

---

# Supplementary materials for An Exact Solver for the Weston-Watkins SVM Subproblem

---

## Contents

<b>A.1</b>	Regarding offsets . . . . .	<b>11</b>
<b>A.2</b>	Proof of Proposition 3.2 . . . . .	<b>11</b>
<b>A.3</b>	Proof of Proposition 3.5 . . . . .	<b>13</b>
<b>A.4</b>	Global linear convergence . . . . .	<b>14</b>
<b>A.5</b>	Proof of Theorem 3.4 . . . . .	<b>18</b>
<b>A.6</b>	Experiments . . . . .	<b>29</b>
<b>A.7</b>	Code availability . . . . .	<b>31</b>

---

### A.1. Regarding offsets

In this section, we review the literature on SVMs in particular with regard to offsets.

For binary kernel SVMs, Steinwart et al. (2011) demonstrates that kernel SVMs without offset achieve comparable classification accuracy as kernel SVMs with offset. Furthermore, they propose algorithms that solve kernel SVMs without offset that are significantly faster than solvers for kernel SVMs with offset.

For binary linear SVMs, Hsieh et al. (2008) introduced coordinate descent for the dual problem associated to linear SVMs without offsets, or with the bias term included in the  $\mathbf{w}$  term. Chiu et al. (2020) studied whether the method of Hsieh et al. (2008) can be extended to allow offsets, but found evidence that the answer is negative. For multiclass linear SVMs, Keerthi et al. (2008) studied block coordinate descent for the CS-SVM and WW-SVM, both without offsets. We are not aware of a multiclass analogue to Chiu et al. (2020) although the situation should be similar.

The previous paragraph discussed coordinate descent in relation to the offset. Including the offset presents challenges to primal methods as well. In Section 6 of Shalev-Shwartz et al. (2011), the authors argue that including an unregularized offset term in the primal objective leads to slower convergence guarantee. Furthermore, Shalev-Shwartz et al. (2011) observed that including an unregularized offset did not significantly change the classification accuracy.

The original Crammer-Singer (CS) SVM was proposed without offsets (Crammer & Singer, 2001). In Section VI of (Hsu & Lin, 2002), the authors show the CS-SVM with offsets do *not* perform better than CS-SVM without offsets. Furthermore, CS-SVM with offsets requires twice as many iterations to converge than without.

### A.2. Proof of Proposition 3.2

Below, let  $i \in [n]$  be arbitrary. First, we note that  $-\boldsymbol{\pi}' = \begin{bmatrix} -\mathbb{1}' \\ \mathbf{I}_{k-1} \end{bmatrix}$  and so

$$\boldsymbol{\pi}'\boldsymbol{\beta}_i = \begin{bmatrix} -\mathbb{1}'\boldsymbol{\beta}_i \\ \boldsymbol{\beta}_i \end{bmatrix}. \quad (14)$$

Now, let  $j \in [k]$ , we have by (3) that

$$[\boldsymbol{\alpha}_i]_j = [-\boldsymbol{\sigma}_{y_i}\boldsymbol{\pi}'\boldsymbol{\beta}_i]_j = [-\boldsymbol{\pi}'\boldsymbol{\beta}_i]_{\boldsymbol{\sigma}_{y_i}(j)}. \quad (15)$$

Note that if  $j \neq y_i$ , then  $\boldsymbol{\sigma}_{y_i}(j) \neq 1$  and so  $[\boldsymbol{\alpha}_i]_j = [-\boldsymbol{\pi}'\boldsymbol{\beta}_i]_{\boldsymbol{\sigma}_{y_i}(j)} = [\boldsymbol{\beta}_i]_{\boldsymbol{\sigma}_{y_i}(j)-1} \in [0, C]$ . On the other hand, if  $j = y_i$ , then  $\boldsymbol{\sigma}_{y_i}(y_i) = 1$  and  $[\boldsymbol{\alpha}_i]_{y_i} = [-\boldsymbol{\pi}'\boldsymbol{\beta}_i]_1 = -\mathbb{1}'\boldsymbol{\beta}_i = -\sum_{t \in [k-1]} [\boldsymbol{\beta}_i]_t = -\sum_{t \in [k]: t \neq y_i} [\boldsymbol{\beta}_i]_{\boldsymbol{\sigma}_{y_i}(t)-1} = -\sum_{t \in [k]: t \neq y_i} [\boldsymbol{\alpha}_i]_t$ . Thus,  $\boldsymbol{\alpha} \in \mathcal{F}$ . This proves that  $\Psi(\mathcal{G}) \subseteq \mathcal{F}$ .

Next, let us define another map  $\Xi : \mathcal{F} \rightarrow \mathbb{R}^{(k-1) \times n}$  as follows: For each  $\alpha \in \mathcal{F}$ , define  $\beta := \Xi(\alpha)$  block-wise by

$$\beta_i := \text{proj}_{2:k}(\sigma_{y_i} \alpha_i) \in \mathbb{R}^{k-1}$$

where

$$\text{proj}_{2:k} = \begin{bmatrix} \mathbf{0} & \mathbf{I}_{k-1} \end{bmatrix} \in \mathbb{R}^{(k-1) \times k}.$$

By construction, we have for each  $j \in [k-1]$  that  $[\beta_i]_j = [\sigma_{y_i} \alpha_i]_{j+1} = [\alpha_i]_{\sigma_{y_i}(j+1)}$ . Since  $j+1 \neq 1$  for any  $j \in [k-1]$ , we have that  $\sigma_{y_i}(j+1) \neq y_i$  for any  $j \in [k-1]$ . Thus,  $[\beta_i]_j = [\alpha_i]_{\sigma_{y_i}(j+1)} \in [0, C]$ . This proves that  $\Xi(\mathcal{F}) \subseteq \mathcal{G}$ .

Next, we prove that for all  $\alpha \in \mathcal{F}$  and  $\beta \in \mathcal{G}$ , we have  $\Xi(\Psi(\beta)) = \beta$  and  $\Psi(\Xi(\alpha)) = \alpha$ .

By construction, the  $i$ -th block of  $\Xi(\Psi(\beta))$  is given by

$$\begin{aligned} \text{proj}_{2:k}(\sigma_{y_i}(-\sigma_{y_i} \pi' \beta_i)) &= -\text{proj}_{2:k}(\sigma_{y_i} \sigma_{y_i} \pi' \beta_i) \\ &= -\text{proj}_{2:k}(\pi' \beta_i) \\ &= -\begin{bmatrix} \mathbf{0} & \mathbf{I}_{k-1} \end{bmatrix} \begin{bmatrix} \mathbf{1}' \\ -\mathbf{I}_{k-1} \end{bmatrix} \beta_i \\ &= \mathbf{I}_{k-1} \beta_i = \beta_i. \end{aligned}$$

For the second equality, we used the fact that  $\sigma_y^2 = \mathbf{I}$  for all  $y \in [k]$ . Thus,  $\Xi(\Psi(\beta)) = \beta$ .

Next, note that the  $i$ -th block of  $\Psi(\Xi(\alpha))$  is, by construction,

$$-\sigma_{y_i} \pi' \text{proj}_{2:k}(\sigma_{y_i} \alpha_i) = -\sigma_{y_i} \pi' \begin{bmatrix} \mathbf{0} & \mathbf{I}_{k-1} \end{bmatrix} \sigma_{y_i} \alpha_i = -\sigma_{y_i} \begin{bmatrix} \mathbf{0} & \pi' \end{bmatrix} \sigma_{y_i} \alpha_i \quad (16)$$

Recall that  $\pi' = \begin{bmatrix} \mathbf{1}' \\ -\mathbf{I}_{k-1} \end{bmatrix}$  and so  $\begin{bmatrix} \mathbf{0} & \pi' \end{bmatrix} = \begin{bmatrix} \mathbf{0} & \mathbf{1}' \\ \mathbf{0} & -\mathbf{I}_{k-1} \end{bmatrix}$ . Therefore,

$$\left[ \begin{bmatrix} \mathbf{0} & \pi' \end{bmatrix} \sigma_{y_i} \alpha_i \right]_1 = \sum_{j=2}^k [\sigma_{y_i} \alpha_i]_j = \sum_{j \in [k]: j \neq y_i} [\alpha_i]_j = -[\alpha_i]_{y_i} = -[\sigma_{y_i} \alpha_i]_1$$

and, for  $j = 2, \dots, k$ ,

$$\left[ \begin{bmatrix} \mathbf{0} & \pi' \end{bmatrix} \sigma_{y_i} \alpha_i \right]_j = -[\sigma_{y_i} \alpha_i]_j.$$

Hence, we have just shown that  $\begin{bmatrix} \mathbf{0} & \pi' \end{bmatrix} \sigma_{y_i} \alpha_i = -\sigma_{y_i} \alpha_i$ . Continuing from (16), we have

$$-\sigma_{y_i} \pi' \text{proj}_{2:k}(\sigma_{y_i} \alpha_i) = -\sigma_{y_i} (-\sigma_{y_i} \alpha_i) = \sigma_{y_i} \sigma_{y_i} \alpha_i = \alpha_i.$$

This proves that  $\Psi(\Xi(\alpha)) = \alpha$ . Thus, we have shown that  $\Psi$  and  $\Xi$  are inverses of one another. This proves that  $\Psi$  is a bijection.

Finally, we prove that

$$f(\Psi(\beta)) = g(\beta).$$

Recall that

$$f(\alpha) := \frac{1}{2} \sum_{i,s \in [n]} x'_s x_i \alpha'_i \alpha_s - \sum_{i \in [k]} \sum_{\substack{j \in [k]: \\ j \neq y_i}} \alpha_{ij}$$

Thus,

$$\alpha'_i \alpha_s = (-\sigma_{y_i} \pi' \beta_i)' (-\sigma_{y_s} \pi' \beta_s) = \beta'_i \pi \sigma_{y_i} \sigma'_{y_s} \pi' \beta_s$$

On the other hand, (3) implies that  $\sigma_{y_i} \alpha_i = -\pi' \beta_i$ . Hence

$$\sum_{j \in [k] \setminus \{y_i\}} \alpha_{ij} = \sum_{j \in [k]: j \neq 1} [\alpha_i]_{\sigma_{y_i}(j)} = \sum_{j \in [k]: j \neq 1} [\sigma_{y_i} \alpha_i]_j = \sum_{j \in [k]: j \neq 1} [-\pi' \beta_i]_j = \sum_{j \in [k-1]} [\beta_i]_j = \mathbf{1}' \beta_i.$$

Thus,

$$f(\alpha) := \frac{1}{2} \sum_{i,s \in [n]} x'_s x_i \alpha'_i \alpha_s - \sum_{i \in [k]} \sum_{\substack{j \in [k]: \\ j \neq y_i}} \alpha_{ij} = \frac{1}{2} \sum_{i,s \in [n]} x'_s x_i \beta'_i \pi \sigma_{y_i} \sigma'_{y_s} \pi' \beta_s - \sum_{i \in [k]} \mathbf{1}' \beta_i = g(\beta)$$

as desired. Finally, we note that  $\sigma_y = \sigma'_y$  for all  $y \in [k]$ . This concludes the proof of Proposition 3.2.  $\square$

### A.3. Proof of Proposition 3.5

We prove the following lemma which essentially unpacks the succinct Proposition 3.5:

**Lemma A.1.** *Recall the situation of Corollary 3.3: Let  $\beta \in \mathcal{G}$  and  $i \in [n]$ . Let  $\alpha = \Psi(\beta)$ . Consider*

$$\min_{\hat{\beta} \in \mathcal{G}} g(\hat{\beta}) \text{ such that } \hat{\beta}_s = \beta_s, \forall s \in [n] \setminus \{i\}. \quad (17)$$

Let  $\mathbf{w}$  be as in (1), i.e.,  $\mathbf{w} = -\sum_{i \in [n]} x_i \alpha'_i$ . Then a solution to (17) is given by  $[\beta_1, \dots, \beta_{i-1}, \tilde{\beta}_i, \beta_{i+1}, \dots, \beta_n]$  where  $\tilde{\beta}_i$  is a minimizer of

$$\min_{\hat{\beta}_i \in \mathbb{R}^{k-1}} \frac{1}{2} \hat{\beta}_i' \Theta \hat{\beta}_i - \hat{\beta}_i' ((1 - \pi \sigma_{y_i} \mathbf{w}' x_i) / \|x_i\|_2^2 + \Theta \beta_i) \text{ such that } 0 \leq \hat{\beta}_i \leq C.$$

Furthermore, the above optimization has a unique minimizer which is equal to the minimizer of (4) where

$$v := (1 - \rho_{y_i} \pi \mathbf{w}' x_i + \Theta \beta_i \|x_i\|_2^2) / \|x_i\|_2^2$$

and  $\mathbf{w}$  is as in (1).

