{i=j}^k \alpha_i := 0$  whenever  
 219  $j > k$ . We denote the uniform distribution on the set  $S$  by  $\mathcal{U}(S)$ , and we denote the  
 220 Bernoulli distribution with probability  $p$  by  $\text{Ber}(p)$ .

221 **2. A Frank-Wolfe method with Taylor-point Updating (TUFW).** In  
 222 this section we present a new type of Frank-Wolfe algorithm whose newness derives  
 223 from the way in which approximate gradients are computed. We call the algorithm  
 224 TUFW for “Taylor-point Updating Frank-Wolfe”. Algorithm TUFW is similar to the  
 225 standard Frank-Wolfe method except that it computes an inexpensive estimate  $g^k$  of  
 226 the gradient  $\nabla F(x^k)$  at iteration  $k$ . And it differs from other Frank-Wolfe methods  
 227 in the way it computes the estimate  $g^k$  in that (i) TUFW explicitly uses second-  
 228 derivatives of the loss functions  $l_i$ , and (ii) it can use very different batch-size rules  
 229 – both stochastic and deterministic. TUFW is described in [Algorithm 2.1](#). We now  
 230 explain the steps of the method.

---

**Algorithm 2.1** (TUFW) Frank-Wolfe Method with Taylor Point Updating

---

```

1: input initial point  $x^0$ , step-size sequence  $\{\gamma_k\}_k$ , and (possibly) number of itera-
   tions  $K$ 
2: initialize Taylor points  $\mathcal{B}_0 \leftarrow [n]$  and  $b_i \leftarrow x^0$  for  $i \in [n]$ 
3: for  $k = 0, 1, 2, \dots, K$ , do
4:   If  $k > 0$ , update Taylor points:
     Use some Rule to construct set  $\mathcal{B}_k \subset [n]$  of indices for Taylor-point updating
     Update Taylor points: for  $i \in \mathcal{B}_k$  update  $b_i \leftarrow x^k$ 
5:   Compute estimates of individual gradients:
      $g_i^k \leftarrow \nabla f_i(b_i) + \nabla^2 f_i(b_i)(x^k - b_i)$  for  $i \in [n]$ 
6:   Compute estimate of full gradient:  $g^k \leftarrow \frac{1}{n} \sum_{i=1}^n g_i^k$ 
7:   Solve LMO:  $s^k \leftarrow \arg \min_{s \in \mathcal{C}} \langle g^k, s \rangle$ 
8:   Update feasible solution:  $x^{k+1} \leftarrow x^k + \gamma_k(s^k - x^k)$ .
9: end for
10: return  $x^{K+1}$ 

```

---

Suppose we are at iteration  $k$ . Then as a prerequisite to solving the LMO (Linear Minimization Oracle) problem in [Line 7](#) of Algorithm TUFW, we first need to compute an (accurate) estimate  $g^k$  of  $\nabla F(x^k)$ , which in turn means computing an (accurate) estimate  $g_i^k$  of  $\nabla f_i(x^k)$  for  $i \in [n]$ . In an exact method we would simply compute  $g_i^k \leftarrow \nabla f_i(x^k)$  and then compute  $g^k \leftarrow \frac{1}{n} \sum_{i=1}^n g_i^k$ , the cost of which grows at least as  $O(np)$  operations. Here we instead consider, for each  $i \in [n]$ , computing the estimate  $g_i^k$  of  $\nabla f_i(x^k)$  by using the second-order Taylor approximation of  $\nabla f_i(x)$  around some point  $b_i$ , resulting in the estimate:

$$g_i^k \leftarrow \nabla f_i(b_i) + \nabla^2 f_i(b_i)(x^k - b_i) \quad \text{for } i \in [n],$$

231 which is [Line 5](#) of TUFW. We refer to  $b_i$  as the “Taylor point” for the Taylor approxi-  
 232 mation of  $\nabla f_i(\cdot)$ , and  $b_1, \dots, b_n$  are the collection of Taylor points. The Taylor points  
 233 for different observations might be different, namely  $b_i \neq b_j$ . Also, the Taylor point  $b_i$   
 234 for observation  $i$  will be *updated* from time to time over the course of the algorithm.  
 235 In [Algorithm 2.1](#) we initially set all Taylor points to be the initial point  $x^0$ , namely



$b_i \leftarrow x^0$  for  $i = 1, 2, \dots, n$  in [Line 2](#). Notice that there is an Taylor point  $b_i$  for each observation index  $i \in [n]$ .

In order for the Taylor points  $b_1, \dots, b_n$  to yield sufficiently accurate gradient estimates  $g_1^k, \dots, g_n^k$  of  $\nabla g_1^k, \dots, \nabla g_n^k$  at iteration  $k$ , we will need to update some/all of the Taylor points to the current feasible point iterate  $x^k$  for some subset  $\mathcal{B}_k \subset [n]$  at iteration  $k$  of Algorithm TUFW. That is, we update  $b_i \leftarrow x^k$  for  $i \in \mathcal{B}_k$ . This update is done in [Line 4](#) of Algorithm TUFW.

The set  $\mathcal{B}_k$  is the set of indices whose Taylor points are updated to  $x^k$  at iteration  $k$ , and  $\mathcal{B}_k$  can be determined by any kind of **Rule**. As a trivial example, if one uses the rule  $\mathcal{B}_k = [n]$  for all  $k$ , then all Taylor points are updated to the current iterate  $x^k$  at each iteration  $k$ , in which case TUFW specializes to an instantiation of the traditional Frank-Wolfe method with exact gradient computation. As another example, one can use the rule that  $\mathcal{B}_k$  is a randomly chosen subset of the indices of cardinality, say, 50, or  $\lfloor n/10 \rfloor$ , or perhaps is iteration dependent such as cardinality  $\lfloor \ln(n)\sqrt{k} \rfloor$ . As a third example, one can use the rule that  $\mathcal{B}_k$  is chosen deterministically in a cyclic fashion with fixed batch-size  $B$ , so that  $\mathcal{B}_1 = \{1, \dots, B\}$ ,  $\mathcal{B}_2 = \{B+1, \dots, 2B\}$ , etc., with modular arithmetic so that  $\mathcal{B}_k \subset [n]$ . Let  $\beta_k := |\mathcal{B}_k|$  be the cardinality of  $\mathcal{B}_k$ . As we will see in succeeding sections, we can attain nearly-optimal complexity by developing batch-size rules for which  $\beta_k$  is a *decreasing* function of  $k$ , both in the convex and non-convex setting. This decreasing batch-size property is in stark contrast to other versions of Frank-Wolfe that use sampling of the observations to compute gradient estimates, wherein the batch-sizes grow with  $k$ , see [\[11, 29\]](#).

Going back to the description [Algorithm 2.1](#), the Taylor-approximated (at Taylor point  $b_i$ ) gradient of  $f_i$  is computed for each observation  $i \in [n]$  in [Line 5](#), and the estimate of the full gradient is computed in [Line 6](#). Finally, the Linear Minimization Oracle (LMO) is called in [Line 7](#) and the new iterate is computed in [Line 8](#) using the step-size  $\gamma_k$  along the cord  $[x^k, s^k]$ .

Inspecting [Algorithm 2.1](#), it appears that the computation of the gradient estimate  $g^k$  in [Line 5](#) and [Line 6](#) requires  $O(np^2)$  flops. However, we now show how to do the computation of  $g^k$  in these steps in only  $O((|\mathcal{B}_k| + 1)p^2)$  flops by taking advantage of the “linear prediction” structure of the loss functions  $f_i(x) := l_i(w_i^\top x)$  for each observation  $i \in [n]$ , and updating  $g^k$  based on the computed value of  $g^{k-1}$ .

**PROPOSITION 2.1.** *Let  $\bar{\beta}_k := |\mathcal{B}_k|$ . Then the updating of the second-order Taylor approximation model requires  $O(\bar{\beta}_k p^2)$  flops and the computation of  $g^k$  can be done via a matrix-vector product, which takes at most  $O(p^2)$  flops.*

*Proof of Proposition 2.1.* Consider iteration  $k - 1$  and let the Taylor points at this iteration be  $b_1, \dots, b_n$ . Let us write out the computation of  $g^{k-1}$  as:

$$\begin{aligned} g^{k-1} &= \frac{1}{n} \sum_{i=1}^n l'_i(w_i^\top b_i) w_i + \frac{1}{n} \sum_{i=1}^n l''_i(w_i^\top b_i) w_i w_i^\top (x^{k-1} - b_i) \\ &= \frac{1}{n} \sum_{i=1}^n (v_i - t_i w_i^\top b_i) w_i + \left[ \frac{1}{n} \sum_{i=1}^n t_i w_i w_i^\top \right] x^{k-1} \\ &= q_{k-1} + H_{k-1} x^{k-1}, \end{aligned}$$

where  $v_i = l'_i(w_i^\top b_i)$ ,  $t_i = l''_i(w_i^\top b_i)$  for  $i \in [n]$ , and  $q_{k-1} = \frac{1}{n} \sum_{i=1}^n (v_i - t_i w_i^\top b_i) w_i$  and  $H_{k-1} = \frac{1}{n} \sum_{i=1}^n t_i w_i w_i^\top$ . We assume that we can store the values  $v_i$ ,  $t_i$  for  $i \in [n]$ , as well as the vector  $q_{k-1}$  and Hessian matrix  $H_{k-1}$ .

At iteration  $k$ , let  $\mathcal{B}_k$  be the set of observation indices  $i$  where we replace the

Taylor point  $b_i$  with  $x^k$ . Using the overbar symbol “ $\bar{\cdot}$ ” to denote an update of a stored value, we have:

$$\begin{aligned} \bar{b}_i &= b_i & , & & \bar{v}_i &= v_i & , & & \bar{t}_i &= t_i & \text{ for } i \notin \mathcal{B}_k , & \text{ and} \\ \bar{b}_i &= x^k & , & & \bar{v}_i &= l'_i(w_i^\top \bar{b}_i) & , & & \bar{t}_i &= l''_i(w_i^\top \bar{b}_i) & \text{ for } i \in \mathcal{B}_k , \end{aligned}$$

and notice that these computations use  $O(\bar{\beta}_k p)$  flops. We then compute  $g^k$  by first updating  $q^{k-1}$  and  $H_{k-1}$  to  $q^k$  and  $H_k$  as follows:

$$\begin{aligned} (2.1) \quad q_k &= q_{k-1} + \frac{1}{n} \sum_{i \in \mathcal{B}_k} (\bar{v}_i - \bar{t}_i w_i^\top \bar{b}_i - v_i + t_i w_i^\top b_i) w_i \\ H_k &= H_{k-1} + \frac{1}{n} \sum_{i \in \mathcal{B}_k} (\bar{t}_i - t_i) w_i w_i^\top \\ g^k &= q_k + H_k x^k , \end{aligned}$$

where the first two steps of updating  $q_k$  and  $H_k$  in the Taylor approximation model require  $O(\bar{\beta}_k p^2)$  flops and the last step of computing  $g^k$  needs only a matrix-vector product, which is at most  $O(p^2)$  flops. After this, we replace the old values of  $b_i$ ,  $v_i$ , and  $t_i$  with their new values  $\bar{b}_i$ ,  $\bar{v}_i$ , and  $\bar{t}_i$  for  $i \in \mathcal{B}_k$ , which requires  $O(\bar{\beta}_k p)$  flops.  $\square$

Moreover, TUFW can be more computationally efficient if utilizing sparsity. Since the Frank-Wolfe often promotes sparse structured solutions [13], the iterate  $x^k$  is usually highly sparse. Besides, for problems such as matrix completion problems [5, 26], the Hessian matrices are also highly sparse [26]. Therefore, the matrix-vector products and the updating of  $H_k$  don't necessarily require  $O(p^2)$  flops if using sparse linear algebra. Furthermore, the  $\mathcal{B}_k$  can frequently become empty in above rules, and in this case the  $q_k$  and  $H_k$  remain unchanged, so instead of making matrix-vector product  $H_k x^k$  when computing  $g^k$ , we can update  $H^k$  by  $H_k x^k \leftarrow (1 - \gamma_{k-1}) H_{k-1} x^{k-1} + \gamma_k H_{k-1} s^{k-1}$ , where  $H_{k-1} x^{k-1}$  is already known when computing  $g^{k-1}$  and computing  $H_{k-1} s^{k-1}$  could utilize the sparsity or low-rank property of  $s^{k-1}$ . For example, if  $\mathcal{C}$  is an  $\ell_1$ -norm ball or a simplex, then  $s^{k-1}$  is a sparse vector in which only one component is nonzero and thus updating  $g^k$  takes at most  $O(p)$  flops when  $|\mathcal{B}_k| = 0$ .

**3. Convergence Guarantees for Convex Loss Functions.** In this section we study convergence guarantees and overall complexity of Algorithm 2.1 for tackling problem (1.2) in the case when the scalar loss functions  $l_i$  are convex for all  $i \in [n]$ , which then implies that  $f_i$  and  $F$  are also convex functions. Later, in section 6, we will present numerical experiments where we compare different versions of TUFW applied to problems with convex loss functions.

Recall that in order to run Algorithm 2.1 we need to specify the step-sizes  $\{\gamma_k\}_k$  and the **Rule** to construct the sets  $\mathcal{B}_k$  of indices for Taylor-point updating in Line 4. Our first rule is designed for  $\mathcal{B}_k$  to be comprised of  $\beta_k := n/\sqrt{k}$  independently drawn samples from  $[n]$  without replacement, for all  $k \geq 1$ . Since  $\beta_k$  is not in general an integer, then rather than ensuring  $|\mathcal{B}_k| = \beta_k$ , we will instead ensure that  $\mathbb{E}|\mathcal{B}_k| = \beta_k$  by proceeding as follows.

**DEFINITION 3.1.** *Rule-SBD $\sqrt{k}$ . Define  $\beta_k := n/\sqrt{k}$  and  $p_k := \beta_k - \lfloor \beta_k \rfloor$ . Sample  $\xi_k$  from the Bernoulli distribution  $\text{Ber}(p_k)$  and then uniformly sample  $\lfloor \beta_k \rfloor + \xi_k$  samples from  $[n]$  without replacement. Then define  $\mathcal{B}_k$  to be the set of these index samples.*

Note that by design of Rule-SBD $\sqrt{k}$  it follows that  $\mathbb{E}_{\xi_k} |\mathcal{B}_k| = \beta_k$ . We call this rule “Rule-SBD $\sqrt{k}$ ” for “Stochastic Batch-size Decreasing at the rate  $\sqrt{k}$ ,” and



we emphasize that the Taylor points are updated less often as  $k$  grows, which is in contrast to many other stochastic Frank-Wolfe methods that compute exact individual gradients more frequently as  $k$  grows, see [11, 29].

We first present a bound on the total number of flops used in the first  $k$  iterations of [Algorithm 2.1](#) using *Rule-SBD* $\sqrt{k}$ . We will presume that computing  $l'_i$  and  $l''_i$  require  $O(1)$  flops, and let  $\text{fLMO}$  denote the number of flops required by the LMO in [Line 7](#) of [Algorithm 2.1](#).

**PROPOSITION 3.2.** *Using *Rule-SBD* $\sqrt{k}$  and  $k \geq 1$ , the expected total number of flops used in the first  $k$  iterations of [Algorithm 2.1](#) is  $O\left(k \cdot (\text{fLMO} + p^2) + \sqrt{k} \cdot np^2\right)$ .*

*Proof of Proposition 3.2.* For the initial iteration of Algorithm TUFW the number of flops is  $O(\text{fLMO} + np^2)$ , and the expected total number of flops in the first  $k$  iterations is

$$\begin{aligned} O\left(k \cdot \text{fLMO} + np^2 + \sum_{i=1}^k (\beta_k + 1)p^2\right) &= O\left(k \cdot \text{fLMO} + kp^2 + np^2 + \sum_{i=1}^k \frac{np^2}{\sqrt{i}}\right) \\ &\leq O\left(k \cdot (\text{fLMO} + p^2) + np^2 + np^2 \cdot 2x^{\frac{1}{2}} \Big|_{x=1}^{k+1}\right) \\ &= O\left(k \cdot (\text{fLMO} + p^2) + np^2\sqrt{k}\right), \end{aligned}$$

where the left-hand side follows from [Proposition 2.1](#) and the inequality uses a standard integral bound.  $\square$

Because the feasibility set  $\mathcal{C}$  is bounded, the following measure of the range of the linear operator  $W^\top$

$$(3.1) \quad D_q := \max_{x, y \in \mathcal{C}} \|W^\top(x - y)\|_q$$

is finite for all  $q \in [1, \infty]$ . We note that (3.1) is a common metric used in the literature of problem (1.2), see [18, 22].

