

INSTITUTE OF POLICY AND PLANNING SCIENCES

Discussion Paper Series

No. 1096

A Linear Classification Model Based on
Conditional Geometric Score

by

Jun-ya Gotoh and Akiko Takeda

September 2004

UNIVERSITY OF TSUKUBA
Tsukuba, Ibaraki 305-8573
JAPAN

A Linear Classification Model Based on Conditional Geometric Score

Jun-ya Gotoh* Akiko Takeda†

August 30, 2004

Revised: September 9, 2004

Abstract

We propose a two-class linear classification model by taking into account the Euclidean distance from each data point to the discriminant hyperplane and introducing a risk measure which is known as the conditional value-at-risk in financial risk management. It is formulated as a nonconvex programming problem and we present a solution method for obtaining either a globally or a locally optimal solution by examining the special structure of the problem. Also, this model is proved to be equivalent to the ν -support vector classification under some parameter setting, and numerical experiments show that the proposed model has better predictive accuracy in general.

The 2000 MSC: 68T10, 90C26, 90C90

Keywords: classification model, discriminant hyperplane, conditional value-at-risk, nonconvex programming

1 Introduction

Mathematical optimization approaches to the classification problem have attracted much attention of many researchers since the 1960s (e.g., [7]). Among such approaches are the support vector classification [15], the robust linear programming method [2], integer programming methods (e.g., [6]), the other kind of linear programming methods (e.g., [4, 5]). For the two-class classification problem in \mathbb{R}^n , these models seek to find a decision function $h : \mathbb{R}^n \rightarrow \{\pm 1\}$ using a set of training data $\{\mathbf{x}^i\} \subset \mathbb{R}^n$ which are labeled with binary values $\{y^i\} \subset \{\pm 1\}$ for $i \in I := \{1, \dots, \ell\}$:

$$(\mathbf{x}^1, y^1), \dots, (\mathbf{x}^\ell, y^\ell) \in \mathbb{R}^n \times \{\pm 1\},$$

so that h will predict as accurately as possible the labels of new data points generated from the same probability distribution with the training data. In particular, the linear classification problem concerned in this paper is reduced to constructing a discriminant hyperplane:

$$H(\mathbf{w}, b) := \{\mathbf{x} \in \mathbb{R}^n \mid \langle \mathbf{w}, \mathbf{x} \rangle + b = 0\}, \quad \mathbf{w} \neq \mathbf{0}, \quad \mathbf{w} \in \mathbb{R}^n, \quad b \in \mathbb{R}, \quad (1)$$

*Graduate School of Systems and Information Engineering, University of Tsukuba, 1-1-1 Tennoudai, Tsukuba, Ibaraki, 305-8573 Japan. e-mail: jgoto@sk.tsukuba.ac.jp

†Department of Mathematical and Computing Sciences, Tokyo Institute of Technology, 2-12-1 Oh-Okayama, Meguro-ku, Tokyo, 152-8552 Japan. e-mail: takeda@is.titech.ac.jp

which corresponds to the decision function $h(\mathbf{x}) = \text{sign}(\langle \mathbf{w}, \mathbf{x} \rangle + b)$ where $\text{sign}(\xi)$ is 1 if $\xi \geq 0$, -1 , otherwise.

To this end, if given set of training data is linearly separable, i.e., there exists (\mathbf{w}, b) such that $\mathbf{w} \neq \mathbf{0}$ and $y^i(\langle \mathbf{w}, \mathbf{x}^i \rangle + b) > 0$ hold for $i \in I$, the hard margin support vector classification (HSVC) [15] provides a most reasonable hyperplane as an optimal solution of

$$\left| \begin{array}{l} \text{maximize} \\ \mathbf{w} \neq \mathbf{0}, b \end{array} \min_{i \in I} g(\mathbf{w}, b; \mathbf{x}^i, y^i), \right. \quad (2)$$

where

$$g(\mathbf{w}, b; \mathbf{x}^i, y^i) := \frac{y^i(\langle \mathbf{w}, \mathbf{x}^i \rangle + b)}{\|\mathbf{w}\|}, \quad \text{for } i \in I. \quad (3)$$

In the paper, we call $g(\mathbf{w}, b; \mathbf{x}^i, y^i)$ the *geometric score* of data (\mathbf{x}^i, y^i) with respect to (\mathbf{w}, b) . Formulation (2) implies that only a portion of training data which have the lowest geometric score contribute to determine the hyperplane.

On the other hand, for linearly inseparable data set, there is no hyperplane which correctly separates the given set of data points, and the criterion based on the lowest geometric score is less persuasive because the other misclassified data points are ignored. For overcoming this drawback, the other models such as the soft margin SVC and the robust linear programming method introduce some other misclassification risk measures to be minimized into their formulation.

In this paper, we propose a classification model which minimizes another misclassification risk. The risk measure applied in our model is based on a conditional expectation of the empirical distribution of the geometric score, which is known as the *conditional value-at-risk* (CVaR) in the financial risk management [12]. The CVaR measure is now highlighted as a tool of perceiving and controlling the financial downside risk since it has preferable properties such as coherency [1] and consistency with stochastic dominance [10, 9] and, in addition, many portfolio problems based on the CVaR minimization can be formulated as a convex programming problem [11].

Contrary to those portfolio problems, our model is basically formulated as a nonconvex quadratic optimization problem due to the nonconvexity of the geometric score g . Under some parameter setting, however, it can be transformed into an equivalent convex quadratic optimization problem. The properties of our model are

- i) under some parameter setting, our model is equivalent to the ν -support vector classification (ν -SVC) model [13], and also equivalent to the HSVM model (2),
- ii) the difficulty of the proposed problem depends on given parameters and data set $\{(\mathbf{x}^i, y^i) : i \in I\}$.

From i), the discussion in this paper provides another interpretation of the ν -SVC, and this similarity opens a gate for the use of the kernel functions in our linear model.

The structure of this paper is as follows. In the next section, we first define the conditional expectation of the geometric score as a risk measure by following Rockafellar & Uryasev [12], and

describe the formulation of the risk minimization. In section 3, we discuss the nonconvex structure of the formulation, and propose an algorithm for obtaining either an optimal solution when the problem is convex, or a locally optimal solution when the problem is essentially nonconvex. In section 4, we discuss a relationship between our model and other classification approaches: the HSVC and the ν -SVC, and show that under some parameter setting, our model is equivalent to them. Section 5 is devoted to some numerical experiments, showing that the nonconvexity in our model, which is eliminated in the ν -SVC, plays an important role in determining the hyperplane. Finally, we conclude the paper by adding some remarks and possible extension.

2 Conditional Geometric Score Optimization

We first introduce the several results of Rockafellar & Uryasev [12] and define the conditional expectation of the geometric score as misclassification risk measure in section 2.1. Then, in section 2.2, we describe the linear classification problem which minimizes the misclassification risk measure.

2.1 A Misclassification Risk Measure

We consider the two-class linear classification problem for a finite number of training data which are assumed to be generated from an unknown probability distribution.

We denote the given labeled data set by $\{(\mathbf{x}^i, y^i) : i \in I\} \subset \mathbb{R}^n \times \{\pm 1\}$ where $I := \{1, \dots, \ell\}$ is the index set of the training data, and define two index sets I^+ and I^- by $I^+ := \{i \in I \mid y^i = +1\}$ and $I^- := \{i \in I \mid y^i = -1\}$, respectively. In addition, we suppose that the empirical probability assigned on the given data is given by $P\{(\mathbf{x}, y) = (\mathbf{x}^i, y^i)\} = p_i > 0$ for $i \in I$, where $\sum_{i \in I} p_i = 1$. One simple choice for p_i is

$$p_i = \frac{1}{|I|} = \frac{1}{\ell} \text{ for } i \in I. \quad (4)$$

Another possibility is

$$p_i = \begin{cases} \frac{1}{2|I^+|} & \text{for } i \in I^+ \\ \frac{1}{2|I^-|} & \text{for } i \in I^- \end{cases} \quad (5)$$

which is often used to balance the two classes in the other classification models [2, 13].

For each data point (\mathbf{x}^i, y^i) , the geometric score is defined as in (3) with some given (\mathbf{w}, b) and we sometimes use the notation $g(\mathbf{w}, b)$ in place of $g(\mathbf{w}, b; \mathbf{x}, y)$ for simplicity. $g(\mathbf{w}, b)$ is the Euclidean distance of a point $\mathbf{x} \in \mathbb{R}^n$ to the hyperplane $H(\mathbf{w}, b)$ if positive, and we can say that a data point with lower geometric score has higher risk to be misclassified. Denoting the minus value of the score g by f , i.e.,

$$f(\mathbf{w}, b) := -g(\mathbf{w}, b),$$

we interpret f as the magnitude of misclassification risk.

