Supplementary Material for Dual Space Gradient Descent for Online Learning

Trung Le, Tu Dinh Nguyen, Vu Nguyen, Dinh Phung Centre for Pattern Recognition and Data Analytics Deakin University, Australia {trung.l, tu.nguyen, v.nguyen, dinh.phung}@deakin.edu.au

1 Suitability of Loss Functions

In this section, we present the suitability of the loss functions for Hinge, smooth Hinge, and Logistic for classification and ℓ_1 , and ε -insensitive for regression. We prove that these losses satisfy the condition: there exists a positive constant A such that $|\nabla_o l(y, o)| \leq A, \forall y, o$. For each loss, we show its two forms used in the paper w.r.t o and w.

Hinge loss

$$
l(y, o) = \max(0, 1 - yo)
$$

$$
l(\mathbf{w}, \mathbf{x}, y) = \max(0, 1 - y\mathbf{w}^\top \Phi(\mathbf{x}))
$$

$$
\nabla_o l(y, o) = -\mathbb{I}_{yo \le 1}y
$$

$$
|\nabla_o l(y, o)| = |\mathbb{I}_{yo \le 1}| \le 1 = A
$$

Logistic loss

$$
l(y, o) = \log (1 + e^{-yo})
$$

$$
l(\mathbf{w}, \mathbf{x}, y) = \log (1 + e^{-y\mathbf{w}^\top \Phi(\mathbf{x})})
$$

$$
\nabla_o l(y, o) = \frac{-ye^{-yo}}{e^{-yo} + 1}
$$

$$
|\nabla_o l(y, o)| = \left| \frac{e^{-yo}}{e^{-yo} + 1} \right| < 1 = A
$$

Smooth Hinge loss [\[4\]](#page-6-0)

$$
l(y, o) = \begin{cases} 0 & \text{if } yo > 1 \\ 1 - yo - \frac{\tau}{2} & \text{if } yo < 1 - \tau \\ \frac{1}{2\tau} (1 - yo)^2 & \text{otherwise} \end{cases}
$$

$$
l(\mathbf{w}, \mathbf{x}, y) = \begin{cases} 0 & \text{if } y\mathbf{w}^\top \Phi(\mathbf{x}) > 1 \\ 1 - y\mathbf{w}^\top \Phi(\mathbf{x}) - \frac{\tau}{2} & \text{if } y\mathbf{w}^\top \Phi(\mathbf{x}) < 1 - \tau \\ \frac{1}{2\tau} (1 - y\mathbf{w}^\top \Phi(\mathbf{x}))^2 & \text{otherwise} \end{cases}
$$

$$
\nabla_o l(y, o) = -\mathbb{I}_{\{yo < 1 - \tau\}} y + \tau^{-1} \mathbb{I}_{1 - \tau \le yo \le 1} (yo - 1) y
$$

$$
|\nabla_o l(y, o)| = |\mathbb{I}_{\{yo < 1 - \tau\}}| + |\tau^{-1} \mathbb{I}_{1 - \tau \le yo \le 1} (yo - 1)|
$$

$$
\le |\mathbb{I}_{\{yo < 1 - \tau\}}| + \tau^{-1} \tau |\mathbb{I}_{1 - \tau \le yo \le 1}| \le 1 = A
$$

30th Conference on Neural Information Processing Systems (NIPS 2016), Barcelona, Spain Acknowledgment: This work is partially supported by the Australian Research Council under the Discovery Project DP160109394.

 ℓ_1 loss

$$
l(y, o) = |y - o|
$$

\n
$$
l(\mathbf{w}, \mathbf{x}, y) = |y - \mathbf{w}^\top \Phi(\mathbf{x})|
$$

\n
$$
\nabla_o l(y, o) = \textbf{sign}(o - y)
$$

\n
$$
|\nabla_o l(y, o)| \le 1 = A
$$

ε-insensitive loss

$$
l(y, o) = \max(0, |y - o| - \varepsilon)
$$

$$
l(\mathbf{w}, \mathbf{x}, y) = \max(0, |y - \mathbf{w}^\top \Phi(\mathbf{x})| - \varepsilon)
$$

$$
\nabla_o l(y, o) = \mathbb{I}_{|y - o| \ge \varepsilon} \text{sign}(o - y)
$$

$$
|\nabla_o l(y, o)| \le 1 = A
$$

We note that \mathbb{I}_A denotes the indicator function which renders 1 if A is true and 0 otherwise.

