
Mixed Optimization for Smooth Functions

Mehrdad Mahdavi Lijun Zhang Rong Jin

Department of Computer Science and Engineering, Michigan State University, MI, USA
 mahdavi, zhanglij, rongjin@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 convex case. The algorithm is a mixed optimization algorithm, which combines the advantages of gradient descent and the proximal method. The algorithm achieves a convergence rate of $O(\ln T)$ for the smooth case and $O(1/T)$ for the strongly convex case.

1 Introduction

Minimizing a smooth function over a convex set is a fundamental problem in optimization. In this paper, we propose a new algorithm for minimizing a smooth function over a convex set.

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

where n is the number of data points, (x_i, y_i) is the i -th data point, $w \in \mathcal{W}$ is the parameter vector, and $g_i(w) = \log(1 + \exp(-y_i w; x_i))$. The function $g_i(w)$ is smooth and convex. The function $\mathcal{G}(w)$ is also smooth and convex. The algorithm achieves a convergence rate of $O(1/T)$ for the smooth case and $O(1/T^2)$ for the strongly convex case.

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

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

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

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

5 Convergence Analysis

$$\begin{aligned} \mathcal{F}_k(\mathbf{w}_*) - \mathcal{F}_k(\tilde{\mathbf{w}}_k) &\leq \frac{1}{2} \|\mathbf{w}_* - \tilde{\mathbf{w}}_k\|^2 \leq \frac{1}{2} \Delta_k^2, \\ \mathcal{F}_k(\tilde{\mathbf{w}}_k) - \min_{\mathbf{w}} \mathcal{F}_k(\mathbf{w}) &\leq \frac{k \Delta_k^2}{2^4}, \\ \mathcal{F}_k(\tilde{\mathbf{w}}_k) - \min_{\mathbf{w}} \mathcal{F}_k(\mathbf{w}) &\leq \frac{k \Delta_k^2}{2^4}. \end{aligned} \quad (3)$$

Theorem 2. Let $\hat{\mathbf{w}}_*^k, \tilde{\mathbf{w}}_{k+1}^k \in \mathbb{B}_R$ and $\hat{\mathbf{w}}_*^k \leq \Delta_k$. Then for $k = 1, 2, \dots$, it holds that

$$\hat{\mathbf{w}}_*^{k+1} \leq \frac{\Delta_k}{2}, \quad \mathcal{F}_k(\tilde{\mathbf{w}}_{k+1}) - \min_{\mathbf{w}} \mathcal{F}_k(\mathbf{w}) \leq \frac{k \Delta_k^2}{2^4}$$

$$\mathcal{F}_k(\tilde{\mathbf{w}}_{k+1}) \leq e^{-9/2}$$

$$T_k \geq \frac{300}{k} \ln \frac{1}{\epsilon}.$$

Proof. We first show that $\hat{\mathbf{w}}_*^{k+1} \leq \frac{\Delta_k}{2}$. Let $\mathbf{w}_*^1 = \mathbf{w}_* \leq R := \Delta_1$.

1. I. Let $\mathbf{w}_*^1 = \mathbf{w}_* \leq R := \Delta_1$. Then $\mathbf{w}_*^1 \in \mathbb{B}_R$.

2. Let $\mathbf{w}_*^m \in \mathbb{B}_R$. Then $\mathcal{F}_m(\tilde{\mathbf{w}}_{m+1}) - \mathcal{F}_m(\hat{\mathbf{w}}_*^m) \leq \frac{1}{2} \|\hat{\mathbf{w}}_*^m - \tilde{\mathbf{w}}_{m+1}\|^2 \leq \frac{1}{2} \Delta_m^2$.

$$\hat{\mathbf{w}}_*^m \leq \frac{\Delta_1}{m-1}; \quad \mathcal{F}_m(\tilde{\mathbf{w}}_{m+1}) - \mathcal{F}_m(\hat{\mathbf{w}}_*^m) \leq \frac{m \Delta_m^2}{2^4} = \frac{1 \Delta_1^2}{2^{3m+1}}$$