Let us define the distribution function $\Phi((\mathbf{w}, b), \alpha)$ of f by

$$\Phi((\mathbf{w}, b), \alpha) := P\{f(\mathbf{w}, b) \leq \alpha\},$$

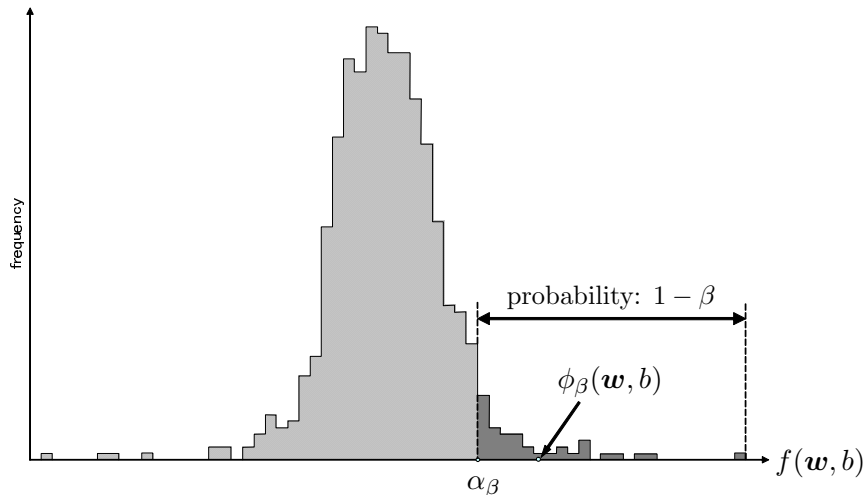


Figure 1: Illustration of the β -tail expectation of f

and a threshold α_β with some confidence level $\beta \in (0, 1)$ by

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

We note that α_β is well-defined because $\Phi((\mathbf{w}, b), \alpha)$ is right continuous and nondecreasing with respect to α . The α_β is known as the *value-at-risk* (VaR) in the context of financial risk management. It is expected that $f(\mathbf{w}, b)$ exceeds α_β only in $(1 - \beta) \times 100\%$.

Following the results of Rockafellar & Uryasev [12], we introduce the β -tail distribution function to focus on the tail part of $\Phi((\mathbf{w}, b), \alpha)$ as

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

Using the expectation operator $E_\beta[\cdot]$ under the β -tail distribution Φ_β , define the β -tail expectation of f as $\phi_\beta(\mathbf{w}, b) := E_\beta[f(\mathbf{w}, b)]$, which is the risk measure known as the *conditional value-at-risk* (CVaR). Denoting the expectation under the original distribution Φ by $E[\cdot]$, the following relation shown in [12]:

$$E[f(\mathbf{w}, b) \mid f(\mathbf{w}, b) \geq \alpha_\beta] \leq \phi_\beta(\mathbf{w}, b) \leq E[f(\mathbf{w}, b) \mid f(\mathbf{w}, b) > \alpha_\beta] \quad (6)$$

implies that $\phi_\beta(\mathbf{w}, b)$ is approximately equal to the conditional expectation of f which exceeds the threshold α_β . In other words, the risk measure $\phi_\beta(\mathbf{w}, b)$ is based on the distribution of the lower geometric score, and a hyperplane $H(\mathbf{w}, b)$ induced from the risk-measure minimization problem:

$$\underset{\mathbf{w} \neq \mathbf{0}, b}{\text{minimize}} \phi_\beta(\mathbf{w}, b) \quad (7)$$

is expected to classify new data points generated from the same distribution correctly.

To solve the problem (7), Rockafellar & Uryasev [12] introduces a simpler auxiliary function $F_\beta : \mathbb{R}^{n+2} \rightarrow \mathbb{R}$:

$$F_\beta(\mathbf{w}, b, \alpha) := \alpha + \frac{1}{1 - \beta} E[[f(\mathbf{w}, b) - \alpha]^+],$$

where $[X]^+ := \max\{X, 0\}$, and provides a shortcut to minimizing $\phi_\beta(\mathbf{w}, b)$ as

$$\underset{(\mathbf{w}, b) \in W}{\text{minimize}} \phi_\beta(\mathbf{w}, b) = \underset{((\mathbf{w}, b), \alpha) \in W \times \mathbb{R}}{\text{minimize}} F_\beta(\mathbf{w}, b, \alpha), \quad (8)$$

where $W := (\mathbb{R}^n \setminus \{\mathbf{0}\}) \times \mathbb{R}$. This relation shows that the minimal value $\phi_\beta(\mathbf{w}, b)$ and an optimal solution can be achieved by minimizing the function $F_\beta(\mathbf{w}, b, \alpha)$ with respect to $(\mathbf{w}, b) \in W$ and $\alpha \in \mathbb{R}$ simultaneously. Furthermore, it is shown in [12] that an optimal α^* of the right-hand side optimization problem is almost equal to α_β .

2.2 Conditional Geometric Score (CGS) Optimization Problem

The optimization of the conditional lower geometric score given in (8) is described as

$$\left| \underset{\mathbf{w} \neq \mathbf{0}, b, \alpha}{\text{minimize}} \alpha + \frac{1}{1-\beta} \sum_{i \in I} p_i \left[-\frac{y^i (\langle \mathbf{w}, \mathbf{x}^i \rangle + b)}{\|\mathbf{w}\|} - \alpha \right]^+ \right., \quad (9)$$

where $\beta \in (0, 1)$ is a given confidence level. We call this problem the *Conditional Geometric Score* (CGS) optimization problem.

We show first that Problem (9) is equivalent to the following nonlinear programming problem with a nonconvex constraint:

$$(Q(\beta)) \left| \begin{array}{l} \underset{\mathbf{w}, b, z, \alpha}{\text{minimize}} \quad \alpha + \frac{1}{1-\beta} \sum_{i \in I} p_i z_i \\ \text{subject to} \quad z_i + y^i (\langle \mathbf{w}, \mathbf{x}^i \rangle + b) + \alpha \geq 0, \quad i \in I, \\ \quad \quad \quad z_i \geq 0, \quad i \in I, \\ \quad \quad \quad \|\mathbf{w}\|^2 = 1. \end{array} \right. \quad (10)$$

The above two problems (9) and (10) have optimal solutions under the following assumption.

Assumption 2.1

- ▷ Each class has at least one data, i.e., $|I^+|, |I^-| > 0$.
- ▷ The parameter β is chosen so that

$$1 - 2 \min\left\{ \sum_{i \in I^+} p_i, \sum_{i \in I^-} p_i \right\} \leq \beta < 1.$$

Lemma 2.2 *Under Assumption 2.1, Problem (10) has an optimal solution.*

Proof. We can find a feasible solution of (10) easily. Indeed, for arbitrary \bar{b} , $\bar{\alpha}$ and $\bar{\mathbf{w}}$ satisfying $\|\bar{\mathbf{w}}\|^2 = 1$, let \bar{z} satisfy

$$\bar{z}_i \geq \max\{-y_i (\langle \bar{\mathbf{w}}, \mathbf{x}^i \rangle + \bar{b}) - \bar{\alpha}, 0\}, \quad \forall i \in I,$$

then $(\bar{\mathbf{w}}, \bar{b}, \bar{z}, \bar{\alpha})$ is a feasible solution of (10).

