

Article

# Kernel-Free Quadratic Surface Minimax Probability Machine for a Binary Classification Problem

Yulan Wang, Zhixia Yang \*  and Xiaomei Yang

College of Mathematics and Systems Science, Xinjiang University, Urumuqi 830046, China; 107551900393@stu.xju.edu.cn (Y.W.); xiaomeiy@xju.edu.cn (X.Y.)

\* Correspondence: yangzhx@xju.edu.cn; Tel.: +86-189-9913-4166

**Abstract:** In this paper, we propose a novel binary classification method called the kernel-free quadratic surface minimax probability machine (QSMPM), that makes use of the kernel-free techniques of the quadratic surface support vector machine (QS SVM) and inherits the advantage of the minimax probability machine (MPM) without any parameters. Specifically, it attempts to find a quadratic hypersurface that separates two classes of samples with maximum probability. However, the optimization problem derived directly was too difficult to solve. Therefore, a nonlinear transformation was introduced to change the quadratic function involved into a linear function. Through such processing, our optimization problem finally became a second-order cone programming problem, which was solved efficiently by an alternate iteration method. It should be pointed out that our method is both kernel-free and parameter-free, making it easy to use. In addition, the quadratic hypersurface obtained by our method was allowed to be any general form of quadratic hypersurface. It has better interpretability than the methods with the kernel function. Finally, in order to demonstrate the geometric interpretation of our QSMPM, five artificial datasets were implemented, including showing the ability to obtain a linear separating hyperplane. Furthermore, numerical experiments on benchmark datasets confirmed that the proposed method had better accuracy and less CPU time than corresponding methods.

**Keywords:** classification; quadratic surface support vector machine; kernel-free; minimax probability machine; second-order cone programming problem



**Citation:** Wang, Y.; Yang, Z.; Yang, X. Kernel-Free Quadratic Surface Minimax Probability Machine for a Binary Classification Problem. *Symmetry* **2021**, *13*, 1378. <https://doi.org/10.3390/sym13081378>

Academic Editor: Jan Awrejcewicz

Received: 1 July 2021  
Accepted: 26 July 2021  
Published: 28 July 2021

**Publisher's Note:** MDPI stays neutral with regard to jurisdictional claims in published maps and institutional affiliations.



**Copyright:** © 2021 by the authors. Licensee MDPI, Basel, Switzerland. This article is an open access article distributed under the terms and conditions of the Creative Commons Attribution (CC BY) license (<https://creativecommons.org/licenses/by/4.0/>).

## 1. Introduction

Machine learning is an important branch in the field of artificial intelligence, which has a wide range of applications in various fields of contemporary science [1]. With the development of machine learning, the classification problem has been widely concerned and studied in the fields of pattern recognition [2], text classification [3], image processing [4], financial time series prediction [5], skin disease [6], intrusion detection systems [7], etc. The classification problem is a vital task in supervised learning that learns a classification rule from a training set with known labels and then uses it to assign a new sample to a class.

At present, there are many famous classification methods. Among these existing methods, Lanckriet et al. [8,9] proposed an excellent classifier, called the minimax probability machine (MPM). For a given binary classification problem, the MPM not only deals with it in the linear case, but also in the nonlinear case by the kernel trick. It is worth noting that the MPM does not have any parameters, which is an important advantage. Therefore, it has been widely used in computer vision [10], engineering technology [11,12], agriculture [13], and novelty detection [14]. Moreover, many researchers have proposed a variety of improved versions of the MPM from different perspectives [14–25]. The representative works can be briefly reviewed as follows. In [15], Thomas and Gregory proposed MPM regression (MPMR), which transformed the regression problem into a classification problem, and then used the classifier MPM to obtain a regression function. To further exploit the structural

information of the training set, Gu et al. [17] proposed the structural MPM (SMPM) by combining the finite mixture models with the MPM. In addition, Yoshiyama et al. [21] proposed the Laplacian MPM (Lap-MPM), which improved the performance of the MPM in semisupervised learning. However, the nonlinear MPM using kernel techniques lacks interpretability and usually depends heavily on the choice of a proper kernel function and the corresponding kernel parameters. Furthermore, choosing the appropriate kernel function and adjusting its parameters may require much computational time and effort. Therefore, it naturally occurs to us that the study of a kernel-free nonlinear MPM is of great significance.

For the first time, Dagher [26] proposed a kernel-free nonlinear classifier, namely the quadratic surface support vector machine (QSSVM), in 2008. It was based on the maximum margin idea, and the training points were separated by a quadratic hypersurface without a kernel function, avoiding the time-consuming process of selecting the appropriate kernel function and its corresponding parameters. Furthermore, in order to improve the classification accuracy and robustness, Luo et al. [27] proposed the soft-margin quadratic surface support vector machine (SQSSVC). After that, Bai et al. [28] proposed the quadratic kernel-free least-squares support vector machine for target diseases' classification. Following these leading works, some scholars performed further studies, e.g., see [29–34] for the classification problem, [35] for the regression problem, and [36] for the cluster problem. The good performance of these methods demonstrates that the quadratic hypersurface is an effective method to flexibly capture the nonlinear structure of data. Thus, it can be seen that it is very interesting to study the kernel-free nonlinear MPM using the above kernel-free technique.

In this paper, for the binary classification problem, a new kernel-free nonlinear method is proposed, which is called the kernel-free quadratic surface minimax probability machine (QSMPM). It was constructed on the basics of the MPM by using the kernel-free techniques of the QSSVM. Specifically, it tries to seek a quadratic hypersurface that separates two classes of samples with maximum probability. However, the optimization problem derived directly was too difficult to solve. Therefore, a nonlinear transformation was introduced to change the quadratic function involved into a linear function. Through such processing, our optimization problem finally became a second-order cone programming problem, which was solved efficiently by an alternate iteration method. It is important to point out that our QSMPM addresses the following key issues. First, our method directly generates a nonlinear (quadratic) hypersurface without the kernel function, so there is no need to select the appropriate kernel. Second, our method does not need to choose any parameters. Third, the quadratic hypersurface obtained by our method has better interpretability than the one by the methods with the kernel function. Fourth, it is rather flexible because the quadratic hypersurface obtained by our method can be any general form of the quadratic hypersurface. In our experiment, the results of five artificial datasets showed that the proposed method can find the general form of the quadratic surface and has also the ability to obtain the linear separating hyperplane. Numerical experiments on 14 benchmark datasets verified that the proposed method was superior to corresponding methods in both accuracy and CPU time. What is more gratifying is that when the number of samples or the dimension is relatively large, our method can obtain good classification performance quickly. In addition, the results of the Friedman test and Nemenyi post-hoc test indicated that our QSMPM was statistically the best one compared to other methods.

The rest of this paper is organized as follows. Section 2 briefly reviews the related works, the QSSVM, and the MPM. Section 3 presents our method QSMPM, gives its algorithm, and analyzes the computational complexity of the QSMPM. In Section 4, we show the interpretability of our method. In Section 5, the results of the numerical experiments on the artificial datasets and benchmark datasets are presented, and a further statistical analysis is performed. Finally, Section 6 gives the conclusion and future work of this paper.

Throughout this paper, we use lower case letters to represent scalars, lower case bold letters to represent vectors, and upper case bold letters to represent matrices.  $\mathbb{R}$  denotes

the set of real numbers.  $\mathbb{R}^d$  denotes the space of  $d$ -dimensional vectors.  $\mathbb{R}^{d \times d}$  denotes the space of  $d \times d$  matrices.  $\mathbb{S}^d$  denotes the set of  $d \times d$  symmetric matrices.  $\mathbb{S}_+^d$  denotes the set of  $d \times d$  symmetric positive semidefinite matrices.  $\mathbf{I}_d$  denotes the  $d \times d$  identity matrix.  $\|\mathbf{x}\|_2$  denotes the two-norm of the vector  $\mathbf{x}$ .

## 2. Related Work

In this section, we briefly introduce the QSSVM and the MPM. For a binary classification problem, the training set is given as:

$$T = \{(\mathbf{x}_1, y_1), (\mathbf{x}_2, y_2), \dots, (\mathbf{x}_{m_+ + m_-}, y_{m_+ + m_-})\}, \quad (1)$$

where  $\mathbf{x}_i \in \mathbb{R}^d$  is the  $i$ -th sample and  $y_i \in \{+1, -1\}$  corresponds to the class label,  $i = 1, 2, \dots, m_+ + m_-$ . The number of samples in class +1 and class -1 is  $m_+$  and  $m_-$ , respectively. For the training set (1), we want to find a hyperplane or quadratic hypersurface:

$$g(\mathbf{x}) = 0, \quad (2)$$

and then use a decision function:

$$f(\mathbf{x}) = \text{sign}(g(\mathbf{x})) \quad (3)$$

to determine whether a new sample  $\mathbf{x} \in \mathbb{R}^d$  is assigned to class +1 or class -1.

### 2.1. Quadratic Surface Support Vector Machine

We first shortly outline the quadratic surface support vector machine (QSSVM) [26]. For the given training set (1), the goal of the QSSVM is to seek a quadratic separating hypersurface:

$$g(\mathbf{x}) = \frac{1}{2} \mathbf{x}^T \mathbf{A} \mathbf{x} + \mathbf{b}^T \mathbf{x} + c = 0, \quad (4)$$

where  $\mathbf{A} \in \mathbb{S}^d$ ,  $\mathbf{b} \in \mathbb{R}^d$ ,  $c \in \mathbb{R}$ , which separates the samples into two classes with the largest margin. In order to obtain the quadratic hypersurface (4), the QSSVM establishes the following optimization problem:

$$\begin{aligned} \min_{\mathbf{A}, \mathbf{b}, c} \quad & \sum_{i=1}^{m_+ + m_-} \|\mathbf{A} \mathbf{x}_i + \mathbf{b}\|_2^2 \\ \text{s.t.} \quad & y_i \left( \frac{1}{2} \mathbf{x}_i^T \mathbf{A} \mathbf{x}_i + \mathbf{b}^T \mathbf{x}_i + c \right) \geq 1, i = 1, \dots, m_+ + m_-. \end{aligned} \quad (5)$$

The optimization problem (5) is a convex quadratic programming problem.