*Proof.* First, we prove a simple identity:

$$\pi \pi' = \begin{bmatrix} \mathbb{1} & -\mathbf{I}_{k-1} \end{bmatrix} \begin{bmatrix} \mathbb{1}' \\ -\mathbf{I}_{k-1} \end{bmatrix} = \mathbf{I} + \mathbf{O} = \Theta. \quad (18)$$

Next, recall that by definition, we have

$$g(\beta) := \left( \frac{1}{2} \sum_{s,t \in [n]} x'_s x_t \beta'_t \pi \sigma_{y_t} \sigma_{y_s} \pi' \beta_s \right) - \left( \sum_{s \in [n]} \mathbb{1}' \beta_s \right).$$

Let us group the terms of  $g(\beta)$  that depends on  $\beta_i$ :

$$\begin{aligned} g(\beta) &= \frac{1}{2} x'_i x_i \beta'_i \pi \sigma_{y_i} \sigma_{y_i} \pi' \beta_i \\ &\quad + \frac{1}{2} \sum_{s \in [n]: s \neq i} x'_s x_i \beta'_i \pi \sigma_{y_i} \sigma_{y_s} \pi' \beta_s \\ &\quad + \frac{1}{2} \sum_{t \in [n]: t \neq i} x'_i x_t \beta'_t \pi \sigma_{y_t} \sigma_{y_i} \pi' \beta_i \\ &\quad + \frac{1}{2} \sum_{s,t \in [n]} x'_s x_t \beta'_t \pi \sigma_{y_t} \sigma_{y_s} \pi' \beta_s - \sum_{s \in [n]} \mathbb{1}' \beta_s \\ &= \frac{1}{2} x'_i x_i \beta'_i \Theta \beta_i \quad \because \sigma_{y_i}^2 = \mathbf{I} \text{ and (18)} \\ &\quad + \sum_{s \in [n]: s \neq i} x'_s x_i \beta'_i \pi \sigma_{y_i} \sigma_{y_s} \pi' \beta_s \\ &\quad - \mathbb{1}' \beta_i \\ &\quad + \underbrace{\frac{1}{2} \sum_{s,t \in [n]} x'_s x_t \beta'_t \pi \sigma_{y_t} \sigma_{y_s} \pi' \beta_s - \sum_{s \in [n]: s \neq i} \mathbb{1}' \beta_s}_{=: C_i} \end{aligned}$$

where  $C_i$  is a scalar quantity which does not depend on  $\beta_i$ . Thus, plugging in  $\hat{\beta}$ , we have

$$g(\hat{\beta}) = \frac{1}{2} \|x_i\|_2^2 \hat{\beta}'_i \Theta \hat{\beta}_i + \sum_{s \in [n]: s \neq i} x'_s x_i \hat{\beta}'_i \pi \sigma_{y_i} \sigma_{y_s} \pi' \beta_s - \mathbb{1}' \hat{\beta}_i + C_i. \quad (19)$$

Furthermore,

$$\begin{aligned}
 \sum_{s \in [n]: s \neq i} x'_s x_i \widehat{\beta}'_i \boldsymbol{\pi} \boldsymbol{\sigma}_{y_i} \boldsymbol{\sigma}_{y_s} \boldsymbol{\pi}' \beta_s &= \sum_{s \in [n]: s \neq i} \widehat{\beta}'_i \boldsymbol{\pi} \boldsymbol{\sigma}_{y_i} \boldsymbol{\sigma}_{y_s} \boldsymbol{\pi}' \beta_s x'_s x_i \\
 &= \widehat{\beta}'_i \boldsymbol{\pi} \boldsymbol{\sigma}_{y_i} \left( \sum_{s \in [n]: s \neq i} \boldsymbol{\sigma}_{y_s} \boldsymbol{\pi}' \beta_s x'_s \right) x_i \\
 &= \widehat{\beta}'_i \boldsymbol{\pi} \boldsymbol{\sigma}_{y_i} \left( -\boldsymbol{\sigma}_{y_i} \boldsymbol{\pi}' \beta_i x'_i + \sum_{s \in [n]} \boldsymbol{\sigma}_{y_s} \boldsymbol{\pi}' \beta_s x'_s \right) x_i \\
 &= \widehat{\beta}'_i \boldsymbol{\pi} \boldsymbol{\sigma}_{y_i} \left( -\boldsymbol{\sigma}_{y_i} \boldsymbol{\pi}' \beta_i x'_i - \sum_{s \in [n]} \alpha_s x'_s \right) x_i \quad \because (3) \\
 &= \widehat{\beta}'_i \boldsymbol{\pi} \boldsymbol{\sigma}_{y_i} (-\boldsymbol{\sigma}_{y_i} \boldsymbol{\pi}' \beta_i x'_i + \mathbf{w}') x_i \quad \because (1) \\
 &= \widehat{\beta}'_i (-\boldsymbol{\pi} \boldsymbol{\sigma}_{y_i} \boldsymbol{\sigma}_{y_i} \boldsymbol{\pi}' \beta_i \|x_i\|_2^2 + \boldsymbol{\pi} \boldsymbol{\sigma}_{y_i} \mathbf{w}' x_i) \\
 &= \widehat{\beta}'_i (\boldsymbol{\pi} \boldsymbol{\sigma}_{y_i} \mathbf{w}' x_i - \boldsymbol{\pi} \boldsymbol{\pi}' \beta_i \|x_i\|_2^2) \quad \because \boldsymbol{\sigma}_{y_i}^2 = \mathbf{I} \\
 &= \widehat{\beta}'_i (\boldsymbol{\pi} \boldsymbol{\sigma}_{y_i} \mathbf{w}' x_i - \boldsymbol{\Theta} \beta_i \|x_i\|_2^2) \quad \because (18)
 \end{aligned}$$

Therefore, we have

$$\begin{aligned}
 g(\widehat{\boldsymbol{\beta}}) &= \frac{1}{2} \|x_i\|_2^2 \widehat{\beta}'_i \boldsymbol{\Theta} \widehat{\boldsymbol{\beta}}_i + \widehat{\beta}'_i (\boldsymbol{\pi} \boldsymbol{\sigma}_{y_i} \mathbf{w}' x_i - \boldsymbol{\Theta} \beta_i \|x_i\|_2^2 - 1) + C_i \\
 &= \frac{1}{2} \|x_i\|_2^2 \widehat{\beta}'_i \boldsymbol{\Theta} \widehat{\boldsymbol{\beta}}_i - \widehat{\beta}'_i (1 - \boldsymbol{\pi} \boldsymbol{\sigma}_{y_i} \mathbf{w}' x_i + \boldsymbol{\Theta} \beta_i \|x_i\|_2^2) + C_i
 \end{aligned}$$

Thus, (17) is equivalent to

$$\begin{aligned}
 \min_{\widehat{\boldsymbol{\beta}} \in \mathcal{G}} \quad & \frac{1}{2} \|x_i\|_2^2 \widehat{\beta}'_i \boldsymbol{\Theta} \widehat{\boldsymbol{\beta}}_i - \widehat{\beta}'_i (1 - \boldsymbol{\pi} \boldsymbol{\sigma}_{y_i} \mathbf{w}' x_i + \boldsymbol{\Theta} \beta_i \|x_i\|_2^2) + C_i \\
 \text{s.t.} \quad & \widehat{\beta}_s = \beta_s, \forall s \in [n] \setminus \{i\}.
 \end{aligned}$$

Dropping the constant  $C_i$  and dividing through by  $\|x_i\|_2^2$  does not change the minimizers. Hence, (17) has the same set of minimizers as

$$\begin{aligned}
 \min_{\widehat{\boldsymbol{\beta}} \in \mathcal{G}} \quad & \frac{1}{2} \widehat{\beta}'_i \boldsymbol{\Theta} \widehat{\boldsymbol{\beta}}_i - \widehat{\beta}'_i ((1 - \boldsymbol{\pi} \boldsymbol{\sigma}_{y_i} \mathbf{w}' x_i) / \|x_i\|_2^2 + \boldsymbol{\Theta} \beta_i) \\
 \text{s.t.} \quad & \widehat{\beta}_s = \beta_s, \forall s \in [n] \setminus \{i\}.
 \end{aligned}$$

Due to the equality constraints, the only free variable is  $\widehat{\beta}_i$ . Note that the above optimization, when restricted to  $\widehat{\beta}_i$ , is equivalent to the optimization (4) with

$$v := (1 - \boldsymbol{\pi} \boldsymbol{\sigma}_{y_i} \mathbf{w}' x_i) / \|x_i\|_2^2 + \boldsymbol{\Theta} \beta_i$$

and  $\mathbf{w}$  is as in (1). The uniqueness of the minimizer is guaranteed by Theorem 3.4.  $\square$

#### A.4. Global linear convergence

Wang & Lin (2014) established the global linear convergence of the so-called *feasible descent method* when applied to a certain class of problems. As an application, they prove global linear convergence for coordinate descent for solving the dual problem of the binary SVM with the hinge loss. Wang & Lin (2014) considered optimization problems of the following form:

$$\min_{x \in \mathcal{X}} f(x) := g(\mathbf{E}x) + b'x \quad (20)$$

where  $f : \mathbb{R}^n \rightarrow \mathbb{R}$  is a function such that  $\nabla f$  is Lipschitz continuous,  $\mathcal{X} \subseteq \mathbb{R}^n$  is a polyhedral set,  $\arg \min_{x \in \mathcal{X}} f(x)$  is nonempty,  $g : \mathbb{R}^m \rightarrow \mathbb{R}$  is a strongly convex function such that  $\nabla g$  is Lipschitz continuous, and  $\mathbf{E} \in \mathbb{R}^{m \times n}$  and  $b \in \mathbb{R}^m$  are fixed matrix and vector, respectively.

Below, let  $\mathcal{P}_{\mathcal{X}} : \mathbb{R}^n \rightarrow \mathcal{X}$  denote the orthogonal projection on  $\mathcal{X}$ .

**Definition A.2.** In the context of (20), an iterative algorithm that produces a sequence  $\{x^0, x^1, x^2, \dots\} \subseteq \mathcal{X}$  is a *feasible descent method* if there exists a sequence  $\{\epsilon^0, \epsilon^1, \epsilon^2, \dots\} \subseteq \mathbb{R}^n$  such that for all  $t \geq 0$

$$x^{t+1} = \mathcal{P}_{\mathcal{X}}(x^t - \nabla f(x^t) + \epsilon^t) \quad (21)$$

$$\|\epsilon^t\| \leq B\|x^t - x^{t+1}\| \quad (22)$$

$$f(x^t) - f(x^{t+1}) \geq \Gamma\|x^t - x^{t+1}\|^2 \quad (23)$$

where  $B, \Gamma > 0$ .

One of the main result of (Wang & Lin, 2014) is

**Theorem A.3** (Theorem 8 from (Wang & Lin, 2014)). *Suppose an optimization problem  $\min_{x \in \mathcal{X}} f(x)$  is of the form (20) and  $\{x^0, x^1, x^2, \dots\} \subseteq \mathcal{X}$  is a sequence generated by a feasible descent method. Let  $f^* := \min_{x \in \mathcal{X}} f(x)$ . Then there exists  $\Delta \in (0, 1)$  such that*

$$f(x^{t+1}) - f^* \leq \Delta(f(x^t) - f^*), \quad \forall t \geq 0.$$

Now, we begin verifying that the WW-SVM dual optimization and the BCD algorithm for WW-SVM satisfies the requirements of Theorem A.3.

Given  $\beta \in \mathbb{R}^{(k-1) \times n}$ , define its vectorization

$$\text{vec}(\beta) = \begin{bmatrix} \beta_1 \\ \vdots \\ \beta_n \end{bmatrix} \in \mathbb{R}^{(k-1)n}.$$

Define the matrix  $\mathbf{P}_{is} = \pi \sigma_{y_i} x'_i x_s \sigma_{y_s} \pi' \in \mathbb{R}^{(k-1) \times (k-1)}$ , and  $\mathbf{Q} \in \mathbb{R}^{(k-1)n \times (k-1)n}$  by

$$\mathbf{Q} = \begin{bmatrix} \mathbf{P}_{11} & \mathbf{P}_{12} & \cdots & \mathbf{P}_{1n} \\ \mathbf{P}_{21} & \mathbf{P}_{22} & \cdots & \mathbf{P}_{2n} \\ \vdots & \vdots & \ddots & \vdots \\ \mathbf{P}_{n1} & \mathbf{P}_{n2} & \cdots & \mathbf{P}_{nn} \end{bmatrix}.$$

Let

$$\mathbf{E} = \begin{bmatrix} x_1 \sigma_{y_1} \pi' \\ x_2 \sigma_{y_2} \pi' \\ \vdots \\ x_n \sigma_{y_n} \pi' \end{bmatrix}.$$

We observe that  $\mathbf{Q} = \mathbf{E}'\mathbf{E}$ . Thus,  $\mathbf{Q}$  is symmetric and positive semi-definite. Let  $\|\mathbf{Q}\|_{op}$  be the operator norm of  $\mathbf{Q}$ .