In the following, we will show that this problem never become unbounded under Assumption 2.1. Introducing the dual variables $(\boldsymbol{\lambda}, \boldsymbol{\mu}, \delta) \in \mathbb{R}_+^\ell \times \mathbb{R}_+^\ell \times \mathbb{R}$ and constructing the Lagrangian function for Problem (10):

$$\begin{aligned} \mathcal{L}(\mathbf{w}, b, \alpha, \mathbf{z}, \boldsymbol{\lambda}, \boldsymbol{\mu}, \delta) \\ := \alpha + \frac{1}{1-\beta} \sum_{i \in I} p_i z_i - \sum_{i \in I} \lambda_i \{z_i + y^i (\langle \mathbf{w}, \mathbf{x}^i \rangle + b) + \alpha\} - \sum_{i \in I} \mu_i z_i + \frac{1}{2} \delta (\|\mathbf{w}\|^2 - 1), \end{aligned}$$

the Lagrangian dual problem is represented as

$$\underset{\boldsymbol{\lambda} \geq \mathbf{0}, \boldsymbol{\mu} \geq \mathbf{0}, \delta}{\text{maximize}} \quad \underset{\mathbf{w}, b, \alpha, \mathbf{z}}{\text{minimize}} \quad \mathcal{L}(\mathbf{w}, b, \alpha, \mathbf{z}, \boldsymbol{\lambda}, \boldsymbol{\mu}, \delta).$$

The stationarity conditions are

$$\sum_{i \in I} \lambda_i y^i \mathbf{x}^i = \delta \mathbf{w} \quad (11)$$

$$\sum_{i \in I} \lambda_i y^i = 0 \quad (12)$$

$$\sum_{i \in I} \lambda_i = 1 \quad (13)$$

$$0 \leq \lambda_i \leq \frac{p_i}{1-\beta}, \quad i \in I, \quad (14)$$

and substituting (11) to (13) into \mathcal{L} , we obtain the Lagrangian function:

$$\mathcal{L} = -\frac{1}{2} \delta (\|\mathbf{w}\|^2 + 1). \quad (15)$$

Therefore, when $\delta \neq 0$, the Lagrangian dual problem is reformulated as

$$\underset{\lambda, \delta}{\text{maximize}} \quad -\frac{1}{2} \left(\delta + \frac{1}{\delta} \left\| \sum_{i \in I} y^i \lambda_i \mathbf{x}^i \right\|^2 \right) \quad \text{subject to} \quad (12) \text{ to } (14),$$

and the constraints (12) to (14) are equivalent to

$$\sum_{i \in I^+} \lambda_i = \sum_{i \in I^-} \lambda_i = \frac{1}{2}, \quad 0 \leq \lambda_i \leq \frac{p_i}{1-\beta}, \quad i \in I.$$

It is easy to verify that the Lagrangian dual is feasible under Assumption 2.1, which proves that the optimal value of (10) is bounded below. Considering that \mathbf{w} is bounded in (10), we find that (10) has an optimal solution. \square

Proposition 2.3 *Problem (9) is equivalent to Problem (10) in the following sense :*

- \triangleright *Let $(\bar{\mathbf{w}}, \bar{b}, \bar{\alpha})$ be an optimal solution of (10). Then, $(\eta \bar{\mathbf{w}}, \eta \bar{b}, \bar{\alpha})$ for any $\eta > 0$ is an optimal solution of (9).*
- \triangleright *Let $(\mathbf{w}^*, b^*, \alpha^*)$ be an optimal solution of (9). Then, $(\frac{\mathbf{w}^*}{\|\mathbf{w}^*\|}, \frac{b^*}{\|\mathbf{w}^*\|}, \alpha^*)$ is an optimal solution of (10).*

Proof. Now we show the proof only for the first statement since the second one can be shown similarly.

For any $\eta > 0$, $(\eta\bar{\mathbf{w}}, \eta\bar{b}, \bar{\alpha})$ is a feasible solution of (9). Suppose that there exists a feasible solution $(\mathbf{w}', b', \alpha')$ of (9) whose objective function value is less than that of $(\eta\bar{\mathbf{w}}, \eta\bar{b}, \bar{\alpha})$, i.e.,

$$\alpha' + \frac{1}{1-\beta} \sum_{i \in I} p_i \left[-\frac{y^i (\langle \mathbf{w}', \mathbf{x}^i \rangle + b')}{\|\mathbf{w}'\|} - \alpha' \right]^+ < \bar{\alpha} + \frac{1}{1-\beta} \sum_{i \in I} p_i \left[-\frac{y^i (\langle \bar{\mathbf{w}}, \mathbf{x}^i \rangle + \bar{b})}{\|\bar{\mathbf{w}}\|} - \bar{\alpha} \right]^+.$$

Then, $(\frac{\mathbf{w}'}{\|\mathbf{w}'\|}, \frac{b'}{\|\mathbf{w}'\|}, \alpha')$ is feasible for (10) and it attains smaller objective function value than $(\bar{\mathbf{w}}, \bar{b}, \bar{\alpha})$, which contradicts the optimality of $(\bar{\mathbf{w}}, \bar{b}, \bar{\alpha})$ to (10). Therefore, for any $\eta > 0$, $(\eta\bar{\mathbf{w}}, \eta\bar{b}, \bar{\alpha})$ is an optimal solution of (9). \square

From Lemma 2.2 and the first statement in Proposition 2.3, we see that Problem (9) has an optimal solution under Assumption 2.1.

3 Algorithm for the CGS Optimization

The conditional geometric score (CGS) optimization problem (10) includes a single nonconvex constraint. In this section, we examine the nonconvex structure of the problem, and then propose an algorithm for achieving either a globally optimal solution when the problem is convex or a locally optimal solution when the problem is essentially nonconvex.

3.1 2-Step Algorithm

The CGS optimization problem (10) is a nonconvex optimization problem and seems to be difficult to deal with in a direct manner. We consider three cases according to the optimal value of Problem (10) and show the way to solve the problem in each case.

Let us define a relaxation problem of Problem (10) as

$$\text{(RQ}(\beta)\text{)} \left\{ \begin{array}{l} \text{minimize}_{\mathbf{w}, b, z, \alpha} \quad \alpha + \frac{1}{1-\beta} \sum_{i \in I} p_i z_i \\ \text{subject to} \quad z_i + y^i (\langle \mathbf{w}, \mathbf{x}^i \rangle + b) + \alpha \geq 0, \quad i \in I, \\ \quad z_i \geq 0, \quad i \in I, \\ \quad \|\mathbf{w}\|^2 \leq 1. \end{array} \right. \quad (16)$$

The existence of an optimal solution in (16) is shown by the similar proof to Lemma 2.2. Denoting the optimal value of Problem (10) by $\text{opt.}(\text{Q}(\beta))$ and that of Problem (16) by $\text{opt.}(\text{RQ}(\beta))$, then we have

$$\text{opt.}(\text{RQ}(\beta)) \leq \text{opt.}(\text{Q}(\beta)). \quad (17)$$

Proposition 3.1 *Let $(\hat{\mathbf{w}}, \hat{b}, \hat{\alpha}, \hat{z})$ be an optimal solution of Problem (16). Then, under Assumption 2.1, the followings hold.*

1. If $\|\hat{\mathbf{w}}\| = 1$ holds, then $(\hat{\mathbf{w}}, \hat{b}, \hat{\alpha}, \hat{\mathbf{z}})$ is an optimal solution of Problem (10). In this case, $\text{opt.}(Q(\beta)) \leq 0$.
2. If $0 < \|\hat{\mathbf{w}}\| < 1$ holds, then $\frac{1}{\|\hat{\mathbf{w}}\|}(\hat{\mathbf{w}}, \hat{b}, \hat{\alpha}, \hat{\mathbf{z}})$ is an optimal solution of Problem (10) and its optimal value is 0.
3. If $\|\hat{\mathbf{w}}\| = 0$ holds, $\text{opt.}(Q(\beta)) \geq 0$.

Proof. Note that $(\mathbf{w}, b, \alpha, \mathbf{z}) = \mathbf{0}$ is a feasible solution of Problem (16), and therefore, $\text{opt.}(\text{RQ}(\beta)) \leq 0$. We use this to prove the statements.

1. We observe that $(\hat{\mathbf{w}}, \hat{b}, \hat{\alpha}, \hat{\mathbf{z}})$ is feasible and, accordingly, optimal to (10). Therefore, $\text{opt.}(Q(\beta)) = \text{opt.}(\text{RQ}(\beta)) \leq 0$ holds.
2. Suppose on contrary that $\text{opt.}(\text{RQ}(\beta)) < 0$, then $\frac{1}{\|\hat{\mathbf{w}}\|}(\hat{\mathbf{w}}, \hat{b}, \hat{\alpha}, \hat{\mathbf{z}})$ is a feasible solution of (16) and $\frac{1}{\|\hat{\mathbf{w}}\|}\text{opt.}(\text{RQ}(\beta)) < \text{opt.}(\text{RQ}(\beta))$. This contradicts that $(\hat{\mathbf{w}}, \hat{b}, \hat{\alpha}, \hat{\mathbf{z}})$ is an optimal solution of (16), and hence, $\text{opt.}(\text{RQ}(\beta)) = 0$ follows. Thus, $\text{opt.}(Q(\beta)) = 0$ and $\frac{1}{\|\hat{\mathbf{w}}\|}(\hat{\mathbf{w}}, \hat{b}, \hat{\alpha}, \hat{\mathbf{z}})$ is an optimal solution of Problem (10).
3. If $\text{opt.}(Q(\beta)) < 0$, the relation (17) indicates that $\text{opt.}(\text{RQ}(\beta)) < 0$. The dual problem of (16) with an additional constraint $\mathbf{w} = \mathbf{0}$ is reduced to

$$\underset{\lambda, \delta}{\text{maximize}} \quad 0 \quad \text{subject to (12) to (14).}$$

Since this problem has a feasible solution under Assumption 2.1, one has $\text{opt.}(\text{RQ}(\beta)) = 0$ which contradicts the assumption. \square

In the third case of Proposition 3.1, any optimal solution of Problem (10) cannot be achieved from that of Problem (16) since the nonconvex constraint is essential. Accordingly, we have to solve a nonconvex programming problem in such a case.