The following theorem presents the computational guarantee for [Algorithm 2.1](#) with *Rule-SBD* $\sqrt{k}$ . The proof of this theorem, as well as the other ensuing results in this section, are presented in [subsection 3.1](#).

**THEOREM 3.3.** *Suppose that  $F$  is convex and [Assumption 1.1](#) holds, and [Algorithm 2.1](#) with *Rule-SBD* $\sqrt{k}$  is applied to the problem (1.2) with step-sizes defined by  $\gamma_k := 2/(k+2)$  for all  $k \geq 0$ . Then for all  $k \geq 1$  we have:*

$$(3.2) \quad \mathbb{E}[F(x^k) - F(x^*)] \leq \frac{2LD_2^2 + 134\hat{L}D_1D_\infty^2}{n(k+1)}.$$

**COROLLARY 3.4.** *Under the hypotheses of [Theorem 3.3](#), it holds that  $\mathbb{E}[F(x^k) - F(x^*)] \leq \varepsilon$  after at most*

$$\frac{2LD_2^2 + 134\hat{L}D_1D_\infty^2}{n\varepsilon}$$

*iterations, and the total number of flops required is at most*

$$O\left((\text{fLMO} + p^2) \left\lceil \frac{LD_2^2 + \hat{L}D_1D_\infty^2}{n\varepsilon} \right\rceil + np^2 \left\lceil \frac{\sqrt{LD_2^2 + \hat{L}D_1D_\infty^2}}{\sqrt{n}\sqrt{\varepsilon}} \right\rceil\right).$$

Let  $\|\cdot\|$  be a given norm on  $\mathbb{R}^p$  and let  $\text{Diam}(\mathcal{C})$  denote the diameter of  $\mathcal{C}$ , namely  $\text{Diam}(\mathcal{C}) := \max_{x,y \in \mathcal{C}} \|x - y\|$ . Suppose that all feature vectors  $w_1, w_2, \dots$  lie in a bounded set  $S \subset \{w \in \mathbb{R}^p : \|w\|_* \leq M\}$  regardless of the number of observations  $n$ , where  $M$  denotes the radius of the (dual norm) ball centered at the origin in  $\mathbb{R}^p$  containing  $S$ . The following corollary follows from (3.2) using  $D_2 \leq \sqrt{n}D_\infty \leq \sqrt{n}M \cdot \text{Diam}(\mathcal{C})$  and  $D_1 \leq nD_\infty \leq nM \cdot \text{Diam}(\mathcal{C})$ .

COROLLARY 3.5. *Under the boundedness of the feature vectors, the bound on the number of flops in Corollary 3.4 to obtain  $\mathbb{E}[F(x^k) - F(x^*)] \leq \varepsilon$  is*

$$O\left((\text{fLMO} + p^2) \cdot \frac{LM^2 \cdot \text{Diam}(\mathcal{C})^2 + \hat{L}M^3 \cdot \text{Diam}(\mathcal{C})^3}{\varepsilon} + np^2 \cdot \frac{\sqrt{LM^2 \cdot \text{Diam}(\mathcal{C})^2 + \hat{L}M^3 \cdot \text{Diam}(\mathcal{C})^3}}{\sqrt{\varepsilon}}\right).$$

Whereas *Rule-SBD* $\sqrt{k}$  is stochastic, we now present the following deterministic rule which achieves nearly identical computational guarantees (to within a constant factor) but whose computational guarantees are deterministic.

DEFINITION 3.6. *Rule-DBD* $\sqrt{k}$ . For any  $k \geq 1$  define:

$$\mathcal{B}_k = \begin{cases} [n] & \text{if } \sqrt{k} \in \mathbb{N} \\ \emptyset & \text{if } \sqrt{k} \notin \mathbb{N}. \end{cases}$$

In *Rule-DBD* $\sqrt{k}$  we do not update any Taylor points unless  $k = 1, 4, 9, 25, \dots$ , and for these values of  $k$  we update all  $n$  Taylor points. We call this approach “*Rule-DBD* $\sqrt{k}$ ” for “Deterministic Batch-frequency Decreasing at the rate  $\sqrt{k}$ ,” and we point out that for *Rule-DBD* $\sqrt{k}$  the Taylor points are updated less often as  $k$  grows (in a different way but with similar effect as in *Rule-SBD* $\sqrt{k}$ ).

Similar to the case of *Rule-SBD* $\sqrt{k}$ , we have:

PROPOSITION 3.7. *Using Rule-DBD* $\sqrt{k}$  and  $k \geq 1$ , the total number of flops used in the first  $k$  iterations of Algorithm 2.1 is  $O(k \cdot (\text{fLMO} + p^2) + \sqrt{k} \cdot np^2)$ .

*Proof of Proposition 3.7.* The proof is almost identical to that of Proposition 3.2. For the initial iteration of Algorithm TUFW the number of flops is  $O(\text{fLMO} + np^2)$ , and the number of flops in the first  $k$  iterations is

$$O\left(k \cdot \text{fLMO} + np^2 + \sum_{i=1}^k (\beta_k + 1)p^2\right) \leq O\left(k \cdot (\text{fLMO} + p^2) + np^2\sqrt{k}\right),$$

where the left-hand side follows from Proposition 2.1 and the inequality follows since  $\sum_{i=1}^k \beta_k \leq n\sqrt{k}$ .  $\square$

THEOREM 3.8. *Suppose that  $F$  is convex and Assumption 1.1 holds, and Algorithm 2.1 with Rule-DBD* $\sqrt{k}$  *is applied to the problem (1.2) with step-sizes defined by  $\gamma_k := 2/(k+2)$  for all  $k \geq 0$ . Then for all  $k \geq 1$  we have:*

$$(3.3) \quad F(x^k) - F(x^*) \leq \frac{2LD_2^2 + 144\hat{L}D_3^3}{n(k+1)}.$$

COROLLARY 3.9. *Under the hypotheses of Theorem 3.8, it holds that  $F(x^k) - F(x^*) \leq \varepsilon$  after at most*

$$\frac{2LD_2^2 + 144\hat{L}D_3^3}{n\varepsilon}$$

*iterations, and the total number of flops required is at most*

$$O\left((\text{fLMO} + p^2) \left[\frac{LD_2^2 + \hat{L}D_3^3}{n\varepsilon}\right] + np^2 \left[\frac{\sqrt{LD_2^2 + \hat{L}D_3^3}}{\sqrt{n}\sqrt{\varepsilon}}\right]\right).$$

*Remark 3.10.* Comparing the bounds in [Theorem 3.3](#) and [Theorem 3.8](#), if we ignore the absolute constants then the bound in [Theorem 3.8](#) is better because  $D_3^3 \leq D_1 D_\infty^2$ .

Also similar to [Corollary 3.5](#), we have the following corollary which shows that the joint dependence on  $n$  and  $\varepsilon$  is  $O(n/\sqrt{\varepsilon})$ :

**COROLLARY 3.11.** *Under the boundedness of the feature vectors, the bound on the number of flops in [Corollary 3.9](#) to obtain  $F(x^k) - F(x^*) \leq \varepsilon$  is*

$$O\left((\text{fLMO} + p^2) \cdot \frac{LM^2 \cdot \text{Diam}(\mathcal{C})^2 + \hat{L}M^3 \cdot \text{Diam}(\mathcal{C})^3}{\varepsilon} + np^2 \cdot \frac{\sqrt{LM^2 \cdot \text{Diam}(\mathcal{C})^2 + \hat{L}M^3 \cdot \text{Diam}(\mathcal{C})^3}}{\sqrt{\varepsilon}}\right).$$

[Table 1](#) shows a comparison of computational guarantees of the standard Frank-Wolfe method, the Constant batch-size Stochastic Frank-Wolfe (CSFW) method developed in [\[22\]](#), and the TUFW method with *Rule-SBD* $\sqrt{k}$  and *Rule-DBD* $\sqrt{k}$ .

TABLE 1

*Complexity bounds for different Frank-Wolfe methods to obtain an  $\varepsilon$ -optimal solution of  $\text{ERM}_\ell$  with convex losses, under the boundedness assumption of the feature vectors. In the table  $\varepsilon_0 := F(x^0) - F(x^*)$ ,  $c_1 := LM^2 \text{Diam}(\mathcal{C})^2$ , and  $c_2 := \hat{L}M^3 \text{Diam}(\mathcal{C})^3$ .*

Method	Optimality Metric	Overall Complexity
<i>Rule-SBD</i> $\sqrt{k}$ ( <a href="#">Cor. 3.5</a> )	$\mathbb{E}[F(x^k) - F(x^*)] \leq \varepsilon$	$O\left((\text{fLMO} + p^2) \cdot \frac{c_1 + c_2}{\varepsilon} + np^2 \cdot \frac{\sqrt{c_1 + c_2}}{\sqrt{\varepsilon}}\right)$
<i>Rule-DBD</i> $\sqrt{k}$ ( <a href="#">Cor. 3.11</a> )	$F(x^k) - F(x^*) \leq \varepsilon$	$O\left((\text{fLMO} + p^2) \cdot \frac{c_1 + c_2}{\varepsilon} + np^2 \cdot \frac{\sqrt{c_1 + c_2}}{\sqrt{\varepsilon}}\right)$
Standard Frank-Wolfe ( <a href="#">[13]</a> )	$F(x^k) - F(x^*) \leq \varepsilon$	$O\left((\text{fLMO} + np) \cdot \frac{c_1}{\varepsilon}\right)$
CSFW ( <a href="#">[22]</a> )	$\mathbb{E}[F(x^k) - F(x^*)] \leq \varepsilon$	$O\left((n \cdot \text{fLMO} + np) \left[\frac{c_1}{\varepsilon} + \frac{\sqrt{\varepsilon_0}}{n\sqrt{\varepsilon}} + \frac{\sqrt{nc_1}}{\sqrt{\varepsilon}}\right]\right)$

Note in particular that with *Rule-SBD* $\sqrt{k}$  or *Rule-DBD* $\sqrt{k}$ , the TUFW's joint dependence on  $n$  and  $\varepsilon$  is  $O(n/\sqrt{\varepsilon})$ , as compared to  $O(n/\varepsilon)$  in [\[13, 22, 18\]](#),  $O(1/\varepsilon^2)$  in [\[11, 29, 30, 9\]](#), or  $O(1/\varepsilon^3)$  in [\[21\]](#). Indeed, in the regime where  $\varepsilon$  is sufficiently small, the dominant term in the TUFW's complexity bound is the left term (which counts operations of the LMO), which is independent of  $n$ . In this regime the overall complexity is nearly optimal: the LMO complexity's dependence on  $\varepsilon$  is  $O(1/\varepsilon)$ , which is essentially the same as in the lower bound complexity result in [\[16\]](#). The LMO complexity of TUFW differs from the lower bound only by the term involving  $\hat{L}$ , which is the Lipschitz constant of the second-order derivatives.

We end this subsection by considering the following rule which does no updates of the Taylor points beyond their initialization in [Line 2](#) of TUFW.

**DEFINITION 3.12.** *Rule- $\emptyset$ .* Define  $\mathcal{B}_0 = [n]$  and  $\mathcal{B}_k := \emptyset$  for  $k \geq 1$ .

*Rule- $\emptyset$*  is quite relevant in the special case when the loss functions  $l_i$  are all quadratic loss functions (as in linear regression, matrix completion, LASSO, etc.) in which case the Hessian  $\nabla^2 F(x) = \nabla^2 F(x^0)$  for any  $x$  and so does not change. Therefore in this case  $g_i^k = \nabla f_i(x^k)$  and  $g^k = \nabla F(x^k)$  for all  $k$ . When using *Rule- $\emptyset$* , the computational cost of computing the initial individual gradients and the Hessian  $\nabla^2 F(x^0)$  is only  $O(np^2)$  flops, and the cost of calculating the (exact) gradient is the cost of a matrix-vector product, which is at most  $O(p^2)$  flops per iteration.

**3.1. Proofs of results in section 3.** We will prove [Theorem 3.3](#) as a consequence of the following sequences of lemmas, corollaries, and propositions.

LEMMA 3.13. Under [Assumption 1.1](#) for problem  $ERM_\ell$ ,

$$(3.4) \quad F(x^{k+1}) \leq F(x^k) + \gamma_k \langle \nabla F(x^k), s^k - x^k \rangle + \gamma_k^2 L D_2^2 / 2n$$

holds for iteration  $k$  of [Algorithm 2.1](#) for all  $k \geq 0$ .

*Proof of Lemma 3.13.* For iteration  $k$  it follows from (1.3) that

$$l_i(w_i^\top x^{k+1}) \leq l_i(w_i^\top x^k) + l'_i(w_i^\top x^k)(w_i^\top x^{k+1} - w_i^\top x^k) + L(w_i^\top x^{k+1} - w_i^\top x^k)^2 / 2.$$

Then since  $F(x^{k+1}) = \frac{1}{n} \sum_{i=1}^n l_i(w_i^\top x^{k+1})$ , we have:

$$\begin{aligned} F(x^{k+1}) &\leq F(x^k) + \gamma_k \langle \nabla F(x^k), s^k - x^k \rangle + \gamma_k^2 L \left\| W^\top s^k - W^\top x^k \right\|_2^2 / 2n \\ &\leq F(x^k) + \gamma_k \langle \nabla F(x^k), s^k - x^k \rangle + \gamma_k^2 L D_2^2 / 2n, \end{aligned}$$

where the last inequality above uses (3.1).  $\square$

For any  $k \geq 0$  and  $i \in [n]$ , we will be interested in the most recent iteration (up through  $k$ ) at which Taylor point  $b_i$  was updated. This is defined as:

$$(3.5) \quad \tau_i^k := \max\{t : t \leq k \text{ and } i \in \mathcal{B}_t\},$$

which means that in the first  $k$  iterations,  $\tau_i^k$  is the latest iteration in which the Taylor point  $b_i$  was updated. In [Algorithm 2.1](#), the relationship between  $\tau_i^k$  and  $\tau_i^{k-1}$  is then:

$$(3.6) \quad \tau_i^k = \begin{cases} k & \text{if } i \in \mathcal{B}_k \\ \tau_i^{k-1} & \text{if } i \notin \mathcal{B}_k. \end{cases}$$

LEMMA 3.14. Under [Assumption 1.1](#) for the problem  $ERM_\ell$ , for any  $u, v \in \mathcal{C}$ ,

$$(3.7) \quad (\nabla F(x^k) - g^k)^\top (u - v) \leq \frac{\hat{L}}{2n} \sum_{i=1}^n \left( \sum_{j=\tau_i^k}^{k-1} \gamma_j \left| w_i^\top (s^j - x^j) \right| \right)^2 \cdot |w_i^\top (u - v)|$$

holds for iteration  $k$  of [Algorithm 2.1](#) for all  $k \geq 1$ .