After obtaining the optimal solution  $\mathbf{A}_*$ ,  $\mathbf{b}_*$ , and  $c_*$  to the optimization problem (5), for a given new sample  $\mathbf{x} \in \mathbb{R}^d$ , its label is assigned to either class +1 or class -1 by the decision function:

$$f(\mathbf{x}) = \text{sgn}\left(\frac{1}{2} \mathbf{x}^T \mathbf{A}_* \mathbf{x} + \mathbf{b}_*^T \mathbf{x} + c_*\right). \quad (6)$$

To allow some samples in the training set (1) to be misclassified, Luo et al. further proposed the soft-margin quadratic surface support vector machine (SQSSVM); please refer to [27].

### 2.2. Minimax Probability Machine

Now, we briefly review the minimax probability machine (MPM) [8,9]. Let us leave the training set (1) aside for a moment and suppose that these samples have some distribution. Specifically, assume that the samples in class +1 are drawn from a distribution with the mean vector  $\boldsymbol{\mu}_+ \in \mathbb{R}^d$  and the covariance matrix  $\boldsymbol{\Sigma}_+ \in \mathbb{S}_+^d$ , without making other specific distributional assumptions. A similar assumption is also given for the samples in class -1 with the mean vector  $\boldsymbol{\mu}_- \in \mathbb{R}^d$  and the covariance matrix  $\boldsymbol{\Sigma}_- \in \mathbb{S}_+^d$ . Denote the two

distributions as  $x_+ \sim (\boldsymbol{\mu}_+, \boldsymbol{\Sigma}_+)$  and  $x_- \sim (\boldsymbol{\mu}_-, \boldsymbol{\Sigma}_-)$ , respectively. Based on the above assumptions, the MPM attempts to obtain a separating hyperplane:

$$g(x) = \boldsymbol{w}^T \boldsymbol{x} - b = 0, \quad (7)$$

where  $\boldsymbol{w} \in \mathbb{R}^d$ ,  $b \in \mathbb{R}$ , which separates the two classes of samples with maximal probability with respect to all distributions having these mean vectors and covariance matrices. This is expressed as:

$$\begin{aligned} \max_{\boldsymbol{w}, b, \alpha} \quad & \alpha \\ \text{s.t.} \quad & \inf_{x_+ \sim (\boldsymbol{\mu}_+, \boldsymbol{\Sigma}_+)} \text{pr}\{\boldsymbol{w}^T \boldsymbol{x}_+ - b \geq 0\} \geq \alpha, \\ & \inf_{x_- \sim (\boldsymbol{\mu}_-, \boldsymbol{\Sigma}_-)} \text{pr}\{\boldsymbol{w}^T \boldsymbol{x}_- - b \leq 0\} \geq \alpha, \end{aligned} \quad (8)$$

where  $\alpha \in (0, 1)$  represents the lower bound of the accuracy for future data, namely the worst-case accuracy. The infimum “inf” is taken over all distributions having these mean vectors  $\boldsymbol{\mu}_\pm \in \mathbb{R}^d$  and covariance matrices  $\boldsymbol{\Sigma}_\pm \in \mathbb{S}_+^d$ .

The constraint condition of the above optimization problem (8) is the probabilistic constraint, which is difficult to solve. In order to convert the probabilistic constraints to easy, tractable constraints, the following lemma [9] is given:

**Lemma 1** ([9]). *Let  $\boldsymbol{x}$  be a  $d$ -dimensional random vector with mean vector  $\boldsymbol{\mu}$  and covariance matrix  $\boldsymbol{\Sigma}$ , where  $\boldsymbol{\Sigma} \in \mathbb{S}_+^d$ . Given  $\boldsymbol{w} \in \mathbb{R}^d$ ,  $b \in \mathbb{R}$ , such that  $\boldsymbol{w}^T \boldsymbol{x} \leq b$  and  $\alpha \in (0, 1)$ , the condition:*

$$\inf_{x \sim (\boldsymbol{\mu}, \boldsymbol{\Sigma})} \text{pr}\{\boldsymbol{w}^T \boldsymbol{x} - b \leq 0\} \geq \alpha \quad (9)$$

holds if and only if:

$$b - \boldsymbol{w}^T \boldsymbol{\mu} \geq \kappa(\alpha) \sqrt{\boldsymbol{w}^T \boldsymbol{\Sigma} \boldsymbol{w}}, \quad (10)$$

where  $\kappa(\alpha) = \sqrt{\frac{\alpha}{1-\alpha}}$ .

Using the above Lemma 1, the optimization problem (8) is equivalent to:

$$\begin{aligned} \max_{\boldsymbol{w}, b, \alpha} \quad & \alpha \\ \text{s.t.} \quad & -b + \boldsymbol{w}^T \boldsymbol{\mu}_+ \geq \kappa(\alpha) \sqrt{\boldsymbol{w}^T \boldsymbol{\Sigma}_+ \boldsymbol{w}}, \\ & b - \boldsymbol{w}^T \boldsymbol{\mu}_- \geq \kappa(\alpha) \sqrt{\boldsymbol{w}^T \boldsymbol{\Sigma}_- \boldsymbol{w}}. \end{aligned} \quad (11)$$

Then, through a series of algebraic operations (see Theorem 2 in [9], for the details), the above optimization problem (11) leads to:

$$\begin{aligned} \min_{\boldsymbol{w}} \quad & \left\| \boldsymbol{\Sigma}_+^{\frac{1}{2}} \boldsymbol{w} \right\|_2 + \left\| \boldsymbol{\Sigma}_-^{\frac{1}{2}} \boldsymbol{w} \right\|_2 \\ \text{s.t.} \quad & \boldsymbol{w}^T (\boldsymbol{\mu}_+ - \boldsymbol{\mu}_-) = 1. \end{aligned} \quad (12)$$

When its optimal solution  $\boldsymbol{w}_*$  is obtained, for the optimization problem (11), the optimal solution with respect to  $b$  is given by:

$$b_* = \boldsymbol{w}_*^T \boldsymbol{\mu}_- + \frac{\left\| \boldsymbol{\Sigma}_-^{\frac{1}{2}} \boldsymbol{w}_* \right\|_2}{\left\| \boldsymbol{\Sigma}_+^{\frac{1}{2}} \boldsymbol{w}_* \right\|_2 + \left\| \boldsymbol{\Sigma}_-^{\frac{1}{2}} \boldsymbol{w}_* \right\|_2}, \quad (13)$$

or:

$$b_* = \mathbf{w}_*^T \boldsymbol{\mu}_+ - \frac{\left\| \boldsymbol{\Sigma}_+^{\frac{1}{2}} \mathbf{w}_* \right\|_2}{\left\| \boldsymbol{\Sigma}_+^{\frac{1}{2}} \mathbf{w}_* \right\|_2 + \left\| \boldsymbol{\Sigma}_-^{\frac{1}{2}} \mathbf{w}_* \right\|_2}. \quad (14)$$

Now, let us return to the training set (1). It is easy to see that these required mean vectors  $\boldsymbol{\mu}_\pm \in \mathbb{R}^d$  and covariance matrices  $\boldsymbol{\Sigma}_\pm \in \mathbb{S}_+^d$  are able to be estimated by the training set (1) as follows:

$$\begin{aligned} \hat{\boldsymbol{\mu}}_\pm &= \frac{1}{m_\pm} \sum_{i=1}^{m_\pm} \mathbf{x}_i \in \mathbb{R}^d, \\ \hat{\boldsymbol{\Sigma}}_\pm &= \frac{1}{m_\pm} \sum_{i=1}^{m_\pm} (\mathbf{x}_i - \hat{\boldsymbol{\mu}}_\pm)(\mathbf{x}_i - \hat{\boldsymbol{\mu}}_\pm)^T \in \mathbb{S}_+^d. \end{aligned} \quad (15)$$

Therefore, in practice, these mean vectors  $\boldsymbol{\mu}_\pm$  and covariance matrices  $\boldsymbol{\Sigma}_\pm$  in (12)–(14) should be replaced by  $\hat{\boldsymbol{\mu}}_\pm$  and  $\hat{\boldsymbol{\Sigma}}_\pm$ , and the optimal solutions of  $\mathbf{w}$  and  $\mathbf{b}$  thus obtained are denoted as  $\hat{\mathbf{w}}_*$  and  $\hat{\mathbf{b}}_*$ . Then, for a given new sample  $\mathbf{x} \in \mathbb{R}^d$ , its label is assigned to either class +1 or class –1 by the decision function:

$$f(\mathbf{x}) = \text{sgn}(\hat{\mathbf{w}}_*^T \mathbf{x} - \hat{\mathbf{b}}_*). \quad (16)$$

In addition, for nonlinear cases and more details, please refer to [8,9].

### 3. Kernel-Free Quadratic Surface Minimax Probability Machine

In this section, we first formulate the kernel-free quadratic surface minimax probability machine (QSMPM). Then, its algorithm is given.

#### 3.1. Optimization Problem

For the binary classification problem with the training set (1), we attempt to find a quadratic separating hypersurface:

$$g(\mathbf{x}) = \frac{1}{2} \mathbf{x}^T \mathbf{A} \mathbf{x} + \mathbf{b}^T \mathbf{x} - c = 0, \quad (17)$$

where  $\mathbf{A} \in \mathbb{S}^d$ ,  $\mathbf{b} \in \mathbb{R}^d$ ,  $c \in \mathbb{R}$ , which separates the two classes of the samples. Inspired by the MPM, we construct the following optimization problem:

$$\begin{aligned} \max_{\mathbf{A}, \mathbf{b}, c, \alpha} \quad & \alpha \\ \text{s.t.} \quad & \inf_{\mathbf{x}_+ \sim (\boldsymbol{\mu}_+, \boldsymbol{\Sigma}_+)} \text{pr}\left\{ \frac{1}{2} \mathbf{x}_+^T \mathbf{A} \mathbf{x}_+ + \mathbf{b}^T \mathbf{x}_+ - c \geq 0 \right\} \geq \alpha, \\ & \inf_{\mathbf{x}_- \sim (\boldsymbol{\mu}_-, \boldsymbol{\Sigma}_-)} \text{pr}\left\{ \frac{1}{2} \mathbf{x}_-^T \mathbf{A} \mathbf{x}_- + \mathbf{b}^T \mathbf{x}_- - c \leq 0 \right\} \geq \alpha, \end{aligned} \quad (18)$$

where  $\alpha \in (0, 1)$  represents the lower bound of the accuracy for future data, namely the worst-case accuracy. The notation  $\mathbf{x}_+ \sim (\boldsymbol{\mu}_+, \boldsymbol{\Sigma}_+)$  refers to the class distribution that has the prescribed mean vector  $\boldsymbol{\mu}_+ \in \mathbb{R}^d$  and covariance matrix  $\boldsymbol{\Sigma}_+ \in \mathbb{S}_+^d$ , but otherwise arbitrary, and likewise for  $\mathbf{x}_-$ .