Corollary 3.2 *Let $(\hat{\mathbf{w}}, \hat{b}, \hat{\alpha}, \hat{\mathbf{z}})$ be an optimal solution of Problem (16). Then, the followings hold.*

- \triangleright If $\text{opt.}(Q(\beta)) < 0$, one has $\|\hat{\mathbf{w}}\| = 1$.
- \triangleright If $\text{opt.}(Q(\beta)) > 0$, one has $\|\hat{\mathbf{w}}\| = 0$.

This corollary is proved by taking contraposition of the statements in Proposition 3.1. The corollary implies that when $\text{opt.}(Q(\beta)) < 0$, Problem (10) turns into the convex problem (16), and when $\text{opt.}(Q(\beta)) \geq 0$, it is equivalent to the nonconvex problem:

$$\left| \begin{array}{l} \underset{\mathbf{w}, b, \mathbf{z}, \alpha}{\text{minimize}} \quad \alpha + \frac{1}{1-\beta} \sum_{i \in I} p_i z_i \\ \text{subject to} \quad z_i + y^i (\langle \mathbf{w}, \mathbf{x}^i \rangle + b) + \alpha \geq 0, \quad i \in I, \\ \quad \quad \quad z_i \geq 0, \quad i \in I, \\ \quad \quad \quad \|\mathbf{w}\|^2 \geq 1. \end{array} \right. \quad (18)$$

Whether the convex constraint $\|\mathbf{w}\|^2 \leq 1$ in (16) is essential depends on the sign of the optimal value $\text{opt.}(Q(\beta))$ of (10) and we can roughly control the difficulty of Problem (10) by adjusting β since $\text{opt.}(Q(\beta))$ is non-decreasing with respect to β .

Now we outline our algorithm for solving the nonconvex program (10) to determine a discriminant hyperplane.

Algorithm 3.3 (*2-step framework*)

Step 1. At first, solve a convex optimization problem (16). For an optimal solution $(\hat{\mathbf{w}}, \hat{b}, \hat{\alpha}, \hat{\mathbf{z}})$,

- ▷ *if $\|\hat{\mathbf{w}}\| = 1$, terminate with an optimal solution $(\hat{\mathbf{w}}, \hat{b}, \hat{\alpha})$ of (10),*
- ▷ *if $0 < \|\hat{\mathbf{w}}\| < 1$, terminate with an optimal solution $(\frac{\hat{\mathbf{w}}}{\|\hat{\mathbf{w}}\|}, \frac{\hat{b}}{\|\hat{\mathbf{w}}\|}, \frac{\hat{\alpha}}{\|\hat{\mathbf{w}}\|})$ of (10),*
- ▷ *if $\|\hat{\mathbf{w}}\| = 0$, go to Step 2.*

Step 2. Solve the nonconvex program (10) approximately via the algorithm proposed in the next subsection.

3.2 Local Optimum Search Method with Approximation Accuracy

In this subsection, we focus on the case that an optimal solution of Problem (16) attains $\|\hat{\mathbf{w}}\| = 0$, that is, the third case of Proposition 3.1 where the optimal value of Problem (10) is nonnegative.

Lemma 3.4 *The following problem provides a lower bound for Problem (10) when the optimal value of (10) is nonnegative, i.e., $\text{opt.}(Q(\beta)) \geq 0$.*

$$(\text{RP}_\infty(\beta)) \left\{ \begin{array}{l} \underset{\mathbf{w}, b, z, \alpha}{\text{minimize}} \quad \alpha + \frac{1}{1-\beta} \sum_{i \in I} p_i z_i \\ \text{subject to} \quad z_i + y^i (\langle \mathbf{w}, \mathbf{x}^i \rangle + b) + \alpha \geq 0, \quad i \in I, \\ z_i \geq 0, \quad i \in I, \\ \|\mathbf{w}\|_\infty = \frac{1}{\sqrt{n}}. \end{array} \right. \quad (19)$$

Proof. Obviously, the optimal value of (19) is nonnegative, that is, $\text{opt.}(\text{RP}_\infty(\beta)) \geq 0$ and, in particular, when $\text{opt.}(Q(\beta)) = 0$, then $\text{opt.}(\text{RP}_\infty(\beta)) = 0$ holds. Hence, we focus on the case that $\text{opt.}(Q(\beta)) > 0$ below. Consider a relaxation problem of (19):

$$(\text{RRP}_\infty(\beta)) \left\{ \begin{array}{l} \underset{\mathbf{w}, b, z, \alpha}{\text{minimize}} \quad \alpha + \frac{1}{1-\beta} \sum_{i \in I} p_i z_i \\ \text{subject to} \quad z_i + y^i (\langle \mathbf{w}, \mathbf{x}^i \rangle + b) + \alpha \geq 0, \quad i \in I, \\ z_i \geq 0, \quad i \in I, \\ \|\mathbf{w}\|_\infty \geq \frac{1}{\sqrt{n}}, \end{array} \right. \quad (20)$$

and note that (20) is also a relaxation problem of (10) because the feasible region of (10) is included in that of (20). Also, we can show that when $\text{opt.}(Q(\beta)) > 0$, any optimal solution of (20) satisfies $\|\mathbf{w}\|_\infty = \frac{1}{\sqrt{n}}$, which implies that (19) is also a relaxation problem for (10). Indeed,

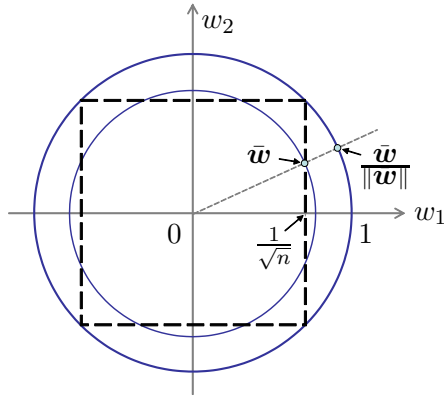


Figure 2: Constraints $\|\mathbf{w}\| = 1$ and $\|\mathbf{w}\|_\infty = 1/\sqrt{n}$

suppose on contrary that an optimal solution $\tilde{\mathbf{w}}$ of (20) satisfies $\|\tilde{\mathbf{w}}\|_\infty > \frac{1}{\sqrt{n}}$, then we obtain a feasible solution $\frac{1}{\sqrt{n}\|\tilde{\mathbf{w}}\|_\infty}(\tilde{\mathbf{w}}, \tilde{b}, \tilde{\alpha}, \tilde{\mathbf{z}})$ of (20) with smaller optimal value, which contradicts the optimality of $\tilde{\mathbf{w}}$. \square

For all $(d, h) \in \{1, \dots, n\} \times \{\pm \frac{1}{\sqrt{n}}\}$, solve the following linear programming problems:

$$(\text{RLP}_\infty(d, h; \beta)) \left\{ \begin{array}{l} \underset{\mathbf{w}, b, \mathbf{z}, \alpha}{\text{minimize}} \quad \alpha + \frac{1}{1-\beta} \sum_{i \in I} p_i z_i \\ \text{subject to} \quad z_i + y^i (\langle \mathbf{w}, \mathbf{x}^i \rangle + b) + \alpha \geq 0, \quad i \in I, \\ z_i \geq 0, \quad i \in I, \\ -\frac{1}{\sqrt{n}} \leq w_j \leq \frac{1}{\sqrt{n}} \text{ for all } j = 1, \dots, n, \\ w_d = h. \end{array} \right. \quad (21)$$

Using the optimal values $\bar{f}(d, h)$ for all $(d, h) \in \{1, \dots, n\} \times \{\pm \frac{1}{\sqrt{n}}\}$ and corresponding solutions, we obtain the optimal value \bar{f} of (19) by setting

$$\bar{f} := \min \{ \bar{f}(d, h) \mid (d, h) \in \{1, \dots, n\} \times \{\pm \frac{1}{\sqrt{n}}\} \},$$

and its optimal solution $(\bar{\mathbf{w}}, \bar{b}, \bar{\alpha}, \bar{\mathbf{z}})$ by defining

$$(\bar{\mathbf{w}}, \bar{b}, \bar{\alpha}, \bar{\mathbf{z}}) : \text{an optimal solution of } (\text{RLP}_\infty(d^*, h^*; \beta)), \quad (22)$$

where $(d^*, h^*) \in \arg \min \{ \bar{f}(d, h) \mid (d, h) \in \{1, \dots, n\} \times \{\pm \frac{1}{\sqrt{n}}\} \}$. Note that $0 < \|\bar{\mathbf{w}}\| < 1$, and $\frac{1}{\|\bar{\mathbf{w}}\|}(\bar{\mathbf{w}}, \bar{b}, \bar{\gamma}, \bar{\mathbf{z}})$ is feasible to (10) whose objective function value is $\frac{1}{\|\bar{\mathbf{w}}\|} \bar{f}$. One then has

$$0 \leq \bar{f} \leq \text{opt.}(Q(\beta)) \leq \frac{1}{\|\bar{\mathbf{w}}\|} \bar{f}.$$