*Proof of Lemma 3.14.* It follows from the definition of  $\tau_i^k$  in (3.5) that

$$\begin{aligned} &(\nabla F(x^k) - g^k)^\top (u - v) \\ &= \frac{1}{n} \sum_{i=1}^n \left( l'_i(w_i^\top x^k) - l'_i(w_i^\top x^{\tau_i^k}) - l''_i(w_i^\top x^{\tau_i^k})(w_i^\top x^k - w_i^\top x^{\tau_i^k}) \right) \cdot w_i^\top (u - v) \\ &\leq \frac{\hat{L}}{2n} \sum_{i=1}^n |w_i^\top x^k - w_i^\top x^{\tau_i^k}|^2 \cdot |w_i^\top (u - v)| \\ &= \frac{\hat{L}}{2n} \sum_{i=1}^n \left( \sum_{j=\tau_i^k}^{k-1} w_i^\top (x^{j+1} - x^j) \right)^2 \cdot |w_i^\top (u - v)| \\ &\leq \frac{\hat{L}}{2n} \sum_{i=1}^n \left( \sum_{j=\tau_i^k}^{k-1} \left| w_i^\top (x^{j+1} - x^j) \right| \right)^2 \cdot |w_i^\top (u - v)| \\ &= \frac{\hat{L}}{2n} \sum_{i=1}^n \left( \sum_{j=\tau_i^k}^{k-1} \left| \gamma_j w_i^\top (s^j - x^j) \right| \right)^2 \cdot |w_i^\top (u - v)| \end{aligned}$$

where the first inequality uses (1.4).  $\square$

402 Because  $|w_i^\top(u - v)| \leq D_\infty$ , [Lemma 3.14](#) implies the following corollary.

COROLLARY 3.15.

$$403 \quad (3.9) \quad \left( \nabla F(x^k) - g^k \right)^\top (u - v) \leq \frac{\hat{L}D_\infty}{2n} \sum_{i=1}^n \left( \sum_{j=\tau_i^k}^{k-1} \gamma_j \left| w_i^\top(s^j - x^j) \right| \right)^2.$$

404 LEMMA 3.16. *For any  $k \in \mathbb{N}$  the following inequality holds:*

$$405 \quad (3.10) \quad \sqrt{\lfloor k/2 \rfloor + 1} \exp(2\sqrt{\lfloor k/2 \rfloor + 1} - 2\sqrt{k+1}) \leq \min\{3(k+2)^{-\frac{1}{2}}, 10(k+2)^{-1}\}.$$

*Proof.* Since  $h(x) := \sqrt{x}$  is a concave function on  $\mathbb{R}_+$ , we have  $h(k+2) \leq h(k+1) + h'(k+1)$ , which yields  $\sqrt{k+1} \geq \sqrt{k+2} - \frac{1}{2\sqrt{k+1}} \geq \sqrt{k+2} - \sqrt{2}/4$  when  $k \geq 1$ . Furthermore  $\sqrt{\lfloor k/2 \rfloor + 1} \leq \frac{\sqrt{2}}{2}\sqrt{k+2}$ , and therefore

$$\sqrt{\lfloor k/2 \rfloor + 1} \exp(2\sqrt{\lfloor k/2 \rfloor + 1} - 2\sqrt{k+1}) \leq \frac{\sqrt{2}e^{\frac{\sqrt{2}}{2}}}{2} (k+2)^{\frac{1}{2}} \exp(-(2 - \sqrt{2})(k+2)^{\frac{1}{2}}).$$

406 Next note that the function  $h(x) := e^{-x}x^t$  on the domain  $x \in (0, \infty)$  for  $t > 0$   
 407 attains its maximal value at  $x = t$ . In particular for  $t = 2$  and  $t = 3$ , respectively, we  
 408 conclude that  $e^{-x} \leq 4e^{-2}x^{-2}$  and  $e^{-x} \leq 27e^{-3}x^{-3}$  (respectively), and thus  $\exp(-(2 - \sqrt{2})(k+2)^{\frac{1}{2}}) \leq 4e^{-2}(2 - \sqrt{2})^{-2}(k+2)^{-1}$  and  $\exp(-(2 - \sqrt{2})(k+2)^{\frac{1}{2}}) \leq 27e^{-3}(2 - \sqrt{2})^{-3}(k+2)^{-3/2}$  (respectively). Therefore, the left-hand side of (3.10) is bounded  
 409 by  $2\sqrt{2}e^{\frac{\sqrt{2}}{2}-2}(2 - \sqrt{2})^{-2}(k+2)^{-\frac{1}{2}}$  and  $\frac{27}{2}\sqrt{2}e^{\frac{\sqrt{2}}{2}-3}(2 - \sqrt{2})^{-3}(k+2)^{-1}$ , respectively,  
 410 which is less than  $3(k+2)^{-\frac{1}{2}}$  and  $10(k+2)^{-1}$ , respectively.  $\square$

413 Let  $F^*$  denote the optimal objective function value of  $\text{ERM}_\ell$  and let  $\varepsilon_k := F(x^k) - F^*$  denote the optimality gap at iteration  $k$ .

415 LEMMA 3.17. *Under [Assumption 1.1](#) for problem  $\text{ERM}_\ell$ , suppose that  $F$  is convex, and let the step-sizes be given by  $\gamma_k = 2/(k+2)$  for all  $k \geq 0$ . Then for all  $k \geq 1$  we have*

$$418 \quad (3.11) \quad \varepsilon_k \leq \frac{2LD_2^2}{n(k+1)} + \frac{2}{(k)(k+1)} \sum_{t=1}^k t (\nabla F(x^{t-1}) - g^{t-1})^\top (s^{t-1} - x^*).$$

419 *Proof of Lemma 3.17.* Utilizing [Lemma 3.13](#) we have

$$\begin{aligned} (3.12) \quad F(x^{k+1}) &\leq F(x^k) + \gamma_k \langle \nabla F(x^k), s^k - x^k \rangle + \gamma_k^2 LD_2^2 / 2n \\ &= F(x^k) + \gamma_k \langle \nabla F(x^k) - g^k, s^k - x^k \rangle + \gamma_k \langle g^k, s^k - x^k \rangle + \gamma_k^2 LD_2^2 / 2n \\ 420 \quad &\leq F(x^k) + \gamma_k \langle \nabla F(x^k) - g^k, s^k - x^k \rangle + \gamma_k \langle g^k, x^* - x^k \rangle + \gamma_k^2 LD_2^2 / 2n \\ &= F(x^k) + \gamma_k \langle \nabla F(x^k) - g^k, s^k - x^* \rangle + \gamma_k \langle \nabla F(x^k), x^* - x^k \rangle + \gamma_k^2 LD_2^2 / 2n \\ &\leq F(x^k) + \gamma_k \langle \nabla F(x^k) - g^k, s^k - x^* \rangle + \gamma_k (F(x^*) - F(x^k)) + \gamma_k^2 LD_2^2 / 2n, \end{aligned}$$

421 where the second inequality uses  $s^k \in \arg \min_{s \in \mathcal{C}} (g^k)^\top s$  from [Line 7](#) of the algorithm,  
 422 and the third inequality uses the gradient inequality applied to  $F$ . Subtracting  $F^*$   
 423 from both sides of the above inequality chain, we arrive at:

$$424 \quad (3.13) \quad \varepsilon_{k+1} \leq (1 - \gamma_k) \varepsilon_k + \gamma_k \left( \nabla F(x^k) - g^k \right)^\top (s^k - x^*) + \gamma_k^2 LD_2^2 / 2n.$$

425 Multiplying both side by  $(k+1)(k+2)$  and telescoping the inequalities yields:

$$426 \quad (3.14) \quad (k+1)(k+2) \varepsilon_{k+1} \leq \frac{2(k+1)LD_2^2}{n} + \sum_{t=0}^k 2(t+1) (\nabla F(x^t) - g^t)^\top (s^t - x^*),$$

427 which yields (3.11) after rearranging and re-indexing the counter  $t$ .  $\square$

LEMMA 3.18. Let the series  $\{d_{ij}\}_{i \in [n], j \in \mathbb{N}}$  be given, and suppose that there exists  $M_1$  and  $M_\infty$  for which  $\sup_j \sum_{i=1}^n d_{ij} \leq M_1$  and  $\sup_{i,j} d_{ij} \leq M_\infty$ . Then for *Algorithm 2.1* using Rule-SBD  $\sqrt{k}$  and step-sizes  $\gamma_k := 2/(k+2)$  for all  $k \geq 0$ , it holds that:

$$(3.15) \quad \mathbb{E} \left[ \sum_{i=1}^n \left( \sum_{j=\tau_i^k}^{k-1} \gamma_j d_{ij} \right) \right] \leq \frac{6M_1}{\sqrt{k+2}} \quad \text{and} \quad \mathbb{E} \left[ \sum_{i=1}^n \left( \sum_{j=\tau_i^k}^{k-1} \gamma_j d_{ij} \right)^2 \right] \leq \frac{134M_1M_\infty}{k+2}.$$

Before proving the above lemma, we first introduce the following proposition.

PROPOSITION 3.19. For integers  $v$  and  $k$  satisfying  $1 \leq v \leq k$  it holds that:

$$(3.16) \quad \sum_{i=1}^v \prod_{j=i}^k \left( 1 - \frac{1}{\sqrt{j}} \right) < (\sqrt{v+1}) \exp(2\sqrt{v+1} - 2\sqrt{k+1}).$$

*Proof of Proposition 3.19.* For convenience let us define  $\rho_j := 1 - \frac{1}{\sqrt{j}}$  for  $j \geq 1$ , and define the function  $h(x) := \exp(2\sqrt{x} - 2\sqrt{k+1})$  for  $x \geq 0$ . We will first prove the following inequality:

$$(3.17) \quad \prod_{j=i}^k \left( 1 - \frac{1}{\sqrt{j}} \right) \leq h(i) \quad \text{for all } i \in [k].$$

For  $i = 1$  the left side above is 0 (because  $\rho_1 = 0$  and hence  $\prod_{j=1}^k \rho_j = 0$ ), and the right side is  $h(1) > 0$ , so (3.17) is true for  $i = 1$ . For  $i \geq 2$  we have:

$$\begin{aligned} \prod_{j=i}^k \rho_j &= \exp \left( \sum_{j=i}^k \ln \left( 1 - \frac{1}{\sqrt{j}} \right) \right) \leq \exp \left( - \sum_{j=i}^k \frac{1}{\sqrt{j}} \right) \\ &\leq \exp \left( - \int_i^{k+1} x^{-\frac{1}{2}} dx \right) = \exp(2\sqrt{i} - 2\sqrt{k+1}) = h(i), \end{aligned}$$

where the first inequality follows from the fact that  $\log(1-x) \leq -x$  for any  $0 \leq x < 1$ , and the second inequality is an integral bound. This establishes (3.17), whereby we have

$$\begin{aligned} \sum_{i=1}^v \prod_{j=i}^k \rho_j &\leq \sum_{i=1}^v h(i) \leq \int_1^{v+1} h(x) dx \\ &= \left( \sqrt{x} - \frac{1}{2} \right) \exp(2\sqrt{x} - 2\sqrt{k+1}) \Big|_1^{v+1} \\ &= \left( \sqrt{v+1} - \frac{1}{2} \right) \exp(2\sqrt{v+1} - 2\sqrt{k+1}) - \left( \frac{1}{2} \right) \exp(2 - 2\sqrt{k+1}) \\ &< \left( \sqrt{v+1} - \frac{1}{2} \right) \exp(2\sqrt{v+1} - 2\sqrt{k+1}), \end{aligned}$$

where the second inequality is an integral bound using the fact that  $h(x)$  is an increasing function.  $\square$

*Proof of Lemma 3.18.* First of all, let us define  $G_k^i := \left( \sum_{j=\tau_i^k}^{k-1} \gamma_j d_{ij} \right)$  and  $Z_k^i := \left( \sum_{j=\tau_i^k}^{k-1} \gamma_j d_{ij} \right)^2$  for  $i \in [n]$ , and also  $G_k := \sum_{i=1}^n G_k^i$  and  $Z_k := \sum_{i=1}^n Z_k^i$ . Notice in  $G_k^i$  and  $Z_k^i$  that the randomness lies only in  $\tau_i^k$ , which in turn depends on the update sets  $\mathcal{B}_0, \mathcal{B}_1, \dots, \mathcal{B}_k$ .

We first analyze  $G_k^i$  and  $G_k$ . We note that  $G_k^i = 0$  when  $i \in \mathcal{B}_k$ , which happens with probability  $\beta_k/n$ . Because of the relationship between  $\{\tau_i^k\}$  and  $\{\tau_i^{k-1}\}$  in (3.6),



we further note that with probability  $1 - \beta_k/n$  that  $i \notin \mathcal{B}_k$  in which case it follows that  $G_k^i = G_{k-1}^i + \gamma_{k-1}d_{i,k-1}$ . Defining  $\rho_k := 1 - \frac{1}{\sqrt{k}} = 1 - \beta_k/n$ , we have

$$\mathbb{E}[G_k] = \sum_{i=1}^n \mathbb{E}[G_k^i] = \rho_k \sum_{i=1}^n \mathbb{E}[G_{k-1}^i + \gamma_{k-1}d_{i,k-1}] \leq \rho_k (\mathbb{E}[G_{k-1}] + \gamma_{k-1}M_1) ,$$

where the inequality uses  $\sup_j \sum_{i=1}^n d_{ij} \leq M_1$ . Recurrently applying (3.18) yields:

$$\mathbb{E}[G_k] \leq \rho_k (\mathbb{E}[G_{k-1}] + \gamma_{k-1}M_1) \leq \dots \leq \sum_{i=1}^k \left( \prod_{j=i}^k \rho_j \right) \gamma_{i-1}M_1 + \left( \prod_{j=1}^k \rho_j \right) \mathbb{E}[G_0] .$$

Noting that  $G_0 = 0$  (since all Taylor points are set to  $x^0$  in Line 2 of Algorithm 2.1) and using  $\gamma_{i-1} = \frac{2}{i+1}$  it follows that:

$$\begin{aligned} \mathbb{E}[G_k] &\leq \sum_{i=1}^{\lfloor k/2 \rfloor} \left( \prod_{j=i}^k \rho_j \right) \frac{2M_1}{i+1} + \sum_{i=\lfloor k/2 \rfloor + 1}^k \left( \prod_{j=i}^k \rho_j \right) \frac{2M_1}{i+1} \\ &\leq \sum_{i=2}^{\lfloor k/2 \rfloor} \left( \prod_{j=i}^k \rho_j \right) \frac{2M_1}{3} + \sum_{i=\lfloor k/2 \rfloor + 1}^k \left( \prod_{j=i}^k \rho_j \right) \frac{2M_1}{\lfloor k/2 \rfloor + 2} , \end{aligned}$$

where the second inequality uses the fact that  $\rho_1 = 0$  and hence  $\prod_{j=1}^k \rho_j = 0$ . By Proposition 3.19 we have:

$$\sum_{i=1}^{\lfloor k/2 \rfloor} \left( \prod_{j=i}^k \rho_j \right) \frac{2M_1}{3} \leq \sqrt{\lfloor k/2 \rfloor + 1} \exp \left( 2\sqrt{\lfloor k/2 \rfloor + 1} - 2\sqrt{k+1} \right) \frac{2M_1}{3} \leq \frac{2M_1}{\sqrt{k+2}} ,$$

where the second inequality above uses Lemma 3.16. Similarly for the second term in the rightmost side of (3.20) we have:

$$\sum_{i=\lfloor k/2 \rfloor + 1}^k \left( \prod_{j=i}^k \rho_j \right) \frac{2M_1}{\lfloor k/2 \rfloor + 2} \leq \sum_{i=1}^k \left( \prod_{j=i}^k \rho_j \right) \frac{2M_1}{\lfloor k/2 \rfloor + 2} \leq \frac{\sqrt{k+1} \cdot 2M_1}{\lfloor k/2 \rfloor + 2} \leq \frac{4M_1}{\sqrt{k+2}} ,$$

where the second inequality uses Proposition 3.19. Substituting (3.21) and (3.22) back into (3.20) yields the first part of (3.15), namely:

$$\mathbb{E}[G_k] \leq 6M_1/\sqrt{k+2} .$$

We employ a similar proof approach to establish the second part of (3.15). We note that  $Z_k^i = 0$  when  $i \in \mathcal{B}_k$ , which happens with probability  $\beta_k/n$ . Because of the relationship between  $\{\tau_i^k\}$  and  $\{\tau_i^{k-1}\}$  in (3.6), we further note that with probability  $1 - \beta_k/n$  that  $i \notin \mathcal{B}_k$  in which case we have

$$\begin{aligned} Z_k^i &= (G_{k-1}^i + \gamma_{k-1}d_{i,k-1})^2 = Z_{k-1}^i + 2G_{k-1}^i\gamma_{k-1}d_{i,k-1} + (\gamma_{k-1}d_{i,k-1})^2 \\ &\leq Z_{k-1}^i + 2M_\infty\gamma_{k-1}G_{k-1}^i + \gamma_{k-1}^2M_\infty d_{i,k-1} . \end{aligned}$$

