# Supplementary Material: $O(\log T)$ Projections for Stochastic Optimization of Smooth and Strongly Convex Functions

**Lijun Zhang**                                                ZHANGLIJ@MSU.EDU

**Tianbao Yang†**                                                TYANG@GE.COM

**Rong Jin**                                                RONGJIN@CSE.MSU.EDU

**Xiaofei He‡**                                                XIAOFEIHE@CAD.ZJU.EDU.CN

Department of Computer Science and Engineering Michigan State University East Lansing MI                A
†GE Global Research San Ramon CA              A
‡ State Key Laboratory of CAD CG College of Computer Science Zhejiang University Hangzhou                China

## A. Proof of Lemma 1

We need the following lemma that characterizes the property of the extra gradient descent

**Lemma 8** Lemma in Nesterov. *Let $\mathcal{Z}$ be a convex compact set in Euclidean space $\mathcal{E}$ with inner product $\langle \cdot, \cdot \rangle$ let $\| \cdot \|$ be a norm on $\mathcal{E}$ and $\| \cdot \|$ be its dual norm, and let $\omega(\mathbf{z}) : \mathcal{Z} \mapsto \mathbb{R}$ be a $\alpha$ strongly convex function with respect to $\| \cdot \|$ The Bregman distance associated with $\omega$ for points $\mathbf{z}, \mathbf{w} \in \mathcal{Z}$ is defined as*

$$B_\omega(\mathbf{z}, \mathbf{w}) = \omega(\mathbf{z}) - \omega(\mathbf{w}) - \langle \mathbf{z} - \mathbf{w}, \nabla \omega(\mathbf{w}) \rangle.$$

*Let $\mathcal{U}$ be a convex and closed subset of $\mathcal{Z}$, and let $\mathbf{z} \in \mathcal{Z}$, let $\boldsymbol{\xi}, \boldsymbol{\eta} \in \mathcal{E}$, and let $\gamma > $ Consider the points*

$$\mathbf{w} = \arg\min_{\mathbf{y} \in}\{\langle \gamma \boldsymbol{\xi} - \nabla \omega(\mathbf{z}), \mathbf{y} \rangle + \omega(\mathbf{y})\},$$

$$\mathbf{z}_+ = \arg\min_{\mathbf{y} \in}\{\langle \gamma \boldsymbol{\eta} - \nabla \omega(\mathbf{z}), \mathbf{y} \rangle + \omega(\mathbf{y})\}.$$

*Then for all $\mathbf{z} \in \mathcal{U}$ one has*

$$\langle \mathbf{w} - \mathbf{z}, \gamma \boldsymbol{\eta} \rangle \le B_\omega(\mathbf{z}, \mathbf{z}) - B_\omega(\mathbf{z}, \mathbf{z}_+) + \frac{\gamma^2}{\alpha}\|\boldsymbol{\eta} - \boldsymbol{\xi}\|^2 - \frac{\alpha}{}\{\|\mathbf{w} - \mathbf{z}\|^2 + \|\mathbf{z}_+ - \mathbf{w}\|^2\}.$$

*Proof of Lemma 1* We first state the inner loop in Algorithm below

**for** $t$ to $M$ **do**

　Compute the average gradient at $\mathbf{w}_t^k$ over $B^k$ calls to the gradient oracle

$$\mathbf{g}_t^k = \frac{}{B^k}\sum_{i=1}^{B^k} \mathbf{g}(\mathbf{w}_t^k, i)$$

　Update

$$\mathbf{z}_t^k = \Pi_{\mathcal{D}}\left(\mathbf{w}_t^k - \eta \mathbf{g}_t^k\right)$$

　Compute the average gradient at $\mathbf{z}_t^k$ over $B^k$ calls to the gradient oracle

$$\mathbf{f}_t^k = \frac{}{B^k}\sum_{i=1}^{B^k} \mathbf{g}(\mathbf{z}_t^k, i)$$

update

$$\mathbf{w}_{t+1}^k \quad \boldsymbol{\Pi}_{\mathcal{D}}\left(\mathbf{w}_t^k - \eta \mathbf{f}_t^k\right)$$