The above optimization problem (18) corresponds to the optimization problem (8), which is used to derive the optimization problem (11). Analogically, the optimization problem (18) should be used to derive the required optimization problem. Unfortunately, it does not have a counterpart when the functions in curly braces in the optimization problem (18) are quadratic because of the lack of the corresponding Lemma 1. In order

to overcome this difficulty, we change the quadratic functions as a linear function by introducing a nonlinear transformation from  $\mathbf{x} = ([x]_1, [x]_2, \dots, [x]_d)^T \in \mathbb{R}^d$  to:

$$\mathbf{z} = \mathbf{z}(\mathbf{x}) = \left( \frac{1}{2}[x]_1^2, [x]_1[x]_2, \dots, [x]_1[x]_d, \frac{1}{2}[x]_2^2, \dots, [x]_2[x]_d, \dots, \frac{1}{2}[x]_d^2, [x]_1, [x]_2, \dots, [x]_d \right)^T \in \mathbb{R}^{\frac{d^2+3d}{2}}. \quad (19)$$

By representing the upper triangular entries of the symmetric matrix:

$$\mathbf{A} = \mathbf{A}^T = (a_{ij})_{d \times d} \in \mathbb{S}^d \quad (20)$$

as a vector:

$$\mathbf{a} = (a_{11}, a_{12}, \dots, a_{1d}, a_{22}, \dots, a_{2d}, \dots, a_{dd})^T \in \mathbb{R}^{\frac{d^2+d}{2}}, \quad (21)$$

and defining:

$$\mathbf{w} = (\mathbf{a}^T, \mathbf{b}^T)^T \in \mathbb{R}^{\frac{d^2+3d}{2}}, \quad (22)$$

the quadratic function (17) of  $\mathbf{x}$  in  $d$ -dimensional space yields the linear function of  $\mathbf{z}$  in  $\frac{d^2+3d}{2}$ -dimensional space as follows:

$$g(\mathbf{x}) = \frac{1}{2}\mathbf{x}^T \mathbf{A} \mathbf{x} + \mathbf{b}^T \mathbf{x} - c = \mathbf{w}^T \mathbf{z} - c. \quad (23)$$

Following the transformation (19), the training set (1) in the  $d$ -dimensional space correspondingly becomes:

$$\tilde{T} = \{(z_1, y_1), (z_2, y_2), \dots, (z_{m_+ + m_-}, y_{m_+ + m_-})\}, \quad (24)$$

where  $z_i = z(\mathbf{x}_i)$ ; in other words,  $z_i \in \mathbb{R}^{\frac{d^2+3d}{2}}$  is obtained by replacing  $\mathbf{x}$  in the formula (19) with  $\mathbf{x}_i = ([x_i]_1, [x_i]_2, \dots, [x_i]_d)^T \in \mathbb{R}^d$ ,  $i = 1, 2, \dots, m_+ + m_-$ . For the training set (24), it is naturally assumed that the samples of the two classes are sampled from  $\mathbf{z}_+ \sim (\boldsymbol{\mu}_{z_+}, \boldsymbol{\Sigma}_{z_+})$  and  $\mathbf{z}_- \sim (\boldsymbol{\mu}_{z_-}, \boldsymbol{\Sigma}_{z_-})$ , respectively, where these mean vectors  $\boldsymbol{\mu}_{z_{\pm}} \in \mathbb{R}^{\frac{d^2+3d}{2}}$  and covariance matrices  $\boldsymbol{\Sigma}_{z_{\pm}} \in \mathbb{S}_+^{\frac{d^2+3d}{2}}$  can be estimated as:

$$\begin{aligned} \hat{\boldsymbol{\mu}}_{z_{\pm}} &= \frac{1}{m_{\pm}} \sum_{i=1}^{m_{\pm}} z_i \in \mathbb{R}^{\frac{d^2+3d}{2}}, \\ \hat{\boldsymbol{\Sigma}}_{z_{\pm}} &= \frac{1}{m_{\pm}} \sum_{i=1}^{m_{\pm}} (z_i - \hat{\boldsymbol{\mu}}_{z_{\pm}})(z_i - \hat{\boldsymbol{\mu}}_{z_{\pm}})^T \in \mathbb{S}_+^{\frac{d^2+3d}{2}}. \end{aligned} \quad (25)$$

Based on the transformation (19), the optimization problem (18) is replaced by:

$$\begin{aligned} \max_{\mathbf{w}, c, \alpha} \quad & \alpha \\ \text{s.t.} \quad & \inf_{\mathbf{z}_+ \sim (\hat{\boldsymbol{\mu}}_{z_+}, \hat{\boldsymbol{\Sigma}}_{z_+})} \text{pr}\{\mathbf{w}^T \mathbf{z}_+ - c \geq 0\} \geq \alpha, \\ & \inf_{\mathbf{z}_- \sim (\hat{\boldsymbol{\mu}}_{z_-}, \hat{\boldsymbol{\Sigma}}_{z_-})} \text{pr}\{\mathbf{w}^T \mathbf{z}_- - c \leq 0\} \geq \alpha. \end{aligned} \quad (26)$$

Now, Lemma 1 [9] is applicable to the optimization problem (26). Thus, we have:

$$\begin{aligned} \max_{\mathbf{w}, c, \alpha} \quad & \alpha \\ \text{s.t.} \quad & -c + \mathbf{w}^T \hat{\boldsymbol{\mu}}_{z_+} \geq \kappa(\alpha) \sqrt{\mathbf{w}^T \hat{\boldsymbol{\Sigma}}_{z_+} \mathbf{w}}, \\ & c - \mathbf{w}^T \hat{\boldsymbol{\mu}}_{z_-} \geq \kappa(\alpha) \sqrt{\mathbf{w}^T \hat{\boldsymbol{\Sigma}}_{z_-} \mathbf{w}}, \end{aligned} \quad (27)$$

where  $\kappa(\alpha) = \sqrt{\frac{\alpha}{1-\alpha}}$ . Moreover, a series of algebraic operation shows that the above optimization problem (27) is equivalent to the following second-order cone programming problem:

$$\begin{aligned} \min_w \quad & \left\| \hat{\Sigma}_{z+}^{\frac{1}{2}} w \right\|_2 + \left\| \hat{\Sigma}_{z-}^{\frac{1}{2}} w \right\|_2 \\ \text{s.t.} \quad & w^T (\hat{\mu}_{z+} - \hat{\mu}_{z-}) = 1. \end{aligned} \quad (28)$$

When its optimal solution  $w_*$  is obtained, for the optimization problem (27), the optimal solution with respect to  $c$  is given by:

$$c_* = w_*^T \hat{\mu}_{z-} + \frac{\left\| \hat{\Sigma}_{z-}^{\frac{1}{2}} w_* \right\|_2}{\left\| \hat{\Sigma}_{z+}^{\frac{1}{2}} w_* \right\|_2 + \left\| \hat{\Sigma}_{z-}^{\frac{1}{2}} w_* \right\|_2}, \quad (29)$$

or:

$$c_* = w_*^T \hat{\mu}_{z+} - \frac{\left\| \hat{\Sigma}_{z+}^{\frac{1}{2}} w_* \right\|_2}{\left\| \hat{\Sigma}_{z+}^{\frac{1}{2}} w_* \right\|_2 + \left\| \hat{\Sigma}_{z-}^{\frac{1}{2}} w_* \right\|_2}. \quad (30)$$

In the next subsection, we show how to solve the optimization problem (28).

### 3.2. Algorithm

Now, we present the solving process of the optimization problem (28), which is achieved by referring to [9]. By constructing an orthogonal matrix  $F \in \mathbb{R}^{\frac{d^2+3d}{2} \times \frac{d^2+3d-2}{2}}$  whose columns span the subspace of vectors orthogonal to  $\hat{\mu}_{z+} - \hat{\mu}_{z-} \in \mathbb{R}^{\frac{d^2+3d}{2}}$ , the unknown variable  $w \in \mathbb{R}^{\frac{d^2+3d}{2}}$  is converted into  $u \in \mathbb{R}^{\frac{d^2+3d-2}{2}}$ . Specifically, let  $w = w_0 + Fu$ , where  $w_0 = \frac{(\hat{\mu}_{z+} - \hat{\mu}_{z-})}{\|\hat{\mu}_{z+} - \hat{\mu}_{z-}\|_2} \in \mathbb{R}^{\frac{d^2+3d}{2}}$ ; the optimization problem (28) is transferred to the unconstrained optimization problem:

$$\min_u \left\| \hat{\Sigma}_{z+}^{\frac{1}{2}} (w_0 + Fu) \right\|_2 + \left\| \hat{\Sigma}_{z-}^{\frac{1}{2}} (w_0 + Fu) \right\|_2, \quad (31)$$

In order to solve the above optimization problem (31), Lanckriet et al. [9] introduced two extra variables  $\beta$  and  $\eta$  and considered the following optimization problem:

$$\min_{u, \beta, \eta} \beta + \frac{1}{\beta} \left\| \hat{\Sigma}_{z+}^{\frac{1}{2}} (w_0 + Fu) \right\|_2^2 + \eta + \frac{1}{\eta} \left\| \hat{\Sigma}_{z-}^{\frac{1}{2}} (w_0 + Fu) \right\|_2^2. \quad (32)$$

