

---

# Mixed Optimization for Smooth Functions

---

Mehrdad Mahdavi    Lijun Zhang    Rong Jin

Department of Computer Science, Michigan State University, East Lansing, MI, USA  
 {fmahdavi, zhanglij, rongjing}@msu.edu

## Abstract

In this paper, we propose a new algorithm for minimizing a smooth function over a convex set. The algorithm achieves a convergence rate of  $O(1/T)$  for the smooth case and  $O(1/T^2)$  for the strongly smooth case. The algorithm is a mixed optimization algorithm, which combines the advantages of the gradient method and the proximal method. The algorithm achieves a convergence rate of  $O(1/T)$  for the smooth case and  $O(1/T^2)$  for the strongly smooth case. The algorithm is a mixed optimization algorithm, which combines the advantages of the gradient method and the proximal method.

## 1 Introduction

Minimizing a smooth function over a convex set is a fundamental problem in optimization. The gradient method is a classic algorithm for minimizing a smooth function over a convex set. The proximal method is another classic algorithm for minimizing a smooth function over a convex set. In this paper, we propose a new algorithm for minimizing a smooth function over a convex set, which combines the advantages of the gradient method and the proximal method.

$$\min_{\mathbf{w} \in \mathcal{W}} G(\mathbf{w}) := \frac{1}{n} \sum_{i=1}^n g_i(\mathbf{w}); \quad (1)$$

where  $n$  is the number of data points,  $\mathbf{x}_i, y_i$  are the data points,  $\mathcal{W}$  is the feasible set,  $g_i(\mathbf{w}) = \log(1 + \exp(\mathbf{w}^\top \mathbf{x}_i - y_i))$  is the loss function, and  $G(\mathbf{w})$  is the average loss function. The algorithm achieves a convergence rate of  $O(1/T)$  for the smooth case and  $O(1/T^2)$  for the strongly smooth case. The algorithm is a mixed optimization algorithm, which combines the advantages of the gradient method and the proximal method.

A simple gradient method achieves a convergence rate of  $O(1/T)$  for the smooth case and  $O(1/T^2)$  for the strongly smooth case. The proximal method achieves a convergence rate of  $O(1/T)$  for the smooth case and  $O(1/T^2)$  for the strongly smooth case. The algorithm achieves a convergence rate of  $O(1/T)$  for the smooth case and  $O(1/T^2)$  for the strongly smooth case. The algorithm is a mixed optimization algorithm, which combines the advantages of the gradient method and the proximal method.

The algorithm achieves a convergence rate of  $O(1/T)$  for the smooth case and  $O(1/T^2)$  for the strongly smooth case. The algorithm is a mixed optimization algorithm, which combines the advantages of the gradient method and the proximal method.



## 2 More Related Work

## Deterministic Smooth Optimization

**Deterministic Smooth Optimization**

GD

5. F

16. A

18, 17. I

GD

15, 18, 16.

$O(1=T)$

$O(1=T^2)$

Stochastic Smooth Optimization 

**Stochastic Smooth Optimization**

### 3 Preliminaries

$$\begin{aligned} & \mathcal{Q} \quad \mathbb{M}^{\dagger} c \quad \mathbb{M} \quad c \quad F \quad \dagger \quad c \quad w; w' \quad 2W, \quad \mathbb{M} \quad hw; w' i \\ & \mathbb{M}^{\dagger} c \quad w \dagger \quad \mathbb{M} w'. \quad \mathcal{Q} \quad \mathcal{Q} \quad \dagger \quad c \quad \mathbb{M} \quad 2- \\ & \dagger \quad \mathcal{Q} \quad \dagger \quad \mathbb{M}^{\dagger} c \quad \mathcal{Q} \quad c \quad G(w) \quad \mathbb{M} \text{ fi} \quad \mathbb{M} \quad (1) \quad \dagger \quad \dagger \quad c \quad n c \quad \mathcal{Q} \quad c \quad . \\ & \mathcal{Q} \quad c \quad \dagger \quad \dagger \quad \mathbb{M} \quad 20, 22. \quad \dagger \quad \mathcal{Q} \quad \dagger \quad \dagger \quad \mathcal{Q} \quad \dagger \quad G(w) \quad \dagger \quad \mathbb{M}^{\dagger} c \quad w_* \quad 2W. \quad - \\ & \mathcal{Q} \quad c \quad \dagger \quad \dagger \quad \mathcal{Q} \quad \dagger \quad \mathcal{Q} \quad \dagger \quad W \quad \mathbb{M}_{R, \dagger} \quad \dagger \quad \dagger \quad \mathcal{Q} \quad R. B \quad \mathbb{M} \quad c \quad \mathbb{M} \quad \mathbb{M}^{\dagger} \\ & \mathcal{Q} \quad c \quad \dagger \quad \dagger \quad \mathcal{Q} \quad \dagger \quad \dagger \quad c \quad g_i(w) \quad \text{-smooth} \quad \dagger \quad \mathbb{M} \text{ fi} \quad \mathbb{M} \quad 16. \end{aligned}$$