3. Let $\mathbf{w}_*^m \in \mathbb{B}_R$. Then $\mathcal{F}_m(\tilde{\mathbf{w}}_{m+1}) - \mathcal{F}_m(\hat{\mathbf{w}}_*^m) \leq \frac{1}{2} \|\hat{\mathbf{w}}_*^m - \tilde{\mathbf{w}}_{m+1}\|^2 \leq \frac{1}{2} \Delta_m^2$.

$$\frac{1}{n} \sum_{i=1}^n g_i(\tilde{\mathbf{w}}_{m+1}) \leq \mathcal{F}_m(\hat{\mathbf{w}}_*^m) + \frac{1 \Delta_1^2}{2^{3m+1}} - \frac{1}{m-1} \|\tilde{\mathbf{w}}_{m+1} - \bar{\mathbf{w}}_m\|^2$$

$$\leq \mathcal{F}_m(\hat{\mathbf{w}}_*^m) + \frac{1 \Delta_1^2}{2^{3m+1}} + \frac{1}{2m-2} \|\tilde{\mathbf{w}}_{m+1} - \bar{\mathbf{w}}_m\|^2$$

4. Let $\hat{\mathbf{w}}_*^{m+1} \leq \Delta_m = \Delta_1^{1-m}$.

$$\bar{\mathbf{w}}_m \leq \sum_{i=1}^m \tilde{\mathbf{w}}_i \leq \sum_{i=1}^m \Delta_i \leq \frac{\Delta_1}{1-m} \leq 2\Delta_1$$

5. Let $\hat{\mathbf{w}}_*^{m+1} \leq \Delta_m = \Delta_1^{1-m}$. Then $\mathcal{F}_m(\tilde{\mathbf{w}}_{m+1}) - \mathcal{F}_m(\hat{\mathbf{w}}_*^m) \leq \frac{1}{2} \|\hat{\mathbf{w}}_*^m - \tilde{\mathbf{w}}_{m+1}\|^2 \leq \frac{1}{2} \Delta_m^2$.

$$\frac{1}{n} \sum_{i=1}^n g_i(\tilde{\mathbf{w}}_{m+1}) \leq \mathcal{F}_m(\hat{\mathbf{w}}_*^m) + \frac{1 \Delta_1^2}{2^{3m+1}} + \frac{2}{2m-2} \|\tilde{\mathbf{w}}_{m+1} - \bar{\mathbf{w}}_m\|^2$$

6. Let $\mathbf{w}_*^m \in \mathbb{B}_R$. Then $\mathcal{F}_m(\tilde{\mathbf{w}}_{m+1}) - \mathcal{F}_m(\hat{\mathbf{w}}_*^m) \leq \frac{1}{2} \|\hat{\mathbf{w}}_*^m - \tilde{\mathbf{w}}_{m+1}\|^2 \leq \frac{1}{2} \Delta_m^2$.

$$\mathcal{F}_m(\mathbf{w}_*^m) \leq \mathcal{F}_m(\mathbf{w}_*) = \frac{1}{n} \sum_{i=1}^n g_i(\mathbf{w}_*) + \frac{1}{2^{m-1}} (\|\mathbf{w}_* - \bar{\mathbf{w}}_m\|^2 + 2 \|\mathbf{w}_* - \bar{\mathbf{w}}_m\| \|\bar{\mathbf{w}}_m\|): \quad (6)$$

7. Let $\mathbf{w}_*^m \in \mathbb{B}_R$. Then $\mathcal{F}_m(\tilde{\mathbf{w}}_{m+1}) - \mathcal{F}_m(\hat{\mathbf{w}}_*^m) \leq \frac{1}{2} \|\hat{\mathbf{w}}_*^m - \tilde{\mathbf{w}}_{m+1}\|^2 \leq \frac{1}{2} \Delta_m^2$.

