

# Supplementary Materials for Coreset Stochastic Variance-Reduced Gradient with Application to Optimal Margin Distribution Machine

Zhi-Hao Tan and Teng Zhang and Wei Wang

National Key Laboratory for Novel Software Technology, Nanjing University, Nanjing 210023, China  
Collaborative Innovation Center of Novel Software Technology and Industrialization, Nanjing 210023, China  
{tanzh, zhangt, wangw}@lamda.nju.edu.cn

In the supplementary material, we will give the detailed proofs of lemma and theorems in main paper. Without loss of generality, we assume that  $\|\phi(x)\| = K(x, x)^{1/2} = 1$ ,  $\forall x \in \mathcal{D}$ . And we only consider the binary classification, so the label  $y$  is either  $-1$  or  $1$  which implies  $|y| = y^2 = 1$ .

## Bounded Approximation Error

In this section, we will present the detailed proofs of Theorem 1, Lemma 2 and Lemma 3, i.e., the intermediate results when we give the upper bound of approximation error.

**Theorem 1.** *If  $\mathbf{w}_{odm}^*$  is the optimal solution of ODM, then there exists a positive constant  $H$  such that  $\|\mathbf{w}_{odm}^*\| \leq H$ , where  $H = \frac{\sqrt{\lambda(1-\theta)^2 + \lambda\mu(1+\theta)^2}}{1-\theta}$ . Moreover, for squared hinge loss SVM, the optimal solution  $\mathbf{w}_{h^2}^*$  satisfies  $\mathbf{w}_{h^2}^* \leq \lambda$ .*

### Proof. 1 For ODM

The optimization problem of ODM can be reformulated as follows:

$$\begin{aligned} \min_{\mathbf{w}, \boldsymbol{\xi}, \boldsymbol{\epsilon}} \quad & \frac{1}{2} \mathbf{w}^\top \mathbf{w} + \frac{\lambda}{m(1-\theta)^2} \boldsymbol{\xi}^\top \boldsymbol{\xi} + \frac{\lambda\mu}{m(1-\theta)^2} \boldsymbol{\epsilon}^\top \boldsymbol{\epsilon} \\ \text{s.t.} \quad & \mathbf{Y} \mathbf{x}^\top \mathbf{w} \geq (1-\theta) \mathbf{e} - \boldsymbol{\xi}, \\ & \mathbf{Y} \mathbf{x}^\top \mathbf{w} \leq (1+\theta) \mathbf{e} + \boldsymbol{\epsilon}. \end{aligned} \quad (1)$$

where  $\mathbf{x}$  is the matrix whose  $i$ -th column is  $\phi(\mathbf{x}_i)$ , i.e.,  $\mathbf{x} = [\phi(\mathbf{x}_1), \dots, \phi(\mathbf{x}_m)]$ ,  $\mathbf{Y}$  is a  $m \times m$  diagonal matrix with  $y_1, \dots, y_m$  as the diagonal elements and  $\mathbf{e}$  stands for the all-one vector.

Introduce the Lagrange multipliers  $\boldsymbol{\zeta} \geq 0$  and  $\boldsymbol{\beta} \geq 0$  for the two constraints respectively, the Lagrangian of (1) leads to

$$\begin{aligned} \mathcal{L}(\mathbf{w}, \boldsymbol{\xi}, \boldsymbol{\epsilon}, \boldsymbol{\alpha}, \boldsymbol{\beta}) \\ = \frac{1}{2} \mathbf{w}^\top \mathbf{w} + \frac{\lambda}{m(1-\theta)^2} \boldsymbol{\xi}^\top \boldsymbol{\xi} + \frac{\lambda\mu}{m(1-\theta)^2} \boldsymbol{\epsilon}^\top \boldsymbol{\epsilon} \\ - \boldsymbol{\zeta}^\top (\mathbf{Y} \mathbf{X}^\top \mathbf{w} - (1-\theta) \mathbf{e} + \boldsymbol{\xi}) \\ + \boldsymbol{\beta}^\top (\mathbf{Y} \mathbf{X}^\top \mathbf{w} - (1+\theta) \mathbf{e} - \boldsymbol{\epsilon}). \end{aligned} \quad (2)$$