**end for**

o si plify the notation we de ne

$$\mathbf{g}_t^k \quad \nabla F\, \mathbf{w}_t^k \ \text{ and } \mathbf{f}_t^k \quad \nabla F\, \mathbf{z}_t^k .$$

Let the two nor s $\|\cdot\|$ and $\|\cdot\|$ n Le a be the vector $\ell_2$ nor Each iteration in the inner oop satis es the conditions in Le a by doing the appings be ow

$$\mathcal{U} \quad \mathcal{Z} \quad \mathcal{E} \leftarrow \mathcal{D},\ \omega\,\mathbf{z} \ \leftarrow -\|\mathbf{z}\|^2,\ \alpha \leftarrow \ ,\ \gamma \leftarrow \eta,\ \mathbf{z} \ \leftarrow \mathbf{w}_t^k,\ \boldsymbol{\xi} \leftarrow \mathbf{g}_t^k,\ \boldsymbol{\eta} \leftarrow \mathbf{f}_t^k,\ \mathbf{w} \leftarrow \mathbf{z}_t^k,\ \mathbf{z}_+ \leftarrow \mathbf{v}_{t+1}^k\, \mathbf{z} \leftarrow \mathbf{w} .$$

Fo owing Le a we have

$$\left\|\mathbf{z}_t^k\right\| -$$

$$\langle \mathbf{z}_t^k - \mathbf{w}\ , \eta \mathbf{f}_t^k \rangle$$

$$\leq \frac{\|\mathbf{w}_t^k - \mathbf{w}\ \|^2}{} - \frac{\|\mathbf{w}_{t+1}^k - \mathbf{w}\ \|^2}{} + \eta^2 \|\mathbf{g}_t^k - \mathbf{f}_t^k\|^2 - -\|\mathbf{w}_t^k - \mathbf{z}_t^k\|^2$$

$$\leq \frac{\|\mathbf{w}_t^k - \mathbf{w}\ \|^2}{} - \frac{\|\mathbf{w}_{t+1}^k - \mathbf{w}\ \|^2}{} + \eta^2 \left(\|\mathbf{g}_t^k - \mathbf{g}_t^k\|^2 + \|\mathbf{f}_t^k - \mathbf{f}_t^k\|^2 + \|\mathbf{g}_t^k - \mathbf{f}_t^k\|^2\right) - -\|\mathbf{w}_t^k - \mathbf{z}_t^k\|^2$$

$$\leq \frac{\|\mathbf{w}_t^k - \mathbf{w}\ \|^2}{} - \frac{\|\mathbf{w}_{t+1}^k - \mathbf{w}\ \|^2}{} + \eta^2 \left(\|\mathbf{g}_t^k - \mathbf{g}_t^k\|^2 + \|\mathbf{f}_t^k - \mathbf{f}_t^k\|^2\right) + \eta^2 \|\mathbf{g}_t^k - \mathbf{f}_t^k\|^2 - -\|\mathbf{w}_t^k - \mathbf{z}_t^k\|^2$$

$$\leq \frac{\|\mathbf{w}_t^k - \mathbf{w}\ \|^2}{} - \frac{\|\mathbf{w}_{t+1}^k - \mathbf{w}\ \|^2}{} + \eta^2 \left(\|\mathbf{g}_t^k - \mathbf{g}_t^k\|^2 + \|\mathbf{f}_t^k - \mathbf{f}_t^k\|^2\right) + \eta^2 L^2 \|\mathbf{w}_t^k - \mathbf{z}_t^k\|^2 - -\|\mathbf{w}_t^k - \mathbf{z}_t^k\|^2$$