**Proposition A.4.** *The optimization (D2) is of the form (20). More precisely, the optimization (D2) can be expressed as*

$$\min_{\beta \in \mathcal{G}} g(\beta) = \varphi(\mathbf{E}\text{vec}(\beta)) - \mathbb{1}'\text{vec}(\beta) \quad (24)$$

where the feasible set  $\mathcal{G}$  is a nonempty polyhedral set (i.e., defined by a system of linear inequalities, hence convex),  $\varphi$  is strongly convex, and  $\nabla g$  is Lipschitz continuous with Lipschitz constant  $L := \|\mathbf{Q}\|_{op}$ . Furthermore, (24) has at least one minimizer.

*Proof.* Observe

$$\begin{aligned}
 g(\boldsymbol{\beta}) &= \frac{1}{2} \sum_{i,s \in [n]} x'_s x_i \beta'_i \boldsymbol{\pi} \boldsymbol{\sigma}_{y_i} \boldsymbol{\sigma}_{y_s} \boldsymbol{\pi}' \beta_s - \sum_{i \in [n]} \mathbb{1}' \beta_i \\
 &= \frac{1}{2} \text{vec}(\boldsymbol{\beta})' \mathbf{Q} \text{vec}(\boldsymbol{\beta}) - \mathbb{1}' \text{vec}(\boldsymbol{\beta}) \\
 &= \frac{1}{2} (\mathbf{E} \text{vec}(\boldsymbol{\beta}))' (\mathbf{E} \text{vec}(\boldsymbol{\beta})) - \mathbb{1}' \text{vec}(\boldsymbol{\beta}) \\
 &= \varphi(\mathbf{E} \text{vec}(\boldsymbol{\beta})) - \mathbb{1}' \text{vec}(\boldsymbol{\beta})
 \end{aligned}$$

where  $\varphi(\bullet) = \frac{1}{2} \|\bullet\|^2$ . Note that  $\text{vec}(\nabla g(\boldsymbol{\beta})) = \mathbf{Q} \text{vec}(\boldsymbol{\beta}) - \mathbb{1}$ . Hence, the Lipschitz constant of  $g$  is  $\|\mathbf{Q}\|_{op}$ . For the ‘‘Furthermore’’ part, note that the above calculation shows that (24) is a quadratic program where the second order term is positive semi-definite and the constraint set is convex. Hence, (24) has at least one minimizer.  $\square$

Let  $B = [0, C]^{k-1}$ . Let  $\boldsymbol{\beta}^t$  be  $\boldsymbol{\beta}$  at the end of the  $t$ -iteration of the outer loop of Algorithm 1. Define

$$\boldsymbol{\beta}^{t,i} := [\beta_1^{t+1}, \dots, \beta_i^{t+1}, \beta_{i+1}^t, \dots, \beta_n^t].$$

By construction, we have

$$\beta_i^{t+1} = \arg \min_{\beta \in B} g([\beta_1^{t+1}, \dots, \beta_{i-1}^{t+1}, \beta, \beta_{i+1}^t, \dots, \beta_n^t]) \quad (25)$$

For each  $i = 1, \dots, n$ , let

$$\nabla_i g(\boldsymbol{\beta}) = \left[ \frac{\partial g}{\partial \beta_{1i}}(\boldsymbol{\beta}), \frac{\partial g}{\partial \beta_{2i}}(\boldsymbol{\beta}), \dots, \frac{\partial g}{\partial \beta_{(k-1)i}}(\boldsymbol{\beta}) \right]'$$

By Lemma 24 (Wang & Lin, 2014), we have

$$\beta_i^{t+1} = \mathcal{P}_B(\beta_i^{t+1} - \nabla_i g(\boldsymbol{\beta}^{t,i}))$$

where  $\mathcal{P}_B$  denotes orthogonal projection on to  $B$ . Now, define  $\boldsymbol{\epsilon}^t \in \mathbb{R}^{(k-1) \times n}$  such that

$$\boldsymbol{\epsilon}_i^t = \beta_i^{t+1} - \beta_i^t - \nabla_i g(\boldsymbol{\beta}^{t,i}) + \nabla_i g(\boldsymbol{\beta}^t).$$

**Proposition A.5.** *The BCD algorithm for the WW-SVM is a feasible descent method. More precisely, the sequence  $\{\boldsymbol{\beta}^0, \boldsymbol{\beta}^1, \dots\}$  satisfies the following conditions:*

$$\boldsymbol{\beta}^{t+1} = \mathcal{P}_{\mathcal{G}}(\boldsymbol{\beta}^t - \nabla g(\boldsymbol{\beta}^t) + \boldsymbol{\epsilon}^t) \quad (26)$$

$$\|\boldsymbol{\epsilon}^t\| \leq (1 + \sqrt{n}L) \|\boldsymbol{\beta}^t - \boldsymbol{\beta}^{t+1}\| \quad (27)$$

$$g(\boldsymbol{\beta}^t) - g(\boldsymbol{\beta}^{t+1}) \geq \Gamma \|\boldsymbol{\beta}^t - \boldsymbol{\beta}^{t+1}\|^2 \quad (28)$$

where  $L$  is as in Proposition A.4,  $\Gamma := \min_{i \in [n]} \frac{\|x_i\|^2}{2}$ ,  $\mathcal{G}$  is the feasible set of (D2), and  $\mathcal{P}_{\mathcal{G}}$  is the orthogonal projection onto  $\mathcal{G}$ .

The proof of Proposition A.5 essentially generalizes Proposition 3.4 of (Luo & Tseng, 1993) to the higher dimensional setting:

*Proof.* Recall that  $\mathcal{G} = B^{\times n} := B \times \dots \times B$ . Note that the  $i$ -th block of  $\boldsymbol{\beta}^t - \nabla g(\boldsymbol{\beta}^t) + \boldsymbol{\epsilon}^t$  is

$$\beta_i^t - \nabla_i g(\boldsymbol{\beta}^t) + \boldsymbol{\epsilon}_i^t = \beta_i^t - \nabla_i g(\boldsymbol{\beta}^t) + (\beta_i^{t+1} - \beta_i^t - \nabla_i g(\boldsymbol{\beta}^{t,i}) + \nabla_i g(\boldsymbol{\beta}^t)) = \beta_i^{t+1} - \nabla_i g(\boldsymbol{\beta}^{t,i}).$$

Thus, the  $i$ -th block of  $\mathcal{P}_{\mathcal{G}}(\boldsymbol{\beta}^t - \nabla g(\boldsymbol{\beta}^t) + \boldsymbol{\epsilon}^t)$  is

$$\mathcal{P}_B(\beta_i^{t+1} - \nabla_i g(\boldsymbol{\beta}^{t,i})) = \beta_i^{t+1}.$$

This is precisely the identity (26).

Next, we have

$$\begin{aligned}\|\epsilon_i^t\| &\leq \|\beta_i^{t+1} - \beta_i^t\| + \|\nabla_i g(\beta^{t,i}) - \nabla_i g(\beta^t)\| \\ &\leq \|\beta_i^{t+1} - \beta_i^t\| + L\|\beta^{t,i} - \beta^t\| \\ &\leq \|\beta_i^{t+1} - \beta_i^t\| + L\|\beta^{t+1} - \beta^t\|.\end{aligned}$$

From this, we get that

$$\begin{aligned}\|\epsilon^t\| &= \sqrt{\sum_{i=1}^n \|\epsilon_i^t\|^2} \\ &\leq \sqrt{\sum_{i=1}^n (\|\beta_i^{t+1} - \beta_i^t\| + L\|\beta^{t+1} - \beta^t\|)^2} \\ &\leq \sqrt{\sum_{i=1}^n \|\beta_i^{t+1} - \beta_i^t\|^2} + \sqrt{\sum_{i=1}^n L^2 \|\beta^{t+1} - \beta^t\|^2} \\ &= \|\beta^{t+1} - \beta^t\| + \sqrt{n}L\|\beta^{t+1} - \beta^t\| \\ &= (1 + \sqrt{n}L)\|\beta^{t+1} - \beta^t\|.\end{aligned}$$

Thus, we conclude that  $\|\epsilon^t\| \leq (1 + \sqrt{n}L)\|\beta^{t+1} - \beta^t\|$  which is (27).

Finally, we show that

$$g(\beta^{t,i-1}) - g(\beta^{t,i}) + \nabla_i g(\beta^{t,i})'(\beta_i^{t+1} - \beta_i^t) \geq \Gamma \|\beta_i^{t+1} - \beta_i^t\|^2$$

where  $\Gamma := \min_{i \in [n]} \frac{\|x_i\|^2}{2}$ .

**Lemma A.6.** *Let  $\beta_1, \dots, \beta_{i-1}, \beta, \beta_{i+1}, \dots, \beta_n \in \mathbb{R}^{k-1}$  be arbitrary. Then there exist  $v \in \mathbb{R}^{k-1}$  and  $C \in \mathbb{R}$  which depend only on  $\beta_1, \dots, \beta_{i-1}, \beta_{i+1}, \dots, \beta_n$ , but not on  $\beta$ , such that*

$$g([\beta_1, \dots, \beta_{i-1}, \beta, \beta_{i+1}, \dots, \beta_n]) = \frac{1}{2}\|x_i\|^2 \beta' \beta - v' \beta - C.$$

In particular, we have

$$\nabla_i g([\beta_1, \dots, \beta_{i-1}, \beta, \beta_{i+1}, \dots, \beta_n]) = \|x_i\|^2 \beta - v.$$

*Proof.* The result follows immediately from the identity (19). □

**Lemma A.7.** *Let  $\beta_1, \dots, \beta_{i-1}, \beta, \eta, \beta_{i+1}, \dots, \beta_n \in \mathbb{R}^{k-1}$  be arbitrary. Then we have*

$$\begin{aligned}&g([\beta_1, \dots, \beta_{i-1}, \eta, \beta_{i+1}, \dots, \beta_n]) - g([\beta_1, \dots, \beta_{i-1}, \beta, \beta_{i+1}, \dots, \beta_n]) \\ &\quad + \nabla_i g([\beta_1, \dots, \beta_{i-1}, \beta, \beta_{i+1}, \dots, \beta_n])'(\beta - \eta) \\ &= \frac{\|x_i\|^2}{2} \|\eta - \beta\|^2\end{aligned}$$

*Proof.* Let  $v, C$  be as in Lemma A.6. We have

$$\begin{aligned}&g([\beta_1, \dots, \beta_{i-1}, \eta, \beta_{i+1}, \dots, \beta_n]) - g([\beta_1, \dots, \beta_{i-1}, \beta, \beta_{i+1}, \dots, \beta_n]) \\ &= \frac{\|x_i\|^2}{2} \|\eta\|^2 - v' \eta - \frac{\|x_i\|^2}{2} \|\beta\|^2 + v' \beta \\ &= \frac{\|x_i\|^2}{2} (\|\eta\|^2 - \|\beta\|^2) + v'(\beta - \eta)\end{aligned}$$

and

$$\nabla_i g([\beta_1, \dots, \beta_{i-1}, \beta, \beta_{i+1}, \dots, \beta_n])'(\beta - \eta) = (\|x_i\|^2 \beta - v)'(\beta - \eta) = \|x_i\|^2 (\|\beta\|^2 - \beta' \eta) - v'(\beta - \eta).$$