This optimization problem (32) is solved by an alternative iteration. The variables are divided into two sets: one is  $\beta$  and  $\eta$ , and the other is  $u$ . At the  $t$ -th iteration, first by fixing  $\beta$  and  $\eta$  to take the derivative of the optimization problem (32) with respect to  $u$ , we have the following updated iteration formula of  $u_t$ :

$$\left( \frac{1}{\beta_t} P + \frac{1}{\eta_t} Q \right) u_t = - \left( \frac{1}{\beta_t} p + \frac{1}{\eta_t} q \right), \quad (33)$$

where  $P = F^T \hat{\Sigma}_{z+} F \in \mathbb{R}^{\frac{d^2+3d-2}{2} \times \frac{d^2+3d-2}{2}}$ ,  $Q = F^T \hat{\Sigma}_{z-} F \in \mathbb{R}^{\frac{d^2+3d-2}{2} \times \frac{d^2+3d-2}{2}}$ ,  $p = F^T \hat{\Sigma}_{z+} w_0 \in \mathbb{R}^{\frac{d^2+3d-2}{2}}$ ,  $q = F^T \hat{\Sigma}_{z-} w_0 \in \mathbb{R}^{\frac{d^2+3d-2}{2}}$ . To ensure the stability, the regularization term  $\delta I_{\frac{d^2+3d-2}{2}}$  ( $\delta > 0$ ) is added. Therefore, the Equation (33) is replaced by:

$$\left( \frac{1}{\beta_t} P + \frac{1}{\eta_t} Q + \delta I \right) u_t = - \left( \frac{1}{\beta_t} p + \frac{1}{\eta_t} q \right). \quad (34)$$

Next, by fixing  $\mathbf{u}$  to take the derivative of the optimization problem (32) with respect to  $\beta$  and  $\eta$ , respectively, we have the following updated iteration formula of  $\beta_t$  and  $\eta_t$ :

$$\beta_t = \left\| \hat{\Sigma}_{g_+}^{\frac{1}{2}} (\mathbf{w}_0 + \mathbf{F}\mathbf{u}_t) \right\|_2, \eta_t = \left\| \hat{\Sigma}_{g_-}^{\frac{1}{2}} (\mathbf{w}_0 + \mathbf{F}\mathbf{u}_t) \right\|_2. \quad (35)$$

When the optimal solution  $\mathbf{u}_*$  is obtained by the above two updated iteration Formulas (34) and (35), the optimal solution  $\mathbf{w}_*$  of the optimization problem (28) is  $\mathbf{w}_* = \mathbf{w}_0 + \mathbf{F}\mathbf{u}_*$ . Then, we summarize the process of finding the optimal solution  $A_*$ ,  $\mathbf{b}_*$ ,  $c_*$  of the optimization problem (18) in Algorithm 1.

---

**Algorithm 1** Kernel-free quadratic surface minimax probability machine (QSMPM).

---

**Input:** Training set (1),  $\delta = 1 \times 10^{-6}$ , number of maximum iterations  $\tau = 100$ .

- 1: Initialize  $\beta_1 = 1, \eta_1 = 1, t = 1$ .
- 2: Obtain  $\mathbf{z}_i$  by (19),  $i = 1, 2, \dots, m_+ + m_-$ ;
- 3: Calculate  $\hat{\mu}_{z_{\pm}}$  and  $\hat{\Sigma}_{z_{\pm}}$  by (25), and calculate  $\mathbf{w}_0 = \frac{(\hat{\mu}_{z_+} - \hat{\mu}_{z_-})}{\left\| \hat{\mu}_{z_+} - \hat{\mu}_{z_-} \right\|_2}$ ;
- 4: Calculate  $\mathbf{P} = \mathbf{F}^T \hat{\Sigma}_{z_+} \mathbf{F}$ ,  $\mathbf{Q} = \mathbf{F}^T \hat{\Sigma}_{z_-} \mathbf{F}$ ,  $\mathbf{p} = \mathbf{F}^T \hat{\Sigma}_{z_+} \mathbf{w}_0$ ,  $\mathbf{q} = \mathbf{F}^T \hat{\Sigma}_{z_-} \mathbf{w}_0$ , where  $\mathbf{F}$  is an orthogonal matrix whose columns span the subspace of vectors orthogonal  $\hat{\mu}_{z_+} - \hat{\mu}_{z_-}$ ;
- 5: **while**  $t < \tau$  **do**
- 6: Given  $\beta_t$  and  $\eta_t$ , update  $\mathbf{u}_t$  by (34);
- 7: Given  $\mathbf{u}_t$ , update  $\beta_t$  and  $\eta_t$  by (35);
- 8:  $t \leftarrow t + 1$ .
- 9: **end**
- 10: Assign  $\mathbf{w}_* = \mathbf{w}_0 + \mathbf{F}\mathbf{u}_t$ , then obtain  $A_*$ ,  $\mathbf{b}_*$  by (20), (21), and (22); further, obtain  $c_*$  by (29) or (30).

**Output:**  $A_*$ ,  $\mathbf{b}_*$ ,  $c_*$ .

---

After obtaining the optimal solution  $A_*$ ,  $\mathbf{b}_*$  and  $c_*$  to the optimization problem (18), for a given new sample  $\mathbf{x} \in \mathbb{R}^d$ , its label is assigned to either class +1 or class -1 by the decision function:

$$f(\mathbf{x}) = \text{sgn}\left(\frac{1}{2}\mathbf{x}^T \mathbf{A}_* \mathbf{x} + \mathbf{b}_*^T \mathbf{x} - c_*\right). \quad (36)$$

It should be pointed out that our QSMPM is kernel-free, which avoids the time-consuming task of selecting the appropriate kernel function and its corresponding parameters. What is more, it does not require any choice of parameter, which makes its use simpler and convenient. Furthermore, from the geometric point of view, the quadratic hypersurface (17) determined by our method is allowed to be any general form of quadratic hypersurface, including hyperplanes, hyperparaboloids, hyperspheres, hyperellipsoids, hyperhyperboloids, and so on, which is shown clearly by five artificial examples in Section 5.

### 3.3. Computational Complexity

Here, we analyze the computational complexity of our QSMPM. Suppose that the number and the dimension of the samples are  $N$  and  $d$ , respectively. Before reformulating the QSMPM as an SOCP problem, all  $d$ -dimensional samples need to be projected into the  $\frac{d^2+3d}{2}$ -dimensional space. Therefore, the total computational complexity of the QSMPM is  $O\left(\left(\frac{d^2+3d}{2}\right)^3 + N\left(\frac{d^2+3d}{2}\right)^2 + Nd^2\right)$ . In addition, we give the computational complexity of the MPM and the SVM. Their complexity is  $O(d^3 + Nd^2)$  [9] and  $O(N^3)$  [19], respectively. Then, by referencing the computational complexity of SVM, we obtain that the computational complexity of the QSSVM is  $O(N^3 + Nd^2)$ . According to the above analysis, assuming that  $N$  is much larger than  $d$ , we can see that the computational complexity of the QSMPM is higher than that of the MPM, but lower than that of the SVM and the QSSVM.

#### 4. The Interpretability

In this section, we discuss the interpretability of our method QSMPM. Suppose we have obtained the optimal solution  $A_*$ ,  $b_*$ ,  $c_*$  to the optimization problem (18), then the quadratic hypersurface (17) has the following component form:

$$g(\mathbf{x}) = \frac{1}{2}\mathbf{x}^T A_* \mathbf{x} + \mathbf{b}_*^T \mathbf{x} - c_* = \frac{1}{2} \sum_{i=1}^d \sum_{j=1}^d a_{ij}^* [x]_i [x]_j + \sum_{i=1}^d b_i^* [x]_i - c_* = 0, \quad (37)$$

where  $[x]_i$  is the  $i$ -th component of the vector  $\mathbf{x} \in \mathbb{R}^d$ ,  $a_{ij}^*$  is the  $i$ -th row and the  $j$ -th column component of the matrix  $A_* \in \mathbb{S}^d$ , and  $b_i^*$  is the  $i$ -th component of the vector  $\mathbf{b}_* \in \mathbb{R}^d$ . Each component of  $\mathbf{x}$  produces the contribution of a quadratic polynomial function. Specifically,  $b_i^*$  is the linear effect coefficient of the  $i$ -th component,  $a_{ii}^*$  ( $i = j$ ) is the quadratic effect coefficient of the  $i$ -th component, and  $a_{ij}^*$  ( $i \neq j$ ) is the interaction coefficient between the  $i$ -th component and the  $j$ -th component. Therefore, for the  $i$ -th component of  $\mathbf{x}$ , consider that the larger  $|a_{ii}^*| + |a_{ij}^*| + |b_i^*|$  ( $j = 1, 2, \dots, d, j \neq i$ ), the greater the contribution of the  $i$ -th component is. In particular, when  $|a_{ii}^*| + |a_{ij}^*| + |b_i^*| = 0$  ( $j = 1, 2, \dots, d, j \neq i$ ), the  $i$ -th component of  $\mathbf{x}$  would not work. Therefore, compared with the methods with the kernel function, the QSMPM has better interpretability.

#### 5. Numerical Experiments

In this section, we provide some numerical experiments to verify the performance of our QSMPM. We compared it with the hard-margin support vector machine (H-SVM), the soft-margin support vector machine (S-SVM), and the MPM with the linear kernel, the quadratic polynomial kernel, and the RBF kernel (H-SVM-L, H-SVM-P, H-SVM-R, S-SVM-L, S-SVM-P, S-SVM-R, MPM-L, MPM-P, and MPM-R, respectively). In addition, we also compared it with the QSSVM and the SQSSVM. In all numerical experiments, the penalty parameter  $C$  in the S-SVM and the kernel parameter  $\sigma$  of the RBF kernel were selected from  $\{2^{-7}, 2^{-6}, \dots, 2^7\}$  by the 10-fold cross-validation method. All numerical experiments were conducted using MATLAB R2016 (b) on a computer equipped with a 2.50 GHz (I5-4210U) CPU, and 4G available memory.

##### 5.1. Artificial Datasets

To show the geometric interpretation of the proposed method QSMPM and to compare it with the original methods, the MPM-L, the MPM-P, and the MPM-R, we performed the following numerical experiments on 4 artificial examples. These 5 artificial examples were all generated in the 2-dimensional space, and each generated 300 samples  $\{\mathbf{x}_i = ([x_i]_1, [x_i]_2)^T\}_{i=1}^{300}$ , of which the first 150 were samples in class +1 and the last 150 samples in class -1. Here, we first illustrate the symbols in all figures. The red "+" represents the samples in class +1, and the blue "o" represents the samples in class -1. The results in the upper right express the accuracy of each method on the artificial example. The curve in bold black represents the hyperplane or quadratic hypersurface. Now, let us introduce the numerical experiments on each artificial example in turn.

