

---

# $\nu$ -Support Vector Machine as Conditional Value-at-Risk Minimization

---

Akiko Takeda

Keio University, 3-14-1 Hiyoshi, Kouhoku, Yokohama, Kanagawa 223-8522, Japan

TAKEDA@AE.KEIO.AC.JP

Masashi Sugiyama

Tokyo Institute of Technology, 2-12-1 O-okayama, Meguro-ku, Tokyo 152-8552, Japan

SUGI@CS.TITECH.AC.JP

## Abstract

The  $\nu$ -support vector classification ( $\nu$ -SVC) algorithm was shown to work well and provide intuitive interpretations, e.g., the parameter  $\nu$  roughly specifies the fraction of support vectors. Although  $\nu$  corresponds to a fraction, it cannot take the entire range between 0 and 1 in its original form. This problem was settled by a non-convex extension of  $\nu$ -SVC and the extended method was experimentally shown to generalize better than original  $\nu$ -SVC. However, its good generalization performance and convergence properties of the optimization algorithm have not been studied yet. In this paper, we provide new theoretical insights into these issues and propose a novel  $\nu$ -SVC algorithm that has guaranteed generalization performance and convergence properties.

## 1. Introduction

*Support vector classification* (SVC) is one of the most successful classification algorithms in modern machine learning (Schölkopf & Smola, 2002). SVC finds a hyperplane that separates training samples in different classes with maximum margin (Boser et al., 1992). The maximum margin hyperplane was shown to minimize an upper bound of the generalization error according to the Vapnik-Chervonenkis theory (Vapnik, 1995). Thus the generalization performance of SVC is theoretically guaranteed.

SVC was extended to be able to deal with non-separable data by trading the margin size with the

data separation error (Cortes & Vapnik, 1995). This soft-margin formulation is commonly referred to as  $C$ -SVC since the trade-off is controlled by the parameter  $C$ .  $C$ -SVC was shown to work very well in a wide range of real-world applications (Schölkopf & Smola, 2002).

An alternative formulation of the soft-margin idea is  $\nu$ -SVC (Schölkopf et al., 2000)—instead of the parameter  $C$ ,  $\nu$ -SVC involves another trade-off parameter  $\nu$  that roughly specifies the fraction of support vectors (or sparseness of the solution). Thus, the  $\nu$ -SVC formulation provides us richer interpretation than the original  $C$ -SVC formulation, which would be potentially useful in real applications.

Since the parameter  $\nu$  corresponds to a fraction, it should be able to be chosen between 0 and 1. However, it was shown that admissible values of  $\nu$  are actually limited (Crisp & Burges, 2000; Chang & Lin, 2001). To cope with this problem, Perez-Cruz et al. (2003) introduced the notion of negative margins and proposed extended  $\nu$ -SVC ( $E\nu$ -SVC) which allows  $\nu$  to take the entire range between 0 and 1. They also experimentally showed that the generalization performance of  $E\nu$ -SVC is often better than that of original  $\nu$ -SVC. Thus the extension contributes not only to elucidating the theoretical property of  $\nu$ -SVC, but also to improving its generalization performance.

However, there remain two open issues in  $E\nu$ -SVC. The first issue is that the reason why a high generalization performance can be obtained by  $E\nu$ -SVC was not completely explained yet. The second issue is that the optimization problem involved in  $E\nu$ -SVC is non-convex and theoretical convergence properties of the  $E\nu$ -SVC optimization algorithm have not been studied yet. The purpose of this paper is to provide new theoretical insights into these two issues.

After reviewing existing SVC methods in Section 2, we

---

Appearing in *Proceedings of the 25<sup>th</sup> International Conference on Machine Learning*, Helsinki, Finland, 2008. Copyright 2008 by the author(s)/owner(s).

elucidate the generalization performance of  $E\nu$ -SVC in Section 3. We first show that the  $E\nu$ -SVC formulation could be interpreted as minimization of the *conditional value-at-risk* (CVaR), which is often used in finance (Rockafellar & Uryasev, 2002; Gotoh & Takeda, 2005). Then we give new generalization error bounds based on the CVaR risk measure. This theoretical result justifies the use of  $E\nu$ -SVC.

In Section 4, we address non-convexity of the  $E\nu$ -SVC optimization problem. We first give a new optimization algorithm that is guaranteed to converge to one of the local optima within a finite number of iterations. Based on this improved algorithm, we further show that the global solution can be actually obtained within finite iterations even though the optimization problem is non-convex.

Finally, in Section 5, we give concluding remarks and future prospects. Proofs of all theorems and lemmas are sketched in Appendix unless mentioned.

## 2. Support Vector Classification

In this section, we formulate the classification problem and briefly review support vector algorithms.

### 2.1. Classification Problem

Let us address the classification problem of learning a decision function  $h$  from  $\mathcal{X}$  ( $\subset \mathbb{R}^n$ ) to  $\{\pm 1\}$  based on training samples  $(\mathbf{x}_i, y_i)$  ( $i \in M := \{1, \dots, m\}$ ). We assume that the training samples are i.i.d. following the unknown probability distribution  $P(\mathbf{x}, y)$  on  $\mathcal{X} \times \{\pm 1\}$ .

The goal of the classification task is to obtain a classifier  $h$  that minimizes the generalization error (or the risk):

$$R[h] := \int \frac{1}{2} |h(\mathbf{x}) - y| dP(\mathbf{x}, y),$$

which corresponds to the misclassification rate for unseen test samples.

For the sake of simplicity, we generally focus on linear classifiers, i.e.,

$$h(\mathbf{x}) = \text{sign}(\langle \mathbf{w}, \mathbf{x} \rangle + b), \quad (1)$$

where  $\mathbf{w}$  ( $\in \mathbb{R}^n$ ) is a non-zero normal vector,  $b$  ( $\in \mathbb{R}$ ) is a bias parameter, and  $\text{sign}(\xi) = 1$  if  $\xi \geq 0$  and  $-1$  otherwise.

Most of the discussions in this paper can be directly applicable to non-linear kernel classifiers (Schölkopf & Smola, 2002). Thus we may not lose generality by restricting ourselves to linear classifiers.

### 2.2. Support Vector Classification

The Vapnik-Chervonenkis theory (Vapnik, 1995) showed that a large margin classifier has a small generalization error. Motivated by this theoretical result, Boser et al. (1992) developed an algorithm for finding the hyperplane  $(\mathbf{w}, b)$  with maximum margin:

$$\min_{\mathbf{w}, b} \frac{1}{2} \|\mathbf{w}\|^2 \quad \text{s.t. } y_i(\langle \mathbf{w}, \mathbf{x}_i \rangle + b) \geq 1, \quad i \in M. \quad (2)$$

This is called (hard-margin) support vector classification (SVC) and valid when the training samples are linearly separable. In the following, we omit “ $i \in M$ ” in the constraint for brevity.