Using the local optimum search technique proposed in [14], the upper bound $\frac{1}{\|\bar{\mathbf{w}}\|} \bar{f}$ might be improved. The procedure for local optimum search exploits the property of Problem (10), whose feasible region is reduced to a polyhedral cone by excluding the nonconvex constraint. We construct a polyhedral set by linearizing the nonconvex constraint as

$$\mathcal{F}(\mathbf{w}_k) := \left\{ (\mathbf{w}, b, \alpha, \mathbf{z}) \left\{ \begin{array}{l} z_i + y_i (\langle \mathbf{w}, \mathbf{x}^i \rangle + b) + \alpha \geq 0, \quad i \in I, \\ z_i \geq 0, \quad i \in I, \\ \langle \mathbf{w}_k, \mathbf{w} \rangle = 1 \end{array} \right. \right\}$$

with some feasible solution \mathbf{w}_k of (10), and then apply a pivoting technique for the linear programming problem:

$$\left| \begin{array}{l} \text{minimize } q(\mathbf{z}, \alpha) := \alpha + \frac{1}{1-\beta} \sum_{i \in I} p_i z_i \\ \text{subject to } (\mathbf{w}, b, \alpha, \mathbf{z}) \in \mathcal{F}(\mathbf{w}_k) \end{array} \right. \quad (23)$$

to find a local optimum of Problem (10). Note that, when the optimal value of Problem (10) is nonnegative, it is equivalent to (18) which includes the constraint $\|\mathbf{w}\|^2 \geq 1$ instead of $\|\mathbf{w}\|^2 = 1$. Since $\mathcal{F}(\mathbf{w}_k)$ is included in the feasible region of (18), the optimal value of (23) is also nonnegative.

Here we introduce some additional definitions for the description of local optimum search algorithm. Let us define a polyhedral cone H by

$$H := \left\{ (\mathbf{w}, b, \alpha, \mathbf{z}) \left| \begin{array}{l} z_i + y^i(\langle \mathbf{w}, \mathbf{x}^i \rangle + b) + \alpha \geq 0, \quad i \in I \\ z_i \geq 0, \quad i \in I \end{array} \right. \right\},$$

and a unit sphere G by

$$G := \{(\mathbf{w}, b, \alpha, \mathbf{z}) \mid \|\mathbf{w}\|^2 = 1\}.$$

Let $\text{ray}H$ denote the set of extreme rays of H , and we call a feasible solution lying on $\text{ray}H \cap G$ a *basic solution*.

Algorithm 3.5 (*Local optimum search*)

Step 0. Choose a feasible basic solution $(\mathbf{w}_0, b_0, \alpha_0, \mathbf{z}_0)$ of (10). Let $k = 0$.

Step 1. At \mathbf{w}_k , construct the polyhedral set $\mathcal{F}(\mathbf{w}_k)$.

Step 2. If there exists an extreme point $(\bar{\mathbf{w}}, \bar{b}, \bar{\alpha}, \bar{\mathbf{z}})$ of $\mathcal{F}(\mathbf{w}_k)$ adjacent to $(\mathbf{w}_k, b_k, \alpha_k, \mathbf{z}_k)$ with $q(\mathbf{z}_k, \alpha_k) \geq q(\bar{\mathbf{z}}, \bar{\alpha})$, then compute the next basic solution

$$(\mathbf{w}_{k+1}, b_{k+1}, \alpha_{k+1}, \mathbf{z}_{k+1}) := \frac{1}{\|\bar{\mathbf{w}}\|}(\bar{\mathbf{w}}, \bar{b}, \bar{\alpha}, \bar{\mathbf{z}}) \quad (24)$$

of (10), let $k = k + 1$ and go to Step 1. Otherwise, terminate with a local minimizer $(\mathbf{w}_k, b_k, \alpha_k, \mathbf{z}_k)$.

In Step 0, we can find a feasible basic solution of (10) easily. After obtaining a feasible solution $\frac{1}{\|\bar{\mathbf{w}}\|}(\bar{\mathbf{w}}, \bar{b}, \bar{\gamma}, \bar{\mathbf{z}})$ from (22), we construct $\mathcal{F}(\frac{\bar{\mathbf{w}}}{\|\bar{\mathbf{w}}\|})$. Then, we find an extreme point $(\mathbf{w}', b', \mathbf{z}', \alpha')$ of $\mathcal{F}(\frac{\bar{\mathbf{w}}}{\|\bar{\mathbf{w}}\|})$ after one pivot procedure, project the extreme point on G by (24), and achieve a feasible basic solution of (10).

The solution $(\mathbf{w}_k, b_k, \alpha_k, \mathbf{z}_k)$ obtained by this algorithm corresponds to a local minimizer of Problem (10), since this solution satisfies the *Karush-Kuhn-Tucker* (KKT) optimality conditions induced from (10). According to [14], it is easily shown that Algorithm 3.5 terminates within finite iterations when the optimal value of Problem (10) is nonnegative, i.e., $\text{opt.}(Q(\beta)) \geq 0$.

4 Relation to the Other Classification Problems

In this section, we discuss a relationship between the conditional geometric score (CGS) optimization problem and other classification problems such as the hard margin support vector classification (HSVC) and the ν -support vector classification (ν -SVC). Under some parameter setting, our model is proved to be equivalent to them.

4.1 The Hard Margin Support Vector Classification

We first explore the relation between the CGS optimization (9) and the hard margin support vector classification (2). Note that (2) can be represented as

$$\text{(HSVC2)} \quad \left| \begin{array}{l} \text{minimize} \\ \mathbf{w} \neq \mathbf{0}, b \end{array} \right. \max_{i \in I} - \frac{y^i (\langle \mathbf{w}, \mathbf{x}^i \rangle + b)}{\|\mathbf{w}\|}.$$

We show below that the formulation (9) coincides with (HSVC2) if $\beta > 1 - \min_i p_i$ holds.

Let $(\mathbf{w}^*, b^*, \alpha^*)$ be an optimal solution of (9) and let

$$k := \arg \max_{i \in I} \left\{ - \frac{y^i (\langle \mathbf{w}^*, \mathbf{x}^i \rangle + b^*)}{\|\mathbf{w}^*\|} \right\}.$$

Suppose that there exists $\gamma > 0$ such that

$$\alpha^* = - \frac{y^k (\langle \mathbf{w}^*, \mathbf{x}^k \rangle + b^*)}{\|\mathbf{w}^*\|} - \gamma.$$

Then,

$$\begin{aligned} \alpha^* + \frac{1}{1-\beta} \sum_{i \in I} p_i \left[- \frac{y^i (\langle \mathbf{w}^*, \mathbf{x}^i \rangle + b^*)}{\|\mathbf{w}^*\|} - \alpha^* \right]^+ &\geq - \frac{y^k (\langle \mathbf{w}^*, \mathbf{x}^k \rangle + b^*)}{\|\mathbf{w}^*\|} + \left(\frac{p_k}{1-\beta} - 1 \right) \gamma \\ &> - \frac{y^k (\langle \mathbf{w}^*, \mathbf{x}^k \rangle + b^*)}{\|\mathbf{w}^*\|}, \end{aligned}$$

which contradicts the optimality of $(\mathbf{w}^*, b^*, \alpha^*)$. Accordingly, the optimal solution satisfies

$$\alpha^* \geq \max_{i \in I} \left\{ - \frac{y^i (\langle \mathbf{w}^*, \mathbf{x}^i \rangle + b^*)}{\|\mathbf{w}^*\|} \right\},$$

and Problem (9) results in (HSVC2).

From this viewpoint, the hard margin support vector classification (HSVC) corresponds to the CGS optimization with sufficiently large β . To attain more accurate classifier, it might be better to solve the CGS optimization problem (10) which takes into consideration the other misclassified data by adjusting the parameter β .