Thus,

$$\begin{aligned} & g([\beta_1, \dots, \beta_{i-1}, \eta, \beta_{i+1}, \dots, \beta_n]) - g([\beta_1, \dots, \beta_{i-1}, \beta, \beta_{i+1}, \dots, \beta_n]) \\ & \quad + \nabla_i g([\beta_1, \dots, \beta_{i-1}, \beta, \beta_{i+1}, \dots, \beta_n])'(\beta - \eta) \\ &= \frac{\|x_i\|^2}{2} (\|\eta\|^2 - \|\beta\|^2) + v'(\beta - \eta) + \|x_i\|^2 (\|\beta\|^2 - \beta' \eta) - v'(\beta - \eta) \\ &= \frac{\|x_i\|^2}{2} (\|\eta\|^2 - \|\beta\|^2) + \|x_i\|^2 (\|\beta\|^2 - \beta' \eta) \\ &= \|x_i\|^2 \left( \frac{1}{2} (\|\eta\|^2 - \|\beta\|^2) + (\|\beta\|^2 - \beta' \eta) \right) \\ &= \|x_i\|^2 \left( \frac{1}{2} (\|\eta\|^2 + \|\beta\|^2) - \beta' \eta \right) \\ &= \frac{\|x_i\|^2}{2} \|\eta - \beta\|^2 \end{aligned}$$

as desired.  $\square$

Applying Lemma A.7, we have

$$g(\beta^{t,i-1}) - g(\beta^{t,i}) + \nabla_i g(\beta^{t,i})'(\beta_i^{t+1} - \beta_i^t) \geq \frac{\|x_i\|^2}{2} \|\beta_i^{t+1} - \beta_i^t\|^2.$$

Since (25) is true, we have by Lemma 24 of (Wang & Lin, 2014) that

$$\nabla_i g(\beta^{t,i})'(\beta_i^t - \beta_i^{t+1}) \geq 0$$

Equivalently,  $\nabla_i g(\beta^{t,i})'(\beta_i^{t+1} - \beta_i^t) \leq 0$ . Thus, we deduce that

$$g(\beta^{t,i-1}) - g(\beta^{t,i}) \geq \frac{\|x_i\|^2}{2} \|\beta_i^{t+1} - \beta_i^t\|^2 \geq \Gamma \|\beta_i^{t+1} - \beta_i^t\|^2$$

Summing the above identity over  $i \in [n]$ , we have

$$g(\beta^{t,0}) - g(\beta^{t,n}) = \sum_{i=1}^n g(\beta^{t,i-1}) - g(\beta^{t,i}) \geq \Gamma \sum_{i=1}^n \|\beta_i^{t+1} - \beta_i^t\|^2 = \Gamma \|\beta^{t+1} - \beta^t\|^2$$

Since  $(\beta^{t,0}) = \beta^t$  and  $\beta^{t,n} = \beta^{t+1}$ , we conclude that  $g(\beta^t) - g(\beta^{t+1}) \geq \Gamma \|\beta^{t+1} - \beta^t\|^2$ .  $\square$

To conclude the proof of Theorem 3.6, we note that Proposition A.5 and Proposition A.4 together imply that the requirements of Theorem 8 from (Wang & Lin, 2014) (restated as Theorem A.3 here) are satisfied for the BCD algorithm for WW-SVM. Hence, we are done.  $\square$

### A.5. Proof of Theorem 3.4

The goal of this section is to prove Theorem 3.4. The time complexity analysis has been carried out at the end of Section 4 of the main article. Below, we focus on the part of the theorem on the correctness of the output. Throughout this section,  $k \geq 2$ ,  $C > 0$  and  $v \in \mathbb{R}^{k-1}$  are assumed to be fixed. Additional variables used are summarized in Table 2.



Table 2. Variables used in Section A.5

VARIABLE(S)	DEFINED IN	NOTA BENE
$t$	ALGORITHM 2	ITERATION INDEX
$\ell, \mathbf{vals}, \delta_t, \gamma_t$	SUBROUTINE 3	$t \in [\ell]$ IS AN ITERATION INDEX
$\mathbf{up}, \mathbf{dn}$	SUBROUTINE 3	SYMBOLS
$\tilde{b}, \tilde{\gamma}, v_{\max}$	LEMMA A.9	
$\langle 1 \rangle, \dots, \langle k-1 \rangle$	ALGORITHM 2	
$n_m^t, n_u^t, S^t, \hat{\gamma}^t, \hat{b}^t$	ALGORITHM 2	$t \in [\ell]$ IS AN ITERATION INDEX
$\llbracket k \rrbracket, I_u^\gamma, I_m^\gamma, n_u^\gamma, n_m^\gamma$	DEFINITION A.10	$\gamma \in \mathbb{R}$ IS A REAL NUMBER
$S^{(n_m, n_u)}, \hat{\gamma}^{(n_m, n_u)}, \hat{b}^{(n_m, n_u)}$	DEFINITION A.13	$(n_m, n_u) \in \llbracket k \rrbracket^2$
$\mathbf{vals}^+$	DEFINITION A.18	
$\mathbf{u}(j), \mathbf{d}(j)$	DEFINITION A.19	$j \in [k-1]$ IS AN INTEGER
$\mathbf{crit}_1, \mathbf{crit}_2$	DEFINITION A.20	
$\text{KKT\_cond}()$	SUBROUTINE 5	

### A.5.1. THE CLIPPING MAP

First, we recall the clipping map:

**Definition A.8.** The *clipping map*  $\text{clip}_C : \mathbb{R}^{k-1} \rightarrow [0, C]^{k-1}$  is the function defined as follows: for  $w \in \mathbb{R}^{k-1}$ ,  $[\text{clip}_C(w)]_i := \max\{0, \min\{C, w_i\}\}$ .

**Lemma A.9.** Let  $v_{\max} = \max_{i \in [k-1]} v_i$ . The optimization (4) has a unique global minimum  $\tilde{b}$  satisfying the following:

1.  $\tilde{b} = \text{clip}_C(v - \tilde{\gamma}\mathbf{1})$  for some  $\tilde{\gamma} \in \mathbb{R}$
2.  $\tilde{\gamma} = \sum_{i=1}^{k-1} \tilde{b}_i$ . In particular,  $\tilde{\gamma} \geq 0$ .
3. If  $v_i \leq 0$ , then  $\tilde{b}_i = 0$ . In particular, if  $v_{\max} \leq 0$ , then  $\tilde{b} = \mathbf{0}$ .
4. If  $v_{\max} > 0$ , then  $0 < \tilde{\gamma} < v_{\max}$ .

*Proof.* We first prove part 1. The optimization (4) is a minimization over a convex domain with strictly convex objective, and hence has a unique global minimum  $\tilde{b}$ . For each  $i \in [k-1]$ , let  $\lambda_i, \mu_i \in \mathbb{R}$  be the dual variables for the constraints  $0 \geq b_i - C$  and  $0 \geq -b_i$ , respectively. The Lagrangian for the optimization (4) is

$$\mathcal{L}(b, \lambda, \mu) = \frac{1}{2} b'(\mathbf{I} + \mathbf{O})b - v'b + (b - C)'\lambda + (-b)'\mu.$$

Thus, the stationarity (or gradient vanishing) condition is

$$0 = \nabla_b \mathcal{L}(b, \lambda, \mu) = (\mathbf{I} + \mathbf{O})b - v + \lambda - \mu.$$

The KKT conditions are as follows:

for all  $i \in [k-1]$ , the following holds:

$$[(\mathbf{I} + \mathbf{O})b]_i + \lambda_i - \mu_i = v_i \quad \text{stationarity} \quad (29)$$

$$C \geq b_i \geq 0 \quad \text{primal feasibility} \quad (30)$$

$$\lambda_i \geq 0 \quad \text{dual feasibility} \quad (31)$$

$$\mu_i \geq 0 \quad \text{"} \quad (32)$$

$$\lambda_i(C - b_i) = 0 \quad \text{complementary slackness} \quad (33)$$

$$\mu_i b_i = 0 \quad \text{"} \quad (34)$$

(29) to (34) are satisfied if and only if  $b = \tilde{b}$  is the global minimum.

Let  $\tilde{\gamma} \in \mathbb{R}$  be such that  $\tilde{\gamma}\mathbb{1} = \mathbf{O}\tilde{b}$ . Note that by definition, part 2 holds. Furthermore, (29) implies

$$\tilde{b} = v - \tilde{\gamma}\mathbb{1} - \lambda + \mu. \quad (35)$$

Below, fix some  $i \in [k-1]$ . Note that  $\lambda_i$  or  $\mu_i$  cannot both be nonzero. Otherwise, (33) and (34) would imply that  $C = \tilde{b}_i = 0$ , a contradiction. We claim the following:

1. If  $v_i - \tilde{\gamma} \in [0, C]$ , then  $\lambda_i = \mu_i = 0$  and  $\tilde{b}_i = v_i - \tilde{\gamma}$ .
2. If  $v_i - \tilde{\gamma} > C$ , then  $\tilde{b}_i = C$ .
3.  $v_i - \tilde{\gamma} < 0$ , then  $\tilde{b}_i = 0$ .

We prove the first claim. To this end, suppose  $v_i - \tilde{\gamma} \in [0, C]$ . We will show  $\lambda_i = \mu_i = 0$  by contradiction. Suppose  $\lambda_i > 0$ . Then we have  $C = \tilde{b}_i$  and  $\mu_i = 0$ . Now, (35) implies that  $C = \tilde{b}_i = v_i - \tilde{\gamma} - \lambda_i$ . However, we now have  $v_i - \tilde{\gamma} - \lambda_i \leq C - \lambda_i < C$ , a contradiction. Thus,  $\lambda_i = 0$ . Similarly, assuming  $\mu_i > 0$  implies

$$0 = \tilde{b}_i = v_i - \lambda + \mu_i \geq 0 + \mu_i > 0,$$

a contradiction. This proves the first claim.

Next, we prove the second claim. Note that

$$C \geq \tilde{b}_i = v_i - \tilde{\gamma} - \lambda_i + \mu_i > C - \lambda_i + \mu_i \implies 0 > -\lambda_i + \mu_i \geq -\lambda_i.$$

In particular, we have  $\lambda_i > 0$  which implies  $C = \tilde{b}_i$  by complementary slackness.

Finally, we prove the third claim. Note that

$$0 \leq \tilde{b}_i = v_i - \tilde{\gamma} - \lambda_i + \mu_i < -\lambda_i + \mu_i \leq \mu_i$$

Thus,  $\mu_i > 0$  and so  $0 = \tilde{b}_i$  by complementary slackness. This proves that  $\tilde{b} = \text{clip}_C(v - \tilde{\gamma}\mathbb{1})$ , which concludes the proof of part 1.

For part 2, note that  $\tilde{\gamma} = \sum_{i=1}^{k-1} \tilde{b}_i$  holds by definition. The ‘‘in particular’’ portion follows immediately from  $\tilde{b} \geq 0$ .

We prove part 3 by contradiction. Suppose there exists  $i \in [k-1]$  such that  $v_i \leq 0$  and  $\tilde{b}_i > 0$ . Thus, by (34), we have  $\mu_i = 0$ . By (29), we have  $b_i + \tilde{\gamma} \leq b_i + \tilde{\gamma} + \lambda_i = v_i \leq 0$ . Thus, we have  $-\tilde{\gamma} \geq b_i > 0$ , or equivalently,  $\tilde{\gamma} < 0$ . However, this contradicts part 2. Thus,  $\tilde{b}_i = 0$  whenever  $v_i \leq 0$ . The ‘‘in particular’’ portion follows immediately from the observation that  $v_{\max} \leq 0$  implies that  $v_i \leq 0$  for all  $i \in [k-1]$ .

For part 4, we first prove that  $\tilde{\gamma} < v_{\max}$  by contradiction. Suppose that  $\tilde{\gamma} \geq v_{\max}$ . Then we have  $v - \tilde{\gamma}\mathbb{1} \leq v - v_{\max}\mathbb{1} \leq 0$ . Thus, by part 1, we have  $\tilde{b} = \text{clip}_C(v - \tilde{\gamma}\mathbb{1}) = \mathbf{0}$ . By part 2, we must have that  $\tilde{\gamma} = \sum_{i=1}^{k-1} \tilde{b}_i = 0$ . However,  $\tilde{\gamma} \geq v_{\max} > 0$ , which is a contradiction.