2 Proofs

Lemma 1. *After the iteration* t*, we have the following representations*

$$
\hat{\mathbf{w}}_t = \sum_{j=1}^t \alpha_j \left(1 - \beta_j\right) \Phi\left(\mathbf{x}_j\right) \tag{1}
$$

$$
\tilde{\mathbf{w}}_t = \sum_{j=1}^t \alpha_j \beta_j \mathbf{z} \left(\mathbf{x}_j \right) \tag{2}
$$

$$
\mathbf{w}_{t} = \sum_{j=1}^{t} \alpha_{j} \Phi\left(\boldsymbol{x}_{j}\right)
$$
\n(3)

where $\alpha_j = -\eta_t \nabla_o l(g_j, f_j^h(\boldsymbol{x}_j))$, $\forall j = 1, \ldots, t$ and $\eta_t = \frac{1}{\lambda t}$.

Proof. Since if $\beta_j = 1$, we perform the budget maintenance procedure and move the current vector to the random-feature space, we have the representations in Eqs. $(1,2,3)$ $(1,2,3)$ $(1,2,3)$. In addition at the iteration j, $\Phi(x_j)$ arrives with the initial coefficient $\alpha_j = -\eta_j \nabla_o l(y_j, \hat{f}_j^h(x_j))$. After the iteration $t > j$, this coefficient becomes

$$
\alpha_j = -\frac{t-1}{t}\frac{t-2}{t-1}...\frac{j}{j+1}\frac{1}{\lambda j}\nabla_o l\left(y_i, f_j^h\left(\boldsymbol{x}_j\right)\right) = -\eta_t \nabla_o l\left(y_j, f_j^h\left(\boldsymbol{x}_j\right)\right)
$$

Theorem 2. With a probability at least $1-2^8\left(\frac{\sigma_\mu A d_\mathcal{X}}{\lambda \varepsilon}\right) \exp\left(-\frac{D\lambda^2\varepsilon^2}{4(M+2)\varepsilon^2}\right)$ $\frac{D\lambda^2\varepsilon^2}{4(M+2)A^2}$ where $d_{\mathcal{X}}$ specifies the *diameter of the compact set* X *, we have*

i)
$$
|f_t(\mathbf{x}) - f_t^h(\mathbf{x})| \le \varepsilon
$$
 for all $t > 0$ and $\mathbf{x} \in \mathcal{X}$.
\n*ii)* $\mathbb{E}[|f_t(\mathbf{x}) - f_t^h(\mathbf{x})|] \le A^{-1} \lambda \varepsilon \sum_{j=1}^t \mathbb{E}[\alpha_j^2]^{1/2} \mu_j^{1/2}$ where $\mu_j = p(\beta_j = 1)$.

Let us define a random map $z : \mathbb{R}^d \to \mathbb{R}^{2D}$ where $z(x) = \frac{1}{D^{1/2}} \left[\cos \left(\omega_i^{\mathsf{T}} x \right), \sin \left(\omega_i^{\mathsf{T}} x \right) \right]_{i=1}^D$ and $\omega_1, ..., \omega_D \stackrel{i.i.d}{\sim} \mathcal{N}(0, \sigma^{-2}I)$ for every $x \in \mathbb{R}^d$. We would like to restate Claim 1 in [\[3\]](#page-6-1).

Let M be a compact subset of \mathbb{R}^d with diameter diam (\mathcal{M}) . Then, for the random mapping \boldsymbol{z} (.), we *have*

$$
\mathbb{P}\left(\sup_{x,x'\in\mathcal{M}}\left|K\left(x,x'\right)-\boldsymbol{z}\left(x\right)^{\mathsf{T}}\boldsymbol{z}\left(x'\right)\right|<\varepsilon\right)\geq 1-2^{8}\left(\frac{\sigma\mathrm{diam}\left(\mathcal{M}\right)}{\varepsilon}\right)\exp\left(\frac{-D\varepsilon^{2}}{4\left(d+2\right)}\right)
$$

where
$$
K\left(x, x^{'}\right) = e^{-\frac{\left\|x - x^{'}\right\|^2}{2\sigma^2}}
$$
.

Proof. We denote

$$
\omega = (\omega_1, ..., \omega_D) \sim p_{\omega}(\omega) = \prod_{i=1}^{D} \mathcal{N}(\omega_i | 0, \sigma^{-2} I)
$$

$$
\tilde{K}(x, x') = z(x)^{\mathsf{T}} z(x') = D^{-1} \sum_{i=1}^{D} \left(\cos \left(\omega_i^{\mathsf{T}} x \right) \cos \left(\omega_i^{\mathsf{T}} x' \right) + \sin \left(\omega_i^{\mathsf{T}} x \right) \sin \left(\omega_i^{\mathsf{T}} x' \right) \right)
$$