##### Example 1.

$$\begin{aligned} [x_i]_2 &= -\frac{1}{2}[x_i]_1 + 2 + \zeta_i, i = 1, \dots, 150, \\ [x_i]_2 &= -\frac{1}{2}[x_i]_1 - 3 + \zeta_i, i = 151, \dots, 300, \end{aligned}$$

where  $[x_i]_1 \sim U[-3, 4]$ ,  $\zeta_i \sim N(0, 1)$ .

Figure 1 illustrates the classification results of the MPM-L, the MPM-P, the MPM-R, and the QSMPM on Example 1, respectively. We can see from Figure 1 that our QSMPM can

obtain classification results as good as the other three methods. In addition, the quadratic hypersurface found by our QSMPM is a straight line, that is a linear separating hyperplane.

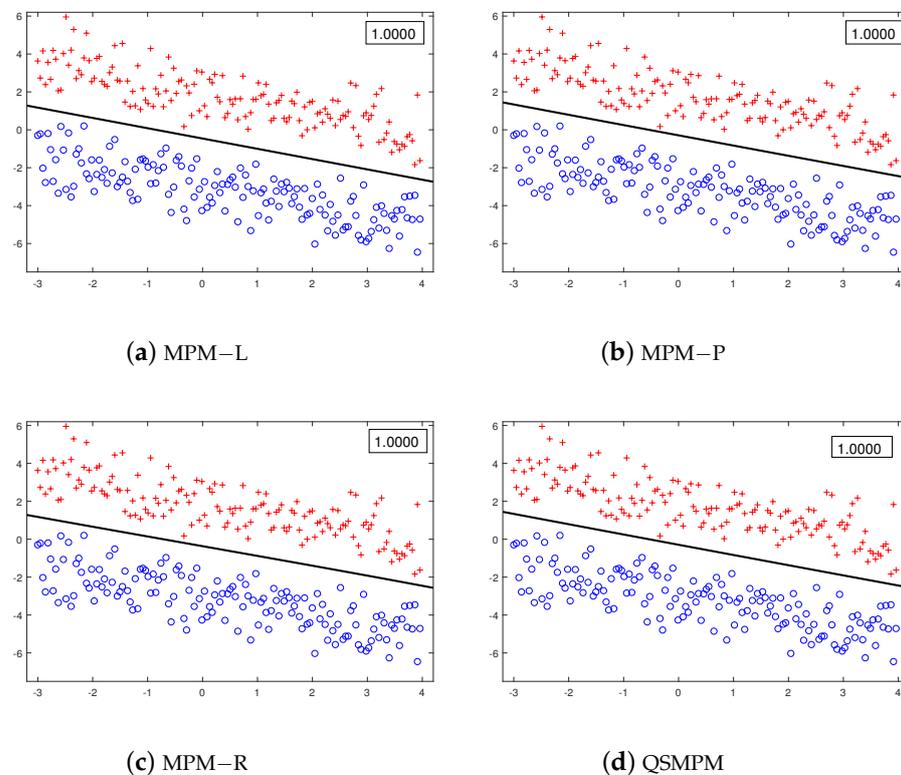


Figure 1. Classification results on Example 1.

### Example 2.

$$[x_i]_2 = \frac{1}{2}[x_i]_1^2 + 2 + \xi_i, i = 1, \dots, 150,$$

$$[x_i]_2 = \frac{1}{2}[x_i]_1^2 - 3 + \xi_i, i = 151, \dots, 300,$$

where  $[x_i]_1 \sim U[-3, 4]$ ,  $\xi_i \sim N(0, 1)$ .

Figure 2 presents the classification results on Example 2. We can observe in Figure 2 that the classification result of our QSMPM is superior to the MPM-L and similar to the MPM-P and the MPM-R. Moreover, the quadratic hypersurface obtained by our QSMPM is a parabola.

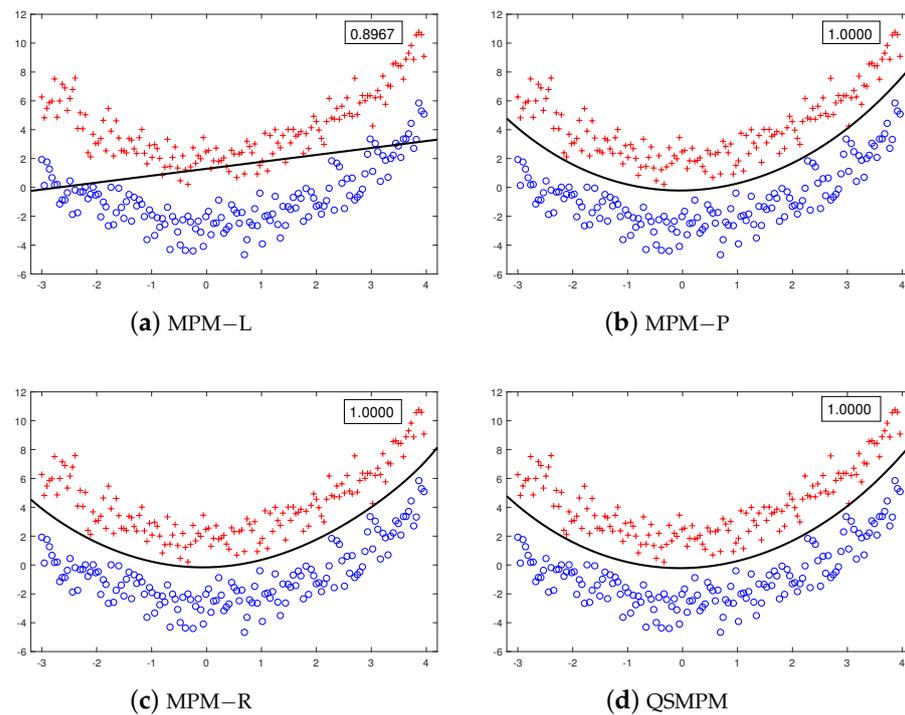
### Example 3.

$$[x_i]_1 = r_i \cos(t_i), [x_i]_2 = r_i \sin(t_i), i = 1, \dots, 150,$$

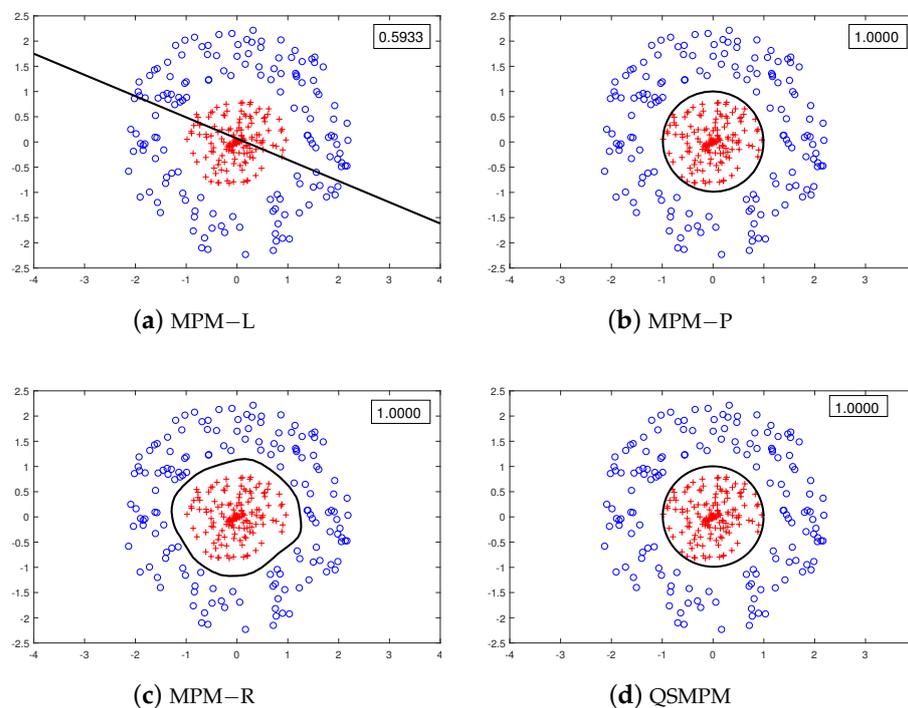
$$[x_i]_1 = 1.3 + r_i \cos(t_i), [x_i]_2 = 1.3 + r_i \sin(t_i), i = 151, \dots, 300,$$

where  $r_i \sim U[0, 1]$ ,  $t_i \sim U[0, 1]$ .

Figure 3 reports the classification results of the MPM-L, the MPM-P, the MPM-R, and the QSMPM on Example 3, respectively. Obviously, in Figure 3, we can see that the classification result of our QSMPM is superior to the MPM-L and is the same as the MPM-P and the MPM-R. Furthermore, the quadratic hypersurface of our QSMPM is a circle.



**Figure 2.** Classification results on Example 2.



**Figure 3.** Classification results on Example 3.

#### Example 4.

$$[x_i]_1 = 2 + a_i \cos(t_i), [x_i]_2 = b_i \sin(t_i), i = 1, \dots, 150,$$

$$[x_i]_1 = 3.5 + a_i \cos(t_i), [x_i]_2 = 1.5 + b_i \cos(t_i), i = 151, \dots, 300,$$

where  $a_i \sim U[0, 1]$ ,  $b_i \sim U[0, 1]$ ,  $t_i \sim U[0, 1]$ .

Figure 4 shows the classification results on Example 4. We can observe in Figure 4 that the QSMPM can obtain the same classification performance as the MPM-P and the MPM-R and is better than the MPM-L. Our QSMPM can find an ellipse.