By setting the partial derivative of  $\mathbf{w}$ ,  $\boldsymbol{\xi}$ ,  $\boldsymbol{\epsilon}$  to zero, we have

$$\begin{aligned} \frac{\partial \mathcal{L}}{\partial \mathbf{w}} = \mathbf{w} - \mathbf{X} \mathbf{Y} \boldsymbol{\zeta} + \mathbf{X} \mathbf{Y} \boldsymbol{\beta} = 0 & \implies \mathbf{w} = \mathbf{X} \mathbf{Y} (\boldsymbol{\zeta} - \boldsymbol{\beta}) \\ \frac{\partial \mathcal{L}}{\partial \boldsymbol{\xi}} = \frac{2\lambda}{m(1-\theta)^2} \boldsymbol{\xi} - \boldsymbol{\zeta} = 0 & \implies \boldsymbol{\xi} = \frac{m(1-\theta)^2}{2\lambda} \boldsymbol{\zeta} \\ \frac{\partial \mathcal{L}}{\partial \boldsymbol{\epsilon}} = \frac{2\lambda\mu}{m(1-\theta)^2} \boldsymbol{\epsilon} - \boldsymbol{\beta} = 0 & \implies \boldsymbol{\epsilon} = \frac{m(1-\theta)^2}{2\lambda\mu} \boldsymbol{\beta} \end{aligned}$$

Substituting the above to (2), we gain the dual form

$$\begin{aligned} \mathcal{G}(\boldsymbol{\alpha}) = -\frac{1}{2} \mathbf{w}^\top \mathbf{w} - \frac{m(1-\theta)^2}{4\lambda} \boldsymbol{\zeta}^\top \boldsymbol{\zeta} - \frac{m(1-\theta)^2}{4\lambda\mu} \boldsymbol{\beta}^\top \boldsymbol{\beta} \\ + (1-\theta) \boldsymbol{\zeta}^\top \mathbf{e} - (1+\theta) \boldsymbol{\beta}^\top \mathbf{e} \end{aligned} \quad (3)$$

Let us denote  $(\mathbf{w}^*, \boldsymbol{\xi}^*, \boldsymbol{\epsilon}^*)$  and  $(\boldsymbol{\zeta}^*, \boldsymbol{\beta}^*)$  be the primal and dual solutions of (1) and (3), respectively. Since the strong duality holds, we have

$$\begin{aligned} \frac{1}{2} \|\mathbf{w}^*\|^2 + \frac{\lambda}{m(1-\theta)^2} \sum_{i=1}^m (\xi_i^{*2} + \mu \epsilon_i^{*2}) \\ = -\frac{1}{2} \|\mathbf{w}^*\|^2 - \frac{m(1-\theta)^2}{4\lambda} \sum_{i=1}^m \zeta_i^{*2} - \frac{m(1-\theta)^2}{4\lambda\mu} \sum_{i=1}^m \beta_i^{*2} \\ + (1-\theta) \sum_{i=1}^m \zeta_i^* - (1+\theta) \sum_{i=1}^m \beta_i^* \end{aligned}$$

Since the loss of ODM is non-negative, we have

$$\begin{aligned} \|\mathbf{w}^*\|^2 \leq -\frac{m(1-\theta)^2}{4\lambda} \sum_{i=1}^m \zeta_i^{*2} - \frac{m(1-\theta)^2}{4\lambda\mu} \sum_{i=1}^m \beta_i^{*2} \\ + (1-\theta) \sum_{i=1}^m \zeta_i^* - (1+\theta) \sum_{i=1}^m \beta_i^* \\ = \sum_{i=1}^m \left( (1-\theta) \zeta_i^* - \frac{m(1-\theta)^2}{4\lambda} \zeta_i^{*2} \right) \\ + \sum_{i=1}^m \left( (1+\theta) \beta_i^* - \frac{m(1-\theta)^2}{4\lambda\mu} \beta_i^{*2} \right) \\ \leq \sum_{i=1}^m \frac{\lambda}{m} + \sum_{i=1}^m \frac{\lambda\mu(1+\theta)^2}{m(1-\theta)^2} \end{aligned}$$

$$= \frac{\lambda(1-\theta)^2 + \lambda\mu(1+\theta)^2}{(1-\theta)^2} = H^2$$

We note that the second inequality uses that the maximum of the quadratic concave function is obtained when  $\zeta_i^* = \frac{2\lambda}{m(1-\theta)}$  and  $\beta_i^* = \frac{2\lambda\mu(1+\theta)}{m(1-\theta)^2}$ . Therefore, for ODM,  $\|\mathbf{w}_{odm}^*\| \leq H$ , where  $H = \frac{\sqrt{\lambda(1-\theta)^2 + \lambda\mu(1+\theta)^2}}{1-\theta}$ .