**Definition 1** (Lipschitz). A differentiable loss function  $f(\mathbf{w})$  is said to be  $L$ -smooth with respect to a norm  $\|\cdot\|_k$ , if it holds that

$$f(\mathbf{w}) = f(\mathbf{w}') + hr f(\mathbf{w}'); \mathbf{w} = \mathbf{w}'i + \frac{1}{2}k\mathbf{w} = \mathbf{w}'k^2; \quad \delta \mathbf{w}; \mathbf{w}' \in W;$$

$$\begin{aligned} & \frac{\partial}{\partial w} r f(w) - L \leq \frac{\partial}{\partial w} h r f(w) + \frac{\partial}{\partial w} r f(w'); w' = w + i k w \\ & \quad G(w), \quad r g_i(w) \\ & \quad T, \quad G(\bar{w}), \quad G(w^*), \quad \text{etc.} \\ & \quad O_f, \quad O(\log T), \quad O(T), \quad O(1=T) \end{aligned}$$
[illegible]

---

**Algorithm 1** MI-EDG-AD
 

---

**Input:**  $\epsilon, m, \eta, \Delta_1, \eta, T_1$  s.t.  $\eta > 1$

```

1:  $\bar{\mathbf{w}}_1 = \mathbf{0}$ 
2: for  $k = 1; \dots; m$  do
3:    $\mathcal{W}_k = \{\mathbf{w} : \mathbf{w} + \bar{\mathbf{w}}_k \in \mathcal{W}; \|\mathbf{w}\| \leq \Delta_k\}$ 
4:    $\mathcal{O}_f = \{r G(\bar{\mathbf{w}}_k)\}$ 
5:    $\mathcal{G}_k = \{\bar{\mathbf{w}}_k + r G(\bar{\mathbf{w}}_k) = \bar{\mathbf{w}}_k + \frac{1}{n} \sum_{i=1}^n r g_i(\bar{\mathbf{w}}_k)\}$ 
6:    $\mathbf{w}_k^1 = \mathbf{0}$ 
7:   for  $t = 1; \dots; T_k$  do
8:      $\mathcal{O}_s = \{g_{i_k}^t(\mathbf{w})\}$ 
9:      $\hat{\mathbf{g}}_k^t = \mathbf{g}_k + r g_{i_k}^t(\mathbf{w}_k^t + \bar{\mathbf{w}}_k) - r g_{i_k}^t(\bar{\mathbf{w}}_k)$ 
10:     $\mathbf{w}_k^{t+1} = \arg \max_{\mathbf{w} \in \mathcal{W}_k} \langle \mathbf{w}, \hat{\mathbf{g}}_k^t \rangle + \frac{1}{2} k \|\mathbf{w} - \mathbf{w}_k^t\|^2$ 
11:   end for
12:    $\tilde{\mathbf{w}}_{k+1} = \frac{1}{T_k+1} \sum_{t=1}^{T_k+1} \mathbf{w}_k^t$ ;  $\bar{\mathbf{w}}_{k+1} = \bar{\mathbf{w}}_k + \tilde{\mathbf{w}}_{k+1}$ 
13:    $\Delta_{k+1} = \Delta_k$ ,  $\eta_{k+1} = \eta$ ,  $T_{k+1} = \eta^2 T_k$ 
14: end for
Return  $\bar{\mathbf{w}}_{m+1}$ 

```

---

**Definition 2** ( $\eta$ -strongly convex). A function  $f(\mathbf{w})$  is said to be  $\eta$ -strongly convex w.r.t a norm  $\|\cdot\|_k$ , if there exists a constant  $\eta > 0$  (often called the modulus of strong convexity) such that it holds

$$f(\mathbf{w}) \leq f(\mathbf{w}') + \langle \mathbf{w} - \mathbf{w}', i \rangle + \frac{1}{2} k \|\mathbf{w} - \mathbf{w}'\|_k^2; \quad \forall \mathbf{w}, \mathbf{w}' \in \mathcal{W}$$