We further denote

$$
g(\omega) = \sup_{x, x' \in \mathcal{M}} \left| K(x, x') - \tilde{K}(x, x') \right|
$$

$$
G_{\varepsilon} = \{ \omega : g(\omega) < A^{-1} \lambda \varepsilon \}
$$

It is certain that $\mathbb{P}_{\omega}(G_{\varepsilon}) \geq 1 - \theta$ where $\theta = 2^8 \left(\frac{\sigma A \text{diam}(\mathcal{M})}{\lambda \varepsilon} \right) exp \left(\frac{-D \lambda^2 \varepsilon^2}{4(d+2)A} \right)$ $\frac{-D\lambda^2\varepsilon^2}{4(d+2)A^2}$ and for every $\omega\in G_\varepsilon$ and $x, x^{'} \in \mathcal{M}$ we have \overline{a} $\overline{}$

$$
\left|K\left(x,x^{'}\right)-\tilde{K}\left(x,x^{'}\right)\right|
$$

We now turn back to Theorem [2.](#page-1-3) It appears that

$$
\left|f_t\left(x\right) - f_t^h\left(x\right)\right| \le \sum_{j=1}^t \beta_j \left|\alpha_j\right| \left|K\left(x_j, x\right) - \tilde{K}\left(x_j, x\right)\right|
$$

Therefore, for every $\omega \in G_{\varepsilon}$ we have

$$
\left|f_t\left(x\right) - f_t^h\left(x\right)\right| \le A^{-1} \lambda \varepsilon \sum_{j=1}^t \beta_j \left| \alpha_j \right|
$$

Let us denote $s = (x_1, y_1), ..., (x_t, y_t)$. Taking expectation of the above inequality w.r.t s, we gain for all $\omega \in G_{\varepsilon}$

$$
\mathbb{E}_{s} \left[\left| f_{t} \left(x \right) - f_{t}^{h} \left(x \right) \right| \right] \leq A^{-1} \lambda \varepsilon \sum_{j=1}^{t} \mathbb{E}_{s} \left[\beta_{j}^{2} \right]^{1/2} \mathbb{E}_{s} \left[\alpha_{j}^{2} \right]^{1/2}
$$

$$
\leq A^{-1} \lambda \varepsilon \sum_{j=1}^{t} \mu_{j} \mathbb{E}_{s} \left[\alpha_{j}^{2} \right]^{1/2}
$$

It means that

$$
\mathbb{P}_{\omega}\left(\mathbb{E}_{s}\left[\left|f_{t}\left(x\right)-f_{t}^{h}\left(x\right)\right|\right]\leq A^{-1}\lambda\varepsilon\sum_{j=1}^{t}\mu_{j}\mathbb{E}_{s}\left[\alpha_{j}^{2}\right]^{1/2}\right)\geq\mathbb{P}_{\omega}\left(G_{\varepsilon}\right)\geq1-\theta
$$

 \Box

Lemma 3. *The following statement holds for all* t

$$
\|\mathbf{w}_t\| \leq \frac{A}{\lambda}
$$

Proof. Using Lemma [1,](#page-1-4) we have

$$
\mathbf{w}_{t} = \sum_{j=1}^{t} \alpha_{j} \Phi\left(x_{j}\right)
$$

where $\alpha_j = -\eta_t \nabla_o l\left(y_j, f_j^h\left(\boldsymbol{x}_j\right)\right)$. It implies that

$$
\|\mathbf{w}_{t}\| \leq \sum_{j=1}^{t} |\alpha_{j}| \|\Phi(\mathbf{x}_{j})\| \leq \sum_{j=1}^{t} |\alpha_{j}| \leq \sum_{j=1}^{t} \frac{A}{\lambda t} = \frac{A}{\lambda}
$$

Lemma 4. *The following statement holds for all* t

$$
||g_t|| \le G = 2A
$$

where we define $g_t = \lambda \mathbf{w}_t + \nabla_{\mathbf{w}} l(\mathbf{w}_t, \mathbf{x}_t, y_t) = \lambda \mathbf{w}_t + \nabla_o l(\mathbf{y}_t, f_t(\mathbf{x}_t)) \Phi(\mathbf{x}_t)$.

Proof. We derive as

$$
||g_t|| \leq \lambda ||\mathbf{w}_t|| + ||\nabla_o l\left(y_t, f_t\left(\mathbf{x}_t\right)\right) \Phi\left(\mathbf{x}_t\right)|| \leq \lambda \frac{A}{\lambda} + A = 2A
$$

Lemma 5. *The following statement holds for all* t

$$
\mathbb{E}\left[\left\|\mathbf{w}_t-\mathbf{w}^{\star}\right\|^2\right] \leq W^2
$$

where $W = \frac{2A(1+\sqrt{5})}{\lambda}$ $\frac{1}{\lambda}$.