### 2.3. $C$ -Support Vector Classification

Cortes and Vapnik (1995) extended the SVC algorithm to non-separable cases and proposed trading the margin size with the data separation error (i.e., “soft-margin”):

$$\begin{aligned} \min_{\mathbf{w}, b, \xi} \quad & \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_{i=1}^m \xi_i \\ \text{s.t. } \quad & y_i(\langle \mathbf{w}, \mathbf{x}_i \rangle + b) \geq 1 - \xi_i, \quad \xi_i \geq 0, \end{aligned}$$

where  $C$  ( $> 0$ ) controls the trade-off. This formulation is usually referred to as  $C$ -SVC, and was shown to work very well in various real-world applications (Schölkopf & Smola, 2002).

### 2.4. $\nu$ -Support Vector Classification

$\nu$ -SVC is another formulation of soft-margin SVC (Schölkopf et al., 2000):

$$\begin{aligned} \min_{\mathbf{w}, b, \xi, \rho} \quad & \frac{1}{2} \|\mathbf{w}\|^2 - \nu \rho + \frac{1}{m} \sum_{i=1}^m \xi_i \\ \text{s.t. } \quad & y_i(\langle \mathbf{w}, \mathbf{x}_i \rangle + b) \geq \rho - \xi_i, \quad \xi_i \geq 0, \quad \rho \geq 0, \end{aligned}$$

where  $\nu$  ( $\in \mathbb{R}$ ) is the trade-off parameter.

Schölkopf et al. (2000) showed that if the  $\nu$ -SVC solution yields  $\rho > 0$ ,  $C$ -SVC with  $C = 1/(m\rho)$  produces the same solution. Thus  $\nu$ -SVC and  $C$ -SVC are equivalent. However,  $\nu$ -SVC has additional intuitive interpretations, e.g.,  $\nu$  is an upper bound on the fraction of margin errors and a lower bound on the fraction of support vectors (i.e., sparseness of the solution). Thus, the  $\nu$ -SVC formulation would be potentially more useful than the  $C$ -SVC formulation in real applications.

### 2.5. $E\nu$ -SVC

Although  $\nu$  has an interpretation as a fraction, it cannot always take its full range between 0 and 1 (Crisp & Burges, 2000; Chang & Lin, 2001).

### 2.5.1. ADMISSIBLE RANGE OF $\nu$

For an optimal solution  $\{\alpha_i^C\}_{i=1}^m$  of dual  $C$ -SVC, let

$$\zeta(C) := \frac{1}{Cm} \sum_{i=1}^m \alpha_i^C,$$

$$\nu_{\min} := \lim_{C \rightarrow \infty} \zeta(C) \quad \text{and} \quad \nu_{\max} := \lim_{C \rightarrow 0} \zeta(C).$$

Then, Chang and Lin (2001) showed that for  $\nu \in (\nu_{\min}, \nu_{\max}]$ , the optimal solution set of  $\nu$ -SVC is the same as that of  $C$ -SVC with some  $C$  (not necessarily unique). In addition, the optimal objective value of  $\nu$ -SVC is strictly negative. However, for  $\nu \in (\nu_{\max}, 1]$ ,  $\nu$ -SVC is unbounded, i.e., there exists no solution; for  $\nu \in [0, \nu_{\min}]$ ,  $\nu$ -SVC is feasible with zero optimal objective value, i.e., we end up with just having a trivial solution ( $\mathbf{w} = \mathbf{0}$  and  $b = 0$ ).

### 2.5.2. INCREASING UPPER ADMISSIBLE RANGE

It was shown by Crisp and Burges (2000) that

$$\nu_{\max} = 2 \min(m_+, m_-) / m,$$

where  $m_+$  and  $m_-$  are the number of positive and negative training samples. Thus, when the training samples are balanced (i.e.,  $m_+ = m_-$ ),  $\nu_{\max} = 1$  and therefore  $\nu$  can reach its upper limit 1. When the training samples are imbalanced (i.e.,  $m_+ \neq m_-$ ), Perez-Cruz et al. (2003) proposed modifying the optimization problem of  $\nu$ -SVC as

$$\min_{\mathbf{w}, b, \xi, \rho} \frac{1}{2} \|\mathbf{w}\|^2 - \nu \rho + \frac{1}{m_+} \sum_{i: y_i=1} \xi_i + \frac{1}{m_-} \sum_{i: y_i=-1} \xi_i$$

s.t.  $y_i(\langle \mathbf{w}, \mathbf{x}_i \rangle + b) \geq \rho - \xi_i, \quad \xi_i \geq 0, \quad \rho \geq 0,$

i.e., the effect of positive and negative samples are balanced. Under this modified formulation,  $\nu_{\max} = 1$  holds even when training samples are imbalanced.

For the sake of simplicity, we assume  $m_+ = m_-$  in the rest of this paper; when  $m_+ \neq m_-$ , all the results can be simply extended in a similar way as above.

### 2.5.3. DECREASING LOWER ADMISSIBLE RANGE

When  $\nu \in [0, \nu_{\min}]$ ,  $\nu$ -SVC produces a trivial solution ( $\mathbf{w} = \mathbf{0}$  and  $b = 0$ ) as shown in Chang and Lin (2001). To prevent this, Perez-Cruz et al. (2003) proposed allowing the margin  $\rho$  to be negative and enforcing the norm of  $\mathbf{w}$  to be unity:

$$\min_{\mathbf{w}, b, \xi, \rho} -\nu \rho + \frac{1}{m} \sum_{i=1}^m \xi_i$$

s.t.  $y_i(\langle \mathbf{w}, \mathbf{x}_i \rangle + b) \geq \rho - \xi_i, \quad \xi_i \geq 0, \quad \|\mathbf{w}\|^2 = 1. \quad (3)$

By this modification, a non-trivial solution can be obtained even for  $\nu \in [0, \nu_{\min}]$ . This modified formulation is called extended  $\nu$ -SVC ( $E\nu$ -SVC).