Therefore

$$\begin{aligned} \mathbb{E}[Z_k] &= \sum_{i=1}^n \mathbb{E}[Z_k^i] \leq \rho_k \sum_{i=1}^n \mathbb{E}[Z_{k-1}^i + 2M_\infty\gamma_{k-1}G_{k-1}^i + \gamma_{k-1}^2M_\infty d_{i,k-1}] \\ &\leq \rho_k (\mathbb{E}[Z_{k-1}] + 2M_\infty\gamma_{k-1}\mathbb{E}[G_{k-1}] + \gamma_{k-1}^2M_\infty M_1) \\ &\leq \rho_k \left( \mathbb{E}[Z_{k-1}] + \frac{24M_\infty M_1}{(k+1)\sqrt{k+1}} + \frac{4M_\infty M_1}{(k+1)^2} \right) \\ &\leq \rho_k \left( \mathbb{E}[Z_{k-1}] + \frac{28M_\infty M_1}{(k+1)\sqrt{k+1}} \right) , \end{aligned}$$

where the second inequality uses  $\sup_j \sum_{i=1}^n d_{ij} \leq M_1$ , and the third inequality uses (3.23). Recurrently applying (3.24) yields:

$$(3.25) \quad \mathbb{E}[Z_k] \leq \rho_k \left( \mathbb{E}[Z_{k-1}] + \frac{28M_1M_\infty}{(k+1)^{3/2}} \right) \leq \cdots \leq \sum_{i=1}^k \left( \prod_{j=i}^k \rho_j \right) \frac{28M_1M_\infty}{(i+1)^{3/2}}.$$

Therefore

$$(3.26) \quad \begin{aligned} \mathbb{E}[Z_k] &\leq \sum_{i=1}^{\lfloor k/2 \rfloor} \left( \prod_{j=i}^k \rho_j \right) \frac{28M_1M_\infty}{(i+1)^{3/2}} + \sum_{i=\lfloor k/2 \rfloor + 1}^k \left( \prod_{j=i}^k \rho_j \right) \frac{28M_1M_\infty}{(i+1)^{3/2}} \\ &\leq \sum_{i=1}^{\lfloor k/2 \rfloor} \left( \prod_{j=i}^k \rho_j \right) \frac{28M_1M_\infty}{3\sqrt{3}} + \sum_{i=\lfloor k/2 \rfloor + 1}^k \left( \prod_{j=i}^k \rho_j \right) \frac{28M_1M_\infty}{(\lfloor k/2 \rfloor + 2)^{3/2}}. \end{aligned}$$

By Proposition 3.19 we have

$$(3.27) \quad \begin{aligned} \sum_{i=1}^{\lfloor k/2 \rfloor} \left( \prod_{j=i}^k \rho_j \right) \frac{28M_1M_\infty}{3\sqrt{3}} &\leq \sqrt{\lfloor k/2 \rfloor + 1} \exp(2\sqrt{\lfloor k/2 \rfloor + 1} - 2\sqrt{k+1}) \cdot \frac{28M_1M_\infty}{(3\sqrt{3})} \\ &\leq \frac{54M_1M_\infty}{k+2}, \end{aligned}$$

where the second inequality above uses Lemma 3.16. Similarly for the second term in the rightmost side of (3.26) we have:

$$(3.28) \quad \begin{aligned} \sum_{i=\lfloor k/2 \rfloor + 1}^k \left( \prod_{j=i}^k \rho_j \right) \frac{28M_1M_\infty}{(\lfloor k/2 \rfloor + 2)^{3/2}} &\leq \sum_{i=1}^k \left( \prod_{j=i}^k \rho_j \right) \frac{28M_1M_\infty}{(\lfloor k/2 \rfloor + 2)^{3/2}} \\ &\leq \sqrt{k+1} \cdot \frac{28M_1M_\infty}{(\lfloor k/2 \rfloor + 2)^{3/2}} \leq \frac{80M_1M_\infty}{k+2}, \end{aligned}$$

where the second inequality uses Proposition 3.19. Substituting (3.27) and (3.28) back into (3.26) yields the second part of (3.15), namely:

$$(3.29) \quad \mathbb{E}[Z_k] \leq \frac{134M_1M_\infty}{k+2}. \quad \square$$

*Proof of Theorem 3.3.* From Lemma 3.17 it holds that:

$$(3.30) \quad \varepsilon_k \leq \frac{2LD_2^2}{n(k+1)} + \frac{2}{k(k+1)} \sum_{t=1}^k t (\nabla F(x^{t-1}) - g^{t-1})^\top (s^{t-1} - x^*),$$

and from Corollary 3.15 it follows that:

$$\varepsilon_k \leq \frac{2LD_2^2}{n(k+1)} + \frac{\hat{L}D_\infty}{nk(k+1)} \sum_{t=1}^k t \sum_{i=1}^n \left( \sum_{j=\tau_i^{t-1}}^{t-2} \gamma_j |w_i^\top (s^j - x^j)| \right)^2.$$

Taking expectations on both sides then yields

$$(3.31) \quad \mathbb{E}[\varepsilon_k] \leq \frac{2LD_2^2}{n(k+1)} + \frac{\hat{L}D_\infty}{nk(k+1)} \sum_{t=1}^k t \mathbb{E} \left[ \sum_{i=1}^n \left( \sum_{j=\tau_i^{t-1}}^{t-2} \gamma_j |w_i^\top (s^j - x^j)| \right)^2 \right].$$

Now define  $d_{ij} := |w_i^\top (s^j - x^j)|$  and  $M_1 := D_1$  and  $M_\infty := D_\infty$ . Then since  $\sup_j \sum_{i=1}^n d_{ij} = \sup_j \sum_{i=1}^n |w_i^\top (s^j - x^j)| \leq D_1 = M_1$  and  $\sup_{i,j} d_{ij}$  is equal to

486  $\sup_{i,j} |w_i^\top (s^j - x^j)|$  which is bounded by  $D_\infty = M_\infty$ , we can apply [Lemma 3.18](#)  
 487 to [\(3.31\)](#) and obtain

$$488 \quad (3.32) \quad \mathbb{E}[\varepsilon_k] \leq \frac{2LD_2^2}{n(k+1)} + \frac{\hat{L}D_\infty}{nk(k+1)} \sum_{t=1}^k t \frac{134D_1D_\infty}{t+1} \leq \frac{2LD_2^2}{n(k+1)} + \frac{134\hat{L}D_1D_\infty^2}{n(k+1)},$$

489 which demonstrates [\(3.2\)](#).  $\square$

490 We next proceed towards the proof of [Theorem 3.8](#). We first prove the following  
 491 lemma.

492 **LEMMA 3.20.** *In any iteration  $k$  of TUFW outlined in [Algorithm 2.1](#) with  $\gamma_k :=$   
 493  $2/(k+2)$  and Rule-DBD $\sqrt{k}$ , for problem setup [\(1.2\)](#), it holds that*

$$494 \quad (3.33) \quad \sum_{t=1}^k t \sum_{i=1}^n \left( \sum_{j=\tau_i^{t-1}}^{t-2} \gamma_j |w_i^\top (z^j - y^j)| \right)^2 |w_i^\top (z^{t-1} - x^*)| \leq 144kD_3^3$$

495 for any  $y^j, z^j \in \mathcal{C}$ ,  $j = 1, \dots, k$ .

496 *Proof of Lemma 3.20.* Notice that in Rule-DBD $\sqrt{k}$ ,  $\tau_i^k = (\lfloor \sqrt{k} \rfloor)^2$ , so  $k - \tau_i^k =$   
 497  $k - (\lfloor \sqrt{k} \rfloor)^2 \leq \lfloor 2\sqrt{k} \rfloor$ , where the latter follows from squaring the inequality  $\lfloor \sqrt{k} \rfloor \geq$   
 498  $\sqrt{k} - 1$  and rearranging. Therefore,

$$\begin{aligned} & \sum_{t=1}^k t \sum_{i=1}^n \left( \sum_{j=\tau_i^{t-1}}^{t-2} \gamma_j |w_i^\top (z^j - y^j)| \right)^2 |w_i^\top (z^{t-1} - x^*)| \\ 499 \quad (3.34) \quad & \leq \sum_{t=1}^k t \sum_{i=1}^n \left( \sum_{j=t-1-\lfloor 2\sqrt{t-1} \rfloor}^{t-2} \frac{2 |w_i^\top (z^j - y^j)|}{t+1-\lfloor 2\sqrt{t-1} \rfloor} \right)^2 |w_i^\top (z^{t-1} - x^*)| \\ & \leq \sum_{t=1}^k \frac{4t \lfloor 2\sqrt{t-1} \rfloor}{(t+1-\lfloor 2\sqrt{t-1} \rfloor)^2} \sum_{j=t-1-\lfloor 2\sqrt{t-1} \rfloor}^{t-2} \sum_{i=1}^n |w_i^\top (z^j - y^j)|^2 |w_i^\top (z^{t-1} - x^*)| \end{aligned}$$

500 where the first inequality is due to  $\tau_i^{t-1} \geq t-1-\lfloor 2\sqrt{t-1} \rfloor$  and the second inequality  
 501 is due to the fact that for any vector  $a \in \mathbb{R}^d$  we have  $\|a\|_1^2 \leq d\|a\|_2^2$ . In the right-hand  
 502 side of [\(3.34\)](#), we have

$$503 \quad (3.35) \quad \frac{4t \lfloor 2\sqrt{t-1} \rfloor}{(t+1-\lfloor 2\sqrt{t-1} \rfloor)^2} \leq \frac{8t^{3/2}}{t^2/9} = 72t^{-1/2},$$

504 where the inequality involving the denominators above follows from verifying  $t+1-\lfloor 2\sqrt{t-1} \rfloor \geq t/3$  via the quadratic formula. Additionally notice that for any vectors  
 505  $a, b \in \mathbb{R}^d$  and letting  $A$  denote the diagonal matrix whose diagonal elements are  
 506  $a_1, a_2, \dots, a_d$ , we have  $\sum_{i=1}^d a_i^2 b_i = e^\top A^2 b \leq \|A^2 e\|_u \|b\|_v = \|a\|_{2u}^2 \|b\|_v$  for any  $u, v$   
 507 such that  $u, v \geq 1$  and  $\frac{1}{u} + \frac{1}{v} = 1$ . Therefore  $\sum_{i=1}^n |w_i^\top (z^j - y^j)|^2 |w_i^\top (z^{t-1} - x^*)| \leq$   
 508  $D_{2u}^2 D_v$  and it then follows that

$$510 \quad (3.36) \quad \sum_{j=t-1-\lfloor 2\sqrt{t-1} \rfloor}^{t-2} \sum_{i=1}^n |w_i^\top (z^j - y^j)|^2 |w_i^\top (z^{t-1} - x^*)| \leq \lfloor 2\sqrt{t-1} \rfloor D_{2u}^2 D_v \leq 2\sqrt{t} D_{2u}^2 D_v.$$

511 The inequality [\(3.33\)](#) then follows by setting  $u = 3/2$ ,  $v = 3$ , and substituting [\(3.35\)](#)  
 512 and [\(3.36\)](#) into [\(3.34\)](#).  $\square$

*Proof of Theorem 3.8.* Similar to the proof of Theorem 3.3, according to Lemmas 3.14 and 3.17,

$$\varepsilon_k \leq \frac{2LD_2^2}{n(k+1)} + \frac{\hat{L}}{nk(k+1)} \sum_{t=1}^k t \sum_{i=1}^n \left( \sum_{j=\tau_i^{t-1}}^{t-2} \gamma_j |w_i^\top (s^j - x^j)| \right)^2 |w_i^\top (s^{t-1} - x^*)|.$$

Next apply Lemma 3.20 to the above inequality, which yields (3.3).  $\square$

**4. Convergence Guarantees for Non-convex Loss Functions.** In this section we study convergence rates and overall complexity of Algorithm 2.1 for tackling problem (1.2) in the case that the scalar loss functions  $l_i$  are not necessarily convex. Later, in section 6, we will present experimental results comparing different versions of TUFW applied to problems with non-convex loss functions.

Recall that the Frank-Wolfe gap  $\mathcal{G}(x)$  at  $x$  was defined in (1.7), and when  $F$  is convex  $\mathcal{G}(x^k)$  is an upper bound on the optimality gap at  $x^k$ , see [13] for example. The Frank-Wolfe gap does not necessarily bound the optimality gap in the non-convex setting; nevertheless  $\mathcal{G}(x)$  is always nonnegative and  $\mathcal{G}(x) = 0$  if and only if  $x$  is a stationary point of  $\text{ERM}_\ell$ , and for this reason  $\mathcal{G}(x)$  is often used as a measure of non-stationarity at  $x$ , see [15, 24, 29, 22]. For an instance of  $\text{ERM}_\ell$  with non-convex objective function  $F$ , we say that a point  $x \in \mathcal{C}$  is an  $\varepsilon$ -stationarity point of  $\text{ERM}_\ell$  if  $\mathcal{G}(x) \leq \varepsilon$ .

Recall also that in order to run Algorithm 2.1, we need to specify the step-sizes  $\{\gamma_k\}_k$  and the Rule for constructing the sets  $\mathcal{B}_k$  of indices for updating the Taylor-points in Line 4 of Algorithm 2.1. Let  $K$  denote the number of iterations that we intend to run. In a somewhat similar spirit as the rule *Rule-SBD* $\sqrt{k}$  for convex losses, we present the following stochastic updating rule that is designed for  $\mathcal{B}_k$  to be comprised of  $\beta_k := n/\sqrt[4]{K}$  independently drawn samples from  $[n]$  without replacement, for all  $k \geq 1$ . If  $\beta_k$  is not integral, we will instead ensure  $\mathbb{E}|\mathcal{B}_k| = \beta_k$  by using a similar approach to *Rule-SBD* $\sqrt{k}$  which uses a Bernoulli random variable. With this in mind, we formally define the following Rule:

**DEFINITION 4.1.** *Rule-SBD* $\sqrt[4]{K}$ . For a fixed value of  $K \geq 1$ , and for all  $k \geq 1$  define  $\beta_k := \beta := n/\sqrt[4]{K}$  and  $p_k := p := \beta_k - \lfloor \beta_k \rfloor$ . Sample  $\xi_k$  from the Bernoulli distribution  $\text{Ber}(p_k)$  and then uniformly sample  $\lfloor \beta_k \rfloor + \xi_k$  samples from  $[n]$  without replacement. Then define  $\mathcal{B}_k$  to be the set of these index samples.

Note that by design of *Rule-SBD* $\sqrt[4]{K}$  it also follows that  $\mathbb{E}_{\xi_k} |\mathcal{B}_k| = \beta_k$ . We call this rule “*Rule-SBD* $\sqrt[4]{K}$ ” for “Stochastic Batch-size Decreasing at the rate  $\sqrt[4]{K}$ ,” because the Taylor points are updated less often as the number of overall iterations grow. Next we present a bound on the total number of flops used in  $K$  iterations of Algorithm 2.1 using *Rule-SBD* $\sqrt[4]{K}$ .

**PROPOSITION 4.2.** *Using Rule-SBD* $\sqrt[4]{K}$  and  $K \geq 1$  iterations, the expected total number of flops used in Algorithm 2.1 is  $O(K \cdot (\text{fLMO} + p^2) + K^{3/4} \cdot np^2)$ .

*Proof of Proposition 4.2.* For the initial iteration of Algorithm 2.1 the number of flops is  $O(\text{fLMO} + np^2)$ , and the expected total number of flops in each iteration thereafter is  $O(\text{fLMO} + p^2 + np^2/\sqrt[4]{K})$ . Summing over the  $K$  iterations yields the result.  $\square$

The following theorem presents a computational guarantee for computing an  $\varepsilon$ -stationarity point of  $\text{ERM}_\ell$  using Algorithm 2.1 with *Rule-SBD* $\sqrt[4]{K}$ . The proof of this theorem, as well as the other ensuing results in this section, is presented in subsection 4.1.