4.2 The ν -Support Vector Classification

We next show that the convex problem (16) solved at the first step of Algorithm 3.3 is equivalent to the ν -SVC when the optimal value of (10) is negative. Besides, it is shown that the ν -SVC results in a meaningless solution when the optimal value of (10) is positive.

Multiplying the objective function of (16) by $1 - \beta > 0$ and introducing the Lagrangian multiplier $\theta \geq 0$ with respect to the convex constraint $\|\mathbf{w}\|^2 \leq 1$, we obtain a Lagrangian relaxation problem of (16):

$$(\text{LRQ}(\theta; \beta)) \left\{ \begin{array}{l} \underset{\mathbf{w}, b, \mathbf{z}, \alpha}{\text{minimize}} \quad \theta(\|\mathbf{w}\|^2 - 1) + (1 - \beta)\alpha + \sum_{i \in I} p_i z_i \\ \text{subject to} \quad z_i + y^i(\langle \mathbf{w}, \mathbf{x}^i \rangle + b) + \alpha \geq 0, \quad i \in I, \\ z_i \geq 0, \quad i \in I. \end{array} \right. \quad (25)$$

Then, the optimal value of $(\text{LRQ}(\theta; \beta))$ with any $\theta \geq 0$ is equal to or less than that of (16).

If we change the variable and parameters in Problem (25) as

$$\alpha = -\rho, \quad \beta = 1 - \nu, \quad p_i = \frac{1}{|I|} \text{ for } i \in I, \quad (26)$$

add nonpositivity on $\alpha \leq 0$ and choose the Lagrangian multiplier as $\theta = \frac{1}{2}$, this problem is led into the primal formulation of the ν -SVC [13]:

$$(\nu\text{-SVC}(\nu)) \left\{ \begin{array}{l} \underset{\mathbf{w}, b, \mathbf{z}, \rho}{\text{minimize}} \quad \frac{1}{2}\|\mathbf{w}\|^2 - \nu\rho + \frac{1}{|I|} \sum_{i \in I} z_i \\ \text{subject to} \quad z_i + y^i(\langle \mathbf{w}, \mathbf{x}^i \rangle + b) - \rho \geq 0, \quad i \in I, \\ z_i \geq 0, \quad i \in I, \\ \rho \geq 0. \end{array} \right. \quad (27)$$

In our classification model, α corresponds to a threshold of the distribution of $f(\mathbf{w}, b)$ with confidence level $\beta \in (0, 1)$ and the nonpositivity is not imposed on α .

As shown in Corollary 3.2, when the optimal value of Problem (10) satisfies $\text{opt.}(Q(\beta)) < 0$, (10) is equivalent to the convex problem (16), and when $\text{opt.}(Q(\beta)) \geq 0$, it is equivalent to the nonconvex problem (18). From this point of view, the sign of Lagrangian multiplier θ should be determined according to the optimal value of (10). In addition, we can show the following theorem.

Theorem 4.1 *The ν -SVC (27) with parameters (26) provides a meaningless solution satisfying $\mathbf{w} = \mathbf{0}$ when the optimal value of (10) is positive.*

Proof. Corollary 3.2 shows that when $\text{opt.}(Q(\beta)) > 0$, Problem (16) provides a meaningless optimal solution satisfying $\mathbf{w} = \mathbf{0}$. In such a case, Problem (16) is equivalent to the following linear programming problem:

$$(\text{RP2}(\beta)) \left\{ \begin{array}{l} \underset{\mathbf{w}, b, \mathbf{z}, \alpha}{\text{minimize}} \quad (1 - \beta)\alpha + \sum_{i \in I} p_i z_i \\ \text{subject to} \quad z_i + y^i(\langle \mathbf{w}, \mathbf{x}^i \rangle + b) + \alpha \geq 0, \quad i \in I, \\ z_i \geq 0, \quad i \in I. \end{array} \right. \quad (28)$$

By duality, the optimal value $\text{opt.}(\text{RP2}(\beta))$ of (28) is 0, and the comparison between the ν -SVC (27) with parameters (26) and Problem (28) implies

$$0 = \text{opt.}(\text{RP2}(\beta)) \leq \text{opt.}(\nu\text{-SVC}(\nu)).$$

Therefore, a feasible solution $(\mathbf{w}, b, \rho, \mathbf{z}) = \mathbf{0}$ is optimal in (27). We can show that there is no optimal solution satisfying $\|\mathbf{w}\| > 0$, since if (27) has such a solution and the optimal value is 0, then $\text{opt.}(\text{RP2}(\beta)) < 0$ holds. Thus, any optimal solution of (27) satisfies $\mathbf{w} = \mathbf{0}$, and the statement follows. \square

Now we introduce a convex quadratic program:

$$(D^+) \left\{ \begin{array}{l} \underset{\lambda}{\text{maximize}} \quad - \sum_{i \in I} \sum_{j \in I} y^i y^j \langle \mathbf{x}^i, \mathbf{x}^j \rangle \lambda_i \lambda_j \\ \text{subject to} \quad \sum_{i \in I} \lambda_i y^i = 0, \\ \sum_{i \in I} \lambda_i = 1, \\ 0 \leq \lambda_i \leq \frac{1}{1 - \beta} p_i, \quad i \in I. \end{array} \right.$$

Lemma 4.2 *When the optimal value of (10) is negative, (D^+) is the Lagrangian dual problem of (10) with zero gap.*

Proof. In the proof of Lemma 2.2, we have provided the Lagrangian dual problem of (10) as

$$\underset{\lambda, \delta}{\text{maximize}} \quad -\frac{1}{2} \left(\delta + \frac{1}{\delta} \left\| \sum_{i \in I} y^i \lambda_i \mathbf{x}^i \right\|^2 \right) \quad \text{subject to (12) to (14),}$$

when $\delta \neq 0$. Since the optimal value of the primal problem (10) is negative, the optimal value of the above Lagrangian dual problem is less than 0, any optimal solution satisfies $\delta^* > 0$ and (11) shows $\left\| \sum_{i \in I} y^i \lambda_i \mathbf{x}^i \right\| > 0$. Then, the relation between the arithmetic and the geometric means implies that

$$-\frac{1}{2} \left(\delta + \frac{1}{\delta} \left\| \sum_{i \in I} y^i \lambda_i \mathbf{x}^i \right\|^2 \right) \leq - \left\| \sum_{i \in I} y^i \lambda_i \mathbf{x}^i \right\| = - \sqrt{\sum_{i \in I} \sum_{j \in I} y^i y^j \langle \mathbf{x}^i, \mathbf{x}^j \rangle \lambda_i \lambda_j},$$

and the Lagrangian dual problem results in (D^+) , where the optimal solution δ^* is achieved at $\delta^* = \frac{1}{2} \sqrt{\sum_{i \in I} \sum_{j \in I} y^i y^j \langle \mathbf{x}^i, \mathbf{x}^j \rangle \lambda_i \lambda_j}$. Since Problem (10) is equivalent to Problem (16) when the optimal value of (10) is negative and both the Lagrangian dual (D^+) and (16) are convex programs, the result follows. \square

Theorem 4.3 *The CGS optimization problem (10) and the ν -SVC (27) with parameters (26) provide the same discriminant hyperplane when the optimal value of (10) is negative.*

Proof. The dual of the ν -SVC (27) is known as

$$(D\nu\text{SVC}) \left\{ \begin{array}{l} \underset{\lambda}{\text{maximize}} \quad - \sum_{i \in I} \sum_{j \in I} y^i y^j \langle \mathbf{x}^i, \mathbf{x}^j \rangle \lambda_i \lambda_j \\ \text{subject to} \quad \sum_{i \in I} \lambda_i y^i = 0, \\ \sum_{i \in I} \lambda_i \geq \nu, \\ 0 \leq \lambda_i \leq \frac{1}{|I|}, \quad i \in I. \end{array} \right.$$

By replacing the variables λ by $\lambda' := \frac{\lambda}{\nu}$, the problem is proved to be equivalent to (D^+) except the inequality $\sum_{i \in I} \lambda'_i \geq 1$ which corresponds to $\rho \geq 0$ in (27). When the optimal value of (10) is negative, $\alpha^* < 0$ holds at any optimal solution α^* . Therefore, these two problems are equivalent when the optimal value of (10) is negative. \square

5 Computational Experiments

In this section, we examine our model by comparing with existing models such as the ν -SVC and the robust linear programming (RLP) approach [2] through numerical experiments. The instances to be used in the experiments are named 1) CANCER, 2) LIVER, and 3) DIABETES, all of which are obtained from the UCI repository of machine learning databases [3]. For the CANCER data, we use only three attributes which are shown to be highly predictable via RLP approach in [8], whereas we use all attributes provided in the database for the latter two instances, i.e., six attributes for LIVER, and eight for DIABETES.