The  $E\nu$ -SVC optimization problem is non-convex due to the equality constraint  $\|\mathbf{w}\|^2 = 1$ . Perez-Cruz et al. (2003) proposed the following iterative algorithm for computing a solution. First, for some initial  $\tilde{\mathbf{w}}$ , solve the problem (3) with  $\|\mathbf{w}\|^2 = 1$  replaced by  $\langle \tilde{\mathbf{w}}, \mathbf{w} \rangle = 1$ . Then, using the optimal solution  $\hat{\mathbf{w}}$ , update  $\tilde{\mathbf{w}}$  by

$$\tilde{\mathbf{w}} \leftarrow \gamma \tilde{\mathbf{w}} + (1 - \gamma) \hat{\mathbf{w}} \quad (4)$$

for  $\gamma = 9/10$ , and iterate this procedure until convergence.

Perez-Cruz et al. (2003) experimentally showed that the generalization performance of  $E\nu$ -SVC with  $\nu \in [0, \nu_{\min}]$  is often better than that with  $\nu \in (\nu_{\min}, \nu_{\max}]$ , implying that  $E\nu$ -SVC is a promising classification algorithm. However, it is not clear how the notion of negative margins influences on the generalization performance and how fast the above iterative algorithm converges. The goal of this paper is to give new theoretical insights into these issues.

## 3. Justification of the $E\nu$ -SVC Criterion

In this section, we give a new interpretation of  $E\nu$ -SVC and theoretically explain why it works well.

### 3.1. New Interpretation of $E\nu$ -SVC as CVaR minimization

Let  $f(\mathbf{w}, b; \mathbf{x}, y)$  be the *margin error* for a sample  $(\mathbf{x}, y)$ :

$$f(\mathbf{w}, b; \mathbf{x}, y) := -\frac{y(\langle \mathbf{w}, \mathbf{x} \rangle + b)}{\|\mathbf{w}\|}.$$

Let us consider the distribution of margin errors over all training samples:

$$\Phi(\alpha | \mathbf{w}, b) := P\{(\mathbf{x}_i, y_i) \mid f(\mathbf{w}, b; \mathbf{x}_i, y_i) \leq \alpha\}.$$

For  $\beta \in [0, 1)$ , let  $\alpha_\beta(\mathbf{w}, b)$  be the  $100\beta$ -percentile of the margin error distribution:

$$\alpha_\beta(\mathbf{w}, b) := \min\{\alpha \mid \Phi(\alpha | \mathbf{w}, b) \geq \beta\}.$$

Thus only the fraction  $(1 - \beta)$  of the margin error  $f(\mathbf{w}, b; \mathbf{x}_i, y_i)$  exceeds the threshold  $\alpha_\beta(\mathbf{w}, b)$  (see Figure 1).  $\alpha_\beta(\mathbf{w}, b)$  is commonly referred to as the *value-at-risk* (VaR) in finance and is often used by security houses or investment banks to measure the market risk of their asset portfolios (Rockafellar & Uryasev, 2002; Gotoh & Takeda, 2005).

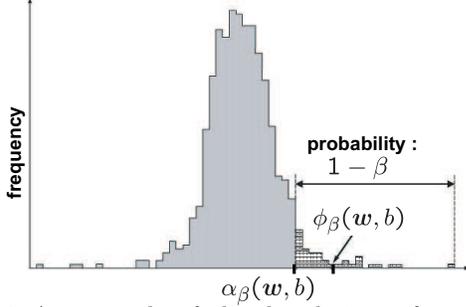


Figure 1. An example of the distribution of margin errors  $f(\mathbf{w}, b; \mathbf{x}_i, y_i)$  over all training samples.  $\alpha_\beta(\mathbf{w}, b)$  is the  $100\beta$ -percentile called the value-at-risk (VaR), and the mean  $\phi_\beta(\mathbf{w}, b)$  of the  $\beta$ -tail distribution is called the conditional VaR (CVaR).

Let us consider the  $\beta$ -tail distribution of  $f(\mathbf{w}, b; \mathbf{x}_i, y_i)$ :

$$\Phi_\beta(\alpha|\mathbf{w}, b) := \begin{cases} 0 & \text{for } \alpha < \alpha_\beta(\mathbf{w}, b), \\ \frac{\Phi(\alpha|\mathbf{w}, b) - \beta}{1 - \beta} & \text{for } \alpha \geq \alpha_\beta(\mathbf{w}, b). \end{cases}$$

Let  $\phi_\beta(\mathbf{w}, b)$  be the mean of the  $\beta$ -tail distribution of  $f(\mathbf{w}, b; \mathbf{x}_i, y_i)$  (see Figure 1 again):

$$\phi_\beta(\mathbf{w}, b) := \mathbf{E}_{\Phi_\beta}[f(\mathbf{w}, b; \mathbf{x}_i, y_i)],$$

where  $\mathbf{E}_{\Phi_\beta}$  denotes the expectation over the distribution  $\Phi_\beta$ .  $\phi_\beta(\mathbf{w}, b)$  is called the *conditional VaR* (CVaR). By definition, the CVaR is always larger than or equal to the VaR:

$$\phi_\beta(\mathbf{w}, b) \geq \alpha_\beta(\mathbf{w}, b). \quad (5)$$

Let us consider the problem of minimizing the CVaR  $\phi_\beta(\mathbf{w}, b)$  (which we refer to as minCVaR):

$$\min_{\mathbf{w}, b} \phi_\beta(\mathbf{w}, b). \quad (6)$$

Then we have the following theorem.

**Theorem 1** *The solution of the minCVaR problem (6) is equivalent to the solution of the  $E\nu$ -SVC problem (3) with*

$$\nu = 1 - \beta.$$

Theorem 1 shows that  $E\nu$ -SVC actually minimizes the CVaR  $\phi_{1-\nu}(\mathbf{w}, b)$ . Thus,  $E\nu$ -SVC could be interpreted as minimizing the mean margin error over a set of “bad” training samples. In contrast, the hard-margin SVC problem (2) can be equivalently expressed in terms of the margin error as

$$\min_{\mathbf{w}, b} \max_{i \in M} f(\mathbf{w}, b; \mathbf{x}_i, y_i).$$

Thus hard-margin SVC minimizes the margin error of the single “worst” training sample. This analysis shows that  $E\nu$ -SVC can be regarded as an extension of hard-margin SVC to be less sensitive to an outlier (i.e., the single “worst” training sample).

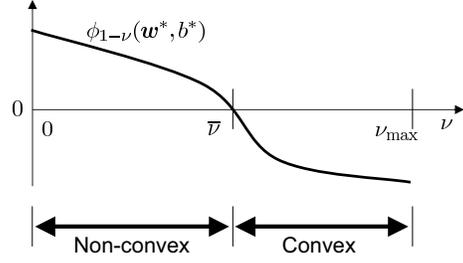


Figure 2. A profile of the CVaR  $\phi_{1-\nu}(\mathbf{w}^*, b^*)$  as a function of  $\nu$ . As shown in Section 4, the  $E\nu$ -SVC optimization problem can be cast as a convex problem if  $\nu \in (\bar{\nu}, \nu_{\max}]$ , while it is essentially non-convex if  $\nu \in (0, \bar{\nu})$ .