$$\leq \frac{\|\mathbf{w}_t^k - \mathbf{w}\ \|^2}{} - \frac{\|\mathbf{w}_{t+1}^k - \mathbf{w}\ \|^2}{} + \eta^2 \left(\|\mathbf{g}_t^k - \mathbf{g}_t^k\|^2 + \|\mathbf{f}_t^k - \mathbf{f}_t^k\|^2\right),$$

wher d f t h e hfth n wes the s oothness

$$\|\mathbf{g}_t^k - \mathbf{f}_t^k\| \quad \|\nabla \mathbf{w}_t^k \mathcal{F}\, \mathbf{z}_t^k \ \mathbf{k}\|\mathbf{w}_t^k - \mathbf{z}_t^k\|.$$

Dividing both sides by $M$ and following Jensen's inequality we have

$$F\left(\frac{1}{M}\sum_{t=1}^{M}\mathbf{z}_t^k\right) - F(\mathbf{w})$$

$$\leq \frac{1}{M}\sum_{t=1}^{M}F(\mathbf{z}_t^k) - F(\mathbf{w})$$

$$\leq \frac{\|\mathbf{w}_1^k - \mathbf{w}\|^2}{M\eta} + \frac{\eta}{M}\left(\sum_{t=1}^{M}\|\mathbf{g}_t^k - \bar{\mathbf{g}}_t^k\|^2 + \sum_{t=1}^{M}\|\mathbf{f}_t^k - \bar{\mathbf{f}}_t^k\|^2\right) + \frac{1}{M}\sum_{t=1}^{M}\langle \mathbf{f}_t^k - \bar{\mathbf{f}}_t^k, \mathbf{z}_t^k - \mathbf{w}\rangle - \frac{\lambda}{M}\sum_{t=1}^{M}\|\mathbf{z}_t^k - \mathbf{w}\|^2.$$

which gives the first inequality in Lemma.

Let $\mathrm{E}_{k-1}[\cdot]$ denote the expectation conditioned on all the randomness up to epoch $k-1$ and $\mathrm{E}_k^{t-1}[\cdot]$ denote the expectation conditioned on all the randomness up to the $t-1$ th iteration in the $k$ th epoch. Taking the conditional expectation of , we have

$$\mathrm{E}_{k-1}\left[F\left(\frac{1}{M}\sum_{t=1}^{M}\mathbf{z}_t^k\right)\right] - F(\mathbf{w})$$

$$\leq \frac{\|\mathbf{w}_1^k - \mathbf{w}\|^2}{M\eta} + \frac{\eta}{M}\left(\sum_{t=1}^{M}\mathrm{E}_{k-1}\left[\|\mathbf{g}_t^k - \bar{\mathbf{g}}_t^k\|^2\right] + \sum_{t=1}^{M}\mathrm{E}_{k-1}\left[\|\mathbf{f}_t^k - \bar{\mathbf{f}}_t^k\|^2\right]\right) + \frac{1}{M}\sum_{t=1}^{M}\mathrm{E}_{k-1}\left[\langle \mathbf{f}_t^k - \bar{\mathbf{f}}_t^k, \mathbf{z}_t^k - \mathbf{w}\rangle\right],$$

where we drop the last term since it is negative. To bound $\mathrm{E}_{k-1}\left[\|\mathbf{g}_t^k - \bar{\mathbf{g}}_t^k\|^2\right]$ we have

$$\mathrm{E}_{k-1}\left[\|\mathbf{g}_t^k - \bar{\mathbf{g}}_t^k\|^2\right] = \mathrm{E}_{k-1}\left[\left\|\frac{1}{B^k}\sum_{i=1}^{B^k}\mathbf{g}(\mathbf{w}_t^k, i) - \bar{\mathbf{g}}_t^k\right\|^2\right] = \mathrm{E}_{k-1}\left[\left\|\frac{1}{B^k}\sum_{i=1}^{B^k}(\mathbf{g}(\mathbf{w}_t^k, i) - \bar{\mathbf{g}}_t^k)\right\|^2\right]$$