Proof. Recall that $g_t = \lambda \mathbf{w}_t + \nabla_{\mathbf{w}} l(\mathbf{w}_t, \mathbf{x}_t, y_t) = \lambda \mathbf{w}_t + \nabla_{\theta} l(\mathbf{y}_t, f_t(\mathbf{x}_t)) \Phi(\mathbf{x}_t)$. It is obvious that g_t satisfies \overline{a}

$$
\mathbb{E}_{(\boldsymbol{x}_t, y_t)}\left[g_t | \mathbf{w}_t\right] = \mathcal{J}^{'}\left(\mathbf{w}_t\right)
$$

We have the following if we denote $\delta g_t = g_t - g_t^h$

$$
\|\mathbf{w}_{t+1} - \mathbf{w}^{\star}\|^2 = \|\mathbf{w}_t - \eta_t g_t^h - \mathbf{w}^{\star}\| = \|\mathbf{w}_t - \eta_t g_t - \mathbf{w}^{\star} + \eta_t \delta g_t\|^2
$$

= $\|\mathbf{w}_t - \mathbf{w}^{\star}\|^2 - 2\eta_t g_t^\top (\mathbf{w}_t - \mathbf{w}^{\star}) + \eta_t^2 \|g_t\|^2 - 2\eta_t^2 g_t^\top \delta g_t + \eta_t^2 \|\delta g_t\|^2 + 2\eta_t (\mathbf{w}_t - \mathbf{w}^{\star})^\top \delta g_t$

It appears that

$$
\delta g_t = \left[\nabla_o l\left(y_t, f_t\left(\boldsymbol{x}_t\right)\right) - \nabla_o l\left(y_t, f_t^h\left(\boldsymbol{x}_t\right)\right)\right] \Phi\left(\boldsymbol{x}_t\right) \|\delta g_t\| = \left|\nabla_o l\left(y_t, f_t\left(\boldsymbol{x}_t\right)\right) - \nabla_o l\left(y_t, f_t^h\left(\boldsymbol{x}_t\right)\right)\right| \leq 2A
$$

Hence, we obtain

$$
\|\mathbf{w}_{t+1} - \mathbf{w}^{\star}\|^2 \leq \|\mathbf{w}_t - \mathbf{w}^{\star}\|^2 - 2\eta_t g_t^\top (\mathbf{w}_t - \mathbf{w}^{\star}) + \eta_t^2 G^2 + 4\eta_t^2 GA + 4\eta_t^2 A^2 + 2\eta_t \|\mathbf{w}_t - \mathbf{w}^{\star}\| \|\delta g_t\|
$$

Taking conditional expectation w.r.t w_t on both sides of the above inequality, we gain

$$
\mathbb{E}\left[\left\|\mathbf{w}_{t+1}-\mathbf{w}^{\star}\right\|^{2}\right] \leq \mathbb{E}\left[\left\|\mathbf{w}_{t}-\mathbf{w}^{\star}\right\|^{2}\right] - 2\eta_{t}\nabla_{\mathbf{w}}\mathcal{J}\left(\mathbf{w}_{t}\right)^{\top}\left(\mathbf{w}_{t}-\mathbf{w}^{\star}\right) + \eta_{t}^{2}G^{2} + 4\eta_{t}^{2}GA
$$

+ 4\eta_{t}^{2}A^{2} + 2\eta_{t}\mathbb{E}\left[\left\|\mathbf{w}_{t}-\mathbf{w}^{\star}\right\|\left\|\delta g_{t}\right\|\right]

$$
\leq \mathbb{E}\left[\left\|\mathbf{w}_{t}-\mathbf{w}^{\star}\right\|^{2}\right] + 16A^{2}\eta_{t}^{2} + 2\eta_{t}\mathbb{E}\left[\left\|\mathbf{w}_{t}-\mathbf{w}^{\star}\right\|\left\|\delta g_{t}\right\|\right] - \frac{1}{t}\left\|\mathbf{w}_{t}-\mathbf{w}^{\star}\right\|
$$

Here we note that we have used

$$
\nabla_{\mathbf{w}} \mathcal{J}(\mathbf{w}_t)^\top (\mathbf{w}_t - \mathbf{w}^\star) \ge \mathcal{J}(\mathbf{w}_t) - \mathcal{J}(\mathbf{w}^\star) + \frac{\lambda}{2} ||\mathbf{w}_t - \mathbf{w}^\star||^2 \ge \frac{\lambda}{2} ||\mathbf{w}_t - \mathbf{w}^\star||^2
$$