### 3.2. Justification of $E\nu$ -SVC

We have shown the equivalence between  $E\nu$ -SVC and minCVaR. Here we derive new bounds of the generalization error based on the notion of CVaR and try to justify the use of  $E\nu$ -SVC.

When training samples are linearly separable, the margin error  $f(\mathbf{w}, b; \mathbf{x}_i, y_i)$  is negative for all samples. Then, at the optimal solution  $(\mathbf{w}^*, b^*)$ , the CVaR  $\phi_{1-\nu}(\mathbf{w}^*, b^*)$  is always negative. However, in non-separable cases,  $\phi_{1-\nu}(\mathbf{w}^*, b^*)$  could be positive particularly when  $\nu$  is close to 0. Regarding the CVaR, we have the following lemma.

**Lemma 2**  *$\phi_{1-\nu}(\mathbf{w}^*, b^*)$  is continuous with respect to  $\nu$  and is strictly decreasing when  $\nu$  is increased.*

Let  $\bar{\nu}$  be such that

$$\phi_{1-\bar{\nu}}(\mathbf{w}^*, b^*) = 0$$

if such  $\bar{\nu}$  exists; we set  $\bar{\nu} = \nu_{\max}$  if  $\phi_{1-\nu}(\mathbf{w}^*, b^*) > 0$  for all  $\nu$  and we set  $\bar{\nu} = 0$  if  $\phi_{1-\nu}(\mathbf{w}^*, b^*) < 0$  for all  $\nu$ . Then we have the following relation (see Figure 2):

$$\begin{aligned} \phi_{1-\nu}(\mathbf{w}^*, b^*) &< 0 \text{ for } \nu \in (\bar{\nu}, \nu_{\max}], \\ \phi_{1-\nu}(\mathbf{w}^*, b^*) &> 0 \text{ for } \nu \in (0, \bar{\nu}). \end{aligned}$$

Below, we analyze the generalization error of  $E\nu$ -SVC depending on the value of  $\nu$ .

#### 3.2.1. JUSTIFICATION WHEN $\nu \in (\bar{\nu}, \nu_{\max}]$

**Theorem 3** *Let  $\nu \in (\bar{\nu}, \nu_{\max}]$ . Suppose that support  $\mathcal{X}$  is in a ball of radius  $R$  around the origin. Then, for all  $(\mathbf{w}, b)$  such that  $\|\mathbf{w}\| = 1$  and  $\phi_{1-\nu}(\mathbf{w}, b) < 0$ , there exists a positive constant  $c$  such that the following bound hold with probability at least  $1 - \delta$ :*

$$R[h] \leq \nu + G(\alpha_{1-\nu}(\mathbf{w}, b)), \quad (7)$$

where

$$G(\gamma) = \sqrt{\frac{2}{m} \left( \frac{4c^2(R^2 + 1)^2}{\gamma^2} \log_2(2m) - 1 + \log \frac{2}{\delta} \right)}.$$

The generalization error bound in (7) is furthermore upper-bounded as

$$\nu + G(\alpha_{1-\nu}(\mathbf{w}, b)) \leq \nu + G(\phi_{1-\nu}(\mathbf{w}, b)).$$

$G(\gamma)$  is monotone decreasing as  $|\gamma|$  increases. Thus, the above theorem shows that when  $\phi_{1-\nu}(\mathbf{w}, b) < 0$ , the upper bound  $\nu + G(\phi_{1-\nu}(\mathbf{w}, b))$  is lowered if the CVaR  $\phi_{1-\nu}(\mathbf{w}, b)$  is reduced. Since  $E\nu$ -SVC minimizes  $\phi_{1-\nu}(\mathbf{w}, b)$  (see Theorem 1), the upper bound of the generalization error is also minimized.

### 3.2.2. JUSTIFICATION WHEN $\nu \in (0, \bar{\nu}]$

Our discussion below depends on the sign of  $\alpha_{1-\nu}(\mathbf{w}, b)$ . When  $\alpha_{1-\nu}(\mathbf{w}, b) < 0$ , we have the following theorem.

**Theorem 4** *Let  $\nu \in (0, \bar{\nu}]$ . Then, for all  $(\mathbf{w}, b)$  such that  $\|\mathbf{w}\| = 1$  and  $\alpha_{1-\nu}(\mathbf{w}, b) < 0$ , there exists a positive constant  $c$  such that the following bound holds with probability at least  $1 - \delta$ :*

$$R[h] \leq \nu + G(\alpha_{1-\nu}(\mathbf{w}, b)).$$

A proof of the above theorem is omitted since the proof follows a similar line to the proof of Theorem 3. This theorem shows that when  $\alpha_{1-\nu}(\mathbf{w}, b) < 0$ , the upper bound  $\nu + G(\alpha_{1-\nu}(\mathbf{w}, b))$  is lowered if  $\alpha_{1-\nu}(\mathbf{w}, b)$  is reduced. On the other hand, Eq.(5) shows that the VaR  $\alpha_{1-\nu}(\mathbf{w}, b)$  is upper-bounded by the CVaR  $\phi_{1-\nu}(\mathbf{w}, b)$ . Therefore, minimizing  $\phi_{1-\nu}(\mathbf{w}, b)$  by  $E\nu$ -SVC may have an effect of lowering the upper bound of the generalization error.

When  $\alpha_{1-\nu}(\mathbf{w}, b) > 0$ , we have the following theorem.

**Theorem 5** *Let  $\nu \in (0, \bar{\nu}]$ . Then, for all  $(\mathbf{w}, b)$  such that  $\|\mathbf{w}\| = 1$  and  $\alpha_{1-\nu}(\mathbf{w}, b) > 0$ , there exists a positive constant  $c$  such that the following bound hold with probability at least  $1 - \delta$ :*

$$R[h] \geq \nu - G(\alpha_{1-\nu}(\mathbf{w}, b)).$$

Moreover, the lower bound of  $R[h]$  is bounded from above as

$$\nu - G(\alpha_{1-\nu}(\mathbf{w}, b)) \leq \nu - G(\phi_{1-\nu}(\mathbf{w}, b)).$$

A proof of the above theorem is also omitted since the proof resembles to Theorem 3. Theorem 5 implies that the lower bound  $\nu - G(\alpha_{1-\nu}(\mathbf{w}, b))$  of the generalization error is upper-bounded by  $\nu - G(\phi_{1-\nu}(\mathbf{w}, b))$ . On the other hand, Eq.(5) and  $\alpha_{1-\nu}(\mathbf{w}, b) > 0$  yields  $\phi_{1-\nu}(\mathbf{w}, b) > 0$ . Thus minimizing  $\phi_{1-\nu}(\mathbf{w}, b)$  by  $E\nu$ -SVC may contribute to lowering the lower bound  $\nu - G(\alpha_{1-\nu}(\mathbf{w}, b))$  of the generalization error.