$$= \frac{1}{B^{k\,2}}\left(\sum_{i=1}^{B^k}\mathrm{E}_{k-1}\left[\|\mathbf{g}(\mathbf{w}_t^k, i) - \bar{\mathbf{g}}_t^k\|^2\right] + \mathrm{E}_{k-1}\left[\sum_{i=j}\langle \mathrm{E}_k^{t-1}[\mathbf{g}(\mathbf{w}_t^k, i) - \bar{\mathbf{g}}_t^k], \mathrm{E}_k^{t-1}[\mathbf{g}(\mathbf{w}_t^k, j) - \bar{\mathbf{g}}_t^k]\rangle\right]\right)$$

$$= \frac{1}{B^{k\,2}}\left(\sum_{i=1}^{B^k}\mathrm{E}_{k-1}\left[\|\mathbf{g}(\mathbf{w}_t^k, i) - \bar{\mathbf{g}}_t^k\|^2\right]\right) \leq \frac{G^2}{B^k},$$

where we make use of the facts $\mathbf{g}(\mathbf{w}_t^k, i)$ and $\mathbf{g}(\mathbf{w}_t^k, j)$ are independent when $i \neq j$ and

$$\mathrm{E}_k^{t-1}[\mathbf{g}(\mathbf{w}_t^k, i) - \bar{\mathbf{g}}_t^k] = , \mathrm{E}_k^{t-1}\left[\|\mathbf{g}(\mathbf{w}_t^k, i) - \bar{\mathbf{g}}_t^k\|^2\right] \leq \mathrm{E}_k^{t-1}\left[\|\mathbf{g}(\mathbf{w}_t^k, i)\|^2\right] \leq G^2, \forall i = , \ldots, B^k.$$

Similarly we also have

$$\mathrm{E}_{k-1}\left[\|\mathbf{f}_t^k - \bar{\mathbf{f}}_t^k\|^2\right] \leq \frac{G^2}{B^k}.$$

Notice that $\mathbf{f}_t^k$ is an unbiased estimate of $\bar{\mathbf{f}}_t^k$ thus

$$\mathrm{E}_{k-1}\left[\langle \mathbf{f}_t^k - \bar{\mathbf{f}}_t^k, \mathbf{z}_t^k - \mathbf{w}\rangle\right] = \mathrm{E}_{k-1}\left[\langle \mathrm{E}_k^{t-1}[\mathbf{f}_t^k - \bar{\mathbf{f}}_t^k], \mathbf{z}_t^k - \mathbf{w}\rangle\right] = .$$

Substituting , and into we get the second inequality in Lemma $\qquad\qquad\square$

## B. Proof of Lemma

Recall that $\mathbf{g}_t^k = \frac{1}{B^k}\sum_{i=1}^{B^k}\mathbf{g}(\mathbf{w}_t^k, i)$ thus

$$\|\mathbf{g}_t^k - \bar{\mathbf{g}}_t^k\| = \left\|\frac{1}{B^k}\sum_{i=1}^{B^k}\mathbf{g}(\mathbf{w}_t^k, i) - \bar{\mathbf{g}}_t^k\right\|.$$

$\text{ınce } \|\mathbf{g} \; \mathbf{w}_t^k, i\; \| \leq G \;$ and $\mathrm{E} \; \mathbf{g} \; \mathbf{w}_t^k, i \qquad \mathbf{g}_t^k \;$ we have with a probab $\imath$ ty at east $\; -\delta$

$$\|\mathbf{g}_t^k -$$

**C.2.** $A > \eta G^2 / \lambda B^k$

Similar to the above proof on event $E_1$ we bound

$$|Z_t^k| \leq \|\mathbf{f}_t^k - \mathbf{f}_t^k\|\|\mathbf{z}_t^k - \mathbf{w}\| \leq \frac{1}{\theta}\|\mathbf{f}_t^k - \mathbf{f}_t^k\|^2 + \frac{\theta}{2}\|\mathbf{z}_t^k - \mathbf{w}\|^2 \leq \frac{C^2}{\theta} + \frac{\theta A}{2},$$