 \Box

 \Box

Taking expectation on both sides again, we obtain

$$
\mathbb{E}\left[\left\|\mathbf{w}_{t+1} - \mathbf{w}^{\star}\right\|^{2}\right] \leq \frac{t-1}{t} \mathbb{E}\left[\left\|\mathbf{w}_{t} - \mathbf{w}^{\star}\right\|^{2}\right] + \frac{16A^{2}}{\lambda^{2}t^{2}} + \frac{4A\mathbb{E}\left[\left\|\mathbf{w}_{t} - \mathbf{w}^{\star}\right\|^{2}\right]^{1/2}}{\lambda t}
$$
\n
$$
\leq \frac{t-1}{t} \mathbb{E}\left[\left\|\mathbf{w}_{t} - \mathbf{w}^{\star}\right\|^{2}\right] + \frac{16A^{2}}{\lambda^{2}t} + \frac{4A\mathbb{E}\left[\left\|\mathbf{w}_{t} - \mathbf{w}^{\star}\right\|^{2}\right]^{1/2}}{\lambda t}
$$
\nChoose $W = \frac{2A(1+\sqrt{5})}{\lambda}$, we have if $\mathbb{E}\left[\left\|\mathbf{w}_{t} - \mathbf{w}^{\star}\right\|^{2}\right] \leq W^{2}$ then $\mathbb{E}\left[\left\|\mathbf{w}_{t+1} - \mathbf{w}^{\star}\right\|^{2}\right] \leq W^{2}$. \square

Theorem 6. *The following statement guarantees for all* T

$$
\mathbb{E}\left[\mathcal{J}\left(\overline{\mathbf{w}}_{T}\right)-\mathcal{J}\left(\mathbf{w}^{\star}\right)\right] \leq \mathbb{E}\left[\frac{1}{T}\sum_{t=1}^{T}\mathcal{J}\left(\mathbf{w}_{t}\right)-\mathcal{J}\left(\mathbf{w}^{\star}\right)\right] \leq \frac{8A^{2}\left(\log\left(T\right)+1\right)}{\lambda T}+\frac{1}{T}W\sum_{t=1}^{T}\mathbb{E}\left[M_{t}^{2}\right]^{1/2}
$$

where $\overline{\mathbf{w}}_T = \frac{1}{T} \sum_{t=1}^T \mathbf{w}_t$, $M_t = \nabla_o l\left(y_t, f_t\left(\boldsymbol{x}_t\right)\right) - \nabla_o l\left(y_t, f_t^h\left(\boldsymbol{x}_t\right)\right)$.

Proof. Recall that $g_t = \lambda \mathbf{w}_t + \nabla_{\mathbf{w}} l(\mathbf{w}_t, \mathbf{x}_t, y_t) = \lambda \mathbf{w}_t + \nabla_o l(\mathbf{y}_t, f_t(\mathbf{x}_t)) \Phi(\mathbf{x}_t)$. It is obvious that g_t satisfies

$$
\mathbb{E}_{(\boldsymbol{x}_t, y_t)}\left[g_t | \mathbf{w}_t\right] = \nabla_{\mathbf{w}} \mathcal{J}\left(\mathbf{w}_t\right)
$$

We have the following if we denote $\delta g_t = g_t - g_t^h$

$$
\|\mathbf{w}_{t+1} - \mathbf{w}^{\star}\|^2 = \|\mathbf{w}_t - \eta_t g_t^h - \mathbf{w}^{\star}\| = \|\mathbf{w}_t - \eta_t g_t - \mathbf{w}^{\star} + \eta_t \delta g_t\|^2
$$

=
$$
\|\mathbf{w}_t - \mathbf{w}^{\star}\|^2 - 2\eta_t g_t^\top (\mathbf{w}_t - \mathbf{w}^{\star}) + \eta_t^2 \|g_t\|^2 - 2\eta_t^2 g_t^\top \delta g_t + \eta_t^2 \|\delta g_t\|^2 + 2\eta_t (\mathbf{w}_t - \mathbf{w}^{\star})^\top \delta g_t
$$

It appears that

$$
\delta g_t = \left[\nabla_o l \left(y_t, f_t \left(\boldsymbol{x}_t \right) \right) - \nabla_o l \left(y_t, f_t^h \left(\boldsymbol{x}_t \right) \right) \right] \Phi \left(x_t \right) \|\delta g_t\| = \left| \nabla_o l \left(y_t, f_t \left(\boldsymbol{x}_t \right) \right) - \nabla_o l \left(y_t, f_t^h \left(\boldsymbol{x}_t \right) \right) \right| \leq 2A
$$

Hence, we obtain

$$
\|\mathbf{w}_{t+1} - \mathbf{w}^{\star}\|^{2} \leq \|\mathbf{w}_{t} - \mathbf{w}^{\star}\|^{2} - 2\eta_{t}g_{t}^{\top}(\mathbf{w}_{t} - \mathbf{w}^{\star}) + \eta_{t}^{2}G^{2} + 4\eta_{t}^{2}GA + 4\eta_{t}^{2}A^{2}
$$