## 4. New Optimization Algorithm

As reviewed in Section 2.5,  $E\nu$ -SVC involves a non-convex optimization problem. In this section, we give a new efficient optimization procedure for  $E\nu$ -SVC. Our proposed procedure involves two optimization algorithms depending on the value of  $\nu$ . We first describe the two algorithms and then show how these two algorithms are chosen for practical use.

### 4.1. Optimization When $\nu \in (\bar{\nu}, \nu_{\max}]$

**Lemma 6** *When  $\nu \in (\bar{\nu}, \nu_{\max}]$ , the  $E\nu$ -SVC problem (3) is equivalent to*

$$\begin{aligned} \min_{\mathbf{w}, b, \xi, \rho} \quad & -\nu\rho + \frac{1}{m} \sum_{i=1}^m \xi_i \\ \text{s.t.} \quad & y_i(\langle \mathbf{w}, \mathbf{x}_i \rangle + b) \geq \rho - \xi_i, \quad \xi_i \geq 0, \quad \|\mathbf{w}\|^2 \leq 1. \end{aligned} \quad (8)$$

This lemma shows that the equality constraint  $\|\mathbf{w}\|^2 = 1$  in the original problem (3) can be replaced by  $\|\mathbf{w}\|^2 \leq 1$  without changing the solution. Due to the convexity of  $\|\mathbf{w}\|^2 \leq 1$ , the above optimization problem is convex and therefore we can easily obtain the global solution by a standard optimization software.

### 4.2. Optimization When $\nu \in (0, \bar{\nu}]$

If  $\nu \in (0, \bar{\nu}]$ , the  $E\nu$ -SVC optimization problem is essentially non-convex and therefore we need a more elaborate algorithm.

#### 4.2.1. LOCAL OPTIMUM SEARCH

Here, we propose the following iterative algorithm for finding a local optimum.

**Algorithm 7 (The  $E\nu$ -SVC local optimum search algorithm for  $\nu \in (0, \bar{\nu}]$ )**

**Step 1:** Initialize  $\tilde{\mathbf{w}}$ .

**Step 2:** Solve the following linear program:

$$\begin{aligned} \min_{\mathbf{w}, b, \xi, \rho} \quad & -\nu\rho + \frac{1}{m} \sum_{i=1}^m \xi_i \\ \text{s.t.} \quad & y_i(\langle \mathbf{w}, \mathbf{x}_i \rangle + b) \geq \rho - \xi_i, \quad \xi_i \geq 0, \quad \langle \tilde{\mathbf{w}}, \mathbf{w} \rangle = 1, \end{aligned} \quad (9)$$

and let the optimal solution be  $(\hat{\mathbf{w}}, \hat{b}, \hat{\xi}, \hat{\rho})$ .

**Step 3:** If  $\tilde{\mathbf{w}} = \hat{\mathbf{w}}$ , terminate and output  $\tilde{\mathbf{w}}$ . Otherwise, update  $\tilde{\mathbf{w}}$  by  $\tilde{\mathbf{w}} \leftarrow \hat{\mathbf{w}} / \|\hat{\mathbf{w}}\|$ .

**Step 4:** Repeat Steps 2–3.

The linear program (9) is the same as the one proposed by Perez-Cruz et al. (2003), i.e., the equality constrained  $\|\mathbf{w}\|^2 = 1$  of the original problem (3) is

replaced by  $\langle \tilde{\mathbf{w}}, \mathbf{w} \rangle = 1$ . The updating rule of  $\tilde{\mathbf{w}}$  in Step 3 is different from the one proposed by Perez-Cruz et al. (2003) (cf. Eq.(4)).

We define a ‘‘corner’’ (or ‘‘0-dimensional face’’) of  $E\nu$ -SVC (3) as the intersection of an edge of the polyhedral cone formed by linear constraints of (3) and  $\|\mathbf{w}\|^2 = 1$ . Under the new update rule, the algorithm visits a corner of  $E\nu$ -SVC (3) in each iteration. Since  $E\nu$ -SVC has finite corners, we can show that Algorithm 7 with the new update rule terminates in a finite number of iterations, i.e., less than or equal to the number of corners of  $E\nu$ -SVC.

**Theorem 8** *Algorithm 7 terminates within a finite number of iterations of Steps 2–3. Furthermore, a solution of the modified  $E\nu$ -SVC algorithm is a local minimizer if it is unique and non-degenerate.*

#### 4.2.2. GLOBAL OPTIMUM SEARCH

Next, we show that the global solution can be actually obtained within finite iterations, despite the non-convexity of the optimization problem.

A naive approach to searching for the global solution is to run the local optimum search algorithm many times with different initial values and choose the best local solution. However, there is no guarantee that this naive approach can find the global solution. Below, we give a more systematic way to find the global solution based on the following lemma.

**Lemma 9** *When  $\nu \in (0, \bar{\nu}]$ , the  $E\nu$ -SVC problem (3) is equivalent to*

$$\min_{\mathbf{w}, b, \xi, \rho} -\nu\rho + \frac{1}{m} \sum_{i=1}^m \xi_i$$

s.t.  $y_i(\langle \mathbf{w}, \mathbf{x}_i \rangle + b) \geq \rho - \xi_i, \xi_i \geq 0, \|\mathbf{w}\|^2 \geq 1. \quad (10)$

Lemma 9 could be proved in a similar way as Lemma 6, so we omit the proof. This lemma shows that the equality constraint  $\|\mathbf{w}\|^2 = 1$  in the original  $E\nu$ -SVC problem (3) can be replaced by  $\|\mathbf{w}\|^2 \geq 1$  without changing the solution if  $\nu \in (0, \bar{\nu}]$ .

The problem (10) is called a *linear reverse convex program* (LRCP), which is a class of non-convex problems consisting of linear constraints and one concave inequality ( $\|\mathbf{w}\|^2 \geq 1$  in the current case). The feasible set of the problem (10) consists of a finite number of *faces*. For LRCPs, Horst and Tuy (1995) showed that the local optimal solutions correspond to 0-dimensional faces (or corners). This implies that all the local optimal solutions of the  $E\nu$ -SVC problem (10) can be traced by checking all the faces.