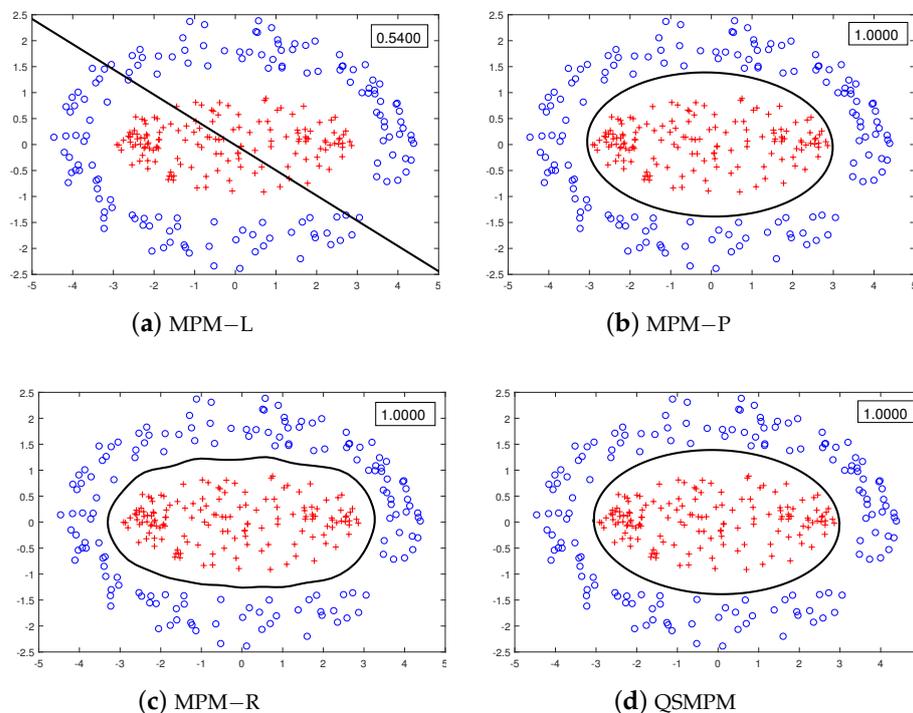


Figure 4. Classification results on Example 4.

#### Example 5.

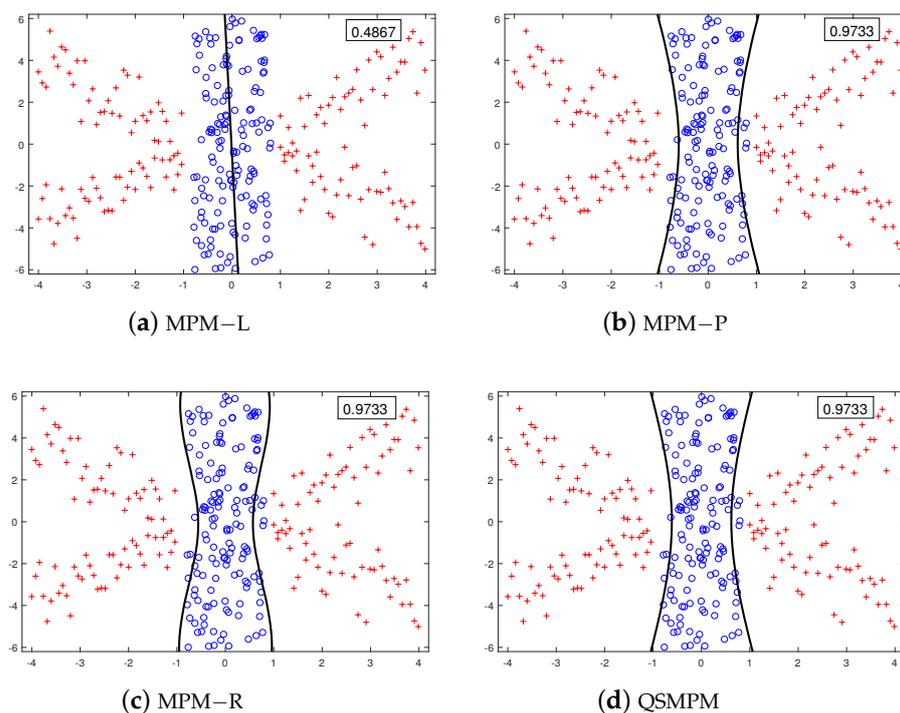
$$[x_i]_2^2 = [x_i]_1^2 - 1 + \xi_i, i = 1, \dots, 150,$$

$$[x_i]_1 \sim U[-0.6, 0.6], [x_i]_2 \sim U[-6, 6], i = 151, \dots, 300,$$

where  $[x_i]_1 \sim U[-4, -1]$  and  $U[1, 4]$ ,  $\xi_i \sim N(0, 1)$ ,  $i = 1, \dots, 150$ .

Figure 5 reports the classification results of the MPM-L, the MPM-P, the MPM-R, and the QSMPM on Example 5, respectively. We can observe in Figure 5 that the classification performance of QSMPM is better than the MPM-L and is similar to the MPM-P and the MPM-R. In addition, our QSMPM finds a hyperbola.

In summary, from Figure 1 to Figure 5, we can see that our QSMPM can find any general form of the quadratic hypersurface, such as the line, parabola, circle, ellipse, and hyperbola found in sequence in the above numerical experiments. Moreover, our method can achieve as good classification performance as the MPM-P and the MPM-R. In addition, it can be seen from Figure 1d that our method can obtain the linear separating hyperplane.



**Figure 5.** Classification results on Example 5.

### 5.2. Benchmark Datasets

To verify the classification performance and computational efficiency of our QSMPM, we performed the following numerical experiments on 14 benchmark datasets. Table 1 summarizes the basic information of the 14 benchmark datasets in the UCI Machine Learning Repository.

**Table 1.** Basic information for the 14 benchmark datasets.

Datasets	Samples	Attributes	Class
lirs	150	4	3
Hepatitis	155	19	2
Wine	178	13	3
Heart	270	13	2
Heart-c	303	14	2
Haberman	306	3	2
Bupa	345	6	2
Pima	768	8	2
QSAR	1055	41	2
Winequality-red	1599	11	6
Wireless	2000	7	4
Image	2310	19	7
Abalone	2649	8	2
Turkiye	5820	32	13

We divided the datasets in Table 1 into two groups for the numerical experiment. The first group was the first seven datasets, and the second group was the last seven datasets. All numerical experimental results were obtained through 10-times 10-fold cross-validation, as well as the numerical experimental results including the mean and standard deviation of accuracy and the CPU time of one experiment. The best results are highlighted in boldface. First of all, Table 2 shows the classification results on the first group.

Table 2. Classification results on the first group.

Methods	Iris	Hepatitis	Wine	Heart	Heart-c	Haberman	Bupa
H-SVM-L	0.6413 ± 0.0245 (4.1293)	0.5835 ± 0.0185 (2.1965)	0.9573 ± 0.0054 (1.8439)	0.6174 ± 0.0287 (7.70490)	<b>1.0000</b> ± 0.0000 (3.9203)	0.5189 ± 0.0024 (15.1649)	0.5573 ± 0.0021 (11.2976)
H-SVM-P	0.9447 ± 0.0055 (1.4087)	0.5150 ± 0.0399 (1.5304)	0.7571 ± 0.0365 (2.6005)	0.6985 ± 0.0126 (3.8797)	0.8678 ± 0.0386 (4.5568)	0.7351 ± 0.0005 (10.0121)	0.5798 ± 0.0003 (13.3459)
H-SVM-R	0.9467 ± 0.0000 (0.8250)	0.5858 ± 0.0280 (0.8734)	0.9341 ± 0.0116 (1.0813)	0.7144 ± 0.0183 (2.8344)	<b>1.0000</b> ± 0.0000 (2.4438)	0.7362 ± 0.0022 (4.0609)	0.5930 ± 0.0192 (3.7438)
S-SVM-L	0.8693 ± 0.0155 (1.4594)	0.6037 ± 0.0206 (1.3416)	0.9582 ± 0.0147 (2.6656)	0.8259 ± 0.0129 (9.2391)	0.9888 ± 0.0156 (8.2672)	0.7241 ± 0.0073 (7.8422)	0.6812 ± 0.0096 (1.7031)
S-SVM-P	0.9653 ± 0.0053 (1.4149)	0.6037 ± 0.0234 (1.4703)	0.7274 ± 0.0235 (1.9407)	<b>0.8296</b> ± 0.0097 (9.0609)	0.8605 ± 0.0087 (11.3328)	0.7228 ± 0.0103 (11.1953)	0.6788 ± 0.0095 (12.3328)
S-SVM-R	0.9547 ± 0.0061 (0.9406)	0.5588 ± 0.0183 (0.8656)	0.8921 ± 0.0131 (1.1641)	0.7989 ± 0.0058 (2.1344)	0.7642 ± 0.0101 (2.6000)	0.7261 ± 0.0080 (2.7719)	0.6826 ± 0.0088 (3.2813)
MPM-L	0.8280 ± 0.0082 (0.3073)	0.6010 ± 0.0114 (0.0244)	0.9731 ± 0.0052 (1.7831)	0.8133 ± 0.0080 (3.1887)	<b>1.0000</b> ± 0.0000 (1.7379)	0.7177 ± 0.0093 (3.4944)	0.6220 ± 0.0079 (0.1201)
MPM-P	0.9747 ± 0.0028 (1.9079)	0.5954 ± 0.0302 (1.7051)	<b>0.9759</b> ± 0.0046 (2.9640)	0.8026 ± 0.0133 (5.6504)	0.9681 ± 0.0066 (6.6425)	0.7159 ± 0.0075 (5.9514)	0.6891 ± 0.0114 (5.5131)
MPM-R	0.9620 ± 0.0063 (4.5256)	0.5382 ± 0.0412 (9.0683)	0.7860 ± 0.0121 (4.3758)	0.6578 ± 0.0105 (10.5238)	0.6981 ± 0.0101 (5.6452)	0.6879 ± 0.0114 (11.9107)	0.7212 ± 0.0088 (8.5723)
QSSVM	0.9533 ± 0.0070 (1.2371)	0.5626 ± 0.0342 (3.7612)	0.9608 ± 0.0052 (2.1109)	0.6989 ± 0.0182 (3.5241)	0.9980 ± 0.0023 (4.2947)	<b>0.7416</b> ± 0.0073 (5.6426)	0.7203 ± 0.0047 (8.0918)
SQSSVM	0.9527 ± 0.0049 (0.9266)	0.5650 ± 0.0117 (4.6098)	0.9622 ± 0.0071 (2.6832)	0.7970 ± 0.0098 (5.2822)	0.9954 ± 0.0024 (7.2993)	0.7296 ± 0.0050 (5.0482)	<b>0.7220</b> ± 0.0074 (7.0653)
QSMPM	<b>0.9767</b> ± 0.0035 (0.3089)	<b>0.6069</b> ± 0.0313 (2.0608)	<b>0.9759</b> ± 0.0053 (1.3884)	0.8293 ± 0.0114 (0.8580)	<b>1.0000</b> ± 0.0000 (2.5179)	0.7205 ± 0.0069 (0.2902)	0.7164 ± 0.0093 (0.2870)

It can be seen from Table 2 that compared with the other methods, our QSMPM obtained better accuracy on the first group of benchmark datasets, among which the accuracy was the best on four benchmark datasets. More specifically, except for Haberman and Bupa, the accuracy of our method was the best compared to the QSSVM and the SQSSVM. The accuracy of our QSMPM was the best compared to the three original kernel versions of the MPM except for Bupa. Furthermore, the accuracy of our method was the best compared to the H-SVM and the S-SVM with three kernel function except for Heart and Haberman. In addition, we can observe that QSMPM had a short CPU time.