+ 2 $\eta_{t} \|\mathbf{w}_{t} - \mathbf{w}^{\star}\| \|\delta g_{t}\|$

$$
g_{t}^{\top}(\mathbf{w}_{t} - \mathbf{w}^{\star}) \leq \frac{\|\mathbf{w}_{t} - \mathbf{w}^{\star}\|^{2} - \|\mathbf{w}_{t+1} - \mathbf{w}^{\star}\|^{2}}{2\eta_{t}} + 8A^{2}\eta_{t} + \|\mathbf{w}_{t} - \mathbf{w}^{\star}\| \|\delta g_{t}\|
$$

Taking conditional expectation w.r.t w_t on both sides, we gain

$$
\nabla_{\mathbf{w}} \mathcal{J}(\mathbf{w}_{t})^{\top} (\mathbf{w}_{t} - \mathbf{w}^{\star}) \leq \mathbb{E} \left[\frac{\|\mathbf{w}_{t} - \mathbf{w}^{\star}\|^{2}}{2\eta_{t}} \right] - \mathbb{E} \left[\frac{\|\mathbf{w}_{t+1} - \mathbf{w}^{\star}\|^{2}}{2\eta_{t}} \right] + 8A^{2}\eta_{t} + \mathbb{E} \left[\|\mathbf{w}_{t} - \mathbf{w}^{\star}\| \|\delta g_{t} \right]
$$

$$
\mathcal{J}(\mathbf{w}_{t}) - \mathcal{J}(\mathbf{w}^{\star}) + \frac{\lambda}{2} \|\mathbf{w}_{t} - \mathbf{w}^{\star}\|^{2} \leq \mathbb{E} \left[\frac{\|\mathbf{w}_{t} - \mathbf{w}^{\star}\|^{2}}{2\eta_{t}} \right] - \mathbb{E} \left[\frac{\|\mathbf{w}_{t+1} - \mathbf{w}^{\star}\|^{2}}{2\eta_{t}} \right]
$$

$$
+ 8A^{2}\eta_{t} + \mathbb{E} \left[\|\mathbf{w}_{t} - \mathbf{w}^{\star}\| \|\delta g_{t} \right]
$$

Taking expectation on both sides once again, we achieve

$$
\mathbb{E}\left[\mathcal{J}\left(\mathbf{w}_{t}\right)-\mathcal{J}\left(\mathbf{w}^{\star}\right)\right] \leq \frac{\lambda}{2}\left(t-1\right)\mathbb{E}\left[\left\|\mathbf{w}_{t}-\mathbf{w}^{\star}\right\|^{2}\right] - \frac{\lambda}{2}t\mathbb{E}\left[\left\|\mathbf{w}_{t+1}-\mathbf{w}^{\star}\right\|^{2}\right] + 8A^{2}\eta_{t} + \mathbb{E}\left[\left\|\mathbf{w}_{t}-\mathbf{w}^{\star}\right\|\left\|\delta g_{t}\right\|\right]
$$

$$
\mathbb{E}\left[\mathcal{J}\left(\mathbf{w}_{t}\right)-\mathcal{J}\left(\mathbf{w}^{\star}\right)\right] \leq \frac{\lambda}{2}\left(t-1\right)\mathbb{E}\left[\left\|\mathbf{w}_{t}-\mathbf{w}^{\star}\right\|^{2}\right] - \frac{\lambda}{2}t\mathbb{E}\left[\left\|\mathbf{w}_{t+1}-\mathbf{w}^{\star}\right\|^{2}\right]
$$

$$
+ 8A^{2}\eta_{t} + \mathbb{E}\left[\left\|\mathbf{w}_{t}-\mathbf{w}^{\star}\right\|^{2}\right]^{1/2}\mathbb{E}\left[\left\|\delta g_{t}\right\|^{2}\right]^{1/2}
$$

Taking sum the above inequality when $t = 1, ..., T$, we obtain

$$
\mathbb{E}\left[\frac{1}{T}\sum_{t=1}^{T} \mathcal{J}\left(\mathbf{w}_{t}\right) - \mathcal{J}\left(\mathbf{w}^{\star}\right)\right] \leq \frac{8A^{2}}{\lambda} \sum_{t=1}^{T} \frac{1}{t} + \frac{1}{T} W \sum_{t=1}^{T} \mathbb{E}\left[M_{t}^{2}\right]^{1/2} \leq \frac{8A^{2} \left(\log T + 1\right)}{\lambda T} + \frac{1}{T} W \sum_{t=1}^{T} \mathbb{E}\left[M_{t}^{2}\right]^{1/2}
$$

Here we note that

$$
\|\delta g_t\| = \left\| \left[\nabla_o l\left(y_t, f_t\left(\boldsymbol{x}_t\right)\right) - \nabla_o l\left(y_t, f_t^h\left(\boldsymbol{x}_t\right)\right) \right] \Phi\left(\boldsymbol{x}_t\right) \right\| = |M_t|
$$

The last conclusion comes from the convexity of the function $J(.)$.