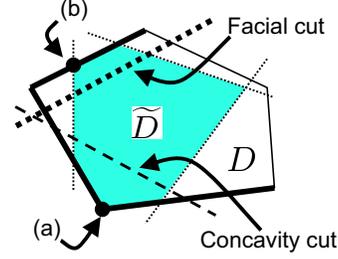


Figure 3. A 0-dimensional face (a) and three proper faces (bold solid lines) of  $D$  are identified in  $\tilde{D}$ . If the corner (a) is found in Step 2, a concavity cut is constructed. If the corner (b) is found, a facial cut is constructed. If these two cuts are added to  $\tilde{D}$ , the remaining area includes no face of  $D$ .

Let  $D$  be the feasible set of  $E\nu$ -SVC (3). Below, we summarize the  $E\nu$ -SVC training algorithm based on the *cutting plane method*, which is an efficient method of tracing faces.

#### Algorithm 10 (The $E\nu$ -SVC global optimum search algorithm for $\nu \in (0, \bar{\nu}]$ )

- Step 1:**  $\tilde{D} \leftarrow D$ .
- Step 2:** Find a local solution by Algorithm 7.
- Step 3:** Identify a face of  $D$  in  $\tilde{D}$  that corresponds the local solution.
- Step 4a:** If the face is a corner, construct a ‘‘concavity cut’’.
- Step 4b:** If the face is a proper face, construct a ‘‘facial cut’’.
- Step 5:** Add the cut to the problem (9) and  $\tilde{D}$ .
- Step 6:** Repeat Steps 2–5 until  $\tilde{D}$  includes no face of  $D$ .
- Step 7:** Output the best local optimal solution as the global solution.

If the local solution obtained in Step 2 is a corner of  $D$  (i.e., the local solution is not on any cutting plane as (a) in Figure 3), a *concavity cut* (Horst & Tuy, 1995) is constructed. The concavity cut has a role of removing the local solution, i.e., a 0-dimensional face of  $D$  and its neighborhood. Otherwise, a *facial cut* (Majthay & Whinston, 1974) is constructed to eliminate the proper face (see (b) in Figure 3).

Since the total number of distinct faces of  $D$  is finite in the current setting and a facial cut or a concavity cut eliminates at least one face at a time, Algorithm 10 is guaranteed to terminate within finite iterations (precisely, less than or equal to the number of all dimensional faces of  $E\nu$ -SVC). Furthermore, since the addition of a concavity cut or a facial cut does not remove local solutions which are better than the best local solution found so far, Algorithm 10 is guaranteed to

trace all *sufficient* local solutions. Thus we can always find a global solution within finite iterations by Algorithm 10. A more detailed discussion on the concavity cut and the facial cut is shown in Horst and Tuy (1995) and Majthay and Whinston (1974), respectively.

### 4.3. Choice of Two Algorithms

We have two convergent algorithms when  $\nu \in (\bar{\nu}, \nu_{\max}]$  and  $\nu \in (0, \bar{\nu}]$ . Thus, choosing a suitable algorithm depending on the value of  $\nu$  would be an ideal procedure. However, the value of the threshold  $\bar{\nu}$  is difficult to explicitly compute since it is defined via the optimal value  $\phi_{1-\bar{\nu}}(\mathbf{w}^*, b^*)$  (see Figure 2). Therefore, it is not straightforward to choose a suitable algorithm for a given  $\nu$ .

When we use  $E\nu$ -SVC in practice, we usually compute the solutions for several different values of  $\nu$  and choose the most promising one based on, e.g., cross-validation. In such scenarios, we can properly switch two algorithms without explicitly knowing the value of  $\bar{\nu}$ —our key idea is that the solution of the problem (8) is non-trivial (i.e.,  $\mathbf{w} \neq \mathbf{0}$ ) if and only if  $\nu \in (\bar{\nu}, \nu_{\max}]$ . Thus if the solutions are computed from large  $\nu$  to small  $\nu$ , the switching point can be identified by checking the triviality of the solution. The proposed algorithm is summarized as follows.

**Algorithm 11 (The  $E\nu$ -SVC algorithm for  $(\nu_{\max} \geq) \nu_1 > \nu_2 > \dots > \nu_k > 0$ )**

**Step 1:**  $i \leftarrow 1$ .

**Step 2:** Compute  $(\mathbf{w}^*, b^*)$  for  $\nu_i$  by solving (8).

**Step 3a:** If  $\mathbf{w}^* \neq \mathbf{0}$ , accept  $(\mathbf{w}^*, b^*)$  as the solution for  $\nu_i$ , increment  $i$ , and go to Step 2.

**Step 3b:** If  $\mathbf{w}^* = \mathbf{0}$ , reject  $(\mathbf{w}^*, b^*)$ .

**Step 4:** Compute  $(\mathbf{w}^*, b^*)$  for  $\nu_i$  by Algorithm 10.

**Step 5:** Accept  $(\mathbf{w}^*, b^*)$  as the solution for  $\nu_i$ , increment  $i$ , and go to Step 4 unless  $i > k$ .

## 5. Conclusions

We characterized the generalization error of  $E\nu$ -SVC in terms of the *conditional value-at-risk* (CVaR, see Figure 1) and showed that a good generalization performance is expected by  $E\nu$ -SVC. We then derived a globally convergent optimization algorithm even though the optimization problem involved in  $E\nu$ -SVC is non-convex.

We introduced the threshold  $\bar{\nu}$  based on the sign of the CVaR (see Figure 2). We can check that the problem (8) is equivalent to  $\nu$ -SVC in the sense that they share the same negative optimal value in  $(\bar{\nu}, \nu_{\max}]$  and  $(\nu_{\min}, \nu_{\max}]$ , respectively (Gotoh & Takeda, 2005). On

the other hand, the problem (8) and  $\nu$ -SVC have the zero optimal value in  $(0, \bar{\nu}]$  and  $[0, \nu_{\min}]$ , respectively. Thus, although the definitions of  $\bar{\nu}$  and  $\nu_{\min}$  are different, they would be essentially the same. We will study the relation between  $\bar{\nu}$  and  $\nu_{\min}$  in more detail in the future work.