#### 4 Mixed Stochastic/Deterministic Gradient Descent

Algorithm 1 (MI-EDG-AD) is a mixed stochastic/deterministic gradient descent algorithm. It consists of two main parts: a deterministic part (lines 1-14) and a stochastic part (lines 15-20). The deterministic part starts with  $\bar{\mathbf{w}}_1 = \mathbf{0}$  and iterates over  $k = 1, \dots, m$ . In each iteration, it defines a set  $\mathcal{W}_k$  and a set of subgradients  $\mathcal{O}_f$ . It then iterates over  $t = 1, \dots, T_k$  to find  $\mathbf{w}_k^{t+1}$  by minimizing a quadratic function. The stochastic part (lines 15-20) involves sampling subgradients  $g_{i_k}^t(\mathbf{w})$  and updating  $\bar{\mathbf{w}}_{k+1}$  and  $\Delta_{k+1}$ .

$$\min_{\substack{\mathbf{w} + \bar{\mathbf{w}}_k \in \mathcal{W} \\ \|\mathbf{w}\| \leq \Delta_k}} \left[ \frac{k}{2} \|\mathbf{w} - \bar{\mathbf{w}}_k\|_k^2 + \frac{1}{n} \sum_{i=1}^n g_i(\mathbf{w} + \bar{\mathbf{w}}_k) \right]; \quad (2)$$

Let  $\Delta_k$  be the radius of  $\mathcal{W}_k$ . By (2), we have  $\mathbf{w}_k^{t+1} = \arg \min_{\mathbf{w} \in \mathcal{W}_k} \left[ \frac{k}{2} \|\mathbf{w} - \bar{\mathbf{w}}_k\|_k^2 + \frac{1}{n} \sum_{i=1}^n g_i(\mathbf{w} + \bar{\mathbf{w}}_k) \right]$ . By (2), we have  $\mathbf{w}_k^{t+1} = \arg \min_{\mathbf{w} \in \mathcal{W}_k} \left[ \frac{k}{2} \|\mathbf{w} - \bar{\mathbf{w}}_k\|_k^2 + \frac{1}{n} \sum_{i=1}^n g_i(\mathbf{w} + \bar{\mathbf{w}}_k) \right]$ . By (2), we have  $\mathbf{w}_k^{t+1} = \arg \min_{\mathbf{w} \in \mathcal{W}_k} \left[ \frac{k}{2} \|\mathbf{w} - \bar{\mathbf{w}}_k\|_k^2 + \frac{1}{n} \sum_{i=1}^n g_i(\mathbf{w} + \bar{\mathbf{w}}_k) \right]$ .

$$\min_{\mathbf{w} \in \mathcal{W}_k} \left[ F_k(\mathbf{w}) = \frac{k}{2} \|\mathbf{w} - \bar{\mathbf{w}}_k\|_k^2 + \frac{1}{n} \sum_{i=1}^n g_i(\mathbf{w} + \bar{\mathbf{w}}_k) \right]; \quad (3)$$

$$\begin{aligned}
W_k &= f\mathbf{w} : \mathbf{w} + \mathbf{w}_k \in \mathcal{W}; \|\mathbf{w}\| \leq \Delta_k \mathcal{G}. \\
F_k(\mathbf{w}) &= \frac{k}{2} \|\mathbf{w}\|^2 + \langle \mathbf{h}(\mathbf{w}; \bar{\mathbf{w}}_k), \mathbf{w} \rangle + \frac{1}{n} \sum_{i=1}^n g_i(\mathbf{w} + \bar{\mathbf{w}}_k)
\end{aligned}$$