Finally, we prove that  $\tilde{\gamma} > 0$  again by contradiction. Suppose that  $\tilde{\gamma} = 0$ . Then part 2 and the fact that  $\tilde{b} \geq \mathbf{0}$  implies that  $\tilde{b} = \mathbf{0}$ . However, by part 1, we have  $\tilde{b} = \text{clip}_C(v)$ . Now, let  $i^*$  be such that  $v_{i^*} = v_{\max}$ . This implies that  $\tilde{b}_{i^*} = \text{clip}_C(v_{\max}) > 0$ , a contradiction.  $\square$

#### A.5.2. RECOVERING $\tilde{\gamma}$ FROM DISCRETE DATA

**Definition A.10.** For  $\gamma \in \mathbb{R}$ , let  $b^\gamma := \text{clip}_C(v - \gamma\mathbb{1}) \in \mathbb{R}^{k-1}$ . Define

$$\begin{aligned} I_u^\gamma &:= \{i \in [k-1] : b_i^\gamma = C\} \\ I_m^\gamma &:= \{i \in [k-1] : b_i^\gamma \in (0, C)\} \\ n_u^\gamma &:= |I_u^\gamma|, \quad \text{and} \quad n_m^\gamma := |I_m^\gamma|. \end{aligned}$$

Let  $\llbracket k \rrbracket := \{0\} \cup [k-1]$ . Note that by definition,  $n_m^\gamma, n_u^\gamma \in \llbracket k \rrbracket$ .

Note that  $I_u^\gamma$  and  $I_m^\gamma$  are determined by their cardinalities. This is because

$$\begin{aligned} I_u^\gamma &= \{\langle 1 \rangle, \langle 2 \rangle, \dots, \langle n_u^\gamma \rangle\} \\ I_m^\gamma &= \{\langle n_u^\gamma + 1 \rangle, \langle n_u^\gamma + 2 \rangle, \dots, \langle n_u^\gamma + n_m^\gamma \rangle\}. \end{aligned}$$

**Definition A.11.** Define

$$\text{disc}^+ := \{v_i : i \in [k-1], v_i > 0\} \cup \{v_i - C : i \in [k-1], v_i - C > 0\} \cup \{0\}.$$

Note that  $\text{disc}^+$  is slightly different from  $\text{disc}$  as defined in the main text.

**Lemma A.12.** Let  $\gamma', \gamma'' \in \text{disc}^+$  be such that  $\gamma \notin \text{disc}^+$  for all  $\gamma \in (\gamma', \gamma'')$ . The functions

$$\begin{aligned} (\gamma', \gamma'') \ni \gamma &\mapsto I_m^\gamma \\ (\gamma', \gamma'') \ni \gamma &\mapsto I_u^\gamma \end{aligned}$$

are constant.

*Proof.* We first prove  $I_m^\lambda = I_m^\rho$ . Let  $\lambda, \rho \in (\gamma', \gamma'')$  be such that  $\lambda < \rho$ . Assume for the sake of contradiction that  $I_m^\lambda \neq I_m^\rho$ . Then either 1)  $i \in [k-1]$  such that  $v_i - \lambda \in (0, C)$  but  $v_i - \rho \notin (0, C)$  or 2)  $i \in [k-1]$  such that  $v_i - \lambda \notin (0, C)$  but  $v_i - \rho \in (0, C)$ . This implies that there exists some  $\gamma \in (\lambda, \rho)$  such that  $v_i - \gamma \in \{0, C\}$ , or equivalently,  $\gamma \in \{v_i, v_i - C\}$ . Hence,  $\gamma \in \text{disc}^+$ , which is a contradiction. Thus, for all  $\lambda, \rho \in (\gamma', \gamma'')$ , we have  $I_m^\lambda = I_m^\rho$ .

Next, we prove  $I_u^\lambda = I_u^\rho$ . Let  $\lambda, \rho \in (\gamma', \gamma'')$  be such that  $\lambda < \rho$ . Assume for the sake of contradiction that  $I_u^\lambda \neq I_u^\rho$ . Then either 1)  $i \in [k-1]$  such that  $v_i - \lambda \geq C$  but  $v_i - \rho < C$  or 2)  $i \in [k-1]$  such that  $v_i - \lambda < C$  but  $v_i - \rho \geq C$ . This implies that there exists some  $\gamma \in (\lambda, \rho)$  such that  $v_i - \gamma = C$ , or equivalently,  $\gamma = v_i - C$ . Hence,  $\gamma \in \text{disc}^+$ , which is a contradiction. Thus, for all  $\lambda, \rho \in (\gamma', \gamma'')$ , we have  $I_u^\lambda = I_u^\rho$ .  $\square$

**Definition A.13.** For  $(n_m, n_u) \in \llbracket k \rrbracket^2$ , define  $S^{(n_m, n_u)}, \hat{\gamma}^{(n_m, n_u)} \in \mathbb{R}$  by

$$\begin{aligned} S^{(n_m, n_u)} &:= \sum_{i=n_u+1}^{n_u+n_m} v_{\langle i \rangle}, \\ \hat{\gamma}^{(n_m, n_u)} &:= \left( C \cdot n_u + S^{(n_m, n_u)} \right) / (n_m + 1). \end{aligned}$$

Furthermore, define  $\hat{b}^{(n_m, n_u)} \in \mathbb{R}^{k-1}$  such that, for  $i \in [k-1]$ , the  $\langle i \rangle$ -th entry is

$$\hat{b}_{\langle i \rangle}^{(n_m, n_u)} := \begin{cases} C & : i \leq n_u \\ v_{\langle i \rangle} - \gamma^{(n_m, n_u)} & : n_u < i \leq n_u + n_m \\ 0 & : n_u + n_m < i. \end{cases}$$

Below, recall  $\ell$  as defined on Subroutine 3-line 2.

**Lemma A.14.** Let  $t \in [\ell]$ . Let  $n_m^t, n_u^t$ , and  $\hat{b}^t$  be as in the for loop of Algorithm 2. Then  $\hat{\gamma}^{(n_m^t, n_u^t)} = \hat{\gamma}^t$  and  $\hat{b}^{(n_m^t, n_u^t)} = \hat{b}^t$ .

*Proof.* It suffices to show that  $S^t = S^{(n_m^t, n_u^t)}$  where the former is defined as in Algorithm 2 and the latter is defined as in Definition A.13. In other words, it suffices to show that

$$S^t = \sum_{j \in [k-1] : n_u^t < j \leq n_u^t + n_m^t} v_{\langle j \rangle}. \quad (36)$$

We prove (36) by induction. The base case  $t = 0$  follows immediately due to the initialization in Algorithm 2-line 4.

Now, suppose that (36) holds for  $S^{t-1}$ :

$$S^{t-1} = \sum_{j \in [k-1] : n_u^{t-1} < j \leq n_u^{t-1} + n_m^{t-1}} v_{\langle j \rangle}. \quad (37)$$

Consider the first case that  $\delta_t = \text{up}$ . Then we have  $n_u^t + n_m^t = n_u^{t-1} + n_m^{t-1}$  and  $n_u^t = n_u^{t-1} + 1$ . Thus, we have

$$\begin{aligned} S^t &= S^{t-1} - v_{\langle n_u^{t-1} \rangle} \quad \because \text{Subroutine 4-line 3,} \\ &= \sum_{j \in [k-1] : n_u^{t-1} + 1 < j \leq n_u^{t-1} + n_m^{t-1}} v_{\langle j \rangle} \quad \because (37) \\ &= \sum_{j \in [k-1] : n_u^t < j \leq n_u^t + n_m^t} v_{\langle j \rangle} \end{aligned}$$

which is exactly the desired identity in (36).

Consider the second case that  $\delta_t = \text{dn}$ . Then we have  $n_u^t + n_m^t = n_u^{t-1} + n_m^{t-1} + 1$  and  $n_u^t = n_u^{t-1}$ . Thus, we have

$$\begin{aligned} S^t &= S^{t-1} + v_{\langle n_u^t + n_m^t \rangle} \quad \because \text{Subroutine 4-line 6,} \\ &= \sum_{j \in [k-1] : n_u^{t-1} + 1 < j \leq n_u^{t-1} + n_m^{t-1} + 1} v_{\langle j \rangle} \quad \because (37) \\ &= \sum_{j \in [k-1] : n_u^t < j \leq n_u^t + n_m^t} v_{\langle j \rangle} \end{aligned}$$

which, again, is exactly the desired identity in (36). □

**Lemma A.15.** *Let  $\tilde{\gamma}$  be as in Lemma A.9. Then we have*

$$\tilde{b} = \widehat{b}^{(n_m^{\tilde{\gamma}}, n_u^{\tilde{\gamma}})} = \text{clip}_C(v - \widehat{\gamma}^{(n_m^{\tilde{\gamma}}, n_u^{\tilde{\gamma}})} \mathbf{1}).$$

*Proof.* It suffices to prove that  $\tilde{\gamma} = \widehat{\gamma}^{(n_m^{\tilde{\gamma}}, n_u^{\tilde{\gamma}})}$ . To this end, let  $i \in [k-1]$ . If  $i \in I_m^{\tilde{\gamma}}$ , then  $\tilde{b}_i = v_i - \tilde{\gamma}$ . If  $i \in I_u^{\tilde{\gamma}}$ , then  $\tilde{b}_i = C$ . Otherwise,  $\tilde{b}_i = 0$ . Thus

$$\tilde{\gamma} = \mathbf{1}^T \tilde{b} = C \cdot n_u^{\tilde{\gamma}} + S^{(n_m^{\tilde{\gamma}}, n_u^{\tilde{\gamma}})} - \tilde{\gamma} \cdot n_m^{\tilde{\gamma}}$$

Solving for  $\tilde{\gamma}$ , we have

$$\tilde{\gamma} = \left( C \cdot n_u^{\tilde{\gamma}} + S^{(n_m^{\tilde{\gamma}}, n_u^{\tilde{\gamma}})} \right) / (n_m^{\tilde{\gamma}} + 1) = \widehat{\gamma}^{(n_m^{\tilde{\gamma}}, n_u^{\tilde{\gamma}})},$$

as desired. □

### A.5.3. CHECKING THE KKT CONDITIONS

**Lemma A.16.** *Let  $(n_m, n_u) \in \llbracket k \rrbracket^2$ . To simplify notation, let  $b := \widehat{b}^{(n_m, n_u)}$ ,  $\gamma := \widehat{\gamma}^{(n_m, n_u)}$ . We have  $\mathbf{O}b = \gamma \mathbf{1}$  and for all  $i \in [k-1]$  that*

$$[(\mathbf{I} + \mathbf{O})b]_{\langle i \rangle} = \begin{cases} C + \gamma & : i \leq n_u \\ v_{\langle i \rangle} & : n_u < i \leq n_u + n_m \\ \gamma & : n_u + n_m < i. \end{cases} \quad (38)$$

Furthermore,  $b$  satisfies the KKT conditions (29) to (34) if and only if, for all  $i \in [k-1]$ ,

$$v_{\langle i \rangle} \begin{cases} \geq C + \gamma & : i \leq n_u \\ \in [\gamma, C + \gamma] & : n_u < i \leq n_u + n_m \\ \leq \gamma & : n_u + n_m < i. \end{cases} \quad (39)$$

*Proof.* First, we prove  $\mathbf{O}b = \gamma \mathbb{1}$  which is equivalent to  $[\mathbf{O}b]_j = \gamma$  for all  $j \in [k-1]$ . This is a straightforward calculation:

$$\begin{aligned}
 [\mathbf{O}b]_j &= \mathbb{1}'b = \sum_{i \in [k-1]} b_{\langle i \rangle} \\
 &= \sum_{i \in [k-1]: i \leq n_u} b_{\langle i \rangle} + \sum_{i \in [k-1]: n_u < i \leq n_u + n_m} b_{\langle i \rangle} + \sum_{i \in [k-1]: n_u + n_m < i} b_{\langle i \rangle} \\
 &= \sum_{i \in [k-1]: i \leq n_u} C + \sum_{i \in [k-1]: n_u < i \leq n_u + n_m} v_{\langle i \rangle} - \gamma \\
 &= C \cdot n_u + \mathcal{S}^{(n_m^t, n_u^t)} - n_m \gamma \\
 &= \gamma.
 \end{aligned}$$

Since  $[(\mathbf{I} + \mathbf{O})b]_i = [\mathbf{I}b]_i + [\mathbf{O}b]_i$ , the identity (38) now follows immediately.