## 2 For Square Hinge Loss SVM

The optimization of square hinge loss SVM is as follows:

$$\begin{aligned} \min_{\mathbf{w}, \boldsymbol{\xi}} \quad & \frac{1}{2} \|\mathbf{w}\|^2 + \frac{\lambda}{m} \sum_{i=1}^m \xi_i^2 \\ \text{s.t.} \quad & y_i \mathbf{w}^\top \phi(\mathbf{x}_i) \geq 1 - \xi_i, \quad i = 1, \dots, m \end{aligned} \quad (4)$$

where  $\boldsymbol{\xi} = [\xi_i]_{i=1}^m$ . The Lagrange function is of the following form

$$\begin{aligned} \mathcal{L}(\mathbf{w}, \boldsymbol{\xi}, \boldsymbol{\alpha}) = & \frac{1}{2} \|\mathbf{w}\|^2 + \frac{\lambda}{m} \sum_{i=1}^m \xi_i^2 \\ & + \sum_{i=1}^m \alpha_i (1 - \xi_i - y_i \mathbf{w}^\top \phi(\mathbf{x}_i)) \end{aligned} \quad (5)$$

By setting the partial derivative of  $\mathbf{w}$ ,  $\boldsymbol{\xi}$  to zero, we have

$$\begin{aligned} \frac{\partial \mathcal{L}}{\partial \mathbf{w}} = \mathbf{w} - \sum_{i=1}^m \alpha_i y_i \phi(\mathbf{x}_i) = 0 & \implies \mathbf{w} = \sum_{i=1}^m \alpha_i y_i \phi(\mathbf{x}_i) \\ \frac{\partial \mathcal{L}}{\partial \xi_i} = \frac{2\lambda}{m} \xi_i - \alpha_i = 0 & \implies \xi_i = \frac{m}{2\lambda} \alpha_i \end{aligned}$$

Substituting the above to (5), we gain the dual form

$$\mathcal{W}(\boldsymbol{\alpha}) = -\frac{1}{2} \|\mathbf{w}\|^2 - \frac{m}{4\lambda} \sum_{i=1}^m \alpha_i^2 + \sum_{i=1}^m \alpha_i \quad (6)$$

Let us denote  $(\mathbf{w}^*, \boldsymbol{\xi}^*)$  and  $\boldsymbol{\alpha}^*$  be the primal and dual solutions of (4) and (5) respectively. Since the strong duality holds, we have

$$\frac{1}{2} \|\mathbf{w}^*\|^2 + \frac{\lambda}{m} \sum_{i=1}^m \xi_i^{*2} = -\frac{1}{2} \|\mathbf{w}^*\|^2 - \frac{m}{4\lambda} \sum_{i=1}^m \alpha_i^{*2} + \sum_{i=1}^m \alpha_i^*$$

Since the loss of square hinge loss SVM is non-negative, we have

$$\begin{aligned} \|\mathbf{w}^*\|^2 &= -\frac{m}{4\lambda} \sum_{i=1}^m \alpha_i^{*2} + \sum_{i=1}^m \alpha_i^* - \frac{\lambda}{m} \sum_{i=1}^m \xi_i^{*2} \\ &\leq \sum_{i=1}^m \left( -\frac{m}{4\lambda} \alpha_i^{*2} + \alpha_i^* \right) \\ &\leq \lambda \end{aligned}$$

Therefore, for square hinge loss SVM, we have  $\mathbf{w}_{h^2}^* \leq \lambda$ . In summary, we gain the conclusion. The proof of Theorem 1 is similar to that of Theorem 1 in (Shalev-Shwartz et al. 2011).  $\square$

**Lemma 2.** For ODM problem, the  $\alpha_t$  satisfies  $\alpha_t^2 \leq A^2$  for all  $t$ , where  $A = \frac{2(H+1+\theta)}{(1-\theta)^2}$ .