## 5 Convergence Analysis

$$\begin{aligned} & \mathbb{E} \left[ \|\mathbf{w}_*^k - \mathbf{w}_*^{k+1}\|^2 \right] \leq \mathbb{E} \left[ \|\mathbf{w}_*^k - \mathbf{w}_*^{k+1}\|^2 \right] + \mathbb{E} \left[ \|\mathbf{w}_*^{k+1} - \mathbf{w}_*^{k+1}\|^2 \right] \\ & \leq \mathbb{E} \left[ \|\mathbf{w}_*^k - \mathbf{w}_*^{k+1}\|^2 \right] + \mathbb{E} \left[ \|\mathbf{w}_*^{k+1} - \mathbf{w}_*^{k+1}\|^2 \right] \\ & \leq \mathbb{E} \left[ \|\mathbf{w}_*^k - \mathbf{w}_*^{k+1}\|^2 \right] + \mathbb{E} \left[ \|\mathbf{w}_*^{k+1} - \mathbf{w}_*^{k+1}\|^2 \right] \\ & \leq \mathbb{E} \left[ \|\mathbf{w}_*^k - \mathbf{w}_*^{k+1}\|^2 \right] + \mathbb{E} \left[ \|\mathbf{w}_*^{k+1} - \mathbf{w}_*^{k+1}\|^2 \right] \end{aligned} \quad (3)$$

**Theorem 2.** Let  $\hat{\mathbf{w}}_*^k$  and  $\hat{\mathbf{w}}_*^{k+1}$  be the optimal solutions that minimize  $F_k(\mathbf{w})$  and  $F_{k+1}(\mathbf{w})$ , respectively, and  $\tilde{\mathbf{w}}_{k+1}$  be the average solution obtained at the end of  $k$ th epoch of MI-EDG AD algorithm. Suppose  $\|\hat{\mathbf{w}}_*^k - \hat{\mathbf{w}}_*^{k+1}\| \leq \Delta_k$ . By setting the step size  $\eta_k = 1/(2^{k/2} \sqrt{3T_k})$ , we have, with a probability  $1 - 2^{-k}$ ,

$$\|\hat{\mathbf{w}}_*^{k+1} - \hat{\mathbf{w}}_*^k\| \leq \Delta_k \quad \text{and} \quad F_k(\tilde{\mathbf{w}}_{k+1}) - \min_{\mathbf{w}} F_k(\mathbf{w}) \leq \frac{k\Delta_k^2}{2^{k/2}}$$

provided that  $\eta_k \leq e^{-9/2}$  and

$$T_k \leq \frac{300 \cdot 8 \cdot 2}{k} \ln \frac{1}{\epsilon}.$$

$$\mathbb{E} \left[ \|\mathbf{w}_*^k - \mathbf{w}_*^{k+1}\|^2 \right] \leq \mathbb{E} \left[ \|\mathbf{w}_*^k - \mathbf{w}_*^{k+1}\|^2 \right] + \mathbb{E} \left[ \|\mathbf{w}_*^{k+1} - \mathbf{w}_*^{k+1}\|^2 \right] \leq 1, \quad \eta_k \leq e^{-9/2}$$

*Proof of Theorem 1.* I  $\mathbf{w}_*^k \in \mathcal{W} \subset \mathbb{B}_R$ ,  $\|\mathbf{w}_*^k - \mathbf{w}_*^{k+1}\| \leq R := \Delta_1$ .

$$\begin{aligned} & \mathbb{E} \left[ \|\mathbf{w}_*^m - \mathbf{w}_*^{m+1}\|^2 \right] \leq \mathbb{E} \left[ \|\mathbf{w}_*^m - \mathbf{w}_*^{m+1}\|^2 \right] + \mathbb{E} \left[ \|\mathbf{w}_*^{m+1} - \mathbf{w}_*^{m+1}\|^2 \right] \\ & \leq \mathbb{E} \left[ \|\mathbf{w}_*^m - \mathbf{w}_*^{m+1}\|^2 \right] + \mathbb{E} \left[ \|\mathbf{w}_*^{m+1} - \mathbf{w}_*^{m+1}\|^2 \right] \\ & \leq \mathbb{E} \left[ \|\mathbf{w}_*^m - \mathbf{w}_*^{m+1}\|^2 \right] + \mathbb{E} \left[ \|\mathbf{w}_*^{m+1} - \mathbf{w}_*^{m+1}\|^2 \right] \end{aligned}$$

$$\|\hat{\mathbf{w}}_*^m - \hat{\mathbf{w}}_*^{m+1}\| \leq \frac{\Delta_1}{m-1}; \quad F_m(\tilde{\mathbf{w}}_{m+1}) - F_m(\hat{\mathbf{w}}_*^m) \leq \frac{m\Delta_m^2}{2^{m/2}} = \frac{1\Delta_1^2}{2^{3m+1}}$$