8. Let $\hat{\mathbf{w}}_*^{m+1} \leq \Delta_m = \Delta_1^{1-m}$. Then $\mathcal{F}_m(\tilde{\mathbf{w}}_{m+1}) - \mathcal{F}_m(\hat{\mathbf{w}}_*^m) \leq \frac{1}{2} \|\hat{\mathbf{w}}_*^m - \tilde{\mathbf{w}}_{m+1}\|^2 \leq \frac{1}{2} \Delta_m^2$.

9. Let $\hat{\mathbf{w}}_*^{m+1} \leq \Delta_m = \Delta_1^{1-m}$. Then $\mathcal{F}_m(\tilde{\mathbf{w}}_{m+1}) - \mathcal{F}_m(\hat{\mathbf{w}}_*^m) \leq \frac{1}{2} \|\hat{\mathbf{w}}_*^m - \tilde{\mathbf{w}}_{m+1}\|^2 \leq \frac{1}{2} \Delta_m^2$.

$$\mathbf{w}_* - \bar{\mathbf{w}}_m \leq \sum_{i=m+1}^{\infty} \tilde{\mathbf{w}}_i \leq \sum_{k=m+1}^{\infty} \Delta_k \leq \frac{\Delta_1}{m(1-\epsilon)} \leq \frac{2\Delta_1}{m}$$

$\mathcal{F}_m(\mathbf{w}_*^m) \leq \frac{1}{n} \sum_{i=1}^n g_i(\mathbf{w}_*) + \frac{1}{2^{m-1}} \left(\frac{4\Delta_1^2}{2m} + \frac{8\Delta_1^2}{m} \right)$
 $= \frac{1}{n} \sum_{i=1}^n g_i(\mathbf{w}_*) + \frac{2}{2^{m-1}} \frac{\Delta_1^2}{2m-1} (2 + \dots) \leq \frac{1}{n} \sum_{i=1}^n g_i(\mathbf{w}_*) + \frac{5}{2^{m-1}} \frac{\Delta_1^2}{2m-1}$ (7)

By (6) and (7),

$$\frac{1}{n} \sum_{i=1}^n g_i(\bar{\mathbf{w}}_{m+1}) - \frac{1}{n} \sum_{i=1}^n g_i(\mathbf{w}_*) \leq \frac{5}{2^{m-2}} \frac{\Delta_1^2}{2m-2} = O(1-2^m),$$

$$T = T_1 \sum_{k=0}^{m-1} 2^k = \frac{T_1 (2^m - 1)}{2 - 1} \leq \frac{T_1}{3} 2^m.$$

□

5.1 Proof of Theorem 2

$\mathcal{F}(\mathbf{w}) = \frac{1}{2} \|\mathbf{w}\|^2 + \langle \mathbf{w}; \bar{\mathbf{w}} \rangle + \frac{1}{n} \sum_{i=1}^n g_i(\mathbf{w} + \bar{\mathbf{w}})$ (8)

$$\mathcal{F}'(\mathbf{w}) = \mathbf{w} + \bar{\mathbf{w}}' + \frac{1}{n} \sum_{i=1}^n g_i'(\mathbf{w} + \bar{\mathbf{w}}')$$
 (9)

$\hat{\mathbf{w}}_*^k = \hat{\mathbf{w}}_*^{k+1} \leq \Delta$, $\hat{\mathbf{w}}_*' \leq \frac{\Delta}{2}$; $\mathcal{F}(\bar{\mathbf{w}}') - \mathcal{F}(\hat{\mathbf{w}}_*) \leq \frac{\Delta^2}{4}$

□

Lemma 1.