*Proof.* We first remind the representation of  $\alpha_t$

$$\begin{aligned} \alpha_t = & \frac{2\lambda}{(1-\theta)^2} \{ (y_t \mathbf{w}^\top \phi(\mathbf{x}_t) + \theta - 1) y_t \mathbb{I}(t \in I_1) \\ & + \mu (y_t \mathbf{w}^\top \phi(\mathbf{x}_t) - \theta - 1) y_t \mathbb{I}(t \in I_2) \} \end{aligned}$$

where  $I_1 \equiv \{i \mid y_t \mathbf{w}^\top \phi(\mathbf{x}_t) < 1 - \theta\}$  and  $I_2 \equiv \{i \mid y_t \mathbf{w}^\top \phi(\mathbf{x}_t) > 1 + \theta\}$ . Then we can obtain

$$\begin{aligned} \|\alpha_t\| &\leq \frac{2\lambda}{(1-\theta)^2} (\|y_t \mathbf{w}_t^\top \phi(\mathbf{x}_t) + \theta - 1\| \mathbb{I}(t \in I_2) \\ &\quad + \mu \|y_t \mathbf{w}_t^\top \phi(\mathbf{x}_t) - \theta - 1\| \mathbb{I}(t \in I_3)) \\ &\leq \frac{2\lambda}{(1-\theta)^2} (\|\mathbf{w}_t^\top \phi(\mathbf{x}_t)\| + \|\theta - 1\| \mathbb{I}(t \in I_2) \\ &\quad + \mu (\|\mathbf{w}_t^\top \phi(\mathbf{x}_t)\| + \|\theta + 1\| \mathbb{I}(t \in I_3))) \\ &\leq \frac{2\lambda}{(1-\theta)^2} (\|\mathbf{w}_t\| \|\phi(\mathbf{x}_t)\| + 1 - \theta) \mathbb{I}(t \in I_2) \\ &\quad + \mu (\|\mathbf{w}_t\| \|\phi(\mathbf{x}_t)\| + 1 + \theta) \mathbb{I}(t \in I_3)) \\ &= \frac{2\lambda}{(1-\theta)^2} (\|\mathbf{w}_t\| + 1 - \theta) \mathbb{I}(t \in I_2) \\ &\quad + \mu (\|\mathbf{w}_t\| + 1 + \theta) \mathbb{I}(t \in I_3)) \\ &\leq \frac{2\lambda}{(1-\theta)^2} (\|\mathbf{w}_t\| + 1 + \theta) \\ &\leq \frac{2\lambda(H+1+\theta)}{(1-\theta)^2} \end{aligned}$$

Note that the first and second inequalities use Minkowski inequality, the third inequality uses Cauchy-Schwarz inequality, and the fourth inequality uses that  $\mu \leq 1$  and  $0 \leq \theta \leq 1$ . Therefore, we have  $\alpha_t^2 \leq A^2$ , where  $A = \frac{2(H+1+\theta)}{(1-\theta)^2}$ .  $\square$

**Lemma 3.** For hinge loss SVM and square hinge loss SVM, the  $\alpha_t$  satisfies  $\alpha_t^2 \leq A^2 = \lambda^2 B^2$  for all  $t$ , where  $B = 1$  and  $B = 2\lambda + 2$  respectively.

*Proof.* For hinge loss SVM, we have

$$\begin{aligned} l(\mathbf{w}; \mathbf{x}, y) &= \max \{0, 1 - y \mathbf{w}^\top \phi(\mathbf{x})\} \\ l'(\mathbf{w}; \mathbf{x}, y) &= -\mathbb{I}_{\{y \mathbf{w}^\top \phi(\mathbf{x}) \leq 1\}} y \phi(\mathbf{x}) \end{aligned}$$

where  $\mathbb{I}_S$  is the indicator function, which equals 1 if the logical statement  $S$  is true and 0 otherwise. Therefore, by taking  $B = 1$ , we have