$$\mathbb{E} \left[ \|\mathbf{w}_*^m - \mathbf{w}_*^{m+1}\|^2 \right] \leq \mathbb{E} \left[ \|\mathbf{w}_*^m - \mathbf{w}_*^{m+1}\|^2 \right] + \mathbb{E} \left[ \|\mathbf{w}_*^{m+1} - \mathbf{w}_*^{m+1}\|^2 \right] \leq \mathbb{E} \left[ \|\mathbf{w}_*^m - \mathbf{w}_*^{m+1}\|^2 \right] + \mathbb{E} \left[ \|\mathbf{w}_*^{m+1} - \mathbf{w}_*^{m+1}\|^2 \right]$$

$$\frac{1}{n} \sum_{i=1}^n g_i(\tilde{\mathbf{w}}_{m+1}) \leq F_m(\hat{\mathbf{w}}_*^m) + \frac{1\Delta_1^2}{2^{3m+1}} - \frac{1}{m-1} \langle \tilde{\mathbf{w}}_{m+1}, \tilde{\mathbf{w}}_m \rangle$$

$$F_m(\hat{\mathbf{w}}_*^m) + \frac{1\Delta_1^2}{2^{3m+1}} + \frac{1k\tilde{\mathbf{w}}_mk\Delta_1}{2^{m-2}}$$

$$\mathbb{E} \left[ \|\mathbf{w}_*^m - \mathbf{w}_*^{m+1}\|^2 \right] \leq \mathbb{E} \left[ \|\mathbf{w}_*^m - \mathbf{w}_*^{m+1}\|^2 \right] + \mathbb{E} \left[ \|\mathbf{w}_*^{m+1} - \mathbf{w}_*^{m+1}\|^2 \right] \leq \mathbb{E} \left[ \|\mathbf{w}_*^m - \mathbf{w}_*^{m+1}\|^2 \right] + \mathbb{E} \left[ \|\mathbf{w}_*^{m+1} - \mathbf{w}_*^{m+1}\|^2 \right]$$

$$k\tilde{\mathbf{w}}_mk \leq \sum_{i=1}^m j\tilde{\mathbf{w}}_ij \leq \sum_{i=1}^m \Delta_i \leq \frac{\Delta_1}{1} \leq 2\Delta_1$$

$$\mathbb{E} \left[ \|\mathbf{w}_*^m - \mathbf{w}_*^{m+1}\|^2 \right] \leq \mathbb{E} \left[ \|\mathbf{w}_*^m - \mathbf{w}_*^{m+1}\|^2 \right] + \mathbb{E} \left[ \|\mathbf{w}_*^{m+1} - \mathbf{w}_*^{m+1}\|^2 \right] \leq \mathbb{E} \left[ \|\mathbf{w}_*^m - \mathbf{w}_*^{m+1}\|^2 \right] + \mathbb{E} \left[ \|\mathbf{w}_*^{m+1} - \mathbf{w}_*^{m+1}\|^2 \right]$$

$$\frac{1}{n} \sum_{i=1}^n g_i(\tilde{\mathbf{w}}_{m+1}) \leq F_m(\hat{\mathbf{w}}_*^m) + \frac{1\Delta_1^2}{2^{3m+1}} + \frac{2 \cdot 1\Delta_1^2}{2^{m-2}}.$$

$$\mathbb{E} \left[ \|\mathbf{w}_*^m - \mathbf{w}_*^{m+1}\|^2 \right] \leq \mathbb{E} \left[ \|\mathbf{w}_*^m - \mathbf{w}_*^{m+1}\|^2 \right] + \mathbb{E} \left[ \|\mathbf{w}_*^{m+1} - \mathbf{w}_*^{m+1}\|^2 \right] \leq \mathbb{E} \left[ \|\mathbf{w}_*^m - \mathbf{w}_*^{m+1}\|^2 \right] + \mathbb{E} \left[ \|\mathbf{w}_*^{m+1} - \mathbf{w}_*^{m+1}\|^2 \right]$$

$$F_m(\mathbf{w}_*^m) - F_m(\mathbf{w}_*) = \frac{1}{n} \sum_{i=1}^n g_i(\mathbf{w}_*) + \frac{1}{2^{m-1}} (k\mathbf{w}_* - \tilde{\mathbf{w}}_mk^2 + 2\langle \mathbf{w}_*, \tilde{\mathbf{w}}_m \rangle) : \quad (6)$$