Next, we prove the ‘‘Furthermore’’ part. First, we prove the ‘‘only if’’ direction. By assumption, we have  $b = \tilde{b}$  and so  $\gamma = \tilde{\gamma}$ . Furthermore, from Lemma A.9 we have  $\tilde{b} = \text{clip}_C(v - \tilde{\gamma} \mathbb{1})$  and so  $b = \text{clip}_C(v - \gamma \mathbb{1})$ . To proceed, recall that by construction, we have

$$b_{\langle i \rangle} = \begin{cases} C & : i \leq n_u \\ v - \gamma & : n_u < i \leq n_u + n_m \\ 0 & : n_u + n_m < i \end{cases}$$

Thus, if  $i \leq n_u$ , then  $C = b_{\langle i \rangle} = [\text{clip}_C(v - \gamma \mathbb{1})]_{\langle i \rangle}$  implies that  $v_{\langle i \rangle} - \gamma \geq C$ . If  $n_u < i \leq n_u + n_m$ , then  $b_{\langle i \rangle} = v_{\langle i \rangle} - \gamma$ . Since  $b_j \in [0, C]$  for all  $j \in [k-1]$ , we have in particular that  $v_{\langle i \rangle} - \gamma \in [0, C]$ . Finally, if  $n_u + n_m < i$ , then  $0 = b_{\langle i \rangle} = [\text{clip}_C(v - \gamma \mathbb{1})]_{\langle i \rangle}$  implies that  $v - \gamma \leq 0$ . In summary,

$$v_{\langle i \rangle} - \gamma \begin{cases} \geq C & : i \leq n_u \\ \in [0, C] & : n_u < i \leq n_u + n_m \\ \leq 0 & : n_u + n_m < i. \end{cases}$$

Note that the above identity immediately implies (39).

Next, we prove the ‘‘if’’ direction. Using (38) and (39), we have

$$[(\mathbf{I} + \mathbf{O})b]_{\langle i \rangle} - v_{\langle i \rangle} \begin{cases} \leq 0 & : i \leq n_u \\ = 0 & : n_u < i \leq n_u + n_m \\ \geq 0 & : n_u + n_m < i. \end{cases}$$

For each  $i \in [k-1]$ , define  $\lambda_i, \mu_i \in \mathbb{R}$  where

$$\lambda_{\langle i \rangle} = \begin{cases} -([(\mathbf{I} + \mathbf{O})b]_{\langle i \rangle} - v_{\langle i \rangle}) & : i \leq n_u \\ 0 & : n_u < i \leq n_u + n_m \\ 0 & : n_u + n_m < i \end{cases}$$

and

$$\mu_{\langle i \rangle} = \begin{cases} 0 & : i \leq n_u \\ 0 & : n_u < i \leq n_u + n_m \\ [(\mathbf{I} + \mathbf{O})b]_{\langle i \rangle} - v_{\langle i \rangle} & : n_u + n_m < i. \end{cases}$$

It is straightforward to verify that all of (29) to (34) are satisfied for all  $i \in [k-1]$ , i.e., the KKT conditions hold at  $b$ .  $\square$

Recall that we use indices with angle brackets  $\langle 1 \rangle, \langle 2 \rangle, \dots, \langle k-1 \rangle$  to denote a fixed permutation of  $[k-1]$  such that

$$v_{\langle 1 \rangle} \geq v_{\langle 2 \rangle} \geq \dots \geq v_{\langle k-1 \rangle}.$$

**Corollary A.17.** Let  $t \in [\ell]$  and  $\tilde{b}$  be the unique global minimum of the optimization (4). Then  $\widehat{b}^t = \tilde{b}$  if and only if `KKT_cond()` returns true during the  $t$ -th iteration of Algorithm 2.

*Proof.* First, by Lemma A.9 we have  $\widehat{b}^t = \tilde{b}$  if and only if  $\widehat{b}^t$  satisfies the KKT conditions (29) to (34). From Lemma A.14, we have  $\widehat{b}^{(n_m^t, n_u^t)} = \widehat{b}^t$  and  $\widehat{\gamma}^{(n_m^t, n_u^t)} = \widehat{\gamma}^t$ . To simplify notation, let  $\gamma = \widehat{\gamma}^{(n_m^t, n_u^t)}$ . By Lemma A.16,  $\widehat{b}^{(n_m^t, n_u^t)}$  satisfies the KKT conditions (29) to (34) if and only if the following are true:

$$v_{\langle i \rangle} \begin{cases} \geq C + \gamma & : i \leq n_u^t \\ \in [\gamma, C + \gamma] & : n_u^t < i \leq n_u^t + n_m^t \\ \leq \gamma & : n_u^t + n_m^t < i. \end{cases}$$

Since  $v_{\langle 1 \rangle} \geq v_{\langle 2 \rangle} \geq \dots$ , the above system of inequalities holds for all  $i \in [k-1]$  if and only if

$$\begin{cases} C + \gamma \leq v_{\langle n_u^t \rangle} & : \text{if } n_u^t > 0. \\ \gamma \leq v_{\langle n_u^t + n_m^t \rangle} \text{ and } v_{\langle n_u^t + 1 \rangle} \leq C + \gamma & : \text{if } n_m^t > 0, \\ v_{\langle n_u^t + n_m^t + 1 \rangle} \leq \gamma & : \text{if } n_u^t + n_m^t < k - 1. \end{cases}$$

Note that the above system holds if and only if `KKT_cond()` returns true. □

#### A.5.4. THE VARIABLES $n_m^t$ AND $n_u^t$

**Definition A.18.** Define the set  $\text{vals}^+ = \{(v_j, \text{dn}, j) : v_j > 0, j = 1, \dots, k-1\} \cup \{(v_j - C, \text{up}, j) : v_j > C, j = 1, \dots, k-1\}$ . Sort the set  $\text{vals}^+ = \{(\gamma_1, \delta_1, j_1), \dots, (\gamma_\ell, \delta_\ell, j_\ell)\}$  so that the ordering of  $\{(\gamma_1, \delta_1), \dots, (\gamma_\ell, \delta_\ell)\}$  is identical to  $\text{vals}$  from Subroutine 3-line 2.

To illustrate the definitions, we consider the following running example

$\langle j \rangle$	=	$\langle 1 \rangle$	$\langle 2 \rangle$	$\langle 3 \rangle$	$\langle 4 \rangle$	$\langle 5 \rangle$	$\langle 6 \rangle$	$\langle 7 \rangle$	$\langle 8 \rangle$	$\langle 9 \rangle$	$\langle 10 \rangle$	$\langle 11 \rangle$	$\langle 12 \rangle$	$\langle 13 \rangle$	$\langle 14 \rangle$
$v_{\langle j \rangle}$	=	1.8	1.4	1.4	1.4	1.2	0.7	0.4	0.4	0.4	0.1	-0.2			
$t$	=	1	2	3	4	5	6	7	8	9	10	11	12	13	14
$\gamma_t$	=	1.8	1.4	1.4	1.4	1.2	0.8	0.7	0.4	0.4	0.4	0.4	0.4	0.2	0.1
$\delta_t$	=	dn	dn	dn	dn	dn	up	dn	up	up	up	dn	dn	up	dn

**Definition A.19.** Define

$$u(j) := \max\{\tau \in [\ell] : v_{\langle j \rangle} - C = \gamma_\tau\}, \quad \text{and} \quad d(j) := \max\{\tau \in [\ell] : v_{\langle j \rangle} = \gamma_\tau\}, \quad (40)$$

where  $\max \emptyset = \ell + 1$ .

Below, we compute  $d(3)$ ,  $d(6)$  and  $u(3)$  for our running example.

				$d(3)$											
				↓											
$t$	=	1	2	3	4	5	6	7	8	9	10	11	12	13	14
$\gamma_t$	=	1.8	1.4	1.4	1.4	1.2	0.8	0.7	0.4	0.4	0.4	0.4	0.4	0.2	0.1
$\delta_t$	=	dn	dn	dn	dn	dn	up	dn	up	up	up	dn	dn	up	dn

**Definition A.20.** Define the following sets

$$\begin{aligned} \text{crit}_1(v) &= \{\tau \in [\ell] : \gamma_\tau > \gamma_{\tau+1}\} \\ \text{crit}_2(v) &= \{\tau \in [\ell] : \gamma_\tau = \gamma_{\tau+1}, \delta_\tau = \text{up}, \delta_{\tau+1} = \text{dn}\} \end{aligned}$$

where  $\gamma_{\ell+1} = 0$ .

Below, we illustrate the definition in our running example. The arrows  $\downarrow$  and  $\Downarrow$  point to elements of  $\text{crit}_1(v)$  and  $\text{crit}_2(v)$ , respectively.

$t$	$=$	$\downarrow$				$\downarrow$	$\downarrow$	$\downarrow$	$\downarrow$		$\Downarrow$		$\downarrow$	$\downarrow$	$\downarrow$
$t$	$=$	1	2	3	4	5	6	7	8	9	10	11	12	13	14
$\gamma_t$	$=$	1.8	1.4	1.4	1.4	1.2	0.8	0.7	0.4	0.4	0.4	0.4	0.4	0.2	0.1
$\delta_t$	$=$	dn	dn	dn	dn	dn	up	dn	up	up	up	dn	dn	up	dn

Later, we will show that Algorithm 2 will halt and output the global optimizer  $\tilde{b}$  on or before the  $t$ -th iteration where  $t \in \text{crit}_1(v) \cup \text{crit}_2(v)$ .

**Lemma A.21.** *Suppose that  $t \in \text{crit}_1(v)$ . Then*

$$\#\{j \in [k-1] : d(j) \leq t\} = \#\{\tau \in [t] : \delta_\tau = \text{dn}\}, \quad \text{and} \quad \#\{j \in [k-1] : u(j) \leq t\} = \#\{\tau \in [t] : \delta_\tau = \text{up}\}.$$

*Proof.* First, we observe that

$$\#\{\tau \in [t] : \delta_\tau = \text{up}\} = \#\{(\gamma, \delta, j') \in \text{vals}^+ : \delta = \text{up}, \gamma \geq \gamma_t\}$$

Next, note that  $j \mapsto (\gamma_{d(j)}, \text{up}, \langle j \rangle)$  is a bijection from  $\{j \in [k-1] : d(j) \leq t\}$  to  $\{(\gamma, \delta, j') \in \text{vals}^+ : \delta = \text{up}, \gamma \geq \gamma_t\}$ . To see this, we view the permutation  $\langle 1 \rangle, \langle 2 \rangle, \dots$  viewed as a bijective mapping  $\langle \cdot \rangle : [k-1] \rightarrow [k-1]$  given by  $j \mapsto \langle j \rangle$ . Denote by  $\rangle \cdot \langle$  the inverse of  $\langle \cdot \rangle$ . Then the (two-sided) inverse to  $j \mapsto (\gamma_{d(j)}, \text{up}, \langle j \rangle)$  is clearly given by  $(\gamma, \text{up}, j') \mapsto j' \langle$ . This proves the first identity of the lemma.

The proof of the second identity is completely analogous.  $\square$

**Lemma A.22.** *The functions  $u$  and  $d : [k-1] \rightarrow [\ell+1]$  are non-decreasing. Furthermore, for all  $j \in [k-1]$ , we have  $u(j) < d(j)$ .*

*Proof.* Let  $j', j'' \in [k-1]$  be such that  $j' < j''$ . By the sorting, we have  $v_{\langle j' \rangle} \geq v_{\langle j'' \rangle}$ . Now, suppose that  $d(j') > d(j'')$ , then by construction we have  $\gamma_{d(j')} < \gamma_{d(j'')}$ . On the other hand, we have

$$\gamma_{d(j')} = v_{\langle j' \rangle} \geq v_{\langle j'' \rangle} = \gamma_{d(j'')}$$

which is a contradiction.