 \Box

Theorem 7. Assume that $l(y, o)$ is a γ -strongly smooth loss function. With a probability at least $1 - \theta$ *, the following statements hold*

$$
\begin{array}{ll}\ni) & \mathbb{E}\left[\mathcal{J}\left(\overline{\mathbf{w}}_{T}\right)-\mathcal{J}\left(\mathbf{w}^{\star}\right)\right] & \leq \qquad \mathbb{E}\left[\frac{1}{T}\sum_{t=1}^{T}\mathcal{J}\left(\mathbf{w}_{t}\right)-\mathcal{J}\left(\mathbf{w}^{\star}\right)\right] & \leq \qquad \frac{8A^{2}(\log T+1)}{\lambda T} + \frac{1}{T}W\gamma\varepsilon\sum_{t=1}^{T}\left(\frac{\sum_{i=1}^{t}\mu_{i}}{t}\right)^{1/2} \\
ii) & \mathbb{E}\left[\mathcal{J}\left(\overline{\mathbf{w}}_{T}\right)-\mathcal{J}\left(\mathbf{w}^{\star}\right)\right] \leq \mathbb{E}\left[\frac{1}{T}\sum_{t=1}^{T}\mathcal{J}\left(\mathbf{w}_{t}\right)-\mathcal{J}\left(\mathbf{w}^{\star}\right)\right] \leq \frac{8A^{2}(\log T+1)}{\lambda T} + W\gamma\varepsilon \\
where \quad \theta = 2^{8}\left(\frac{\sigma_{\mu}A d_{\mathcal{X}}}{\lambda\varepsilon}\right)\exp\left(-\frac{D\lambda^{2}\varepsilon^{2}}{4(M+2)A^{2}}\right).\n\end{array}
$$

Proof. From the smoothness of the loss function, we have

$$
\left|\nabla_{o}l\left(y_{t},f_{t}\left(\boldsymbol{x}_{t}\right)\right)-\nabla_{o}l\left(y_{t},f_{t}^{h}\left(\boldsymbol{x}_{t}\right)\right)\right|\leq\gamma\left|f_{t}\left(\boldsymbol{x}_{t}\right)-f_{t}^{h}\left(\boldsymbol{x}_{t}\right)\right|
$$

Referring to Lemma [2,](#page-1-3) with a probability at least $1-2^8\left(\frac{\sigma_\mu A d_\mathcal{X}}{\lambda \varepsilon}\right) \exp\left(-\frac{D\lambda^2 \varepsilon^2}{4(M+2)\varepsilon^2}\right)$ $\frac{D\lambda^2\varepsilon^2}{4(M+2)A^2}\Big) = 1 - \theta$ we have

$$
|M_t| \le \gamma A^{-1} \lambda \varepsilon \sum_{j=1}^t |\alpha_j| \beta_i \le \gamma A^{-1} \lambda \varepsilon \sum_{j=1}^t \frac{A}{\lambda t} \beta_j = \frac{\gamma \varepsilon}{t} \sum_{j=1}^t \beta_j
$$

$$
M_t^2 \le \frac{\gamma^2 \varepsilon^2}{t^2} \left(\sum_{j=1}^t \beta_j\right)^2 \le \frac{\gamma^2 \varepsilon^2}{t} \sum_{j=1}^t \beta_j^2 = \frac{\gamma^2 \varepsilon^2}{t} \sum_{j=1}^t \beta_j \quad \text{(since } \beta_i = 0 \text{ or } 1\text{)}
$$

$$
\mathbb{E}\left[M_t^2\right] \le \frac{\gamma^2 \varepsilon^2}{t} \left(\sum_{j=1}^t \mu_j\right)
$$

and $|M_t| \leq \gamma \varepsilon$. Therefore, with a probability at least $1 - \theta$ we achieve