$$\mathbb{E} \left[ \|\mathbf{w}_*^m - \mathbf{w}_*^{m+1}\|^2 \right] \leq \mathbb{E} \left[ \|\mathbf{w}_*^m - \mathbf{w}_*^{m+1}\|^2 \right] + \mathbb{E} \left[ \|\mathbf{w}_*^{m+1} - \mathbf{w}_*^{m+1}\|^2 \right] \leq \mathbb{E} \left[ \|\mathbf{w}_*^m - \mathbf{w}_*^{m+1}\|^2 \right] + \mathbb{E} \left[ \|\mathbf{w}_*^{m+1} - \mathbf{w}_*^{m+1}\|^2 \right]$$

$$\begin{aligned} & \mathbb{E} \left[ \|\mathbf{w}_*^m - \mathbf{w}_*^{m+1}\|^2 \right] \leq \mathbb{E} \left[ \|\mathbf{w}_*^m - \mathbf{w}_*^{m+1}\|^2 \right] + \mathbb{E} \left[ \|\mathbf{w}_*^{m+1} - \mathbf{w}_*^{m+1}\|^2 \right] \\ & \leq \mathbb{E} \left[ \|\mathbf{w}_*^m - \mathbf{w}_*^{m+1}\|^2 \right] + \mathbb{E} \left[ \|\mathbf{w}_*^{m+1} - \mathbf{w}_*^{m+1}\|^2 \right] \\ & \leq \mathbb{E} \left[ \|\mathbf{w}_*^m - \mathbf{w}_*^{m+1}\|^2 \right] + \mathbb{E} \left[ \|\mathbf{w}_*^{m+1} - \mathbf{w}_*^{m+1}\|^2 \right] \end{aligned}$$

$$k\mathbf{w}_* - \tilde{\mathbf{w}}_mk \leq \sum_{i=m+1}^{\infty} j\tilde{\mathbf{w}}_ij \leq \sum_{k=m+1}^{\infty} \Delta_k \leq \frac{\Delta_1}{m(1-1)} \leq \frac{2\Delta_1}{m}$$



$$\begin{aligned}
& F(\mathbf{w}_*) - F(\mathbf{w}_*) + \frac{\beta}{2} k \mathbf{w}_*^2 \leq \max(k \mathbf{w}_*^2; k \mathbf{w}_{T+1}^2) \leq \Delta, \\
& \mathbb{E} \sum_{t=1}^{T+1} F(\mathbf{w}_t) - F(\widehat{\mathbf{w}}_*) \leq \Delta^2 \left( \frac{1}{2} + \frac{1}{2} \right) + \frac{1}{2} A_T + B_T + C_T. \tag{10}
\end{aligned}$$

**Lemma 2.** For  $A_T$  defined above we have  $A_T \leq 6 \Delta^2 T$ .

**Lemma 3.** With a probability  $1 - 2^{-\beta}$ , we have

$$B_T \leq \Delta^2 \left( \ln \frac{1}{\delta} + \sqrt{2T \ln \frac{1}{\delta}} \right) \text{ and } C_T \leq 2 \Delta^2 \left( \ln \frac{1}{\delta} + \sqrt{2T \ln \frac{1}{\delta}} \right).$$

$$\begin{aligned}
& \sum_{t=1}^{T+1} F(\mathbf{w}_t) - F(\widehat{\mathbf{w}}_*) \leq \Delta^2 \left( \frac{1}{2} + \frac{1}{2} + 6 \Delta^2 T + 3 \ln \frac{1}{\delta} + 3 \sqrt{2T \ln \frac{1}{\delta}} \right) \\
& = 1 + 2 \Delta^2 T, \\
& \sum_{t=1}^{T+1} F(\mathbf{w}_t) - F(\widehat{\mathbf{w}}_*) \leq \Delta^2 \left( 2 \Delta^2 T + 3 \ln \frac{1}{\delta} + 3 \sqrt{2T \ln \frac{1}{\delta}} \right) \\
& \leq \Delta^2 \frac{5 \sqrt{3 \ln[1/\delta]}}{T+1}; \quad \widehat{\Delta}^2 = k \widehat{\mathbf{w}}_*^2 \leq \Delta^2 \frac{5 \sqrt{3 \ln[1/\delta]}}{T+1}. \\
& T \leq [300 \delta^{-2} \ln \frac{1}{\delta}] = 2, \quad \widehat{\Delta}^2 \leq \frac{\Delta^2}{4}; \quad \mathbb{E} F(\widehat{\mathbf{w}}) - F(\widehat{\mathbf{w}}_*) \leq \frac{1}{2} \Delta^2. \tag{11}
\end{aligned}$$

**Lemma 4.** We have  $k \widehat{\mathbf{w}}_*' k \leq k \widehat{\mathbf{w}}_* k$ .

C  $\leq \Delta^2$  (11)  $\leq 4 \Delta^2$ ,  $k \widehat{\mathbf{w}}_*' k \leq \Delta^2$ .