## References

- Boser, B. E., Guyon, I. M., & Vapnik, V. N. (1992). A training algorithm for optimal margin classifiers. *COLT* (pp. 144–152). ACM Press.
- Chang, C.-C., & Lin, C.-J. (2001). Training  $\nu$ -support vector classifiers: Theory and algorithms. *Neural Computation*, *13*, 2119–2147.
- Cortes, C., & Vapnik, V. (1995). Support-vector networks. *Machine Learning*, *20*, 273–297.
- Crisp, D. J., & Burges, C. J. C. (2000). A geometric interpretation of  $\nu$ -SVM classifiers. *NIPS 12* (pp. 244–250). MIT Press.
- Gotoh, J., & Takeda, A. (2005). A linear classification model based on conditional geometric score. *Pacific Journal of Optimization*, *1*, 277–296.
- Horst, R., & Tuy, H. (1995). *Global optimization: Deterministic approaches*. Berlin: Springer-Verlag.
- Majthay, A., & Whinston, A. (1974). Quasi-concave minimization subject to linear constraints. *Discrete Mathematics*, *9*, 35–59.
- Perez-Cruz, F., Weston, J., Hermann, D. J. L., & Schölkopf, B. (2003). Extension of the  $\nu$ -SVM range for classification. *Advances in Learning Theory: Methods, Models and Applications 190* (pp. 179–196). Amsterdam: IOS Press.
- Rockafellar, R. T., & Uryasev, S. (2002). Conditional value-at-risk for general loss distributions. *Journal of Banking & Finance*, *26*, 1443–1472.
- Schölkopf, B., Smola, A., Williamson, R., & Bartlett, P. (2000). New support vector algorithms. *Neural Computation*, *12*, 1207–1245.
- Schölkopf, B., & Smola, A. J. (2002). *Learning with kernels*. Cambridge, MA: MIT Press.
- Vapnik, V. N. (1995). *The nature of statistical learning theory*. Berlin: Springer-Verlag.

## A. Sketch of Proof of Theorem 1

Let  $(\mathbf{w}^*, b^*, \alpha^*)$  be the optimal solution of

$$\min_{\mathbf{w}, b, \alpha} F_{\beta}(\mathbf{w}, b, \alpha), \quad (11)$$

where, for  $[X]^+ := \max\{X, 0\}$ ,

$$F_{\beta}(\mathbf{w}, b, \alpha) := \alpha + \frac{\sum_{i \in M} [f(\mathbf{w}, b; \mathbf{x}_i, y_i) - \alpha]^+}{(1 - \beta)m}. \quad (12)$$

Then Rockafellar and Uryasev (2002) showed that

$$F_\beta(\mathbf{w}^*, b^*, \alpha^*) = \phi_\beta(\mathbf{w}^*, b^*) = \min_{\mathbf{w}, b} \phi_\beta(\mathbf{w}, b), \quad (13)$$

i.e., the problems (6) and (11) are equivalent.

Introducing slack variables  $\xi_i$ , imposing  $\|\mathbf{w}\|^2 = 1$  (which does not change the solution essentially; only the scale is changed), and letting  $\nu = 1 - \beta$  and  $\rho = -\alpha$  in Eq.(11), we establish the theorem.

## B. Sketch of Proof of Lemma 2

Since Eq.(11) only involves continuous functions, continuity of  $F_\beta(\mathbf{w}^*, b^*, \alpha^*)$  with respect to  $\beta$  is clear. From Eq.(13),  $\phi_\beta(\mathbf{w}^*, b^*)$  is also continuous. Let  $(\mathbf{w}_{\beta_1}^*, b_{\beta_1}^*, \alpha_{\beta_1}^*)$  be the optimal solutions of Eq.(11) for  $0 < \beta_1 < \beta_2 < 1$ . Then we have

$$\begin{aligned} \phi_{\beta_1}(\mathbf{w}_{\beta_1}^*, b_{\beta_1}^*) &= F_{\beta_1}(\mathbf{w}_{\beta_1}^*, b_{\beta_1}^*, \alpha_{\beta_1}^*) \leq F_{\beta_1}(\mathbf{w}_{\beta_2}^*, b_{\beta_2}^*, \alpha_{\beta_2}^*) \\ &< F_{\beta_2}(\mathbf{w}_{\beta_2}^*, b_{\beta_2}^*, \alpha_{\beta_2}^*) = \phi_{\beta_2}(\mathbf{w}_{\beta_2}^*, b_{\beta_2}^*), \end{aligned}$$

where the first inequality is due to optimality of  $(\mathbf{w}_{\beta_1}^*, b_{\beta_1}^*, \alpha_{\beta_1}^*)$  and the second strict inequality is clear from Eq.(12). Thus  $\phi_\beta(\mathbf{w}^*, b^*)$  is strictly increasing with respect to  $\beta$ , implying that  $\phi_{1-\nu}(\mathbf{w}^*, b^*)$  is strictly decreasing with respect to  $\nu$ .

## C. Sketch of Proof of Theorem 3

For a homogeneous classifier  $h(\tilde{\mathbf{x}}) = \text{sign}(\langle \tilde{\mathbf{w}}, \tilde{\mathbf{x}} \rangle)$ , the following lemma holds:

**Lemma 12** (Schölkopf et al., 2000) *Suppose that support  $\mathcal{X}$  of  $\tilde{\mathbf{x}}$  is in a ball of radius  $\tilde{R}$  around the origin. Then, for all  $\tilde{\mathbf{w}}$  such that  $\|\tilde{\mathbf{w}}\| = 1$ , there exists a positive constant  $c$  such that the following bound holds with probability at least  $1 - \delta$ :*

$$\begin{aligned} R[h] &\leq \frac{|\{i \mid y_i \langle \tilde{\mathbf{w}}, \tilde{\mathbf{x}}_i \rangle < \tilde{\gamma}\}|}{m} \\ &+ \sqrt{\frac{2}{m} \left( \frac{4c^2 \tilde{R}^2}{\tilde{\gamma}^2} \log_2(2m) - 1 + \log \frac{2}{\delta} \right)}. \end{aligned}$$

Let  $\tilde{\mathbf{w}} = \frac{(\mathbf{w}^\top, b)^\top}{\sqrt{1+b^2}}$  and  $\tilde{\mathbf{x}} = (\mathbf{x}^\top, 1)^\top$ . Then our classifier (1) can be regarded as homogeneous. The assumption that all the data points  $\mathbf{x}$  live in a centered ball of radius  $R$  implies that all the data points  $\tilde{\mathbf{x}}$  live in a centered ball of radius

$$\tilde{R} = \sqrt{R^2 + 1}.$$

The assumption  $\|\mathbf{w}\| = 1$  implies  $\|\tilde{\mathbf{w}}\| = 1$ . Then we can apply Lemma 12 to the current setting. The condition  $y_i \langle \tilde{\mathbf{w}}, \tilde{\mathbf{x}}_i \rangle < \tilde{\gamma}$  results in

$$y_i \langle \mathbf{w}, \mathbf{x}_i \rangle + b < \tilde{\gamma} \sqrt{1 + b^2} := \gamma.$$