Then, the classification results on the second group are reported in Table 3. The symbol “-” indicates that the corresponding method cannot obtain the classification results, because it cannot choose the optimal parameter in a limited amount of time or because the dimension and the number of the dataset are relatively large, resulting in insufficient memory.

Table 3. Classification results on the second group.

Methods	Pima	QSAR	Winequality-Red	Wireless	Image	Abalone	Turkiye
H-SVM-L	0.6799 ± 0.0035 (128.9255)	0.3675 ± 0.0014 (223.9172)	0.6156 ± 0.0014 (813.7250)	0.7280 ± 0.0094 (3956.2000)	0.6470 ± 0.0011 (947.0836)	0.7357 ± 0.0040 (2825.5000)	0.5051 ± 0.0005 (3013.1000)
H-SVM-P	0.5710 ± 0.0155 (75.1160)	0.8133 ± 0.0081 (451.2469)	0.4653 ± 0.0000 (475.2321)	0.8463 ± 0.0234 (4449.2000)	0.6979 ± 0.0020 (996.2086)	0.4934 ± 0.0000 (1429.9000)	0.5033 ± 0.0038 (23897.0000)
H-SVM-R	0.6919 ± 0.0117 (14.9234)	0.8171 ± 0.0069 (110.5313)	<b>0.7628</b> ± 0.0058 (3707.5000)	<b>0.9877</b> ± 0.0010 (404.7438)	0.9698 ± 0.0019 (567.4688)	0.8067 ± 0.0062 (693.4125)	-
S-SVM-L	<b>0.7669</b> ± 0.0076 (82.0609)	0.8340 ± 0.0152 (211.0188)	0.7291 ± 0.0023 (640.9000)	0.9139 ± 0.0044 (1793.9000)	0.7531 ± 0.0304 (2595.7000)	0.8108 ± 0.0009 (1189.7000)	-
S-SVM-P	0.7585 ± 0.0066 (27.8571)	<b>0.8677</b> ± 0.0058 (259.1922)	0.7492 ± 0.0021 (1123.8000)	0.9799 ± 0.0017 (3120.7000)	0.9602 ± 0.0015 (4820.8000)	0.8298 ± 0.0018 (3120.7000)	-
S-SVM-R	0.6941 ± 0.0078 (16.5313)	0.8503 ± 0.0031 (103.9438)	0.7352 ± 0.0024 (294.0688)	0.9853 ± 0.0054 (717.2984)	<b>0.9715</b> ± 0.0020 (506.2984)	0.8242 ± 0.0008 (984.2250)	-
MPM-L	0.7405 ± 0.0048 (0.1014)	0.8296 ± 0.0039 (0.8377)	0.7409 ± 0.0025 (0.2683)	0.9108 ± 0.0008 (0.1560)	0.8555 ± 0.0012 (0.2969)	0.8144 ± 0.0009 (0.1266)	0.5779 ± 0.0026 (0.3167)
MPM-P	0.7442 ± 0.0035 (32.9178)	0.8225 ± 0.0061 (68.1350)	0.7326 ± 0.0020 (122.7891)	0.9361 ± 0.0013 (187.2266)	0.8499 ± 0.0011 (257.7266)	0.8109 ± 0.0009 (413.9861)	0.5780 ± 0.0014 (2669.2000)
MPM-R	0.7356 ± 0.0054 (37.2874)	0.8373 ± 0.0037 (83.8022)	0.7234 ± 0.0020 (139.5734)	0.9850 ± 0.0006 (248.8594)	0.9413 ± 0.0018 (402.6281)	0.8264 ± 0.0014 (423.6125)	0.5689 ± 0.0016 (2386.7000)
QSSVM	0.7663 ± 0.0049 (70.7730)	-	0.4722 ± 0.0032 (27.7094)	0.6315 ± 0.0078 (14.7250)	0.5714 ± 0.0000 (59.1266)	0.5125 ± 0.0012 (29.5344)	-
SQSSVM	0.7589 ± 0.0036 (43.0610)	-	0.7452 ± 0.0036 (36.8234)	0.9782 ± 0.0012 (45.3036)	0.5714 ± 0.0000 (63.7109)	0.8280 ± 0.0015 (50.1641)	-
QSMPM	0.7530 ± 0.0049 (0.5585)	0.8482 ± 0.0047 (16.6328)	0.7470 ± 0.0027 (1.3525)	0.9427 ± 0.0012 (1.0608)	0.8731 ± 0.0013 (3.6953)	<b>0.8299</b> ± 0.0011 (1.1984)	<b>0.5852</b> ± 0.0018 (14.6360)

From Table 3, we can see that our QSMPM had good classification results on the second group of benchmark datasets. Especially on QSAR and Turkiye, the H-SVM-R, the three kernel versions of the S-SVM, the QSSVM, and the SQSSVM could not obtain the corresponding classification results, but our QSMPM could obtain good classification performance. Here, we mention the reason for this situation. According to the computational complexity of each method, we know that when the sample dimension and the number of samples are relatively large, the SVM and the QSSVM need a larger memory space. In addition, our QSMPM had the fastest running time except the MPM-L, and it ran quite fast when the number of samples or the dimension was large.

### 5.3. Statistical Analysis

To further compare the performance of the above 12 methods, the Friedman test and the post-hoc test were performed. The ranks of the 12 methods on all benchmark datasets is shown in Table 4.

First, the Friedman test was used to compare the average ranks of different methods. The null hypothesis states that all methods have the same performance, that is their average ranks are the same. Based on the average ranks displayed in Table 4, we can calculate the Friedman statistic  $\tau_F$  by the following formula:

$$\tau_{\chi^2} = \frac{12N}{k(k+1)} \left( \sum_{i=1}^k r_i - \frac{k(k+1)^2}{4} \right),$$

$$\tau_F = \frac{(N-1)\tau_{\chi^2}}{N(k-1)-\tau_{\chi^2}},$$
(38)

where  $N$  and  $k$  are the number of datasets and methods, respectively.  $r_i$  is the average rank of the  $i$ -th method. According to the formula (38),  $\tau_F = 4.1825$ . For  $\alpha = 0.05$ , we can obtain  $F_{\alpha} = 1.8526$ . Since  $\tau_F > F_{\alpha}$ , we rejected the null hypothesis. Then, we proceeded with a post-hoc test (the Nemenyi test) to find out which methods significantly differed. To be more specific, the performance of two methods was considered to be significantly different if the difference of their average ranks was larger than the critical difference (CD). The CD can be calculated by:

$$CD = q_{\alpha} \sqrt{\frac{k(k+1)}{6N}}.$$
(39)

For  $\alpha=0.05$ , we know  $q_{\alpha} = 3.2680$ . Thus, we obtained  $CD = 4.4535$  by the formula (39).

Figure 6 visually displays the results of the Friedman test and Nemenyi post-hoc test, where the average ranks of each method are marked along an axis. The axis is turned so that the lowest (best) ranks are to the right. Groups of methods that are not significantly different are linked by a red line. In Figure 6, we can see that our QSMPM was the best one statistically among the compared methods. Furthermore, there was no significant difference in performance between the QSMPM and the MPM-R.

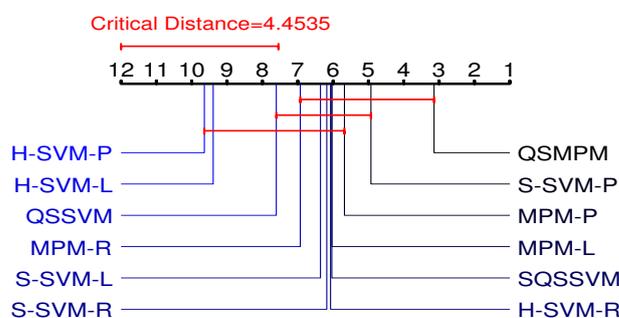


Figure 6. Results of the Friedman test and the Nemenyi post-hoc test.

Table 4. Ranks of accuracy.

Datasets	H-SVM-L	H-SVM-P	H-SVM-R	S-SVM-L	S-SVM-P	S-SVM-R	MPM-L	MPM-P	MPM-R	QSSVM	SQSSVM	QSMPM
Iris	12	9	8	10	3	5	11	2	4	6	7	1
Hepatitis	7	12	6	2.5	2.5	10	4	5	11	9	8	1
Wine	7	11	8	6	12	9	3	1.5	10	5	4	1.5
Heart	12	10	8	3	1	6	4	5	11	9	7	2
Heart-c	2.5	9	2.5	7	10	11	2.5	8	12	5	6	2.5
Haberman	12	3	2	6	7	5	9	10	11	1	4	8
Bupa Liver	12	11	10	7	8	6	9	5	2	3	1	4
Pima	11	12	10	1	4	9	7	6	8	2	3	5
QSAR	10	9	8	5	1	2	6	7	4	11.5	11.5	3
Winequality-red	10	12	1	8	2	6	5	7	9	11	4	3
Wireless	11	10	1	8	4	2	9	7	3	12	5	6
Image Segmentation	10	9	2	8	3	1	6	7	4	11.5	11.5	5
Abalone	10	12	9	8	2	5	6	7	4	11	3	1
Turkiye	5	6	9.5	9.5	9.5	9.5	3	2	4	9.5	9.5	1
Average ranks	9.3929	9.6429	6.0714	6.3571	4.9286	6.1786	6.0357	5.6786	6.9286	7.6071	6.0357	3.1429

## 6. Conclusions

For the binary classification problem, a new classifier, called the kernel-free quadratic surface minimax probability machine (QSMPM), was proposed by using the kernel-free techniques of the QSSVM and the classification idea of the MPM. Specifically, our goal was to find a quadratic hypersurface that separates two classes of samples with maximum probability. However, the optimization problem derived directly was too difficult to solve. Therefore, a nonlinear transformation was introduced to change the quadratic function involved into a linear function. Through such processing, our optimization problem finally became a second-order cone programming problem, which was solved efficiently by an alternate iteration method. Here, we clarify the main contributions of this paper. Unlike the methods realizing nonlinear separation, our method was kernel-free and had better interpretability. Then, our method was easy to use because it did not have any parameters. Furthermore, numerical experiments on five artificial datasets showed that the quadratic hypersurfaces found by our method were rather general, including that it could obtain the linear separating hyperplane. In addition, numerical experiments on benchmark datasets confirmed that the proposed method was superior to some relevant methods in both accuracy and computational time. Especially when the number of samples or dimension was relatively large, our method could also quickly obtain good classification performance. Finally, the results of the statistical analysis showed that our QSMPM was statistically the best one compared with the corresponding methods. Our QSMPM focuses on the standard binary classification problem, which we will extend to the multiclassification problem.