where $\theta$ can be any nonnegative real number. Denote the sum of conditional variances by

$$\sigma_M^2 = \sum_{t=1}^{M} E_k^{t-1}\left[\,Z_t^{k\,2}\,\right] \leq C^2 \sum_{t=1}^{M} \|\mathbf{z}_t - \mathbf{w}\|^2 \leq C^2 A,$$

where $E_k^{t-1}[\cdot]$ denote the expectation conditioned on all the randomness up to the $t-1$ th iteration in the $k$ th epoch.

Notice that $A$ in the upper bound for $|Z_t^k|$ and $\sigma_M^2$ is a random variable thus we cannot directly apply Theorem. To address this challenge we make use of the peeling technique described in Bartlett et al. and have

$$\Pr\left(\sum_{t=1}^{M} Z_t^k \geq \sqrt{C^2 A \tau} + \left(\frac{C^2}{\theta} + \frac{\theta A}{2}\right)\tau\right)$$

$$\Pr\left(\sum_{t=1}^{M} Z_t^k \geq \sqrt{C^2 A \tau} + \left(\frac{C^2}{\theta} + \frac{\theta A}{2}\right)\tau, \frac{\eta G^2}{\lambda B^k} < A \leq \frac{MG^2}{\lambda^2}\right)$$

$$\Pr\left(\sum_{t=1}^{M} Z_t^k \geq \sqrt{C^2 A \tau} + \left(\frac{C^2}{\theta} + \frac{\theta A}{2}\right)\tau, \max_t|Z_t^k| \leq \frac{C^2}{\theta} + \frac{\theta A}{2}, \sigma_M^2 \leq C^2 A, \frac{\eta G^2}{\lambda B^k} < A \leq \frac{MG^2}{\lambda^2}\right)$$

$$\leq \sum_{i=1}^{n} \Pr\left(\sum_{t=1}^{M} Z_t^k \geq \sqrt{C^2 A \tau} + \left(\frac{C^2}{\theta} + \frac{\theta A}{2}\right)\tau, \max_t|Z_t^k| \leq \frac{C^2}{\theta} + \frac{\theta A}{2}, \sigma_M^2 \leq C^2 A, \frac{\eta G^2}{\lambda B^k}2^{i-1} < A \leq \frac{\eta G^2}{\lambda B^k}2^{i}\right)$$

$$\leq \sum_{i=1}^{n} \Pr\left(\sum_{t=1}^{M} Z_t^k \geq \sqrt{\left(C^2 \frac{\eta G^2}{\lambda B^k}2^{i-1}\right)\tau} + \left(\frac{C^2}{\theta} + \frac{\theta}{2}\frac{\eta G^2}{\lambda B^k}2^{i-1}\right)\tau, \max_t|Z_t^k| \leq \frac{C^2}{\theta} + \frac{\theta}{2}\frac{\eta G^2}{\lambda B^k}2^{i}, \sigma_M^2 \leq C^2 \frac{\eta G^2}{\lambda B^k}2^{i}\right)$$

$$\leq \sum_{i=1}^{n} \Pr\left(\sum_{t=1}^{M} Z_t^k \geq \sqrt{\left(C^2 \frac{\eta G^2}{\lambda B^k}2^{i}\right)\tau} + \left(\frac{C^2}{\theta} + \frac{\theta}{2}\frac{\eta G^2}{\lambda B^k}2^{i}\right)\tau, \max_t|Z_t^k| \leq \frac{C^2}{\theta} + \frac{\theta}{2}\frac{\eta G^2}{\lambda B^k}2^{i}, \sigma_M^2 \leq C^2 \frac{\eta G^2}{\lambda B^k}2^{i}\right)$$