For the ‘‘Furthermore’’ part, suppose the contrary that  $u(j) \geq d(j)$ . Then we have  $\gamma_{u(j)} \leq \gamma_{d(j)}$ . However, by definition, we have  $\gamma_{u(j)} = v_{\langle j \rangle} > v_{\langle j \rangle} - C = \gamma_{d(j)}$ . This is a contradiction.  $\square$

**Lemma A.23.** *Let  $t \in \text{crit}_1(v)$ . Then  $n_u^t = \#\{j \in [k-1] : u(j) \leq t\}$ . Furthermore,  $[n_u^t] = \{j \in [k-1] : u(j) \leq t\}$ . Equivalently, for each  $j \in [k-1]$ , we have  $j \leq n_u^t$  if and only if  $u(j) \leq t$ .*

*Proof.* First, we note that

$$\begin{aligned} n_u^t &= \#\{\tau \in [t] : \delta_\tau = \text{up}\} \quad \because \text{Subroutine 4-line 2} \\ &= \#\{j \in [k-1] : u(j) \leq t\} \quad \because \text{Lemma A.21} \end{aligned}$$

This proves the first part. For the ‘‘Furthermore’’ part, let  $N := \#\{j \in [k-1] : u(j) \leq t\}$ . Since  $u$  is monotonic non-decreasing (Lemma A.22), we have  $\{j \in [k-1] : u(j) \leq t\} = [N]$ . Since  $N = n_u^t$  by the first part, we are done.  $\square$

**Lemma A.24.** *Let  $\hat{t}, \check{t} \in \text{crit}_1(v)$  be such that there exists  $t \in [\ell]$  where*

$$n_m^t = \#\{j \in [k-1] : d(j) \leq \check{t}\} - \#\{j \in [k-1] : u(j) \leq \hat{t}\}. \quad (41)$$

*Then  $d(j) \leq \check{t}$  and  $\hat{t} < u(j)$  if and only if  $n_u^{\hat{t}} < j \leq n_u^{\hat{t}} + n_m^t$ .*

*Proof.* By Lemma A.23 and (41), we have  $\#\{j \in [k-1] : d(j) \leq \check{t}\} = n_u^{\check{t}} + n_m^{\check{t}}$ . By Lemma A.22,  $d$  is monotonic non-decreasing and so  $[n_u^{\check{t}} + n_m^{\check{t}}] = \{j \in [k-1] : d(j) \leq \check{t}\}$ . Now,

$$\begin{aligned} & \{j \in [k-1] : d(j) \leq \check{t}, \hat{t} < u(j)\} \\ &= \{j \in [k-1] : d(j) \leq \check{t}\} \cap \{j \in [k-1] : \hat{t} < u(j)\} \\ &= \{j \in [k-1] : d(j) \leq \check{t}\} \setminus \{j \in [k-1] : u(j) \leq \hat{t}\} \\ &= [n_u^{\check{t}} + n_m^{\check{t}}] \setminus [n_u^{\hat{t}}], \end{aligned}$$

where in the last equality, we used Lemma A.23. □

**Corollary A.25.** *Let  $t \in \text{crit}_1(v)$ . Then  $d(j) \leq t$  and  $t < u(j)$  if and only if  $n_u^t < j \leq n_u^t + n_m^t$ .*

*Proof.* We apply Lemma A.24 with  $t = \hat{t} = \check{t}$ , which requires checking that

$$n_m^t = \#\{j \in [k-1] : d(j) \leq t\} - \#\{j \in [k-1] : u(j) \leq t\}.$$

This is true because from Subroutine 4-line 2 and 5, we have

$$n_m^t = \#\{\tau \in [t] : \delta_\tau = \text{dn}\} - \#\{\tau \in [t] : \delta_\tau = \text{up}\}.$$

Applying Lemma A.21, we are done. □

**Lemma A.26.** *Let  $t \in \text{crit}_1(v)$ . Let  $\varepsilon > 0$  be such that for all  $\tau, \tau' \in \text{crit}_1(v)$  where  $\tau' < \tau$ , we have  $\gamma_{\tau'} - \varepsilon > \gamma_\tau$ . Then  $(n_m^t, n_u^t) = (n_m^{\gamma_t - \varepsilon}, n_u^{\gamma_t - \varepsilon})$ .*

*Proof.* We claim that

$$v_{(j)} - \gamma_t + \varepsilon \begin{cases} < 0 & : t < d(j) \\ \in (0, C) & : d(j) \leq t < u(j) \\ > C & : u(j) \leq t. \end{cases} \quad (42)$$

To prove the  $t < d(j)$  case of (42), we have

$$\begin{aligned} v_{(j)} - \gamma_t + \varepsilon &= \gamma_{d(j)} - \gamma_t + \varepsilon \quad \because (40) \\ &< -\varepsilon + \varepsilon = 0 \quad \because t < d(j) \text{ implies that } \gamma_t - \varepsilon > \gamma_{d(j)}. \end{aligned}$$

To prove the  $d(j) \leq t < u(j)$  case of (42), we note that

$$\begin{aligned} v_{(j)} - \gamma_t + \varepsilon &= \gamma_{d(j)} - \gamma_t + \varepsilon \quad (40) \\ &\geq \varepsilon > 0 \quad \because d(j) \leq t \text{ implies } \gamma_{d(j)} \geq \gamma_t. \end{aligned}$$

For the other inequality,

$$\begin{aligned} v_{(j)} - \gamma_t + \varepsilon &= \gamma_{u(j)} + C - \gamma_t + \varepsilon \quad \because (40) \\ &< -\varepsilon + C + \varepsilon = C \quad \because t < u(j) \text{ implies } \gamma_t - \varepsilon > \gamma_{u(j)}. \end{aligned}$$

Finally, we prove the  $u(j) \leq t$  case of (42). Note that

$$\begin{aligned} v_{(j)} - \gamma_t + \varepsilon &= \gamma_{u(j)} + C - \gamma_t + \varepsilon \quad \because (40) \\ &\geq C + \varepsilon > C \quad \because u(j) \leq t \text{ implies that } \gamma_{u(j)} \geq \gamma_t. \end{aligned}$$

Thus, we have proven (42). By Lemma A.23 and Corollary A.25, (42) can be rewritten as

$$v_{(j)} - \gamma_t + \varepsilon \begin{cases} < 0 & : n_u^t + n_m^t < j, \\ \in (0, C) & : n_u^t < j \leq n_u^t + n_m^t, \\ > C & : j \leq n_u^t. \end{cases} \quad (43)$$

Thus, we have  $I_u^{\gamma_t - \varepsilon} = \{\langle 1 \rangle, \dots, \langle n_u^t \rangle\}$  and  $I_m^{\gamma_t - \varepsilon} = \{\langle n_u^t + 1 \rangle, \dots, \langle n_u^t + n_m^t \rangle\}$ . By the definitions of  $n_u^{\gamma_t - \varepsilon}$  and  $n_m^{\gamma_t - \varepsilon}$ , we are done. □



							$\check{t}$		$t$		$\hat{t}$		$d(9)$
	1	2	3	4	5	6	$\downarrow$	8	$\downarrow$	11	$\downarrow$	13	$\downarrow$
$\gamma_t =$	1.8	1.4	1.4	1.4	1.2	0.8	0.7	0.4	0.4	0.4	0.4	0.2	0.1
$\delta_t =$	dn	dn	dn	dn	dn	up	dn	up	up	up	dn	dn	up

Figure 3. Example of a critical iterate type 2. The first case that  $\check{t} < d(j)$  where  $j = 9$ .

**Lemma A.27.** Let  $t \in \text{crit}_2(v)$ . Then  $(n_m^t, n_u^t) = (n_m^{\check{t}}, n_u^{\check{t}})$ .

*Proof.* Let  $\hat{t} \in \text{crit}_1(v)$  be such that  $\gamma_{\hat{t}} = \gamma_t$ , and  $\check{t} = \max\{\tau \in \text{crit}_1(v) : \gamma_\tau > \gamma_t\}$ . We claim that

$$v_{\langle j \rangle} - \gamma_{\check{t}} \begin{cases} \leq 0 & : \check{t} < d(j), \\ \in (0, C) & : d(j) \leq \check{t}, \hat{t} < u(j), \\ \geq C & : u(j) \leq \hat{t}. \end{cases} \quad (44)$$

Note that by definition, we have  $\gamma_{\check{t}} > \gamma_{\hat{t}}$ , which implies that  $\check{t} < \hat{t}$ .

Consider the first case of (44) that  $\check{t} < d(j)$ . See the running example Figure 3. We have by construction that  $v_{\langle j \rangle} = \gamma_{d(j)}$  and so  $v_{\langle j \rangle} - \gamma_{\check{t}} = \gamma_{d(j)} - \gamma_{\check{t}} \leq 0$ .

Next, consider the case when  $d(j) \leq \check{t}$  and  $\hat{t} < u(j)$ . Thus,

$$\begin{aligned} v_{\langle j \rangle} - \gamma_{\check{t}} &> v_{\langle j \rangle} - \gamma_{\hat{t}} && \because \gamma_{\check{t}} > \gamma_{\hat{t}} \\ &= \gamma_{d(j)} - \gamma_{\hat{t}} && \because \text{definition of } d(j) \\ &\geq 0 && \because d(j) \leq \check{t} \implies \gamma_{d(j)} \geq \gamma_{\hat{t}}. \end{aligned}$$

On the other hand

$$\begin{aligned} v_{\langle j \rangle} - \gamma_{\check{t}} &= \gamma_{u(j)} + C - \gamma_{\check{t}} && \because \text{definition of } u(j) \\ &< C && \because \hat{t} < u(j) \implies \gamma_{\check{t}} > \gamma_{u(j)} \end{aligned}$$

Thus, we've shown that in the second case, we have  $v_{\langle j \rangle} - \gamma_{\check{t}} \in (0, C)$ .

We consider the final case that  $u(j) \leq \hat{t}$ . We have

$$\begin{aligned} v_{\langle j \rangle} - \gamma_{\check{t}} &= \gamma_{u(j)} + C - \gamma_{\check{t}} && \because \text{definition of } t \\ &\geq C && \because u(j) \leq \hat{t} \implies \gamma_{u(j)} \geq \gamma_{\hat{t}}. \end{aligned}$$

Thus, we have proven (44).

Next, we claim that  $t, \hat{t}, \check{t}$  satisfy the condition (41) of Lemma A.24, i.e.,

$$n_m^t = \#\{j \in [k-1] : d(j) \leq \check{t}\} - \#\{j \in [k-1] : u(j) \leq \hat{t}\}.$$

To this end, we first recall that

$$n_m^t = \#\{\tau \in [t] : \delta_\tau = \text{dn}\} - \#\{\tau \in [t] : \delta_\tau = \text{up}\}.$$

By assumption on  $t$ , for all  $\tau$  such that  $\check{t} < \tau \leq t$ , we have  $\delta_\tau = \text{up}$ . Thus,

$$\#\{\tau \in [t] : \delta_\tau = \text{dn}\} = \#\{\tau \in [\check{t}] : \delta_\tau = \text{dn}\} = \#\{j \in [k-1] : d(j) \leq \check{t}\}$$

where for the last equality, we used Lemma A.21. Similarly, for all  $\tau$  such that  $t < \tau \leq \hat{t}$ , we have  $\delta_\tau = \text{dn}$ . Thus, we get that analogous result

$$n_u^t = \#\{\tau \in [t] : \delta_\tau = \text{up}\} = \#\{\tau \in [\hat{t}] : \delta_\tau = \text{up}\} = \#\{j \in [k-1] : u(j) \leq \hat{t}\} = n_u^{\hat{t}}. \quad (45)$$

Thus, we have verified the condition (41) of Lemma A.24. Now, applying Lemma A.23 and Lemma A.24, we get

$$v_{\langle j \rangle} - \gamma_t \begin{cases} \leq 0 & : n_u^t + n_m^t < j, \\ \in (0, C) & : n_u^t < j \leq n_u^t + n_m^t \\ \geq C & : j \leq n_u^t. \end{cases} \quad (46)$$

By (45) and that  $\gamma_t = \gamma_t$ , the above reduces to

$$v_{\langle j \rangle} - \gamma_t \begin{cases} \leq 0 & : n_u^t + n_m^t < j, \\ \in (0, C) & : n_u^t < j \leq n_u^t + n_m^t \\ \geq C & : j \leq n_u^t. \end{cases} \quad (47)$$

Thus,  $I_u^{\gamma_t} = \{\langle 1 \rangle, \dots, \langle n_u^t \rangle\}$  and  $I_m^{\gamma_t} = \{\langle n_u^t + 1 \rangle, \dots, \langle n_u^t + n_m^t \rangle\}$ . By the definitions of  $n_u^{\gamma_t}$  and  $n_m^{\gamma_t}$ , we are done.  $\square$

#### A.5.5. PUTTING IT ALL TOGETHER

If  $v_{\max} \leq 0$ , then Algorithm 2 returns  $\emptyset$ .

Otherwise, by Lemma A.9, we have  $\tilde{\gamma} \in (0, v_{\max})$ .

**Lemma A.28.** *Let  $t \in [\ell]$  be such that  $(n_m^t, n_u^t) = (n_m^{\tilde{\gamma}}, n_u^{\tilde{\gamma}})$ . Then during the  $t$ -th loop of Algorithm 2 we have  $\tilde{b} = \hat{b}^t$  and  $\text{KKT\_cond}()$  returns true. Consequently, Algorithm 2 returns the optimizer  $\tilde{b}$  on or before the  $t$ -th iteration.*

*Proof.* We have

$$\begin{aligned} \tilde{b} &= \hat{b}^{(n_m^{\tilde{\gamma}}, n_u^{\tilde{\gamma}})} \quad \because \text{Lemma A.15} \\ &= \hat{b}^{(n_m^t, n_u^t)} \quad \because \text{Assumption} \\ &= \hat{b}^t \quad \because \text{Lemma A.14.} \end{aligned}$$

Thus, by Corollary A.17  $\text{KKT\_cond}()$  returns true on the  $t$ -th iteration. This means that Algorithm 2 halts on or before iteration  $t$ . Let  $\tau \in [\ell]$  be the iteration where Algorithm 2 halts and outputs  $\hat{b}^\tau$ . Then  $\tau \leq t$ . Furthermore, by Corollary A.17,  $\hat{b}^\tau = \tilde{b}$ , which proves the ‘‘Consequently’’ part of the lemma.  $\square$

By Lemma A.28, it suffices to show that  $(n_m^t, n_u^t) = (n_m^{\tilde{\gamma}}, n_u^{\tilde{\gamma}})$  for some  $t \in [\ell]$ .