$$\mathcal{F}(\mathbf{w}_t) - \mathcal{F}(\hat{\mathbf{w}}_*) \leq \frac{\|\mathbf{w}_t - \hat{\mathbf{w}}_*\|^2}{2} - \frac{\|\mathbf{w}_{t+1} - \hat{\mathbf{w}}_*\|^2}{2} + \frac{1}{2} \langle \nabla \hat{g}_t(\mathbf{w}_t) + \mathbf{w}_t; \mathbf{w}_t - \mathbf{w}_{t+1} \rangle$$

$$+ \langle \nabla \hat{\mathcal{F}}(\hat{\mathbf{w}}_*) - \nabla \hat{g}_t(\hat{\mathbf{w}}_*); \mathbf{w}_t - \hat{\mathbf{w}}_* \rangle + \langle -\nabla \hat{g}_t(\mathbf{w}_t) + \nabla \hat{g}_t(\hat{\mathbf{w}}_*) - \nabla \hat{\mathcal{F}}(\hat{\mathbf{w}}_*) + \nabla \hat{\mathcal{F}}(\mathbf{w}_t); \mathbf{w}_t - \hat{\mathbf{w}}_* \rangle$$

$\sum_{t=1}^T \mathcal{F}(\mathbf{w}_t) - \mathcal{F}(\hat{\mathbf{w}}_*) \leq \frac{\|\hat{\mathbf{w}}_*\|^2}{2} - \frac{\|\mathbf{w}_{T+1} - \hat{\mathbf{w}}_*\|^2}{2} - \langle \mathbf{g}; \mathbf{w}_{T+1} \rangle$

$$+ \underbrace{\frac{1}{2} \sum_{t=1}^T \langle \nabla \hat{g}_t(\mathbf{w}_t) + \mathbf{w}_t; \mathbf{w}_t - \mathbf{w}_{t+1} \rangle}_{:= A_T} + \underbrace{\sum_{t=1}^T \langle \nabla \hat{\mathcal{F}}(\hat{\mathbf{w}}_*) - \nabla \hat{g}_t(\hat{\mathbf{w}}_*); \mathbf{w}_t - \hat{\mathbf{w}}_* \rangle}_{:= B_T}$$

$$+ \underbrace{\sum_{t=1}^T \langle -\nabla \hat{g}_t(\mathbf{w}_t) + \nabla \hat{g}_t(\hat{\mathbf{w}}_*) - \nabla \hat{\mathcal{F}}(\hat{\mathbf{w}}_*) + \nabla \hat{\mathcal{F}}(\mathbf{w}_t); \mathbf{w}_t - \hat{\mathbf{w}}_* \rangle}_{:= C_T}$$

$$\mathbf{g} = \nabla \mathcal{F}(\mathbf{0})$$

$$\mathcal{F}(\mathbf{w}_{T+1}) - \mathcal{F}(\mathbf{0}) \leq \langle \nabla \mathcal{F}(\mathbf{0}); \mathbf{w}_{T+1} \rangle + \frac{1}{2} \|\mathbf{w}_{T+1}\|^2 = \langle \mathbf{g}; \mathbf{w}_{T+1} \rangle + \frac{1}{2} \|\mathbf{w}_{T+1}\|^2$$

$$\mathcal{F}(\mathbf{0}) \leq \mathcal{F}(\mathbf{w}_*) + \frac{\beta}{2} \|\mathbf{w}_*\|^2, \quad \max(\|\mathbf{w}_*\|; \|\mathbf{w}_{T+1}\|) \leq \Delta,$$

$$-\mathbf{g}; \mathbf{w}_{T+1} \leq \mathcal{F}(\mathbf{0}) - \mathcal{F}(\mathbf{w}_{T+1}) + \frac{1}{2} \Delta^2 \leq \Delta^2 - (\mathcal{F}(\mathbf{w}_{T+1}) - \mathcal{F}(\widehat{\mathbf{w}}_*))$$

$$\sum_{t=1}^{T+1} \mathcal{F}(\mathbf{w}_t) - \mathcal{F}(\widehat{\mathbf{w}}_*) \leq \Delta^2 \left(\frac{1}{2} + \dots \right) + \frac{1}{2} A_T + B_T + C_T: \quad (10)$$

Lemma 2. $A_T \leq 6 \Delta^2 T$.

Lemma 3. $B_T \leq \Delta^2 \left(\ln \frac{1}{\delta} + \sqrt{2T \ln \frac{1}{\delta}} \right)$, $C_T \leq 2 \Delta^2 \left(\ln \frac{1}{\delta} + \sqrt{2T \ln \frac{1}{\delta}} \right)$.