$$\|l'(\mathbf{w}; \mathbf{x}, y)\| \leq \|\phi(\mathbf{x})\| \leq 1 = B$$

For squared hinge loss SVM, we have

$$\begin{aligned} l(\mathbf{w}; \mathbf{x}, y) &= \max \{0, 1 - y \mathbf{w}^\top \phi(\mathbf{x})\}^2 \\ l'(\mathbf{w}; \mathbf{x}, y) &= -\mathbb{I}_{\{y \mathbf{w}^\top \phi(\mathbf{x}) \leq 1\}} 2y \phi(\mathbf{x}) (1 - y \mathbf{w}^\top \phi(\mathbf{x})) \end{aligned}$$

Then by taking  $B = 2\lambda + 2$  we can obtain

$$\begin{aligned} \|l'(\mathbf{w}; \mathbf{x}, y)\| &\leq \|2y\phi(\mathbf{x})(1 - y\mathbf{w}^\top\phi(\mathbf{x}))\| \\ &\leq 2\|1 - y\mathbf{w}^\top\phi(\mathbf{x})\|\|\phi(\mathbf{x})\|\|y\| \\ &\leq 2\|y\mathbf{w}^\top\phi(\mathbf{x})\| + 2 \\ &\leq 2\|\phi(\mathbf{x})\|\|\mathbf{w}\|\|y\| + 2 \\ &= 2\|\mathbf{w}\| + 2 \\ &\leq 2\lambda + 2 = B \end{aligned}$$

The last inequality uses Theorem 1.

Therefore, there exists a positive constant  $B$  such that  $\|l'(\mathbf{w}; \mathbf{x}, y)\| \leq B$  for hinge loss SVM and square hinge loss SVM. Based on this, we have the following proof.

Note that  $\|\phi(\mathbf{x})\| = K(x, x) = 1$ , by taking  $A = \lambda B$ , we have

$$\alpha_t^2 = \alpha_t^2 K(\mathbf{x}_t, \mathbf{x}_t) = \lambda^2 \|l'(\mathbf{w}_t; \mathbf{x}_t, y_t)\|^2 \leq \lambda^2 B^2 = A^2$$

Therefore we gain the conclusion  $\alpha_t^2 \leq A^2$ .  $\square$

Based on the above results, we obtain the Theorem 4, whose proof is in the main paper.

**Theorem 4.** Assume that the p.s.d. and isotropic kernel  $K(\mathbf{x}_i, \mathbf{x}_j) = k(\|\mathbf{x}_i - \mathbf{x}_j\|^2)$  is used, where  $k(\cdot)$  is a monotonically continuous decreasing function with  $k(0) = 1$ . And let  $\delta$  be the diameter of coresets coverage. Then for the approximation error  $\Delta_t$  as indicated in (10), we have  $\|\Delta_t\| \leq \frac{3}{2}A\delta_\phi$ , where  $\delta_\phi = 2\sqrt{2(1 - \kappa(\delta^2/4))}$ .

## Convergence Analysis

In this section, we will present the proof of intermediate conclusions when analysing the convergence of CSVRG. The proofs of Lemma 6,7 are similar to the proof of Lemma 21, Lemma 22 in (Le et al. 2017).

**Lemma 5.** When using CSVRG to train ODM, we have  $\mathbb{E}[\|\mathbf{w}_t - \mathbf{w}^*\|^2] \leq W^2$  for all  $t$ , where  $W = 2H$ .

*Proof.* According to Theorem 1, by taking  $W = 2H$ , we have the following

$$\begin{aligned} \mathbb{E}[\|\mathbf{w}_t - \mathbf{w}^*\|^2] &\leq 2\mathbb{E}[\|\mathbf{w}_t\|^2] + 2\mathbb{E}[\|\mathbf{w}^*\|^2] \\ &\leq 4H^2 = W^2 \end{aligned}$$

$\square$

**Lemma 6.** When using CSVRG to train hinge loss SVM or square hinge loss SVM, there exists a positive constant  $P$  such that  $\mathbb{E}[\|\mathbf{w}_t\|^2] \leq P^2$  for all  $t$ , where  $P = 2A + \frac{3}{2}A\delta_\phi$ .