Distribution of Optimal Geometric Score Figure 3 shows histograms and statistical characteristics for the optimal distribution of the geometric score $f(\mathbf{w}^*, b^*; \mathbf{x}, y) = -g(\mathbf{w}^*, b^*; \mathbf{x}, y)$, where $(\mathbf{w}^*, b^*, \alpha^*, \mathbf{z}^*)$ is an optimal solution of Problem (10) when the empirical probability p_i is given by (4). Figures a) and b) are the results of the CANCER data under the different parameter β . Figures a), c) and d) show the results of CANCER, DIABETES and LIVER, respectively, under the parameter β chosen so that the number of misclassification becomes very small. It is worth noting that the optimal values ϕ_β in Figures c) and d) are positive and the nonconvexity is essential. From these figures, we see that the shape of the optimal distribution depends on β and the given data sets.

Misclassification and Misprediction To show the potential superiority of our model over the ν -SVC, we performed ten-fold cross-validation for the three data sets. Figure 4 shows the misclassification (training error) and the misprediction (testing error) rates obtained by Problem (10) with $\beta = 0.4$ to 0.55 when the DIABETES data set is applied. The points highlighted by ‘o’ and ‘x’ in the figure correspond to $\beta \geq 0.47$. For such β , at least one of the ten trials through the cross-validation proceeded to the Step 2 of Algorithm 3.3. It implies that when $\beta \geq 0.47$, that is, $\nu \leq 0.53$, the ν -SVC fails to attain any hyperplane at least one of the ten trials while the numerical results of our model and the ν -SVC are exactly the same when $\beta \leq 0.46$. Consequently, it is possible to find more predictable hyperplane by our model for larger β so that the ν -SVC returns a meaningless solution, i.e., $\mathbf{w}^* = \mathbf{0}$. This indicates that the nonconvexity plays an important role in the classification problem. Similar results are observed with the other data sets.

Tables 1 to 3 summarize the results of misclassification and misprediction of three data sets by using three models: the CGS, the ν -SVC, and the RLP. For the CGS and the ν -SVC, we compute discriminant hyperplanes under several values of β and report the best results. Each cell shows

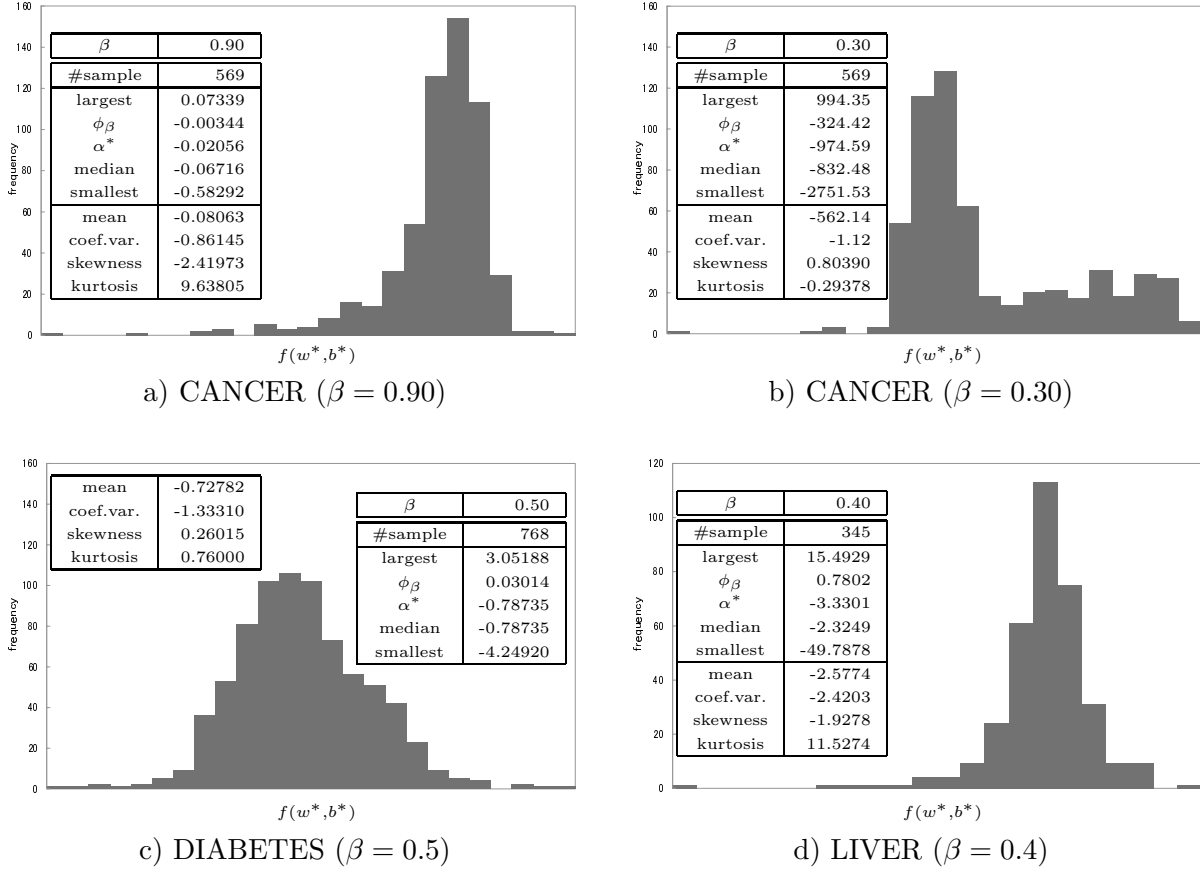


Figure 3: Histograms of f

the sample average and the sample standard deviation of the misclassification or misprediction rates. The column named ‘+1’ shows the rates for the ‘+1’-labeled data set while the ‘-1’ shows those for the other data set. The rows named ‘CGS1’ show results achieved by the model proposed in the paper when the empirical probability is given by (4), whereas ‘CGS2’ means the same model when the probability is given by (5). Similarly, ‘RLP1’ means the minimization of equally weighted misclassification error, i.e., $\frac{1}{\ell} \sum_{i \in I} \xi_i$, whereas ‘RLP2’ minimizes weighted error terms defined by $\frac{1}{|I^+|} \sum_{i \in I^+} \xi_i + \frac{1}{|I^-|} \sum_{i \in I^-} \xi_i$, where $\xi_i = [1 - y_i(\langle \mathbf{w}, \mathbf{x}^i \rangle + b)]^+$. Correspondingly, ‘ ν -SVC1’ means (27), whereas ‘ ν -SVC2’ means a variation of (27) whose objective function is replaced by $\frac{1}{2} \|\mathbf{w}\|^2 - \nu \rho + \frac{1}{2|I^+|} \sum_{i \in I^+} z_i + \frac{1}{2|I^-|} \sum_{i \in I^-} z_i$.

From these tables, our model achieves the best predictive accuracy in terms of the sample mean through ten-fold cross-validation. For the CANCER data, our model performs as well as the RLP2 which is shown to be highly predictable in [8] when the weighted empirical probability (5) is applied. However, the weighted probability did not always attain the better total predictability, but the misprediction rates of the two labels are more balanced than the equally weighted ones.

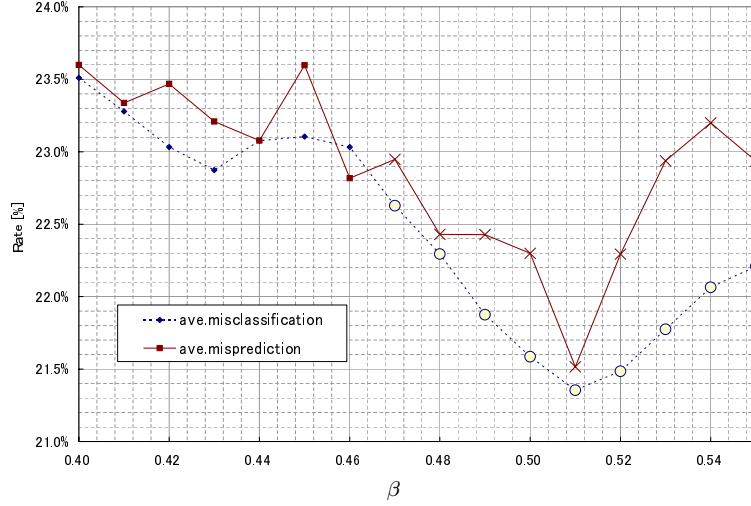


Figure 4: Misclassification vs. misprediction through Algorithm 3.3 (the DIABETES data)