Lemma 4. $\widehat{\mathbf{w}}_*' \leq \widehat{\mathbf{w}}_* + \Delta$.

$$\sum_{t=1}^{T+1} \mathcal{F}(\mathbf{w}_t) - \mathcal{F}(\widehat{\mathbf{w}}_*) \leq \Delta^2 \left(\frac{1}{2} + 6 \Delta^2 T + 3 \ln \frac{1}{\delta} + 3 \sqrt{2T \ln \frac{1}{\delta}} \right)$$

$$\sum_{t=1}^{T+1} \mathcal{F}(\mathbf{w}_t) - \mathcal{F}(\widehat{\mathbf{w}}_*) \leq \Delta^2 \left(2 \ln \frac{1}{\delta} + 3 \ln \frac{1}{\delta} + 3 \sqrt{2T \ln \frac{1}{\delta}} \right)$$

$$\widehat{\mathbf{w}} = \sum_{i=1}^{T+1} \mathbf{w}_i / (T+1), \quad \mathcal{F}(\widehat{\mathbf{w}}) - \mathcal{F}(\widehat{\mathbf{w}}_*) \leq \Delta^2 \frac{5 \sqrt{3 \ln \frac{1}{\delta}}}{T+1}, \quad \widehat{\Delta}^2 = \|\widehat{\mathbf{w}} - \widehat{\mathbf{w}}_*\|^2 \leq \Delta^2 \frac{5 \sqrt{3 \ln \frac{1}{\delta}}}{T+1}.$$

$$T \geq \lceil 300 \delta^{-2} \ln \frac{1}{\delta} \rceil, \quad \widehat{\Delta}^2 \leq \frac{\Delta^2}{4}; \quad \mathcal{F}(\widehat{\mathbf{w}}) - \mathcal{F}(\widehat{\mathbf{w}}_*) \leq \frac{1}{2} \Delta^2: \quad (11)$$

Lemma 4. $\widehat{\mathbf{w}}_*' \leq \widehat{\mathbf{w}}_* + \Delta$.

By (11) and Lemma 4, $\widehat{\mathbf{w}}_*' \leq \Delta$.

6 Conclusions and Open Questions

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

Acknowledgments. ON A N000141210431 N F (II -1251031).

References

1. A. A. , P. L. B. , P. D. , M. J. , I. F. , 58(5):3235–3249, 2012.
2. A. B. , M. , M. , 31(3):167–175, 2003.
3. L. B. , O. B. , I. , 161–168, 2008.
4. , G. L. , O. B. , C. F. , I. A. , 208–240, 2003.
5. , L. O. , 2004.
6. H. B. , G. M. C. , J. N. F. , 134(1):127–155, 2012.
7. A. C. F. , O. , N. , K. , B. , 1647–1655, 2011.
8. O. D. , G. -B. , O. , L. , F. O. , 13:165–202, 2012.
9. M. P. F. , M. , H. , 34(3):A1380–A1405, 2012.
10. E. H. , A. A. , K. , L. F. , 69(2-3):169–192, 2007.
11. E. H. , K. , B. , 19:421–436, 2011.
12. , C. , J. P. , A. , 1008.5204, 2010.
13. A. N. , A. J. , G. L. , A. , 19:1574–1609, 2009.
14. A. , N. , D. B. , P. F. , 1983.
15. , N. , A. , (1/ 2). I. , 27, 372–376, 1983.
16. , N. , P. , 2004.
17. , N. , Ex. , 16(1):235–249, 2005.
18. , N. , 103(1):127–152, 2005.
19. A. , O. , K. , M. , 2012.
20. N. L. , F. x, M. , F. B. , A. , 2672–2680, 2012.
21. , N. , F. P. : P. , 807–814, 2007.
22. , 14:567599, 2013.
23. O. , 2013.
24. L. , J. , H. O(F.) , 2013.