*Proof.* We prove by induction that  $\mathbb{E}[\|\mathbf{w}_t\|^2] \leq P^2$  where  $P = 2A + \frac{3}{2}A\delta_\phi$  for all  $t$ . Assume that the claim is holding for  $t - 1$ , we have the detailed proof in Figure 1.

The third inequality uses  $\sqrt{\mathbb{E}\|\mathbf{a} + \mathbf{b}\|^2} \leq \sqrt{\mathbb{E}\|\mathbf{a}\|^2} + \sqrt{\mathbb{E}\|\mathbf{b}\|^2}$ . The fourth inequality uses Lemma 3. And the fifth inequality uses  $\mathbb{E}\|\xi - \mathbb{E}\xi\|^2 = \mathbb{E}\|\xi\|^2 - \|\mathbb{E}\xi\|^2 \leq \mathbb{E}\|\xi\|^2$  for any random vector  $\xi$ . Finally, the sixth inequality uses Theorem 4.  $\square$

**Lemma 7.** Assume that  $f(\mathbf{w})$  is  $\nu$ -strongly convex, when using CSVRG to train hinge loss SVM or square hinge loss SVM, we have  $\mathbb{E}[\|\mathbf{w}_t - \mathbf{w}^*\|^2] \leq W^2$  for all  $t$ , where  $W = \frac{3A\delta_\phi + \sqrt{9A^2\delta_\phi^2 + 16(1-\eta\nu)P^2}}{2\nu}$ .

*Proof.* We first remind the definitions of the relevant variables

$$\begin{aligned} \mathbf{g} &= \mathbf{w} + \lambda l'(\mathbf{w}; \mathbf{x}_t, y_t) = \mathbf{w} + \alpha_t \phi(\mathbf{x}_t) \\ \mathbf{v}_t &= \nabla\psi_t(\mathbf{w}_{t-1}) - \nabla\psi_t(\tilde{\mathbf{w}}) + \tilde{\boldsymbol{\mu}} \\ &= \mathbf{w}_{t-1} + (\alpha_t - \tilde{\alpha}_t)\phi(\mathbf{x}_t) + \frac{1}{m} \sum_{i=1}^m \tilde{\alpha}_i \phi(\mathbf{x}_i) \\ \Delta_t &= (\alpha_t - \tilde{\alpha}_t)[\phi(\mathbf{c}_t) - \phi(\mathbf{x}_t)] \\ &\quad + \frac{1}{m} \sum_{i=1}^m \tilde{\alpha}_i [\phi(\mathbf{c}_i) - \phi(\mathbf{x}_i)] \\ \mathbf{h}_t &= \mathbf{v}_t + \Delta_t \end{aligned}$$

Let  $\mathbf{d}_t = \mathbf{h}_t - \mathbf{w}_{t-1}$ , according to the proof of lemma 2, we have

$$\begin{aligned} (1 - \eta) \sqrt{\mathbb{E}\|\mathbf{w}_{t-1}\|^2} + \eta \sqrt{\mathbb{E}\|\mathbf{d}_t\|^2} &\leq P \\ \implies \sqrt{\mathbb{E}\|\mathbf{d}_t\|^2} &\leq P \end{aligned}$$

Therefore, we can obtain

$$\sqrt{\mathbb{E}\|\mathbf{h}_t\|^2} \leq \sqrt{\mathbb{E}\|\mathbf{w}_{t-1}\|^2} + \sqrt{\mathbb{E}\|\mathbf{d}_t\|^2} \leq 2P \quad (7)$$

Conditioned on  $\mathbf{w}_{t-1}$ , we have  $\mathbb{E}[\mathbf{v}_t] = \nabla f(\mathbf{w}_{t-1})$ . Then we can obtain