THEOREM 4.3. Suppose that [Assumption 1.1](#) holds and  $F$  is not necessarily convex, and [Algorithm 2.1](#) with  $\text{Rule-SBD}\sqrt[4]{K}$  is applied to the problem (1.2) with step-sizes defined by  $\gamma_k := \gamma := 1/\sqrt{K+1}$  for all  $k \geq 0$ , where  $K \geq 1$  is given. Then:

$$(4.1) \quad \mathbb{E}[\min_{k \in \{0, \dots, K\}} \mathcal{G}(x^k)] \leq \sum_{k=0}^K \frac{\mathbb{E}[\mathcal{G}(x^k)]}{K+1} \leq \frac{F(x^0) - F(x^*)}{\sqrt{K+1}} + \frac{3\hat{L}D_1D_\infty^2 + LD_2^2}{2n\sqrt{K+1}}.$$

The following corollary presents an alternative version of [Theorem 4.3](#) by randomly choosing an iteration index  $\hat{k} \in [K]$ .

COROLLARY 4.4. Let

$$K \geq \left\lceil \frac{\left(2n(F(x^0) - F(x^*)) + 3\hat{L}D_1D_\infty^2 + LD_2^2\right)^2}{(2n\varepsilon)^2} \right\rceil,$$

and let the iteration index  $\hat{k}$  be chosen uniformly from  $[K]$ , namely,  $\hat{k} \sim \mathcal{U}(\{1, \dots, K\})$ . Then  $\mathbb{E}_{\hat{k} \sim \mathcal{U}([K])} \mathbb{E}[\mathcal{G}(x^{\hat{k}})] \leq \varepsilon$ , and the total number of flops required is at most

$$O\left((\text{fLMO} + p^2) \left( \frac{n(F(x^0) - F(x^*)) + \hat{L}D_1D_\infty^2 + LD_2^2}{n\varepsilon} \right)^2 + np^2 \left( \frac{n(F(x^0) - F(x^*)) + \hat{L}D_1D_\infty^2 + LD_2^2}{n\varepsilon} \right)^{3/2} \right).$$

Also similar to [Corollary 3.5](#), we have the following corollary which shows that the joint dependence on  $n$  and  $\varepsilon$  is  $O(n/\varepsilon^{3/2})$  under the hypothesis that all feature vectors  $w_1, w_2, \dots$  lie in a bounded set  $S \subset \{w \in \mathbb{R}^p : \|w\|_* \leq M\}$ .

COROLLARY 4.5. Under the boundedness of the feature vectors, the bound on the number of flops in [Corollary 4.4](#) to obtain  $\mathbb{E}_{\hat{k} \sim \mathcal{U}([K])} \mathbb{E}[\mathcal{G}(x^{\hat{k}})] \leq \varepsilon$  is

$$O\left((\text{fLMO} + p^2) \left[ \frac{\left((F(x^0) - F(x^*)) + LM^2 \cdot \text{Diam}(C)^2 + \hat{L}M^3 \cdot \text{Diam}(C)^3\right)^2}{\varepsilon^2} \right] + np^2 \left[ \frac{\left((F(x^0) - F(x^*)) + LM^2 \cdot \text{Diam}(C)^2 + \hat{L}M^3 \cdot \text{Diam}(C)^3\right)^{3/2}}{\varepsilon^{3/2}} \right] \right).$$

In the expectation  $\mathbb{E}_{\hat{k} \sim \mathcal{U}([K])} \mathbb{E}[\mathcal{G}(x^{\hat{k}})]$  in [Corollaries 4.4](#) and [4.5](#), the randomness comes from both the random sampling at each iteration  $k$  to construct the update batches  $\mathcal{B}_k$  in  $\text{Rule-SBD}\sqrt[4]{K}$  as well as the random sampling of the iteration index  $\hat{k}$  from the uniform distribution on  $[K]$ .

In a similar spirit as  $\text{Rule-DBD}\sqrt{k}$ , we also present the following deterministic rule which achieves nearly identical computational guarantees (to within a constant factor) as  $\text{Rule-SBD}\sqrt[4]{K}$  but whose computational guarantee is deterministic.

DEFINITION 4.6.  $\text{Rule-DBD}\sqrt[4]{K}$ . For a fixed value of  $K \geq 1$ , and for any  $k \geq 1$ , define:

$$\mathcal{B}_k = \begin{cases} [n] & \text{if } k/\lfloor \sqrt[4]{K} \rfloor \in \mathbb{N} \\ \emptyset & \text{if } k/\lfloor \sqrt[4]{K} \rfloor \notin \mathbb{N}. \end{cases}$$

In  $\text{Rule-DBD}\sqrt[4]{K}$  we do not update any Taylor points unless  $k$  is integer times of  $\lfloor \sqrt[4]{K} \rfloor$ , and for these values of  $k$  we update all  $n$  Taylor points. We point out

that for *Rule-DBD*  $\sqrt[4]{K}$  the Taylor points are updated less often as  $K$  grows (in a different way but with similar effect as in *Rule-SBD*  $\sqrt[4]{K}$ ). Similar to the case of *Rule-SBD*  $\sqrt[4]{K}$ , we have:

PROPOSITION 4.7. *Using Rule-DBD  $\sqrt[4]{K}$  and  $K \geq 1$  iterations, the total number of flops used in Algorithm 2.1 is  $O(K \cdot (\text{fLMO} + p^2) + K^{3/4} \cdot np^2)$ .*

(We omit the proof as it is nearly identical to that of Proposition 4.2.)

THEOREM 4.8. *Suppose that Assumption 1.1 holds and  $F$  is not necessarily convex, and Algorithm 2.1 with Rule-DBD  $\sqrt[4]{K}$  is applied to the problem (1.2) with step-sizes defined by  $\gamma_k := \gamma := 1/\sqrt{K+1}$  for all  $k \geq 0$ , where  $K \geq 1$  is given. Then:*

$$(4.2) \quad \min_{k \in \{0, \dots, K\}} \mathcal{G}(x^k) \leq \sum_{k=0}^K \frac{\mathcal{G}(x^k)}{K+1} \leq \frac{F(x^0) - F(x^*)}{\sqrt{K+1}} + \frac{\hat{L}D_3^3 + LD_2^2}{2n\sqrt{K+1}}.$$

COROLLARY 4.9. *Let*

$$K \geq \left\lceil \frac{\left(2n(F(x^0) - F(x^*)) + \hat{L}D_3^3 + LD_2^2\right)^2}{(2n\varepsilon)^2} \right\rceil,$$

and let the iteration index  $\hat{k}$  be chosen uniformly from  $[K]$ , namely,  $\hat{k} \sim \mathcal{U}(\{1, \dots, K\})$ .  $\blacksquare$

Then  $\mathbb{E}_{\hat{k} \sim \mathcal{U}([K])}[\mathcal{G}(x^{\hat{k}})] \leq \varepsilon$ , and the total number of flops required is at most

$$O\left((\text{fLMO} + p^2) \left( \frac{n(F(x^0) - F(x^*)) + \hat{L}D_3^3 + LD_2^2}{n\varepsilon} \right)^2 + np^2 \left( \frac{n(F(x^0) - F(x^*)) + \hat{L}D_3^3 + LD_2^2}{n\varepsilon} \right)^{3/2} \right).$$

Also similar to Corollary 4.5, the following corollary shows that the joint dependence on  $n$  and  $\varepsilon$  is  $O(n/\varepsilon^{3/2})$  under the hypothesis that all feature vectors  $w_1, w_2, \dots$  lie in a bounded set  $S \subset \{w \in \mathbb{R}^p : \|w\|_* \leq M\}$ .

COROLLARY 4.10. *Under the boundedness of the feature vectors, the bound on the number of flops in Corollary 4.9 to obtain  $\mathbb{E}_{\hat{k} \sim \mathcal{U}([K])}[\mathcal{G}(x^{\hat{k}})] \leq \varepsilon$  is*

$$O\left((\text{fLMO} + p^2) \left[ \frac{\left((F(x^0) - F(x^*)) + LM^2 \cdot \text{Diam}(C)^2 + \hat{L}M^3 \cdot \text{Diam}(C)^3\right)^2}{\varepsilon^2} \right] + np^2 \left[ \frac{\left((F(x^0) - F(x^*)) + LM^2 \cdot \text{Diam}(C)^2 + \hat{L}M^3 \cdot \text{Diam}(C)^3\right)^{3/2}}{\varepsilon^{3/2}} \right] \right).$$

Table 2 shows a comparison of the computational guarantees of the standard Frank-Wolfe method and TUFW with *Rule-SBD*  $\sqrt[4]{K}$  and *Rule-DBD*  $\sqrt[4]{K}$ .

Note in particular that with *Rule-SBD*  $\sqrt[4]{K}$  or *Rule-DBD*  $\sqrt[4]{K}$ , the joint dependence on  $n$  and  $\varepsilon$  is  $O(n/\varepsilon^{3/2})$ , as compared to  $O(n^{2/3}/\varepsilon^2)$  in [24],  $O(n^{1/2}/\varepsilon^2)$  in [29], and  $O(1/\varepsilon^3)$  in [30, 9]. In the regime where  $\varepsilon$  is sufficiently small, the dominant term in the complexity bound is the left-most term which counts the operations of the LMO, which is  $O(1/\varepsilon^2)$  (independent of  $n$ ), and is essentially the same as the lower bound complexity in [16]. Because of this, we say that TUFW's overall complexity is nearly optimal. We also mention that the joint dependence on  $n$  and  $\varepsilon$  for CASPIDERRG in



TABLE 2

Complexity bounds for different Frank-Wolfe methods to obtain an  $\varepsilon$ -stationary solution of  $ERM_\ell$  with non-convex losses, under the boundedness assumption of the feature vectors. In the table  $\varepsilon_0 := F(x^0) - F(x^*)$ ,  $c_1 := LM^2 \text{Diam}(\mathcal{C})^2$ , and  $c_2 := \hat{L}M^3 \text{Diam}(\mathcal{C})^3$ .

Method	Optimality Metric	Overall Complexity
<i>Rule-SBD</i> $\sqrt[4]{K}$ (Cor. 4.5)	$\mathbb{E}_{\hat{k} \sim \mathcal{U}([K])} \mathbb{E}[\mathcal{G}(x^{\hat{k}})] \leq \varepsilon$	$O\left((\text{fLMO} + p^2) \cdot \frac{(\varepsilon_0 + c_1 + c_2)^2}{\varepsilon^2} + np^2 \cdot \frac{(\varepsilon_0 + c_1 + c_2)^{3/2}}{\varepsilon^{3/2}}\right)$
<i>Rule-DBD</i> $\sqrt[4]{K}$ (Cor. 4.10)	$\mathbb{E}_{\hat{k} \sim \mathcal{U}([K])} [\mathcal{G}(x^{\hat{k}})] \leq \varepsilon$	$O\left((\text{fLMO} + p^2) \cdot \frac{(\varepsilon_0 + c_1 + c_2)^2}{\varepsilon^2} + np^2 \cdot \frac{(\varepsilon_0 + c_1 + c_2)^{3/2}}{\varepsilon^{3/2}}\right)$
Standard Frank-Wolfe ([15])	$\mathbb{E}_{\hat{k} \sim \mathcal{U}([K])} [\mathcal{G}(x^{\hat{k}})] \leq \varepsilon$	$O\left((\text{fLMO} + np) \cdot \frac{(\varepsilon_0 + c_1)^2}{\varepsilon^2}\right)$

[26] is  $O(n^{3/4}/\varepsilon^{3/2})$ ; however CASPIDERRG requires access to the exact full gradients periodically, and in our computational experience we found that CASPIDERRG was not computationally viable relative to other stochastic Frank-Wolfe methods.

Finally, we point out that when the loss functions  $l_i$  are all quadratic loss functions, the Hessian  $\nabla^2 F(x)$  is constant for all  $x$ , and we can use *Rule- $\emptyset$*  which is the rule that does not update any of the Taylor points.

**4.1. Proofs of results in section 4.** Our proofs will use Lemmas 3.13 and 3.14 since they do not assume convexity of  $F$ . In addition our proofs will rely on the following sequences of lemmas. We will use  $\varepsilon_0$  to denote  $F(x^0) - F(x^*)$  in the rest of this section.

LEMMA 4.11. Under Assumption 1.1 for problem (1.2), if we use the step-size  $\gamma_k := \gamma := 1/\sqrt{K+1}$  in Algorithm 2.1, then

$$(4.3) \quad \sum_{k=0}^K \frac{\mathcal{G}(x^k)}{K+1} \leq \frac{\varepsilon_0}{\sqrt{K+1}} + \frac{1}{K+1} \sum_{k=0}^K (\nabla F(x^k) - g^k)^\top (s^k - \bar{s}^k) + \frac{LD_2^2}{2n\sqrt{K+1}}.$$

*Proof of Lemma 4.11.* For any  $k \geq 0$  we have:

$$(4.4) \quad \begin{aligned} F(x^{k+1}) &\leq F(x^k) + \gamma_k (\nabla F(x^k))^\top (s^k - x^k) + \gamma_k^2 LD_2^2 / 2n \\ &= F(x^k) + \gamma_k (g^k)^\top (s^k - x^k) + \gamma_k (\nabla F(x^k) - g^k)^\top (s^k - x^k) + \gamma_k^2 LD_2^2 / 2n \\ &\leq F(x^k) + \gamma_k (g^k)^\top (\bar{s}^k - x^k) + \gamma_k (\nabla F(x^k) - g^k)^\top (s^k - x^k) + \gamma_k^2 LD_2^2 / 2n \\ &= F(x^k) - \gamma_k \mathcal{G}(x^k) + \gamma_k (g^k - \nabla F(x^k))^\top (\bar{s}^k - s^k) + \gamma_k^2 LD_2^2 / 2n, \end{aligned}$$

where  $\bar{s}^k \in \arg \min_{s \in \mathcal{C}} (\nabla F(x^k))^\top s$ , and the first inequality uses Lemma 3.13, and the second inequality uses the fact that  $s^k \in \arg \min_{s \in \mathcal{C}} (g^k)^\top s$ . Substituting the value of  $\gamma_k$  in the statement of the lemma yields after rearranging terms:

$$(4.5) \quad \frac{\mathcal{G}(x^k)}{\sqrt{K+1}} \leq F(x^k) - F(x^{k+1}) + \frac{1}{\sqrt{K+1}} (\nabla F(x^k) - g^k)^\top (s^k - \bar{s}^k) + \frac{LD_2^2}{2n(K+1)}.$$

The inequality (4.3) then follows by first summing the inequalities (4.5) for  $k \in \{0, \dots, K\}$  and dividing by  $\sqrt{K+1}$ , and noting that  $F(x^0) - F(x^{K+1}) \leq \varepsilon_0$ .  $\square$

LEMMA 4.12. Let the series  $\{d_{ij}\}_{i \in [n], j \in \mathbb{N}}$  be given, and suppose that there exists  $M_1$  and  $M_\infty$  for which  $\sup_j \sum_{i=1}^n d_{ij} \leq M_1$  and  $\sup_{i,j} d_{ij} \leq M_\infty$ . Then for Algorithm 2.1 using Rule-SBD  $\sqrt[4]{K}$  with step-sizes  $\gamma_k \equiv \gamma := 1/\sqrt{K+1}$ , it holds that:

$$(4.6) \quad \mathbb{E}\left[\sum_{i=1}^n \sum_{j=\tau_i^k}^{k-1} \gamma_j d_{ij}\right] \leq \frac{M_1}{(K+1)^{1/4}} \quad \text{and} \quad \mathbb{E}\left[\sum_{i=1}^n \left(\sum_{j=\tau_i^k}^{k-1} \gamma_j d_{ij}\right)^2\right] \leq \frac{3M_1 M_\infty}{\sqrt{K+1}}.$$

647 *Proof of Lemma 4.12.* Similar to the proof of Lemma 3.18, let us first define  
 648  $G_k^i := \left( \sum_{j=\tau_i^k}^{k-1} \gamma_j d_{ij} \right)$  and  $Z_k^i := \left( \sum_{j=\tau_i^k}^{k-1} \gamma_j d_{ij} \right)^2$  for  $i \in [n]$ , and also  $G_k := \sum_{i=1}^n G_k^i$   
 649 and  $Z_k := \sum_{i=1}^n Z_k^i$ . First notice in  $G_k^i$  and  $Z_k^i$  that the randomness lies only in  $\tau_i^k$ ,  
 650 which in turn depends on the update sets  $\mathcal{B}_0, \mathcal{B}_1, \dots, \mathcal{B}_k$ . We note that  $G_k^i = 0$  when  
 651  $i \in \mathcal{B}_k$ , which happens with probability  $\beta/n = 1/\sqrt[4]{K}$ . Because of the relationship  
 652 between  $\{\tau_i^k\}$  and  $\{\tau_i^{k-1}\}$  in (3.6), we further note that with probability  $1 - \beta/n$   
 653 that  $i \notin \mathcal{B}_k$  in which case it follows that  $G_k^i = G_{k-1}^i + \gamma_{k-1} d_{i,k-1} = G_{k-1}^i + \gamma d_{i,k-1}$ .  
 654 Defining  $\rho := 1 - \beta/n = 1 - 1/\sqrt[4]{K}$ , we have

$$655 \quad (4.7) \quad \mathbb{E}[G_k] = \sum_{i=1}^n \mathbb{E}[G_k^i] = \rho \sum_{i=1}^n \mathbb{E}[G_{k-1}^i + \gamma d_{i,k-1}] \leq \rho (\mathbb{E}[G_{k-1}] + \gamma M_1) ,$$

656 where the inequality uses  $\sup_j \sum_{i=1}^n d_{ij} \leq M_1$ . Now (4.7) can be rearranged to yield:

$$657 \quad (4.8) \quad \mathbb{E}[G_k] - \frac{\rho \gamma M_1}{1 - \rho} \leq \rho \left( \mathbb{E}[G_{k-1}] - \frac{\rho \gamma M_1}{1 - \rho} \right) \leq \dots \leq \rho^k \left( \mathbb{E}[G_0] - \frac{\rho \gamma M_1}{1 - \rho} \right) \leq 0 ,$$

658 where the last inequality uses  $G_0 = 0$  (since all Taylor points are set to  $x^0$  in Line 2  
 659 of Algorithm 2.1). We therefore have:

$$660 \quad (4.9) \quad \mathbb{E}[G_k] \leq \frac{\rho \gamma M_1}{1 - \rho} = \frac{(1 - K^{-1/4})(K + 1)^{-1/2} M_1}{K^{-1/4}} \leq \frac{M_1}{(K + 1)^{1/4}} ,$$

661 which proves the first inequality of (4.6).