$$
\mathbb{E}\left[\mathcal{J}\left(\overline{\mathbf{w}}_{T}\right) - \mathcal{J}\left(\mathbf{w}^{\star}\right)\right] \leq \mathbb{E}\left[\frac{1}{T}\sum_{t=1}^{T}\mathcal{J}\left(\mathbf{w}_{t}\right) - \mathcal{J}\left(\mathbf{w}^{\star}\right)\right]
$$
\n
$$
\leq \frac{8A^{2}\left(\log T + 1\right)}{\lambda T} + \frac{1}{T}W\gamma\varepsilon\sum_{t=1}^{T}\frac{\left(\sum_{j=1}^{t}\mu_{j}\right)^{1/2}}{t^{1/2}}
$$
\n
$$
\leq \frac{8A^{2}\left(\log T + 1\right)}{\lambda T} + \frac{1}{T}W\gamma\varepsilon\sum_{t=1}^{T}\left(\frac{\sum_{j=1}^{t}\mu_{j}}{t}\right)^{1/2}
$$

and

$$
\mathbb{E}\left[\mathcal{J}\left(\overline{\mathbf{w}}_{T}\right) - \mathcal{J}\left(\mathbf{w}^{\star}\right)\right] \leq \mathbb{E}\left[\frac{1}{T}\sum_{t=1}^{T}\mathcal{J}\left(\mathbf{w}_{t}\right) - \mathcal{J}\left(\mathbf{w}^{\star}\right)\right] \leq \frac{8A^{2}\left(\log T + 1\right)}{\lambda T} + \frac{1}{T}W\sum_{t=1}^{T}\gamma\varepsilon \leq \frac{8A^{2}\left(\log T + 1\right)}{\lambda T} + W\gamma\varepsilon
$$

 \Box

3 Computational Complexities of DualSGD and FOGD

We compare the computational complexities of our proposed DualSGD and Fourier Online Gradient Descent (FOGD) [\[2\]](#page-6-2). Recall that M and D denote the dimensions of input space and feature space, and B the budget size. There are four operators: (i) random feature mapping; (ii) kernel function; (iii) sorting coefficients of support vectors and (iv) prediction. The random feature mapping first projects the input data vector to random feature space with $O(MD)$ computational complexity, and then compute *sin*, *cos* on the random feature dimension with $O(2D * 2^{\log n} n \log^2 n)$ where *n* is the number of bits accuracy [\[1\]](#page-6-3). The kernel function, sorting coefficients and prediction operate in $\mathcal{O}(MB)$, $\mathcal{O}(B \log B)$ and $\mathcal{O}(D)$ complexity, respectively. The FOGD performs random feature mapping and prediction whilst the DualSGD performs all four operators.

Let D_1 and D_2 denote the number of random features of FOGD and DualSGD. The computational complexities of FOGD and DualSGD reads

$$
\mathcal{O}_{\text{FOGD}} = \mathcal{O}\left(MD_1 + 2D_1 * 2^{\log n} n \log^2 n + D_1\right) = U\left(MD_1 + 2D_1 * 2^{\log n} n \log^2 n + D_1\right)
$$

$$
\mathcal{O}_{\text{DualSGD}} = \mathcal{O}\left(MD_2 + 2D_2 * 2^{\log n} n \log^2 n + D_2 + MB + B \log B\right)
$$

$$
= V\left(MD_2 + 2D_2 * 2^{\log n} n \log^2 n + D_2 + MB + B \log B\right)
$$

where U, V are the number of iterations.

Taking the subtraction of $\mathcal{O}_{\text{FOGD}}$ and $\mathcal{O}_{\text{DualSGD}}$, we obtain:

$$
\hat{\mathcal{O}} = \mathcal{O}_{\text{FOGD}} - \mathcal{O}_{\text{DuatsGD}}
$$

= $M (UD_1 - VD_2 - B) + (UD_1 - VD_2) (2 * 2^{\log n} n \log^2 n + 1) - B \log B$

According Fig. 1 in the introduction section, $D_1 \gg D_2$ and $D_1 \gg B$, thus $D_1 - D_2 \gg B$. In addition, we assume that $U = V$ and normally use double-precision floating-point with $n = 64$ (bits) for storing and computing real number, thus $2 * 2^{\log n} n \log^2 n + 1 > \log B$. Finally, we can see that $\mathcal{O} \gg 0$, thus the computational complexity of DualSGD, in practice, is significantly lower than that of FOGD.

References

- [1] R. P Brent and P. Zimmermann. *Modern computer arithmetic*, volume 18. Cambridge University Press, 2010.
- [2] J. Lu, S. C.H. Hoi, J. Wang, P. Zhao, and Z.-Y. Liu. Large scale online kernel learning. *J. Mach. Learn. Res.*, 2015.
- [3] A. Rahimi and B. Recht. Random features for large-scale kernel machines. In *Advances in Neural Infomration Processing Systems*, 2007.
- [4] S. Shalev-Shwartz and T. Zhang. Stochastic dual coordinate ascent methods for regularized loss. *Journal of Machine Learning Research*, 14(1):567–599, 2013.