$$\begin{aligned} \mathbb{E}[\|\mathbf{w}_t - \mathbf{w}^*\|^2] &= \mathbb{E}\left[\left\|\prod_S (\mathbf{w}_{t-1} - \eta\mathbf{h}_t) - \mathbf{w}^*\right\|^2\right] \\ &\leq \mathbb{E}[\|\mathbf{w}_{t-1} - \eta\mathbf{h}_t - \mathbf{w}^*\|^2] \\ &= \|\mathbf{w}_{t-1} - \mathbf{w}^*\|^2 + \eta^2 \mathbb{E}[\|\mathbf{h}_t\|^2] - 2\eta \mathbb{E}[\langle \mathbf{w}_{t-1} - \mathbf{w}^*, \mathbf{v}_t \rangle] \\ &\quad - 2\eta \mathbb{E}[\langle \mathbf{w}_{t-1} - \mathbf{w}^*, \Delta_t \rangle] \\ &= \|\mathbf{w}_{t-1} - \mathbf{w}^*\|^2 + \eta^2 \mathbb{E}[\|\mathbf{h}_t\|^2] - 2\eta (\mathbf{w}_{t-1} - \mathbf{w}^*)^\top \mathbb{E}[\mathbf{v}_t] \\ &\quad + 2\eta \mathbb{E}[\|\mathbf{w}_{t-1} - \mathbf{w}^*\|^2 \|\Delta_t\|^2]^{1/2} \\ &\leq \|\mathbf{w}_{t-1} - \mathbf{w}^*\|^2 + \eta^2 \mathbb{E}[\|\mathbf{h}_t\|^2] - \eta\nu \|\mathbf{w}_{t-1} - \mathbf{w}^*\|^2 \\ &\quad + 2\eta \mathbb{E}[\|\mathbf{w}_{t-1} - \mathbf{w}^*\|^2 \|\Delta_t\|^2]^{1/2} \end{aligned}$$

where  $S = \mathcal{B}(\mathbf{0}, \lambda)$ . The last inequality uses the  $\nu$ -strongly convexity of  $f(\mathbf{w})$  as follows

$$\begin{aligned} &(\mathbf{w}_t - \mathbf{w}^*)^\top \nabla f(\mathbf{w}_t) \\ &\geq f(\mathbf{w}_t) - f(\mathbf{w}^*) + \frac{\lambda}{2} \|\mathbf{w}_t - \mathbf{w}^*\|^2 \\ &\geq \frac{\lambda}{2} \|\mathbf{w}_t - \mathbf{w}^*\|^2 \end{aligned}$$

$$\begin{aligned}
\sqrt{\mathbb{E} [\|\mathbf{w}_t\|^2]} &\leq \sqrt{\mathbb{E} \left[ \left\| \prod_S (\mathbf{w}_{t-1} - \eta \mathbf{h}_t) \right\|^2 \right]} \leq \sqrt{\mathbb{E} [\|\mathbf{w}_{t-1} - \eta \mathbf{h}_t\|^2]} \\
&= \sqrt{\mathbb{E} [\|\mathbf{w}_{t-1} - \eta (\mathbf{w}_{t-1} + \lambda \ell'(\mathbf{w}_{t-1}; x_t, y_t) - \tilde{\mathbf{g}}_t + \tilde{\boldsymbol{\mu}} + \Delta_t)\|^2]} \\
&\leq (1 - \eta) \sqrt{\mathbb{E} [\|\mathbf{w}_{t-1}\|^2]} + \eta \lambda \sqrt{\mathbb{E} [\|\ell'(\mathbf{w}_{t-1}; x_t, y_t)\|^2]} + \eta \sqrt{\mathbb{E} [\|\tilde{\boldsymbol{\mu}} - \tilde{\mathbf{g}}_t\|^2]} \\
&\quad + \eta \sqrt{\mathbb{E} [\|\Delta_t\|^2]} \\
&\leq (1 - \eta) \sqrt{\mathbb{E} [\|\mathbf{w}_{t-1}\|^2]} + \eta A + \eta \sqrt{\mathbb{E} \left\| \tilde{\boldsymbol{\alpha}}_t \phi(\mathbf{x}_t) - \frac{1}{m} \sum_{i=1}^m \tilde{\boldsymbol{\alpha}}_i \phi(\mathbf{x}_t) \right\|^2} \\
&\quad + \eta \sqrt{\mathbb{E} [\|\Delta_t\|^2]} \\
&\leq (1 - \eta) \sqrt{\mathbb{E} [\|\mathbf{w}_{t-1}\|^2]} + \eta A + \eta \sqrt{\mathbb{E} \|\tilde{\boldsymbol{\alpha}}_t \phi(\mathbf{x}_t)\|^2} + \eta \sqrt{\mathbb{E} [\|\Delta_t\|^2]} \\
&\leq (1 - \eta) \sqrt{\mathbb{E} [\|\mathbf{w}_{t-1}\|^2]} + 2\eta A + \frac{3}{2} \eta A \delta_\phi \\
&\leq (1 - \eta) P + 2\eta A + \frac{3}{2} \eta A \delta_\phi = P
\end{aligned}$$