662 We proceed similarly to establish the second inequality of (4.6). We note that  
 663  $Z_k^i = 0$  when  $i \in \mathcal{B}_k$ , which happens with probability  $\beta/n = 1/\sqrt[4]{K}$ . Because of the  
 664 relationship between  $\{\tau_i^k\}$  and  $\{\tau_i^{k-1}\}$  in (3.6), we further note that with probability  
 665  $\rho = 1 - \beta/n$  that  $i \notin \mathcal{B}_k$  in which case it follows that  $Z_k^i = (G_{k-1}^i + \gamma d_{i,k-1})^2$ . We  
 666 thus have

$$\begin{aligned} \mathbb{E}[Z_k] &= \sum_{i=1}^n \mathbb{E}[Z_k^i] = \rho \sum_{i=1}^n \mathbb{E}[(G_{k-1}^i + \gamma d_{i,k-1})^2] \\ &= \rho \sum_{i=1}^n \left( \mathbb{E}[Z_{k-1}^i] + 2\gamma \mathbb{E}[G_{k-1}^i] d_{i,k-1} + \gamma^2 (d_{i,k-1})^2 \right) \\ 667 \quad (4.10) \quad &\leq \rho \left( \mathbb{E}[Z_{k-1}] + 2\gamma \mathbb{E}[G_{k-1}] M_\infty + \gamma^2 (M_\infty M_1) \right) \\ &\leq \rho \left( \mathbb{E}[Z_{k-1}] + 2\gamma \frac{M_1 M_\infty}{(K + 1)^{1/4}} + \gamma^2 M_\infty M_1 \right) \\ &= \rho \left( \mathbb{E}[Z_{k-1}] + 2 \frac{M_1 M_\infty}{(K + 1)^{3/4}} + \frac{M_\infty M_1}{K + 1} \right) \\ &\leq \rho \left( \mathbb{E}[Z_{k-1}] + \frac{3M_1 M_\infty}{(K + 1)^{3/4}} \right) . \end{aligned}$$

668 Now (4.10) can be rearranged to yield:

$$\begin{aligned} 669 \quad (4.11) \quad \mathbb{E}[Z_k] - \frac{3M_1 M_\infty}{(K + 1)^{3/4}} \cdot \frac{\rho}{1 - \rho} &\leq \rho \left( \mathbb{E}[Z_{k-1}] - \frac{3M_1 M_\infty}{(K + 1)^{3/4}} \cdot \frac{\rho}{1 - \rho} \right) \\ &\leq \dots \leq \rho^k \left( \mathbb{E}[Z_0] - \frac{3M_1 M_\infty}{(K + 1)^{3/4}} \cdot \frac{\rho}{1 - \rho} \right) \leq 0 , \end{aligned}$$

670 where the last inequality uses  $Z_0 = 0$  (since all Taylor points are set to  $x^0$  in Line 2  
 671 of Algorithm 2.1), from which it follows that

$$672 \quad (4.12) \quad \mathbb{E}[Z_k] \leq \frac{3M_1 M_\infty}{(K + 1)^{3/4}} \cdot \frac{\rho}{1 - \rho} \leq \frac{3M_1 M_\infty}{\sqrt{K + 1}} ,$$

673 which proves the second inequality of (4.6).  $\square$

*Proof of Theorem 4.3.* Notice that the first inequality of (4.1) follows from the chain

$$\mathbb{E}[\min_{k \in \{0, \dots, K\}} \mathcal{G}(x^k)] \leq \min_{k \in \{0, \dots, K\}} \mathbb{E}[\mathcal{G}(x^k)] \leq \sum_{k=0}^K \frac{\mathbb{E}[\mathcal{G}(x^k)]}{K+1}.$$

674 From Lemma 4.11 we have:

$$675 \quad (4.13) \quad \sum_{k=0}^K \frac{\mathcal{G}(x^k)}{K+1} \leq \frac{2n\varepsilon_0 + LD_2^2}{2n\sqrt{K+1}} + \frac{1}{K+1} \sum_{k=0}^K (\nabla F(x^k) - g^k)^\top (s^k - \bar{s}^k),$$

676 and applying Corollary 3.15 we obtain:

$$677 \quad (4.14) \quad \sum_{k=0}^K \frac{\mathcal{G}(x^k)}{K+1} \leq \frac{2n\varepsilon_0 + LD_2^2}{2n\sqrt{K+1}} + \frac{\hat{L}D_\infty}{2n(K+1)} \sum_{k=0}^K \sum_{i=1}^n \left( \sum_{j=\tau_i^k}^{k-1} \gamma_j |w_i^\top (s^j - x^j)| \right)^2.$$

678 Taking expectations on both sides, it follows that:

$$679 \quad (4.15) \quad \sum_{k=0}^K \frac{\mathbb{E}[\mathcal{G}(x^k)]}{K+1} \leq \frac{2n\varepsilon_0 + LD_2^2}{2n\sqrt{K+1}} + \frac{\hat{L}D_\infty}{2n(K+1)} \sum_{k=0}^K \mathbb{E} \left[ \sum_{i=1}^n \left( \sum_{j=\tau_i^k}^{k-1} \gamma_j |w_i^\top (s^j - x^j)| \right)^2 \right].$$

680 Now define  $d_{ij} := |w_i^\top (s^j - x^j)|$  and  $M_1 := D_1$  and  $M_\infty := D_\infty$ . Then since  
 681  $\sup_j \sum_{i=1}^n d_{ij}$  is equal to  $\sup_j \sum_{i=1}^n |w_i^\top (s^j - x^j)|$  which is at most  $D_1 = M_1$ , and  
 682  $\sup_{i,j} d_{ij}$  is equal to  $\sup_{i,j} |w_i^\top (s^j - x^j)|$  which is at most  $D_\infty = M_\infty$ , we can apply  
 683 the second inequality of Lemma 4.12 to (4.15) and we obtain the second inequality of  
 684 (4.1).  $\square$

685 We next proceed towards the proof of Theorem 4.8. Let us first prove the following  
 686 lemma.

687 **LEMMA 4.13.** *In any iteration  $k$  of Algorithm 2.1 with Rule-DBD  $\sqrt[4]{K}$  and  $\gamma_k :=$   
 688  $\gamma := 1/\sqrt{K+1}$ , for problem setup (1.2), it holds that*

$$689 \quad (4.16) \quad \sum_{i=1}^n \left( \sum_{j=\tau_i^k}^{k-1} |w_i^\top (z^j - y^j)| \right)^2 \cdot |w_i^\top (z^k - \bar{z}^k)| \leq K^{1/2} D_3^3$$

690 holds for any  $y^j, z^j, \bar{z}^j \in \mathcal{C}$ ,  $j = 1, \dots, k$ .

691 *Proof of Lemma 4.13.* Notice in Rule-DBD  $\sqrt[4]{K}$  that  $k - \tau_i^k \leq \lfloor K^{1/4} \rfloor - 1$ . Now  
 692 let  $u, v \geq 1$  satisfy  $\frac{1}{u} + \frac{1}{v} = 1$ . Then it holds that

$$\begin{aligned} & \sum_{i=1}^n \left( \sum_{j=\tau_i^k}^{k-1} |w_i^\top (z^j - y^j)| \right)^2 \cdot |w_i^\top (z^k - \bar{z}^k)| \\ 693 \quad (4.17) \quad & \leq (\lfloor K^{1/4} \rfloor - 1) \sum_{i=1}^n \sum_{j=\max\{0, k+1-\lfloor K^{1/4} \rfloor\}}^{k-1} |w_i^\top (z^j - y^j)|^2 \cdot |w_i^\top (z^k - \bar{z}^k)| \\ & \leq (\lfloor K^{1/4} \rfloor - 1) \sum_{j=\max\{0, k+1-\lfloor K^{1/4} \rfloor\}}^{k-1} D_{2u}^2 D_v \leq (\lfloor K^{1/4} \rfloor - 1)^2 D_{2u}^2 D_v, \end{aligned}$$

694 where the first inequality uses  $\|a\|_1^2 \leq d\|a\|_2^2$  for  $a \in \mathbb{R}^d$ , and the second inequality  
 695 follows similar to the equivalent part of proof of Lemma 3.20 that uses  $\sum_{i=1}^n a_i^2 b_i \leq$   
 696  $\|a\|_{2u}^2 \|b\|_v$  for  $a, b \in \mathbb{R}^n$  and  $u, v \geq 1$  and  $\frac{1}{u} + \frac{1}{v} = 1$ . Setting  $u = \frac{3}{2}$  and  $v = 3$ , the  
 697 right-hand side of (4.17) is bounded above by  $K^{1/2} D_3^3$ .  $\square$

698 *Proof of Theorem 4.8.* The first inequality in (4.2) is obvious. For the second  
 699 inequality, note from Lemma 4.11 that:

$$700 \quad (4.18) \quad \sum_{k=0}^K \frac{\mathcal{G}(x^k)}{K+1} \leq \frac{2n\varepsilon_0 + LD_2^2}{2n\sqrt{K+1}} + \frac{1}{K+1} \sum_{k=0}^K (\nabla F(x^k) - g^k)^\top (s^k - \bar{s}^k).$$

701 Let  $Q_K$  denote the second term of the right-hand side of (4.18). We have

$$702 \quad (4.19) \quad Q_K \leq \frac{\hat{L}}{2n(K+1)^2} \sum_{k=0}^K \sum_{i=1}^n \left( \sum_{j=\tau_i^k}^{k-1} |w_i^\top (s^j - x^j)| \right)^2 \cdot |w_i^\top (s^k - \bar{s}^k)| \leq \frac{\hat{L}D_3^3}{2n\sqrt{K+1}},$$

703 where the first inequality is due to Lemma 3.14, and the second inequality is due to  
 704 Lemma 4.13. Substituting this bound back to (4.18) yields the second inequality of  
 705 (4.2).  $\square$

706 **5. Extensions.** We also propose an adaptive step-size for TUFW. Besides, TUFW  
 707 can also be easily implemented in the general ERM problems, and the theoretical re-  
 708 sults also have the corresponding version in the general ERM problems.

709 **5.1. Adaptive Step-size.** In addition to the standard step-size  $\gamma_k$ , which is  $\frac{2}{k+2}$   
 710 for convex loss functions and  $\frac{1}{\sqrt{K+1}}$  for non-convex loss functions, we also develop a  
 711 kind of adaptive step-size  $\tilde{\gamma}_k$  as follows:

$$712 \quad (5.1) \quad \tilde{\gamma}_k := \begin{cases} \min \left\{ \gamma_k, \frac{(g^k)^\top (x^k - s^k)}{(s^k - x^k)^\top H_k (s^k - x^k)} \right\} & \text{when } (s^k - x^k)^\top H_k (s^k - x^k) > 0, \\ \gamma_k & \text{when } (s^k - x^k)^\top H_k (s^k - x^k) \leq 0, \end{cases}$$

713 where  $H_k$  is defined in (2.1) and  $\gamma_k$  is the standard step-size.

714 This adaptive step-size  $\tilde{\gamma}_k$  can approximately minimize the quadratic approxima-  
 715 tion of the objective function in the range of  $[0, \gamma_k]$ . Let  $x(\lambda) := x^k + \gamma(s^k - x^k)$  and  
 716 then

$$(5.2) \quad \begin{aligned} F(x(\gamma)) &= F(x^k) + \gamma(g^k)^\top (s^k - x^k) + \frac{\gamma^2}{2} (s^k - x^k)^\top H_k (s^k - x^k) \\ &+ \frac{1}{n} \sum_{i=1}^n \int_{t=0}^{\gamma} \left( \nabla f_i(x^k + t(s^k - x^k)) - \nabla f_i(b_i) - \nabla^2 f_i(b_i)(x^k + t(s^k - x^k) - b_i) \right)^\top (s^k - x^k) dt. \end{aligned}$$

It could be further proven that when  $\gamma \in [0, \gamma_k]$ , the first three terms of the right-hand side dominates and therefore  $\tilde{\gamma}_k$  defined in (5.1), which is also the closed-form solution of

$$\arg \min_{\gamma \in [0, \gamma_k]} F(x^k) + \gamma(g^k)^\top (s^k - x^k) + \frac{\gamma^2}{2} (s^k - x^k)^\top H_k (s^k - x^k),$$

718 can be approximately regarded as  $\arg \min_{\gamma \in [0, \gamma_k]} F(x(\gamma))$ , which yields more decrease  
 719 of the objective value than the standard step-size. See the technical report for the  
 720 proof of the convergence rates of using adaptive step-sizes.

721 **5.2. TUFW for General ERM Problems.** In this subsection, we extend  
 722 TUFW to the general ERM problems. Instead of Assumption 1.1, the new assumption  
 723 for (1.1) states that, a) the diameter of the  $\mathcal{C}$  is bounded, b) for each  $i \in [n]$ , the  
 724 gradient  $\nabla f_i(\cdot)$  is  $L$ -Lipschitz continuous on  $\mathcal{C}$ , c) and the Hessian  $\nabla^2 f_i(x)$  is  $\hat{L}$ -  
 725 Lipschitz continuous on  $\mathcal{C}$ . If still using the rules we introduced before, the overall  
 726 complexiy is as follows:

COROLLARY 5.1. *Suppose that the above assumption holds, then*

1. *when  $F$  is convex, and  $\gamma_k := \frac{2}{k+2}$ , and Algorithm 2.1 with Rule-SBD $\sqrt{k}$  (or Rule-DBD $\sqrt{k}$ ) is applied to the problem (1.1), in order to obtain  $\mathbb{E}[F(x^k) - F(x^*)] \leq \varepsilon$  (or  $F(x^k) - F(x^*) \leq \varepsilon$ ), the bound of the number of flops is*

$$O\left((\text{fLMO} + p^2) \cdot \frac{L \cdot \text{Diam}(\mathcal{C})^2 + \hat{L} \cdot \text{Diam}(\mathcal{C})^3}{\varepsilon} + np^2 \cdot \frac{\sqrt{L \cdot \text{Diam}(\mathcal{C})^2 + \hat{L} \cdot \text{Diam}(\mathcal{C})^3}}{\sqrt{\varepsilon}}\right),$$

2. *and if  $\gamma_k := \frac{1}{\sqrt[k]{K+1}}$ , and Algorithm 2.1 with Rule-SBD $\sqrt[k]{K}$  (or Rule-DBD $\sqrt[k]{K}$ ) is applied to the problem (1.1), in order to obtain  $\mathbb{E}_{\hat{k} \sim \mathcal{U}([K])} \mathbb{E}[\mathcal{G}(x^{\hat{k}})] \leq \varepsilon$  (or  $\mathbb{E}_{\hat{k} \sim \mathcal{U}([K])} [\mathcal{G}(x^{\hat{k}})] \leq \varepsilon$ ), the bound of the number of flops is*

$$O\left((\text{fLMO} + p^2) \left[ \frac{\left((F(x^0) - F(x^*)) + L \cdot \text{Diam}(\mathcal{C})^2 + \hat{L} \cdot \text{Diam}(\mathcal{C})^3\right)^2}{\varepsilon^2} \right] + np^2 \left[ \frac{\left((F(x^0) - F(x^*)) + L \cdot \text{Diam}(\mathcal{C})^2 + \hat{L} \cdot \text{Diam}(\mathcal{C})^3\right)^{3/2}}{\varepsilon^{3/2}} \right] \right).$$

See the technical report for the proof of the convergence of TUFW methods on ERM problems.