## 6 Conclusions and Open Questions

MI EDG AD  $O(1=T)$   $O(T)$   $O(\log T)$   $smoothness$   $O(\ln T)$   $O(\log T)$   $O(1=T^2)$

**Acknowledgments.** ON A  $\Delta$  N000141210431  $\Delta$  F (II -1251031).



## References

- 1 A. A. Agresti, P. L. Bartlett, P. D. Bickel, and M. J. Wainwright. *IEEE Transactions on Information Theory*, 58(5):3235–3249, 2012.
- 2 A. B. Berger, M. J. Wainwright, and M. J. Wainwright. *Oper. Res. Lett.*, 31(3):167–175, 2003.
- 3 L. B. Berger, O. B. Berger, and M. J. Wainwright. *NIPS*, 161–168, 2008.
- 4 B. Berger, G. L. Berger, O. B. Berger, C. Berger, and M. J. Wainwright. *Advanced Lectures on Machine Learning*, 208–240, 2003.
- 5 B. Berger, L. O. Berger, and M. J. Wainwright. *Convex Optimization*. Cambridge University Press, 2004.
- 6 H. B. Berger, G. M. C. Berger, J. N. C. Berger, and M. J. Wainwright. *Mathematical programming*, 134(1):127–155, 2012.
- 7 A. C. Berger, O. Berger, N. Berger, K. Berger, B. Berger, and I. Berger. *NIPS*, 1647–1655, 2011.
- 8 O. D. Berger, G. L. Berger, O. Berger, J. L. Berger, O. Berger, and M. J. Wainwright. *The Journal of Machine Learning Research*, 13:165–202, 2012.
- 9 M. P. F. Berger, M. J. Wainwright, H. Berger, and M. J. Wainwright. *SIAM Journal on Scientific Computing*, 34(3):A1380–A1405, 2012.
- 10 E. H. Berger, A. A. Berger, J. K. Berger, L. Berger, and M. J. Wainwright. *Machine Learning*, 69(2-3):169–192, 2007.
- 11 E. H. Berger, J. K. Berger, B. Berger, and M. J. Wainwright. *Journal of Machine Learning Research - Proceedings Track*, 19:421–436, 2011.
- 12 L. Berger, C. Berger, J. P. Berger, A. Berger, and M. J. Wainwright. *arXiv preprint arXiv:1008.5204*, 2010.
- 13 A. N. Berger, A. Berger, G. L. Berger, A. Berger, and M. J. Wainwright. *SIAM J. on Optimization*, 19:1574–1609, 2009.
- 14 A. N. Berger, D. B. Berger, P. Berger, and M. J. Wainwright. *Math. Program.*, 1983.
- 15 N. Berger, A. Berger, and M. J. Wainwright. *Soviet Mathematics Doklady*, 27, 372–376, 1983.
- 16 N. Berger. *Introductory Lectures on Convex Optimization: A Basic Course*. Kluwer Academic Publishers, 2004.
- 17 N. Berger, E. Berger, and M. J. Wainwright. *SIAM Journal on Optimization*, 16(1):235–249, 2005.
- 18 N. Berger. *Math. Program.*, 103(1):127–152, 2005.
- 19 A. Berger, O. Berger, K. Berger, M. Berger, and M. J. Wainwright. *ICML*, 2012.
- 20 N. L. Berger, M. Berger, J. F. Berger, A. Berger, and M. J. Wainwright. *NIPS*, 2672–2680, 2012.
- 21 N. Berger, J. N. Berger, P. Berger, and M. J. Wainwright. *ICML*, 807–814, 2007.
- 22 N. Berger, J. N. Berger, J. F. Berger, and M. J. Wainwright. *JMLR*, 14:567599, 2013.
- 23 O. Berger, J. N. Berger, J. F. Berger, and M. J. Wainwright. *ICML*, 2013.
- 24 L. Berger, J. N. Berger, J. F. Berger, H. Berger, and M. J. Wainwright. *ICML*, 2013.