Figure 1: Proof of Lemma 6

The last inequality uses  $f(\mathbf{w}_t) - f(\mathbf{w}^*) \geq 0$ . Then by taking expectation again and substituting (4), we have

$$\begin{aligned}
&\mathbb{E} [\|\mathbf{w}_t - \mathbf{w}^*\|^2] \\
&\leq (1 - \eta\nu) \mathbb{E} [\|\mathbf{w}_{t-1} - \mathbf{w}^*\|^2] + 4\eta^2 P^2 \\
&\quad + 2\eta \mathbb{E} [\|\mathbf{w}_{t-1} - \mathbf{w}^*\|^2 \|\Delta_t\|^2]^{1/2} \\
&\leq (1 - \eta\nu) \mathbb{E} [\|\mathbf{w}_{t-1} - \mathbf{w}^*\|^2] + 4\eta^2 P^2 \\
&\quad + 3\eta A \delta_\phi \mathbb{E} [\|\mathbf{w}_{t-1} - \mathbf{w}^*\|^2]^{1/2}
\end{aligned}$$

We prove by induction in  $t$ . Choosing  $W = \frac{3A\delta_\phi + \sqrt{9A^2\delta_\phi^2 + 16(1-\eta\nu)P^2}}{2\nu}$ , which is the solution of following equation

$$(1 - \eta\nu) W^2 + 4\eta^2 P^2 + 3\eta A \delta_\phi W = W^2$$

Therefore, assuming that  $\mathbb{E} [\|\mathbf{w}_{t-1} - \mathbf{w}^*\|^2] \leq W^2$ , we obtain

$$\mathbb{E} [\|\mathbf{w}_t - \mathbf{w}^*\|^2] \leq (1 - \eta\nu) W^2 + 4\eta^2 P^2 + 3\eta A \delta_\phi W = W^2$$

Therefore, we gain the conclusion  $\mathbb{E} [\|\mathbf{w}_t - \mathbf{w}^*\|^2] \leq W^2$  for all  $t$ .  $\square$

According to the Lemma 5, 6, 7, we can obtain the convergence rate as follows, the proof of which is in the main paper.

**Theorem 8.** Consider CSVRG in Algorithm 1 with option II and use it to solve SVMs and ODM. Assume that all  $\psi_i(\mathbf{w})$  are convex and  $L$ -smooth,  $f(\mathbf{w})$  is  $\nu$ -strongly convex. Let  $\mathbf{w}^* = \operatorname{argmin}_{\mathbf{w}} f(\mathbf{w})$ . Assume that  $T$  is sufficiently large so that

$$\rho = \frac{1}{\nu\eta(1 - 4L\eta)T} + \frac{4L\eta}{1 - 4L\eta} < 1$$

then we have linear convergence in expectation for CSVRG:

$$\mathbb{E}[f(\tilde{\mathbf{w}}_s) - f(\mathbf{w}^*)] \leq \rho^s \mathbb{E}[f(\tilde{\mathbf{w}}_0) - f(\mathbf{w}^*)] + \frac{1 - \rho^s}{1 - \rho} \Omega$$

where  $\Omega$  is a constant gap caused by coresset approximation, and  $\Omega \rightarrow 0$  when the radius of coverage approaches 0.

## References

- Le, T.; Nguyen, T. D.; Nguyen, V.; and Phung, D. 2017. Approximation vector machines for large-scale online learning. *The Journal of Machine Learning Research* 18(1):3962–4016.
- Shalev-Shwartz, S.; Singer, Y.; Srebro, N.; and Cotter, A. 2011. Pegasos: Primal estimated sub-gradient solver for svm. *Mathematical programming* 127(1):3–30.