**6. Computational Experiments.** In this section we present the results of computational experiments where we compare Algorithm TUFW (with the four different batch construction Rules developed earlier) and five other methods, namely the standard Frank-Wolfe method [6, 13] plus related methods in the recent literature, as follows:

- (FW) – the standard Frank-Wolfe method, see [6, 13]
- (FW-ada) – the standard Frank-Wolfe method with adaptive step-size [3], which is defined by  $\gamma_k := \min\{1, \mathcal{G}(x^k)/(L\|x^k - s^k\|^2)\}$ , where  $L$  is the Lipschitz constant
- (SPIDER-FW) – the Frank-Wolfe method with stochastic path-integrated differential estimator technique developed in [29]
- (CSFW) – the constant batch-size stochastic Frank-Wolfe method developed in [22], and
- (CASPIDERG) – curvature-aided stochastic path-integrated differential estimator gradient developed in [26].

We tested these methods on both convex and non-convex instances of  $\text{ERM}_\ell$ . We first conducted computational tests on problems with convex losses, using instances of the  $\ell_1$ -constrained logistic regression problem (LR):

$$(6.1) \quad \text{LR:} \quad \min_x F(x) := \frac{1}{n} \sum_{i=1}^n \ln(1 + \exp(-y_i w_i^\top x)) \quad \text{s.t.} \quad \|x\|_1 \leq \lambda,$$

where the  $\ell_1$ -ball constraint is designed to induce sparsity and/or regularization. We chose instances of LR from LIBSVM [2] for which the number of observations  $n \gg 0$  but not so large as to render the problem intractable. This yielded 12 instances from LIBSVM, namely **a1a**, **a2a**, **a8a**, **a9a**, **w1a**, **w2a**, **w7a**, **w8a**, **svmguide3**, **phishing**, **ijcnn1**, and **covtype**, whose dimensions  $n$  and  $p$  are displayed in the third and forth columns of Table 3. The value of  $\lambda$  was specified by 5-fold cross-validation of each instance. However, we used the same cross-validated  $\lambda$  for **a1a**, **a2a**, **a8a**, and **a9a**,

and similarly for w1a, w2a, w7a, for w8a. (Separately, we also tested the methods on larger feasible regions by enlarging each instance’s feasible region using  $\lambda' := 100\lambda$ .) All methods were implemented and run in Python on the MIT Engaging Cluster. To ensure fairness, our implementations did not use sparse linear algebra for all methods, including the techniques in the last paragraph of section 2. **Zikai: I have made the changes. Please also refer to the last paragraph of section 2**

We measured the performance of methods on problems with convex losses using the Frank-Wolfe gap  $\mathcal{G}(x^k)$  defined in (1.7), since  $\mathcal{G}(x^k)$  is a (conservative) upper bound on the optimality gap, namely  $F(x^k) - F^* \leq \mathcal{G}(x^k)$  (see [13]), and the optimal value  $F^*$  is not known. We computed  $\mathcal{G}(x^k)$  separately and offline after-the-fact.

Note that *Rule-SBD* $\sqrt{k}$  and *Rule-SBD* $\sqrt[4]{K}$  for TUFW require uniform sampling of  $\tilde{\beta}_k := |\mathcal{B}_k|$  indices from  $[n]$  without replacement. These indices are then used to perform matrix and vector operations on the corresponding indexed columns of  $W$ . Because matrix and vector operations involving random indices are inherently inefficient in practice, we instead implemented *Rule-SBD* $\sqrt{k}$  and *Rule-SBD* $\sqrt[4]{K}$  by uniformly sampling a starting index  $\hat{i} \sim \mathcal{U}([n])$  and then assigning the other  $|\tilde{\beta}_k| - 1$  indices successively to construct  $\mathcal{B}_k := \{\hat{i}, \hat{i} + 1, \dots, \hat{i} + \tilde{\beta}_k - 1\}$  (modulo  $n$ ). This modification does not affect any of the theoretical convergence guarantees because (3.18) and (3.24) still hold and thus Lemma 3.18 is still true, and similarly (4.7) and (4.10) still hold and thus Lemma 4.12 is still true.

Figure 1 shows the performance of TUFW with *Rule-SBD* $\sqrt{k}$  and *Rule-DBD* $\sqrt{k}$ , and the standard Frank-Wolfe method, on the dataset a9a (which we found to be representative of typical performance). The methods were run both with and without adaptive step-sizes. Notice that the TUFW methods significantly outperform the Frank-Wolfe method. *Rule-DBD* $\sqrt{k}$  using adaptive step-sizes significantly outperforms the other methods in terms of overall runtime and has similar LMO dependence as *Rule-SBD* $\sqrt{k}$ . The second-best performance is from TUFW with *Rule-SBD* $\sqrt{k}$  with adaptive step-sizes. As these results are typical, we will focus on *Rule-SBD* $\sqrt{k}$  and *Rule-DBD* $\sqrt{k}$  with adaptive step-sizes in further comparisons with other Frank-Wolfe methods.

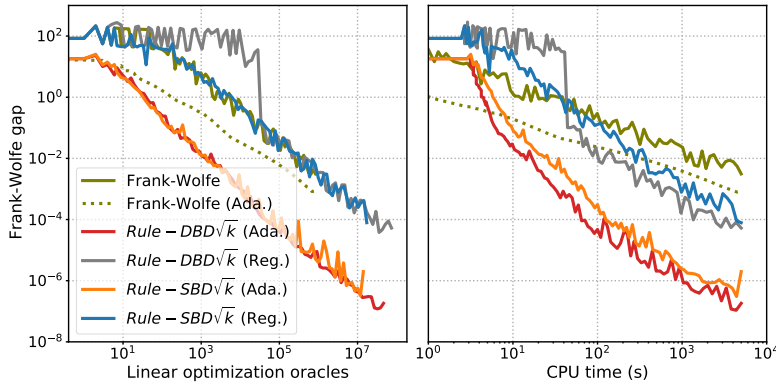


FIG. 1. Performance of TUFW with *Rule-SBD* $\sqrt{k}$  and *Rule-DBD* $\sqrt{k}$  and the standard Frank-Wolfe method, with and without adaptive step-sizes, on the logistic regression problem (6.1), on the dataset a9a.

In Table 3 we compare the CPU runtimes (in seconds) of the TUFW methods using *Rule-SBD* $\sqrt{k}$  and *Rule-DBD* $\sqrt{k}$ , with FW, FW-ada, SPIDER-FW, and CSFW.



(We did not include CASPIDERG in these tests as it is not designed for convex objective functions.) The batch-size of CSFW was set to  $n/100$ , which is exactly as in the implementation in [22]. For each method and each dataset we ran 10 trials, and we report the average of these trials. We only report average runtimes that did not exceed 5000 seconds, otherwise the entry is blank. In each row of the table the bold number highlights the best runtime among the six methods. The last column of the table reports the speed-up of TUFW (the best of  $Rule-SBD\sqrt{k}$  and  $Rule-DBD\sqrt{k}$ ) over the best of the other four methods. Speed-up values larger than 1 are shown in bold. We observe from Table 3 that the TUFW method – using either  $Rule-SBD\sqrt{k}$  or  $Rule-DBD\sqrt{k}$ , significantly outperforms the other methods on almost all problem instances and almost all target optimality tolerances. Furthermore, as the optimality tolerance  $\varepsilon$  decreases and/or as the number of observations  $n$  increase, the advantage of TUFW methods over the other methods becomes even more evident.

In order to test our TUFW methods on problems with larger feasible regions, we increased the size of the feasibility set in (6.1) by inflating the value of  $\lambda$  to  $\lambda' = 100\lambda$  (where recall  $\lambda$  was determined by cross validation). With the larger feasible regions, the advantage of the TUFW methods is even more pronounced. The details of this test is discussed in the Technical Report [?].

TABLE 3

Comparison of average CPU runtimes (in seconds) required to achieve  $\mathcal{G}(x^k) \leq \varepsilon$  for methods on the logistic regression problem (6.1) with  $\lambda$  chosen by cross-validation. (A blank indicates the method used more than 5000 seconds.)

$\varepsilon$	dataset	$n$	$p$	$Rule-SBD\sqrt{k}$	$Rule-DBD\sqrt{k}$	FW	FW-ada	SPIDER-FW	CSFW	Speed-up
1e-1	a1a	1605	123	0.73	<b>0.37</b>	2.02	1.18	348.58	5.03	<b>3.20</b>
1e-3	a1a	1605	123	14.40	<b>5.14</b>	138.87	380.10		160.02	<b>27.02</b>
1e-5	a1a	1605	123	3029.61	<b>1034.66</b>	2956.50				<b>2.86</b>
1e-1	a2a	2265	123	1.02	<b>0.52</b>	3.03	1.55	358.72	5.47	<b>2.97</b>
1e-3	a2a	2265	123	21.24	<b>8.16</b>	197.80	472.72		167.66	<b>20.54</b>
1e-5	a2a	2265	123	2039.75	<b>631.23</b>					
1e-1	a8a	22696	123	7.75	<b>4.66</b>	78.18	12.83	443.85	15.22	<b>2.75</b>
1e-3	a8a	22696	123	51.29	<b>23.71</b>		3178.02		1250.12	<b>52.72</b>
1e-5	a8a	22696	123	420.64	<b>174.87</b>					
1e-1	a9a	32561	123	11.26	<b>6.80</b>	94.50	21.00	440.73	19.82	<b>2.91</b>
1e-3	a9a	32561	123	65.25	<b>31.22</b>		3844.13		1660.13	<b>53.17</b>
1e-5	a9a	32561	123	490.36	<b>197.15</b>					
1e-1	w1a	2477	300	10.45	<b>5.71</b>	17.51	331.22	3566.83	44.39	<b>3.07</b>
1e-3	w1a	2477	300	51.51	<b>21.93</b>	327.38			1219.03	<b>14.93</b>
1e-5	w1a	2477	300	1743.08	<b>541.71</b>					
1e-1	w2a	3470	300	13.11	<b>6.80</b>	37.34	574.42		54.06	<b>5.50</b>
1e-3	w2a	3470	300	128.14	<b>50.77</b>	369.82			1207.01	<b>7.28</b>
1e-5	w2a	3470	300		<b>1422.93</b>					
1e-1	w7a	24692	300	88.86	<b>55.43</b>	541.10			129.19	<b>2.33</b>
1e-3	w7a	24692	300	320.17	<b>178.87</b>				4413.77	<b>24.68</b>
1e-5	w7a	24692	300	3106.98	<b>1516.83</b>					
1e-1	w8a	49749	300	174.89	<b>114.66</b>	1198.35			212.80	<b>1.86</b>
1e-3	w8a	49749	300	553.94	<b>355.23</b>					
1e-5	w8a	49749	300		<b>2707.92</b>					
1e-1	svmguide3	1243	22	0.09	<b>0.02</b>	1.35	0.31	21.46	2.06	<b>13.43</b>
1e-3	svmguide3	1243	22	8.01	<b>2.64</b>	53.64	175.49		48.85	<b>18.53</b>
1e-5	svmguide3	1243	22	1132.50	<b>410.63</b>	629.39			1484.23	<b>1.53</b>
1e-7	svmguide3	1243	22							
1e-1	phishing	11055	68	0.47	0.38	0.01	0.71	0.01	<b>0.01</b>	0.02
1e-3	phishing	11055	68	2.81	1.24	0.67	201.65	<b>0.32</b>	0.53	0.26
1e-5	phishing	11055	68	45.36	<b>16.10</b>	47.01		1391.81	43.70	<b>2.71</b>
1e-7	phishing	11055	68	2478.92	<b>812.38</b>	3370.12				<b>4.15</b>
1e-1	ijcnn1	49990	22	0.85	0.24	1.21	9.02	<b>0.22</b>	0.92	0.91
1e-3	ijcnn1	49990	22	4.15	<b>0.86</b>	54.72	565.14	193.34	71.65	<b>63.40</b>
1e-5	ijcnn1	49990	22	7.67	<b>1.66</b>		2322.76			<b>1402.31</b>
1e-7	ijcnn1	49990	22	10.20	<b>2.32</b>		3910.71			<b>1688.65</b>
1e-1	covtype	581012	54	54.62	29.76	143.55	123.79	<b>5.47</b>	18.91	0.18
1e-3	covtype	581012	54	552.47	<b>231.70</b>	2964.49		843.60	690.78	<b>2.98</b>
1e-5	covtype	581012	54		<b>3777.45</b>					

We next tested the methods on non-convex instances of  $\text{ERM}_\ell$ , using instances of the  $\ell_1$ -constrained binary classification problem with least-squares thresholding using the sigmoid function, namely:

$$(6.2) \quad \text{BCP:} \quad \min_x \frac{1}{n} \sum_{i=1}^n \left( y_i - \frac{1}{1 + \exp(-w_i^\top x)} \right)^2 \quad \text{s.t.} \quad \|x\|_1 \leq \lambda,$$

using the same 12 datasets from LIBSVM as described earlier. Recall that the Frank-Wolfe gap  $\mathcal{G}(x)$  is a measure of non-stationarity, see [15, 24, 29, 22]. In the non-convex setting the design of most Frank-Wolfe methods involves setting the number of iterations  $K$  in advance and using a fixed step-size (usually  $1/\sqrt{K+1}$ ), and the convergence rate for the Frank-Wolfe gap is typically measured using the expectation of  $\mathcal{G}(x^{\hat{k}})$  where  $\hat{k}$  is uniformly sampled over  $\{0, \dots, K\}$ . In our experiments we therefore compute the average Frank-Wolfe gap by computing  $\sum_{i=0}^K \mathcal{G}(x^i)/(K+1)$  and then averaging this value over 5 independent trials. In order to save on computation when  $K \geq 100n$ , we only compute and record the Frank-Wolfe gaps every 200 iterations. Just as in the convex case, we computed the Frank-Wolfe gaps separately and offline after-the-fact.

For each  $K \in \times\{10, 20, 10 \times 2^2, \dots, 10 \times 2^{20}\}$ , we ran 5 independent trials of the methods on the 12 datasets. Figure 2 shows the performance of TUFW with *Rule-SBD*  $\sqrt[4]{K}$  and *Rule-DBD*  $\sqrt[4]{K}$ , and the standard Frank-Wolfe method, on the dataset a9a, both with and without adaptive step-sizes. Each dot in the right-most subfigure corresponds to the average Frank-Wolfe gap, and the average runtime (over 5 trials) using a particular value of  $K$  as described above. Similar to the convex case, the TUFW methods significantly outperform the Frank-Wolfe method. *Rule-DBD*  $\sqrt[4]{K}$  using adaptive step-sizes exhibits the best performance over the other methods in terms of overall runtime and has nearly identical LMO dependence as *Rule-SBD*  $\sqrt[4]{K}$ . The second-best performance is from TUFW with *Rule-SBD*  $\sqrt[4]{K}$  with adaptive step-sizes. We focus on *Rule-SBD*  $\sqrt[4]{K}$  and *Rule-DBD*  $\sqrt[4]{K}$  with adaptive step-sizes in further comparisons with other Frank-Wolfe methods.