$$\leq n e^{-\tau},$$

where

$$n = \left\lceil \log_2 \frac{MB^k}{\eta \lambda} \right\rceil,$$

and the last step follows the Bernstein's inequality for martingales in Theorem, setting

$$\theta = \frac{\lambda}{\tau},$$

$$\tau = \log\frac{n}{\delta},$$

with a probability at least $1 - \delta/?$ we have

$$\sum_{t=1}^{M} Z_t^k$$

$$\leq \sqrt{C^2 A \tau} + \left(\frac{C^2}{\theta} + \frac{\theta A}{2}\right)\tau = \sqrt{C^2 A \tau} + \frac{C^2}{\lambda}\tau^2 + \frac{\lambda A}{2}$$

$$\leq \frac{1}{\lambda}C^2\tau + \frac{\lambda A}{2} + \frac{C^2}{\lambda}\tau^2 + \frac{\lambda A}{2} = \frac{C^2}{\lambda}\left(\log\frac{n}{\delta} + \log^2\frac{n}{\delta}\right) + \frac{\lambda A}{2}.$$

We complete the proof by combining and

## D. Proof of Lemma 7

We fo ow the ogic used in the proof of Le a

It is straightforward to check that

$$B^k \quad \alpha\eta\lambda^{k-1} \quad \frac{\alpha\eta G^2}{V_k}.$$

Then $k$ with a probability $-\delta^{1-1}$ we have

$$_1 \quad F \ \mathbf{w}_1^1 \ - F \ \mathbf{w} \quad \overset{(1)}{\leq} \ \frac{G^2}{\lambda} \quad \frac{G^2}{\lambda^{1-2}} \quad V_1.$$
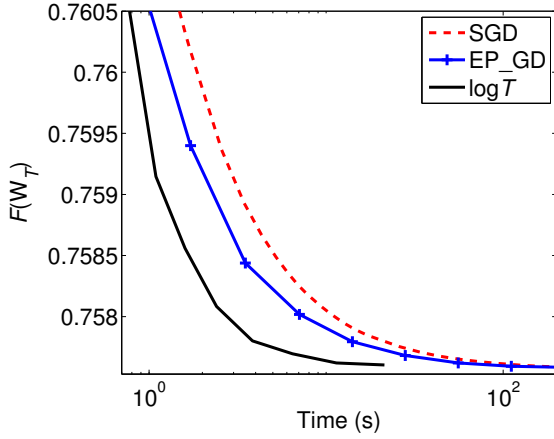
Assu e that with a probability at east $-\delta^{k-1}$ $_k \leq V_k$ for so e $k \geq$ We now prove the case for $k$. Notice that $N$ de ned in is arger than $n$ de ned in Fro Le a with a probability at east $-\delta$ we have

$$_{k+1} \quad F \ \mathbf{w}_1^{k+1} \ - F \ \mathbf{w}$$
$$\leq \frac{\|\mathbf{w}_1^k - \mathbf{w}\|^2}{M\eta} + \frac{G^2\eta}{B^k} \text{ og}^2 \frac{M}{\delta} + \frac{G^2}{\lambda B^k M} \left[ + \text{ og}^2 \frac{M}{\delta} \left( \text{og} \frac{N}{\delta} + - \text{og}^2 \frac{N}{\delta} \right) \right]$$
$$\leq \frac{k}{\alpha} + \frac{}{\alpha} \text{ og}^2 \frac{M\,V_k}{\delta} + \frac{}{\alpha} \left[ + \text{ og}^2 \frac{M}{\delta} \left( \text{og} \frac{N}{\delta} + - \text{og}^2 \frac{N}{\delta} \right) \right] \frac{V_k}{}.$$

Using the de nition of $\alpha$ in with a probability at east $-\delta^k$ we have

$$_{k+1} \leq -V_k + -V_k + -V_k \quad -V_k \quad V_{k+1}.$$

## E. More Results for the Regularized Distance Metric Learning



(a) Mushrooms