When all the data points  $\mathbf{x}$  live in a centered ball of radius  $R$ , we can assume without loss of generality that  $|b| \leq R$ . Then we have

$$\frac{1}{\tilde{\gamma}^2} = \frac{1 + b^2}{\gamma^2} \leq \frac{1 + R^2}{\gamma^2}.$$

Now let us set

$$\gamma = -\alpha_{1-\nu}(\mathbf{w}, b).$$

Then we can show that

$$\frac{1}{m} |\{i \mid y_i \langle \mathbf{w}, \mathbf{x}_i \rangle + b < -\alpha_{1-\nu}(\mathbf{w}, b)\}| \leq \nu.$$

We omit its proof due to lack of space. Then we obtain the upper bound  $\nu + G(\alpha_{1-\nu}(\mathbf{w}, b))$ ; the upper bound  $\nu + G(\phi_{1-\nu}(\mathbf{w}, b))$  is clear from Eq.(5).

## D. Sketch of Proof of Lemma 6

Since the difference between the problems (3) and (8) is only the norm constraint of  $\mathbf{w}$ , it is enough to show that for  $\nu \in (\bar{\nu}, \nu_{\max}]$ ,  $\|\mathbf{w}^*\|^2 = 1$  holds at the optimal solution  $(\mathbf{w}^*, b^*, \boldsymbol{\xi}^*, \rho^*)$  of the problem (8). For such  $\nu$ ,  $\phi_{1-\nu}(\mathbf{w}^*, b^*) < 0$  holds, i.e., the optimal value of  $E\nu$ -SVC is negative. If we suppose  $\|\mathbf{w}^*\|^2 < 1$ , another feasible solution  $(\mathbf{w}^*, b^*, \boldsymbol{\xi}^*, \rho^*)/\|\mathbf{w}^*\|$  achieves a smaller optimal value than  $(\mathbf{w}^*, b^*, \boldsymbol{\xi}^*, \rho^*)$ . This contradicts to the optimality of (8), and hence  $\|\mathbf{w}^*\|^2 = 1$  is proved.

## E. Sketch of Proof of Theorem 8

Let  $(\hat{\mathbf{w}}_k, \hat{b}_k, \hat{\boldsymbol{\xi}}_k, \hat{\rho}_k)$  be an optimal solution of the linear program (9) in the  $k$ -th iteration. Then, a feasible solution of  $E\nu$ -SVC (3) is given by

$$(\tilde{\mathbf{w}}_k, \tilde{b}_k, \tilde{\boldsymbol{\xi}}_k, \tilde{\rho}_k) = (\hat{\mathbf{w}}_k, \hat{b}_k, \hat{\boldsymbol{\xi}}_k, \hat{\rho}_k)/\|\hat{\mathbf{w}}_k\|.$$

Since  $(\hat{\mathbf{w}}_k, \hat{b}_k, \hat{\boldsymbol{\xi}}_k, \hat{\rho}_k)$  is at a corner of the feasible set of the linear program (9),  $(\tilde{\mathbf{w}}_k, \tilde{b}_k, \tilde{\boldsymbol{\xi}}_k, \tilde{\rho}_k)$  is also a corner of the feasible set of  $E\nu$ -SVC (3).

Let  $q(\cdot)$  be the objective function of  $E\nu$ -SVC (3), which is also the objective function of the linear program (9). Then we have

$$q(\tilde{\boldsymbol{\xi}}_{k-1}, \tilde{\rho}_{k-1}) > q(\hat{\boldsymbol{\xi}}_k, \hat{\rho}_k) \geq q(\tilde{\boldsymbol{\xi}}_k, \tilde{\rho}_k) = q(\hat{\boldsymbol{\xi}}_k, \hat{\rho}_k)/\|\hat{\mathbf{w}}_k\|,$$

where the first inequality comes from the optimality of  $(\hat{\boldsymbol{\xi}}_k, \hat{\rho}_k)$  of the linear program (9). The second inequality comes from  $\|\hat{\mathbf{w}}_k\| > 1$ , which is ensured by  $\langle \tilde{\mathbf{w}}_{k-1}, \hat{\mathbf{w}}_k \rangle = 1$ . Thus the algorithm finds a distinct corner of  $E\nu$ -SVC (3) in each iteration. Since the number of corners of  $E\nu$ -SVC (3) is finite, the algorithm terminates within finite iterations.

Let  $\Delta \mathbf{d} = (\Delta \mathbf{w}^\top \ \Delta b^\top \ \Delta \rho^\top \ \Delta \boldsymbol{\xi}^\top)^\top$  be a perturbation from the solution  $\mathbf{d}^* = (\mathbf{w}^*, b^*, \rho^*, \boldsymbol{\xi}^*)$  of Algorithm 7. Note that  $\mathbf{d}^*$  is an optimal solution of the linear program (9) with  $\tilde{\mathbf{w}} = \mathbf{w}^*$ . Using the Karush-Kuhn-Tucker (KKT) optimality conditions, we can express the increase  $\Delta q$  of the objective value as

$$\begin{aligned} \Delta q &:= -\nu \Delta \rho + \frac{1}{m} \sum_{i \in M} \Delta \xi_i \\ &= \Delta \mathbf{d}^\top \begin{pmatrix} y_1 \mathbf{x}_1 & \dots & y_m \mathbf{x}_m & \mathbf{0} \\ y_1 & \dots & y_m & \mathbf{0} \\ \mathbf{1} & \dots & \mathbf{1} & \mathbf{I} \end{pmatrix} \begin{pmatrix} \mathbf{0} \\ \mathbf{0} \\ \mathbf{0} \\ \mathbf{I} \end{pmatrix} \begin{pmatrix} \boldsymbol{\lambda}^* \\ \boldsymbol{\mu}^* \end{pmatrix} - \delta^* \Delta \mathbf{w}^\top \mathbf{w}^*, \end{aligned}$$

where  $\boldsymbol{\lambda}^* \in \mathbb{R}_+^m$ ,  $\boldsymbol{\mu}^* \in \mathbb{R}_+^m$ , and  $\delta^* \leq 0$  are KKT multipliers. If  $\Delta \mathbf{d}$  is a feasible perturbation (i.e.,  $\mathbf{d}^* + \Delta \mathbf{d}$  is feasible), we can show that  $\Delta q > 0$  (we omit its proof due to lack of space), which implies that  $\mathbf{d}^*$  is locally optimal.