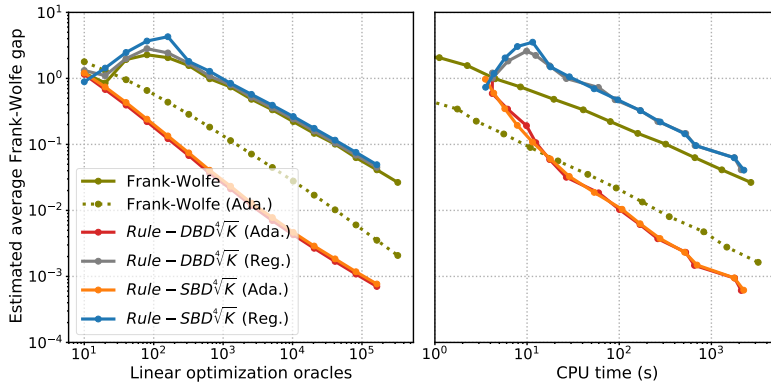


FIG. 2. Performance of TUFW with *Rule-SBD*  $\sqrt[4]{K}$  and *Rule-DBD*  $\sqrt[4]{K}$  and the standard Frank-Wolfe method, with and without adaptive step-sizes, on the binary classification problem (6.2), on the dataset a9a.

In Table 4 we compare the CPU runtimes (in seconds) of the TUFW methods using *Rule-SBD*  $\sqrt[4]{K}$  and *Rule-DBD*  $\sqrt[4]{K}$ , with FW, FW-ada, SPIDER-FW, and

CASPIDERG. CASPIDERG calculates full gradients every  $n^{1/4}K^{1/4}$  iterations and its batch-size is set to  $\max\{n^{3/4}K^{-1/4}, 1\}$ . (We did not include CSFW in these tests as it is not designed for non-convex objective functions.) All methods were run for each of the 21 different values of  $K$  and were terminated if they exceeded 5000 seconds. The numbers in parentheses in the table are the number of iterations  $K$  at which the given average Frank-Wolfe gap was attained. In each row of the table the bold number highlights the best runtime among the six methods. The last column of the table reports the speed-up of TUFW (the best of  $Rule-SBD \sqrt[4]{K}$  and  $Rule-DBD \sqrt[4]{K}$ ) over the best of the other four methods. Speed-up values larger than 1 are shown in bold. Notice from Table 4 that the TUFW method – using either  $Rule-SBD \sqrt[4]{K}$  or  $Rule-DBD \sqrt[4]{K}$  – significantly outperforms the other methods on the datasets a1a, a2a, a8a, and a9a, for sufficiently small values of  $\varepsilon$ . For the datasets w1a, w2a, w7a, and w8a, the TUFW methods are generally dominated by one of the other methods, with FW being the best on these datasets when  $\varepsilon$  is smaller. On the datasets svmguide3, phishing, ijcnn1, and covtype the TUFW method outperforms the other methods, all the more as the required  $\varepsilon$  gets smaller.

TABLE 4

Comparison of average CPU runtimes (in seconds) required to achieve  $\frac{1}{K+1} \sum_{k=0}^K \mathcal{G}(x^k) \leq \varepsilon$  for methods on the non-convex binary classification problem (6.2). The numbers in parentheses are the number of iterations  $K$  at which the given average Frank-Wolfe gap was attained. (A blank indicates the method used more than 5000 seconds.)

$\mathcal{G}(x^k)$	dataset	$n$	$p$	$Rule-SBD \sqrt[4]{K}$	$Rule-DBD \sqrt[4]{K}$	FW	FW-ada	SPIDER-FW	CASPIDERG	Speed-up
1e-2	a1a	1605	119	6.05(1.8e4)	<b>4.44(2.3e4)</b>	10.46(3.7e6)	9.00(5.3e6)	15.77(2.5e7)		<b>2.03</b>
1e-3	a1a	1605	119	90.34(1.6e5)	<b>65.11(1.6e5)</b>	818.40(4.2e7)	237.93(9.2e6)			<b>3.65</b>
1e-4	a1a	1605	119	2225.61(6.3e6)	<b>1453.87(6.3e6)</b>		4953.53(2.5e7)			<b>3.41</b>
1e-2	a2a	2265	119	6.47(2.3e4)	<b>4.94(1.6e4)</b>	12.62(5.2e6)	10.55(2.1e7)	13.54(2.1e7)		<b>2.14</b>
1e-3	a2a	2265	119	93.17(2.0e5)	<b>70.69(2.0e5)</b>	781.05(4.2e7)	303.67(2.1e7)			<b>4.30</b>
1e-4	a2a	2265	119	2168.72(6.3e6)	<b>1389.71(9.4e6)</b>					-
1e-2	a8a	22696	123	46.68(2.5e4)	41.92(2.5e4)	243.81(8.9e6)	140.10(8.7e6)	<b>35.70(4.2e7)</b>		0.85
1e-3	a8a	22696	123	668.35(1.8e5)	<b>603.88(2.5e5)</b>		3317.18(2.6e6)			<b>5.49</b>
1e-4	a8a	22696	123							-
1e-2	a9a	32561	123	83.24(1.4e4)	79.91(1.2e4)	358.08(6.3e6)	240.37(2.6e6)	<b>46.05(4.2e7)</b>		0.58
1e-3	a9a	32561	123	1165.35(1.3e5)	<b>1106.16(1.3e5)</b>					-
1e-4	a9a	32561	123							-
5e-2	w1a	2477	300	8.99(2.5e4)	8.82(2.9e4)	<b>0.44(1.8e4)</b>	12.66(1.2e6)	0.96(6.6e5)	101.83(4.2e7)	0.05
1e-2	w1a	2477	300	74.94(3.3e5)	102.18(5.2e5)	<b>4.69(8.2e4)</b>	157.90(5.1e6)	19.96(1.3e7)		0.06
2e-3	w1a	2477	300	1278.96(1.5e7)	4063.80(4.2e7)	<b>109.23(1.4e6)</b>	1768.32(1.2e7)			0.09
5e-2	w2a	3470	300	12.64(1.5e4)	11.90(1.0e4)	<b>0.50(1.0e4)</b>	17.41(7.0e6)	1.01(7.9e5)	120.51(3.8e7)	0.04
1e-2	w2a	3470	300	93.47(2.8e5)	122.01(2.9e5)	<b>7.44(8.2e4)</b>	226.70(8.0e6)	21.08(1.7e7)		0.08
2e-3	w2a	3470	300	900.09(3.1e6)	2545.77(4.2e7)	<b>152.15(2.4e6)</b>	2415.59(1.4e7)			0.17
5e-2	w7a	24692	300	95.86(1.4e4)	90.12(1.6e4)	7.79(1.2e4)	271.26(2.3e6)	<b>2.16(6.6e5)</b>	595.38(2.5e7)	0.02
1e-2	w7a	24692	300	544.57(1.6e5)	561.90(4.9e4)	80.96(4.9e4)	3596.39(2.0e7)	<b>29.75(8.4e6)</b>		0.05
2e-3	w7a	24692	300	4078.15(1.0e6)		<b>1631.47(6.6e5)</b>		2990.71(4.2e7)		0.40
5e-2	w8a	49749	300	222.71(2.3e4)	216.21(1.0e4)	16.41(1.4e4)	534.20(8.5e6)	<b>3.25(5.2e5)</b>	1122.68(2.7e7)	0.02
1e-2	w8a	49749	300	1292.84(4.1e4)	1303.19(4.1e4)	176.70(9.8e4)		<b>41.53(1.0e7)</b>		0.03
2e-3	w8a	49749	300			<b>3268.55(7.9e5)</b>				0.47
1e-1	svmguide3	1243	22	0.47(5.1e4)	<b>0.15(3.7e4)</b>	2.84(4.2e7)	1.96(6.3e6)			<b>13.39</b>
1e-2	svmguide3	1243	22	7.95(1.4e5)	<b>2.41(6.6e4)</b>		74.73(6.3e6)			<b>30.98</b>
1e-3	svmguide3	1243	22	122.45(1.0e6)	<b>33.68(1.0e6)</b>		2946.23(3.4e7)			<b>87.49</b>
1e-4	svmguide3	1243	22	2418.83(1.0e7)	<b>804.17(2.1e7)</b>					-
1e-1	phishing	11055	68	1.59(1.6e4)	1.17(1.6e4)	2.36(4.6e5)	102.48(1.4e7)	<b>1.17(4.7e6)</b>		0.99
1e-2	phishing	11055	68	17.53(2.9e4)	<b>12.18(3.5e4)</b>	158.49(1.5e7)	2506.03(2.2e7)			<b>13.01</b>
1e-3	phishing	11055	68	216.38(3.9e5)	<b>154.79(3.9e5)</b>					-
1e-4	phishing	11055	68	5049.56(1.0e7)	<b>3558.37(1.0e7)</b>					-
1e-1	ijcnn1	49990	22	2.32(8.2e3)	<b>0.96(9.2e3)</b>	10.88(6.6e5)	103.03(8.7e6)	1.58(3.7e6)	368.77(4.2e7)	<b>1.64</b>
1e-2	ijcnn1	49990	22	24.26(2.5e4)	<b>10.00(1.3e4)</b>	728.57(2.7e7)	2315.87(4.8e6)			<b>72.85</b>
1e-3	ijcnn1	49990	22	298.45(1.6e5)	<b>179.40(1.3e6)</b>					-
1e-4	ijcnn1	49990	22		<b>2750.09(5.2e6)</b>					-
5e-2	covtype	581012	54	156.41(4.5e6)	<b>110.28(5.8e6)</b>	2152.50(4.2e7)	2894.39(4.5e6)			<b>19.52</b>
1e-2	covtype	581012	54	785.43(4.7e6)	<b>572.97(5.8e6)</b>					-
2e-3	covtype	581012	54	4292.90(5.8e6)	<b>3142.24(5.8e6)</b>					-

## REFERENCES

- [1] K. BALASUBRAMANIAN AND S. GHADIMI, *Zeroth-order (non)-convex stochastic optimization via conditional gradient and gradient updates*, in Proceedings of the 32nd International Conference on Neural Information Processing Systems, 2018, pp. 3459–3468.
- [2] C.-C. CHANG AND C.-J. LIN, *Libsvm: a library for support vector machines*, ACM transactions on intelligent systems and technology (TIST), 2 (2011), p. 27.
- [3] V. DEMYANOV AND A. RUBINOV, *The minimization of a smooth convex functional on a convex set*, SIAM Journal on Control, (1967).
- [4] C. FANG, C. J. LI, Z. LIN, AND T. ZHANG, *Spider: near-optimal non-convex optimization via stochastic path integrated differential estimator*, Advances in neural information processing systems, 32 (2018), pp. 687–697.
- [5] M. FAZEL, *Matrix rank minimization with applications*, PhD thesis, PhD thesis, Stanford University, 2002.
- [6] M. FRANK, P. WOLFE, ET AL., *An algorithm for quadratic programming*, Naval research logistics quarterly, 3 (1956), pp. 95–110.
- [7] R. M. FREUND, P. GRIGAS, AND R. MAZUMDER, *An extended Frank–Wolfe method with “in-face” directions, and its application to low-rank matrix completion*, SIAM Journal on Optimization, 27 (2017), pp. 319–346.
- [8] H. GAO AND H. HUANG, *Can stochastic zeroth-order frank-wolfe method converge faster for non-convex problems?*, in International Conference on Machine Learning, PMLR, 2020, pp. 3377–3386.
- [9] H. HASSANI, A. KARBASI, A. MOKHTARI, AND Z. SHEN, *Stochastic conditional gradient++:(non) convex minimization and continuous submodular maximization*, SIAM Journal on Optimization, 30 (2020), pp. 3315–3344.
- [10] T. HASTIE, R. TIBSHIRANI, J. H. FRIEDMAN, AND J. H. FRIEDMAN, *The elements of statistical learning: data mining, inference, and prediction*, vol. 2, Springer, 2009.
- [11] E. HAZAN AND H. LUO, *Variance-reduced and projection-free stochastic optimization*, in International Conference on Machine Learning, PMLR, 2016, pp. 1263–1271.
- [12] F. HUANG, L. TAO, AND S. CHEN, *Accelerated stochastic gradient-free and projection-free methods*, in International Conference on Machine Learning, PMLR, 2020, pp. 4519–4530.
- [13] M. JAGGI, *Revisiting frank-wolfe: Projection-free sparse convex optimization*, in International Conference on Machine Learning, PMLR, 2013, pp. 427–435.
- [14] R. JOHNSON AND T. ZHANG, *Accelerating stochastic gradient descent using predictive variance reduction*, Advances in neural information processing systems, 26 (2013), pp. 315–323.
- [15] S. LACOSTE-JULIEN, *Convergence rate of frank-wolfe for non-convex objectives*, arXiv preprint arXiv:1607.00345, (2016).
- [16] G. LAN, *The complexity of large-scale convex programming under a linear optimization oracle*, arXiv preprint arXiv:1309.5550, (2013).
- [17] G. LAN AND Y. ZHOU, *Conditional gradient sliding for convex optimization*, SIAM Journal on Optimization, 26 (2016), pp. 1379–1409.
- [18] H. LU AND R. M. FREUND, *Generalized stochastic frank–wolfe algorithm with stochastic “substitute” gradient for structured convex optimization*, Mathematical Programming, 187 (2021), pp. 317–349.
- [19] J. MAIRAL, R. JENATTON, G. OBOZINSKI, AND F. BACH, *Convex and network flow optimization for structured sparsity.*, Journal of Machine Learning Research, 12 (2011).
- [20] S. MEI, Y. BAI, AND A. MONTANARI, *The landscape of empirical risk for nonconvex losses*, The Annals of Statistics, 46 (2018), pp. 2747–2774.
- [21] A. MOKHTARI, H. HASSANI, AND A. KARBASI, *Stochastic conditional gradient methods: From convex minimization to submodular maximization*, Journal of machine learning research, (2020).
- [22] G. NEGIAR, G. DRESDNER, A. TSAI, L. E. GHAOUI, F. LOCATELLO, R. FREUND, AND F. PEDREGOSA, *Stochastic frank-Wolfe for constrained finite-sum minimization*, in Proceedings of the 37th International Conference on Machine Learning, vol. 119 of Proceedings of Machine Learning Research, PMLR, 13–18 Jul 2020, pp. 7253–7262.
- [23] C. QU, Y. LI, AND H. XU, *Non-convex conditional gradient sliding*, in International Conference on Machine Learning, PMLR, 2018, pp. 4208–4217.
- [24] S. J. REDDI, S. SRA, B. PÓCZOS, AND A. SMOLA, *Stochastic frank-wolfe methods for nonconvex optimization*, in 2016 54th Annual Allerton Conference on Communication, Control, and Computing (Allerton), IEEE, 2016, pp. 1244–1251.
- [25] A. K. SAHU, M. ZAHEER, AND S. KAR, *Towards gradient free and projection free stochastic optimization*, in The 22nd International Conference on Artificial Intelligence and Statistics,

- PMLR, 2019, pp. 3468–3477.
- [26] Z. SHEN, C. FANG, P. ZHAO, J. HUANG, AND H. QIAN, *Complexities in projection-free stochastic non-convex minimization*, in The 22nd International Conference on Artificial Intelligence and Statistics, PMLR, 2019, pp. 2868–2876.
- [27] R. TIBSHIRANI, *Regression shrinkage and selection via the lasso*, Journal of the Royal Statistical Society: Series B (Methodological), 58 (1996), pp. 267–288.
- [28] V. VAPNIK, *Principles of risk minimization for learning theory*, Advances in neural information processing systems, 4 (1991).
- [29] A. YURTSEVER, S. SRA, AND V. CEVHER, *Conditional gradient methods via stochastic path-integrated differential estimator*, in International Conference on Machine Learning, PMLR, 2019, pp. 7282–7291.
- [30] M. ZHANG, Z. SHEN, A. MOKHTARI, H. HASSANI, AND A. KARBASI, *One sample stochastic frank-wolfe*, in International Conference on Artificial Intelligence and Statistics, PMLR, 2020, pp. 4012–4023.