In our future work, there will be some issues to be addressed to extend our QSMPM. For example, we need to investigate further how to add appropriate regularization terms to our method. Meanwhile, we need to consider that the worst-case accuracies for two classes are not the same, and that will be interesting. Furthermore, we will pay attention to how the QSMPM achieves the dual purpose of feature selection and classification simultaneously. In addition, we can apply our method to practical problems in many fields in the future, especially image recognition in the medical field.

**Author Contributions:** Conceptualization, X.Y.; Data curation, Y.W.; Formal analysis, X.Y.; Methodology, Y.W. and Z.Y.; Software, Y.W.; Supervision, Z.Y.; Writing—original draft, Y.W.; Writing—review and editing, Z.Y. and X.Y. All authors have read and agreed to the published version of the manuscript.

**Funding:** This work is supported by the Xinjiang Provincial Natural Science Foundation of China (No. 2020D01C028) and the National Natural Science Foundation of China (No. 12061071).

**Data Availability Statement:** All of the benchmark datasets used in our numerical experiments are from the UCI Machine Learning Repository, which are available at <http://archive.ics.uci.edu/ml/> (accessed on 10 June 2021).

**Conflicts of Interest:** The authors declare that they have no conflict of interest regarding this work.

## References

1. Langley, P.; Simon, H.A. Applications of machine learning and rule induction. *Commun. ACM* **1995**, *38*, 54–64. [[CrossRef](#)]
2. Vapnik, V.; Sterin, A. On structural risk minimization or overall risk in a problem of pattern recognition. *Autom. Remote Control* **1977**, *10*, 1495–1503.
3. Forman, G. An extensive empirical study of feature selection metrics for text classification. *J. Mach. Learn. Res.* **2003**, *3*, 1289–1305.
4. Ying, S.H.; Qiao, H. Lie group method: A new approach to image matching with arbitrary orientations. *Int. J. Imaging Syst. Technol.* **2010**, *20*, 245–252. [[CrossRef](#)]
5. Cao, L.J.; Tay, F. Support vector machine with adaptive parameters in financial time series forecasting. *IEEE Trans. Neural. Netw.* **2003**, *14*, 1506–1518. [[CrossRef](#)] [[PubMed](#)]
6. Srinivasu, P.N.; Sivasai, J.G.; Ijaz, M.F.; Bhoi, A.K.; Kim, W.; Kang, J.J. Classification of skin disease using deep learning neural networks with MobileNet V2 and LSTM. *Sensors* **2021**, *21*, 2852. [[CrossRef](#)] [[PubMed](#)]
7. Panigrahi, R.; Borah, S.; Bhoi, A.K.; Ijaz, M.F.; Pramanik, M.; Jhaveri, R.H.; Chowdhary, C.L. Performance assessment of supervised classifiers for designing intrusion detection systems: A comprehensive review and recommendations for future research. *Mathematics* **2021**, *9*, 690. [[CrossRef](#)]
8. Lanckriet, G.R.G.; Ghaoui, L.E.; Bhattacharyya, C.; Jordan, M.I. Minimax probability machine. *Adv. Neural Inf. Process. Syst.* **2001**, *37*, 192–200.

9. Lanckriet, G.R.G.; Ghaoui, L.E.; Bhattacharyya, C.; Jordan, M.I. A robust minimax approach to classification. *J. Mach. Learn. Res.* **2003**, *3*, 555–582.
10. Johnny, K.C.; Zhong, Y.; Yang, S. A comparative study of minimax probability machine-based approaches for face recognition. *Pattern Recognit. Lett.* **2007**, *28*, 1995–2002.
11. Deng, Z.; Wang, S.; Chung, F.L. A minimax probabilistic approach to feature transformation for multi-class data. *Appl. Soft Comput.* **2013**, *13*, 116–127. [[CrossRef](#)]
12. Jiang, B.; Guo, Z.; Zhu, Q.; Huang, G. Dynamic minimax probability machine based approach for fault diagnosis using pairwise discriminate analysis. *IEEE Trans. Control Syst. Technol.* **2017**, *27*, 806–813. [[CrossRef](#)]
13. Yang, L.; Gao, Y.; Sun, Q. A new minimax probabilistic approach and its application in recognition the purity of hybrid seeds. *Comput. Model. Eng. Sci.* **2015**, *104*, 493–506.
14. Kwok, J.T.; Tsang, W.H.; Zurada, J.M. A Class of Single-Class Minimax Probability Machines for Novelty Detection. *IEEE. Trans. Neural Netw.* **2007**, *18*, 778–785. [[CrossRef](#)]
15. Strohmann, T.; Grudic, G.Z. A formulation for minimax probability machine regression. *Adv. Neural Inf. Process. Syst.* **2003**, *76*, 9–776.
16. Huang, k.; Yang, H.; King, I.; Lyu, M.R.; Chan, L. The minimum error minimax probability machine. *J. Mach. Learn. Res.* **2004**, *5*, 1253–1286.
17. Gu, B.; Sun, X.; Sheng, V.S. Structural minimax probability machine. *IEEE Trans. Neural Netw. Learn. Syst.* **2017**, *1*, 1646–1656. [[CrossRef](#)] [[PubMed](#)]
18. Maldonado, S.; Carrasco, M.; Lopez, J. Regularized minimax probability machine. *Knowl. Based Syst.* **2019**, *177*, 127–135. [[CrossRef](#)]
19. Yang, L.M.; Wen, Y.K.; Zhag, M.; Wang, X. Twin minimax probability machine for pattern classification. *Neural Netw.* **2020**, *131*, 201–214. [[CrossRef](#)]
20. Cousins, S.; Shawe-Taylor, J. High-probability minimax probability machines. *Mach. Learn.* **2017**, *106*, 863–886. [[CrossRef](#)]
21. Yoshiyama, K.; Sakurai, A. Laplacian minimax probability machine. *Pattern Recognit. Lett.* **2014**, *37*, 192–200. [[CrossRef](#)]
22. He, K.X.; Zhong, M.Y.; Du, W.L. Weighted incremental minimax probability machine-based method for quality prediction in gasoline blending process. *Chemometr. Intell. Lab. Syst.* **2019**, *196*, 103909. [[CrossRef](#)]
23. Ma, J.; Yang, L.M.; Wen, Y.K.; Sun, Q. Twin minimax probability extreme learning machine for pattern recognition. *Knowl. Based Syst.* **2020**, *187*, 104806. [[CrossRef](#)]
24. Ma, J.; Shen J.M. A novel twin minimax probability machine for classification and regression. *Knowl. Based Syst.* **2020**, *196*, 105703. [[CrossRef](#)]
25. Deng, Z.H.; Chen, J.Y.; Zhang, T.; Cao, L.B.; Wang, S.T. Generalized hidden-mapping minimax probability machine for the training and reliability learning of several classical intelligent models. *Inf. Sci.* **2018**, *436–437*, 302–319. [[CrossRef](#)]
26. Dagher, I. Quadratic kernel-free non-linear support vector machine. *J. Glob. Optim.* **2008**, *41*, 15–30. [[CrossRef](#)]
27. Luo, J.; Fang, S.C.; Deng, Z.B.; Guo, X.L. Soft quadratic surface support vector machine for binary classification. *Asia Pac. J. Oper. Res.* **2016**, *33*, 1650046. [[CrossRef](#)]
28. Bai, Y.Q.; Han, X.; Chen, T.; Yu, H. Quadratic kernel-free least squares support vector machine for target diseases classification. *J. Comb. Optim.* **2015**, *30*, 850–870. [[CrossRef](#)]
29. Gao, Q.Q.; Bai, Y.R.; Zhan, Y.R. Quadratic kernel-free least square twin support vector machine for binary classification problems. *J. Oper. Res. Soc. China.* **2019**, *7*, 539–559. [[CrossRef](#)]
30. Mousavi, A.; Gao, Z.M.; Han, L.S.; Lim, A. Quadratic surface support vector machine with L1 norm regularization. *J. Ind. Manag. Optim.* **2019**. [[CrossRef](#)]
31. Gao, Z.M.; Fang, S.C.; Luo, J.; Medhin, N. A kernel-free double well potential support vector machine with applications. *Eur. J. Oper. Re.* **2020**, *290*, 248–262. [[CrossRef](#)]
32. Luo, J.; Fang, S.C.; Bai, Y.Q.; Deng, Z.B. Fuzzy quadratic surface support vector machine based on fisher discriminant analysis. *J. Ind. Mangn. Optim.* **2015**, *12*, 357–373. [[CrossRef](#)]
33. Yan, X.; Bai, Y.Q.; Fang, S.C.; Luo, J. A kernel-free quadratic surface support vector machine for semi-supervised learning. *J. Oper. Res. Soc.* **2016**, *67*, 1001–1011. [[CrossRef](#)]
34. Tian, Y.; Yong, Z.Y.; Luo, J. A new approach for reject inference in credit scoring using kernel-free fuzzy quadratic surface support vector machines. *Appl. Soft Comput.* **2018**, *73*, 96–105. [[CrossRef](#)]
35. Zhai, Q.R.; Tian, Y.; Zhou, J.Y. Linear twin quadratic surface support vector regression. *Math. Probl. Eng.* **2020**, *2020*, 1–18. [[CrossRef](#)]
36. Luo, J.; Tian, Y.; Yan, X. Clustering via fuzzy one-class quadratic surface support vector machine. *Soft Comput.* **2017**, *21*, 5859–5865. [[CrossRef](#)]