Table 1: Misclassification and misprediction rates for the CANCER data

		misclassification			misprediction		
		total	+1	-1	total	+1	-1
CGS1 ($\beta = .97$)	av.	3.03%	4.09%	2.40%	3.16%	4.26%	2.51%
	s.d.	(0.53%)	(0.62%)	(0.49%)	(1.99%)	(4.29%)	(1.62%)
CGS2 ($\beta = .92$)	av.	2.60%	3.30%	2.18%	2.63%	3.46%	2.21%
	s.d.	(0.29%)	(0.37%)	(0.32%)	(1.70%)	(3.32%)	(1.73%)
ν -SVC1 ($\beta = .88$)	av.	3.12%	6.29%	1.24%	3.16%	6.82%	1.10%
	s.d.	(0.37%)	(0.88%)	(0.20%)	(2.45%)	(5.67%)	(1.92%)
ν -SVC2 ($\beta = .88$)	av.	2.75%	4.19%	1.90%	2.81%	4.01%	2.25%
	s.d.	(0.34%)	(0.92%)	(0.44%)	(1.69%)	(4.09%)	(1.76%)
RLP1	av.	3.14%	5.66%	1.65%	3.33%	6.32%	1.65%
	s.d.	(0.42%)	(0.92%)	(0.30%)	(2.54%)	(6.06%)	(1.92%)
RLP2	av.	2.54%	3.51%	1.96%	2.63%	3.46%	2.21%
	s.d.	(0.26%)	(0.56%)	(0.21%)	(1.70%)	(3.32%)	(1.73%)

Table 2: Misclassification and misprediction rates for the DIABETES data

		misclassification			misprediction		
		total	+1	-1	total	+1	-1
CGS1 ($\beta = .51$)	av.	21.35%	11.11%	40.48%	21.52%	10.97%	40.36%
	s.d.	(0.51%)	(0.68%)	(1.29%)	(7.43%)	(5.18%)	(10.74%)
CGS2 ($\beta = .42$)	av.	24.07%	21.14%	29.57%	24.25%	21.77%	28.71%
	s.d.	(0.45%)	(0.94%)	(1.19%)	(7.92%)	(5.98%)	(13.84%)
ν -SVC1 ($\beta = .46$)	av.	23.03%	11.83%	43.95%	22.82%	11.80%	43.08%
	s.d.	(0.97%)	(1.19%)	(1.04%)	(7.18%)	(4.89%)	(11.22%)
ν -SVC2 ($\beta = .42$)	av.	24.07%	21.14%	29.57%	24.25%	21.77%	28.71%
	s.d.	(0.45%)	(0.94%)	(1.19%)	(7.92%)	(5.98%)	(13.84%)
RLP1	av.	22.15%	11.09%	42.79%	22.56%	11.61%	42.77%
	s.d.	(0.84%)	(0.56%)	(1.43%)	(7.12%)	(4.89%)	(10.83%)
RLP2	av.	23.35%	21.03%	27.69%	23.20%	20.72%	27.55%
	s.d.	(0.81%)	(0.77%)	(1.44%)	(7.57%)	(5.60%)	(12.48%)

Table 3: Misclassification and Misprediction for the LIVER data

		misclassification			misprediction		
		total	+1	-1	total	+1	-1
CGS1 ($\beta = .41$)	av.	29.66%	43.38%	19.73%	28.64%	42.88%	18.51%
	s.d.	(1.85%)	(2.16%)	(2.28%)	(7.07%)	(13.77%)	(6.69%)
CGS2 ($\beta = .31$)	av.	31.46%	34.33%	29.40%	31.55%	36.07%	28.51%
	s.d.	(1.79%)	(2.26%)	(2.64%)	(7.72%)	(12.25%)	(8.30%)
ν -SVC1 ($\beta = .25$)	av.	29.44%	49.09%	15.23%	30.69%	51.88%	15.46%
	s.d.	(0.93%)	(3.54%)	(1.92%)	(5.63%)	(9.94%)	(5.75%)
ν -SVC2 ($\beta = .23$)	av.	32.27%	30.73%	33.40%	33.91%	34.17%	34.03%
	s.d.	(1.75%)	(2.10%)	(2.57%)	(6.64%)	(11.30%)	(6.04%)
RLP1	av.	31.69%	31.80%	31.62%	33.03%	34.80%	32.01%
	s.d.	(2.03%)	(1.95%)	(2.74%)	(7.87%)	(13.45%)	(6.32%)
RLP2	av.	28.86%	44.92%	17.23%	30.39%	48.46%	17.46%
	s.d.	(0.78%)	(1.96%)	(1.34%)	(6.07%)	(9.84%)	(6.58%)

6 Concluding Remarks

In this paper, we proposed a classification model based on the conditional expectation of the geometric score which is a simple extension of the Euclidean distance of data point to the discriminant hyperplane. The model is formulated as a nonconvex optimization problem, and we present a two-step algorithm for obtaining either a globally or a locally optimal solution of the problem. At the first step of the algorithm, a convex problem (16) is solved which is shown to be almost the same as the ν -SVC formulation. When the first step results in a meaningless one, i.e., the zero normal vector, the algorithm proceeds to the second step at which a local solution of Problem (10) is provided by solving a finite number of additional linear programming problems.

The numerical experiments indicate that the nonconvexity in our model plays an important role in achieving higher predictability. In fact, for any instance of CANCER, LIVER and DIABETES, the most predictable discriminant hyperplane was induced from an approximate solution of the nonconvex problem (10). In addition, our model with β satisfying Assumption 2.1 always provides a discriminant hyperplane while the ν -SVC results in a meaningless solution for large β .

One possible extension is to construct a nonlinear discriminant surface by incorporating the kernel functions into the dual problem (D^+) of Problem (16). Since (D^+) is equivalent to the ν -SVC formulation under some condition, it is easy to introduce the kernel functions into (D^+) in a similar manner to the ν -SVC. On the other hand, we need more discussion to incorporate the positive definite kernel functions into the nonconvex problem (18). The detailed exploration remains to be done in the future work.

Acknowledgement

Research of the first author is supported by MEXT Grant-in-Aid for Young Scientists (B) 14780343. Research of the second author is supported by MEXT Grant-in-Aid for Young Scientists (B) 16710110.

References

- [1] P. Artzner, F. Delbaen, J.M. Eber and D. Heath, Coherent measures of risk, *Math. Finance* 9 (1999) 203-228.
- [2] K.P. Bennett and O.L. Mangasarian, Robust linear programming discrimination of two linearly inseparable sets, *Optim. Methods Softw.* 1 (1992) 23-34.
- [3] C.L. Blake and C.J. Merz, UCI Repository of machine learning databases [<http://www.ics.uci.edu/mlearn/MLRepository.html>]. Irvine, CA: University of California, Department of Information and Computer Science, 1998.

- [4] N. Freed and F. Glover, A linear programming approach to the discriminant problem, *Decis. Sci.* 12 (1981) 68-74.
- [5] D.J. Hand, *Discrimination and classification*, John Wiley, Chichester U.K., 1981.
- [6] G.J. Koehler and S.S. Erenguc, Minimizing misclassification in linear discriminant analysis, *Decis. Sci.* 21 (1990) 63-85.
- [7] O.L. Mangasarian, Multisurface method of pattern separation, *IEEE Trans. Inform. Theory* IT-14 (1968) 801-807.
- [8] O.L. Mangasarian, W.N. Street and W.H. Wolberg, Breast cancer diagnosis and prognosis via linear programming, *Oper. Res.* 43 (1995) 570-577.
- [9] W. Ogryczak and A. Ruszczyński, Dual stochastic dominance and related mean-risk models, *SIAM J. Optim.* 13 (2002) 60-78.
- [10] G.Ch. Pflug, Some remarks on the value-at-risk and the conditional value-at-risk, in *Probabilistic Constrained Optimization: Methodology and Applications*, S. Uryasev (ed.), Kluwer Academic Publishers, Dordrecht, 2000, 272-281.
- [11] R.T. Rockafellar and S. Uryasev, Optimization of conditional value-at-risk, *J. Risk* 2 (2000) 21-41.
- [12] R.T. Rockafellar and S. Uryasev, Conditional value-at-risk for general loss distributions, *J. Bank. Finance* 26 (2002) 1443-1471.
- [13] B. Schölkopf, A. J. Smola, R. C. Williamson and P. L. Bartlett, New support vector algorithms, *Neural Comput.* 12 (2000) 1207-1245.
- [14] A. Takeda and H. Nishino, On measuring the inefficiency with the inner-product norm in date envelopment analysis, *European J. Oper. Res.* 133 (2001) 377-393.
- [15] V. Vapnik, *The nature of statistical learning theory*, Springer Verlag, New York, 1995.