We first consider the case when  $\tilde{\gamma} \neq \gamma_t$  for any  $t \in \text{crit}_1(v)$ . Thus, there exists  $t \in \text{crit}_1(v)$  such that  $\gamma_{t+1} < \tilde{\gamma} < \gamma_t$ , where we recall that  $\gamma_{\ell+1} := 0$ .

Now, we return to the proof of Theorem 3.4.

$$\begin{aligned} (n_m^t, n_u^t) &= (n_m^{\gamma_t - \varepsilon}, n_u^{\gamma_t - \varepsilon}) \quad \because \text{Lemma A.26} \\ &= (n_m^{\tilde{\gamma}}, n_u^{\tilde{\gamma}}) \quad \because \text{Lemma A.12, and that both } \tilde{\gamma} \text{ and } \gamma_t - \varepsilon \in (\gamma_{t+1}, \gamma_t). \end{aligned}$$

Thus, Lemma A.28 implies the result of Theorem 3.4 under the assumption that  $\tilde{\gamma} \neq \gamma_t$  for any  $t \in \text{crit}_1(v)$ .

Next, we consider when  $\tilde{\gamma} = \gamma_t$  for some  $t \in \text{crit}_1(v)$ . There are three possibilities:

1. There does not exist  $j \in [k-1]$  such that  $v_{\langle j \rangle} = \gamma_t$ ,
2. There does not exist  $j \in [k-1]$  such that  $v_{\langle j \rangle} - C = \gamma_t$ ,
3. There exist  $j_1, j_2 \in [k-1]$  such that  $v_{\langle j_1 \rangle} = \gamma_t$  and  $v_{\langle j_2 \rangle} - C = \gamma_t$ .

First, we consider case 1. We claim that

$$(n_m^{\gamma_t}, n_u^{\gamma_t}) = (n_m^{\gamma_t - \varepsilon'}, n_u^{\gamma_t - \varepsilon'}) \quad \text{for all } \varepsilon' > 0 \text{ sufficiently small.} \quad (48)$$

We first note that  $n_u^{\gamma_t} = n_u^{\gamma_t - \varepsilon'}$  for all  $\varepsilon' > 0$  sufficiently small. To see this, let  $i \in [k-1]$  be arbitrary. Note that

$$\begin{aligned} i \in I_u^{\gamma_t} &\iff v_i - \gamma_t \geq C \iff v_i - \gamma_t + \varepsilon' \geq C, \forall \varepsilon' > 0, \text{ sufficiently small} \\ &\iff i \in I_u^{\gamma_t - \varepsilon'}, \forall \varepsilon' > 0, \text{ sufficiently small.} \end{aligned}$$

Next, we show that  $n_m^{\gamma_t} = n_m^{\gamma_t - \varepsilon'}$  for all  $\varepsilon' > 0$  sufficiently small. To see this, let  $i \in [k-1]$  be arbitrary. Note that

$$\begin{aligned} i \in I_m^{\gamma_t} &\iff v_i - \gamma_t \in (0, C) \stackrel{\dagger}{\iff} v_i - \gamma_t + \varepsilon' \in (0, C), \forall \varepsilon' > 0, \text{ sufficiently small} \\ &\iff i \in I_m^{\gamma_t - \varepsilon'}, \forall \varepsilon' > 0, \text{ sufficiently small} \end{aligned}$$

where at “ $\stackrel{\dagger}{\iff}$ ”, we used the fact that  $v_i - \gamma_t \neq 0$  for any  $i \in [k-1]$ . Thus, we have proven (48). Taking  $\varepsilon' > 0$  so small so that both (48) and the condition in Lemma A.26 hold, we have

$$(n_m^t, n_u^t) = (n_m^{\gamma_t - \varepsilon'}, n_u^{\gamma_t - \varepsilon'}) = (n_m^{\gamma_t}, n_u^{\gamma_t}) = (n_m^{\tilde{\gamma}}, n_u^{\tilde{\gamma}}).$$

This proves Theorem 3.4 under case 1.

Next, we consider case 2. We claim that

$$(n_m^{\gamma_t}, n_u^{\gamma_t}) = (n_m^{\gamma_t + \varepsilon''}, n_u^{\gamma_t + \varepsilon''}) \quad \text{for all } \varepsilon'' > 0 \text{ sufficiently small.} \quad (49)$$

We first note that  $n_u^{\gamma_t} = n_u^{\gamma_t - \varepsilon''}$  for all  $\varepsilon'' > 0$  sufficiently small. To see this, let  $i \in [k-1]$  be arbitrary. Note that

$$\begin{aligned} i \in I_u^{\gamma_t} &\iff v_i - \gamma_t \geq C \stackrel{\ddagger}{\iff} v_i - \gamma_t - \varepsilon'' \geq C, \forall \varepsilon'' > 0, \text{ sufficiently small} \\ &\iff i \in I_u^{\gamma_t + \varepsilon''}, \forall \varepsilon'' > 0, \text{ sufficiently small.} \end{aligned}$$

where at “ $\stackrel{\ddagger}{\iff}$ ”, we used the fact that  $v_i - \gamma_t \neq C$  for any  $i \in [k-1]$ . Next, we show that  $n_m^{\gamma_t} = n_m^{\gamma_t - \varepsilon''}$  for all  $\varepsilon'' > 0$  sufficiently small. To see this, let  $i \in [k-1]$  be arbitrary. Note that

$$\begin{aligned} i \in I_m^{\gamma_t} &\iff v_i - \gamma_t \in (0, C) \stackrel{\ddagger}{\iff} v_i - \gamma_t - \varepsilon'' \in (0, C), \forall \varepsilon'' > 0, \text{ sufficiently small} \\ &\iff i \in I_m^{\gamma_t + \varepsilon''}, \forall \varepsilon'' > 0, \text{ sufficiently small} \end{aligned}$$

where again at “ $\stackrel{\ddagger}{\iff}$ ”, we used the fact that  $v_i - \gamma_t \neq C$  for any  $i \in [k-1]$ . Thus, we have proven (49). Since  $\tilde{\gamma} = \gamma_t \in (0, v_{\max})$  and  $\gamma_1 = v_{\max}$ , we have in particular that  $\gamma_t < \gamma_1$ . Thus, there exists  $\tau \in \text{crit}_1(v)$  such that  $\tau < t$  and  $\gamma_t < \gamma_\tau$ . Furthermore, we can choose  $\tau$  such that for all  $\gamma \in (\gamma_t, \gamma_\tau)$ ,  $\gamma \notin \text{crit}_1(v)$ . Let  $\varepsilon'' > 0$  be so small that  $\gamma_t + \varepsilon'', \gamma_\tau - \varepsilon'' \in (\gamma_t, \gamma_\tau)$ , and furthermore both (49) and the condition in Lemma A.26 hold. We have

$$\begin{aligned} (n_m^\tau, n_u^\tau) &= (n_m^{\gamma_\tau - \varepsilon''}, n_u^{\gamma_\tau - \varepsilon''}) \quad \because \text{Lemma A.26} \\ &= (n_m^{\gamma_t + \varepsilon''}, n_u^{\gamma_t + \varepsilon''}) \quad \because \text{Lemma A.12 and } \gamma_t + \varepsilon'', \gamma_\tau - \varepsilon'' \in (\gamma_t, \gamma_\tau) \\ &= (n_m^{\gamma_t}, n_u^{\gamma_t}) \quad \because (49) \\ &= (n_m^{\tilde{\gamma}}, n_u^{\tilde{\gamma}}) \quad \because \text{Assumption.} \end{aligned}$$

This proves Theorem 3.4 under case 2.

Finally, we consider the last case. Under the assumptions, we have  $t \in \text{crit}_2(v)$ . Then Lemma A.27  $(n_m^t, n_u^t) = (n_m^{\gamma_t}, n_u^{\gamma_t}) = (n_m^{\tilde{\gamma}}, n_u^{\tilde{\gamma}})$ . Thus, we have proven Theorem 3.4 under case 3.  $\square$

## A.6. Experiments

The Walrus solver is available at:

<https://github.com/YutongWangUMich/liblinear>

The actual implementation is in the file `linear.cpp` in the class `Solver_MCSVM_LW`.

All code for downloading the datasets used, generating the train/test split, running the experiments and generating the figures are included. See the `README.md` file for more information.

All experiments are run on a single machine with the following specifications:

Operating system and kernel:

4.15.0-122-generic #124-Ubuntu SMP Thu Oct 15 13:03:05 UTC 2020 x86\_64 GNU/Linux

Processor:

Intel(R) Core(TM) i7-6850K CPU @ 3.60GH

Memory:

31GiB System memory

#### A.6.1. ON SHARKS LINEAR WW-SVM SOLVER

Shark’s linear WW-SVM solver is publicly available in the GitHub repository <https://github.com/Shark-ML>. Specifically, the C++ code is in `Algorithms/QP/QpMcLinear.h` in the class `QpMcLinearWW`. Our reimplementation follows their implementation with two major differences. In our implementations, neither Shark nor Walrus use the shrinking heuristic. Furthermore, we use a stopping criterion based on duality gap, following (Steinwart et al., 2011).

We also remark that Shark solves the following variant of the WW-SVM which is equivalent to ours after a change of variables. Let  $0 < A \in \mathbb{R}$  be a hyperparameter.

$$\min_{\mathbf{u} \in \mathbb{R}^{d \times k}} F_A(\mathbf{u}) := \frac{1}{2} \|\mathbf{u}\|_F^2 + A \sum_{i=1}^n \sum_{\substack{j \in [k]: \\ j \neq y_i}} \text{hinge}((u'_{y_i} x_i - u'_j x_i)/2). \quad (50)$$

Recall the formulation (P) that we consider in this work, which we repeat here:

$$\min_{\mathbf{w} \in \mathbb{R}^{d \times k}} G_C(\mathbf{w}) := \frac{1}{2} \|\mathbf{w}\|_F^2 + C \sum_{i=1}^n \sum_{\substack{j \in [k]: \\ j \neq y_i}} \text{hinge}(w'_{y_i} x_i - w'_j x_i). \quad (51)$$

The formulation (50) is used by Weston & Watkins (1999), while the formulation (51) is used by Vapnik (1998). These two formulations are equivalent under the change of variables  $\mathbf{w} = \mathbf{u}/2$  and  $A = 4C$ . To see this, note that

$$\begin{aligned} G_C(\mathbf{w}) &= G_C(\mathbf{u}/2) \\ &= \frac{1}{2} \|\mathbf{u}/2\|_F^2 + C \sum_{i=1}^n \sum_{\substack{j \in [k]: \\ j \neq y_i}} \text{hinge}((u'_{y_i} x_i - u'_j x_i)/2) \\ &= \frac{1}{8} \|\mathbf{u}\|_F^2 + C \sum_{i=1}^n \sum_{\substack{j \in [k]: \\ j \neq y_i}} \text{hinge}((u'_{y_i} x_i - u'_j x_i)/2) \\ &= \frac{1}{4} \left( \frac{1}{2} \|\mathbf{u}\|_F^2 + 4C \sum_{i=1}^n \sum_{\substack{j \in [k]: \\ j \neq y_i}} \text{hinge}((u'_{y_i} x_i - u'_j x_i)/2) \right) \\ &= \frac{1}{4} F_{4C}(\mathbf{u}) = \frac{1}{4} F_A(\mathbf{u}). \end{aligned}$$

Thus, we have proven

**Proposition A.29.** *Let  $C > 0$  and  $\mathbf{u} \in \mathbb{R}^{d \times k}$ . Then  $\mathbf{u}$  is a minimizer of  $F_{4C}$  if and only if  $\mathbf{u}/2$  is a minimizer of  $G_C$ .*

In our experiments, we use the above proposition to rescale the variant formulation to the standard formulation.