

A new parallel data geometry analysis algorithm to select training data for support vector machine

Yunfeng Shi^a, Shu Lv^{a,*}, Kaibo Shi^b

^a*School of Mathematical Sciences, University of Electronic Science and Technology of China, Chengdu 611731, China*

^b*School of Information Science and Engineering, Chengdu University, Chengdu, 610106, China*

Abstract

Support vector machine (SVM) is one of the most powerful technologies of machine learning, which has been widely concerned because of its remarkable performance. However, when dealing with the classification problem of large-scale datasets, the high complexity of SVM model leads to low efficiency or become impractical. Due to the sparsity of SVM in the sample space, this paper presents a new parallel data geometry analysis (PDGA) algorithm to reduce the training set of SVM, which helps to improve the efficiency of SVM training. The PDGA introduces mahalanobis distance to measure the distance from each sample to its centroid, and based on this, defines hyperellipsoid spatial density to help remove dense redundant data. When further reducing the training set, cosine angle distance analysis method is proposed to determine whether the samples are redundant data, to ensure that the valuable data are not removed. Different from the previous data geometry analysis methods, the PDGA algorithm is implemented in parallel, which leads to substantial saving in the computational cost. Experimental results on artificial dataset and 6 real datasets show that the algorithm can adapt to different sample distributions, significantly reduce the training time and memory requirements without sacrificing the classification accuracy, and its performance is significantly better than the other 4 competitive algorithms.

Keywords: Support vector machine, Sample reduction, Geometry analysis, Parallel

1. Introduction

In the era of big data, information presents explosive growth. How to find valuable information from large-scale data shows great importance. In recent years, the research on classifier of large-scale data has become a hot research field [1, 2]. Many classifiers, such as Naive Bayes [3, 4], Artificial Neural Network [5], Decision Tree and Random Forest [6, 7], Logical Regression [8], K-Nearest Neighbor [9], Support Vector Machine [10] have been used for the classification of big data. Among these classifiers, support vector machine (SVM) [11, 12, 13] proposed by Vapnik has been widely used due to its solid theoretical basis. Based on VC dimension theory and minimum structural risk function, SVM aims at maximizing the margin between the two classes, and thus has strong generalization ability. Bhowmik et al. [14, 15, 16] mentioned that SVM achieved better classification accuracy than other supervised learning methods. The research results of Tang et al [17, 18, 19] showed that the average performance of SVM is even equal to that of deep learning. This is why SVM has been widely concerned in recent years.

However, it is impractical to use SVM for the classification of large-scale data, because the key to training SVM is to solve a quadratic programming problem (QP), which is dependent on the size of training dataset greatly. Theoretically, time complexity of training SVM is $O(m^3)$ and space complexity is $O(m^2)$ [20]. This results in slow training of SVM even on some medium-size datasets [13]. This brings great challenges to the classification of large-scale data by SVM [21]. In recent years, a number of methods have been developed to reduce the computational complexity of SVM. These algorithms can be classified into three categories including techniques which (i) decomposition of QP problems, (ii) parallel training and (iii) sample reduction.

*Corresponding author. Shu Lv

Email address: lvshu@uestc.edu.cn (Shu Lv)

- 20 (i) Decomposition of QP problems [22]: These methods decompose large QP problem into a series of small QP problems by making use of the sparsely characteristic [23] of SVM solutions. By solving the subproblem iteratively, the calculation of whole kernel matrix can be avoided, thus requirements of memory space and solution time are reduced. Among these methods, the most famous one is the Sequence Minimum Optimization (SMO) algorithm [24, 25], which only needs to optimize two variables in each iteration and has a very fast convergence rate. LIBSVM [26], one of the most popular SVM implementations, is proposed based on SMO algorithm. With LIBSVM, the time complexity of SVM training was reduced to $O(m^2)$.
- (ii) Parallel training: The basic idea of most SVM parallel training methods is to divide training set into a series of independent subsets, and then train SVM simultaneously on different processors, thus to filter out the support vectors (SVs) of each subset. Finally use these SVs to train the real SVM decision plane, as shown in [27, 28]. The parallelization training idea of Scholkopf et al. [29] is different. It approximates kernel matrix of SVM to a block diagonal matrix, so that the original optimization problem can be decomposed into hundreds of sub-problems that are easy to be solved in parallel.
- 30 (iii) Sample reduction: In most cases, the solution of SVM is determined by a small subset of training data, which are called support vectors (SVs) [30, 31], and the number of SVs is much less than samples in whole training set. Sample reduction methods [32] reduce time and memory requirement of SVM by eliminating redundant data that are unlikely to be SVs and retaining only samples that may be SVs.

Among the three types of techniques, decomposition of QP problems are applied to either reduce the complexity of the underlying optimization problem, or to handle the optimization process more efficiently. However, these approaches still induce the problem of high memory complexity in big data problems [33]. Parallel training methods [34] require multiple iterations in the process of filtering SVs, which introduce additional memory burden and reduce classification accuracy. Therefore, these two methods are not suitable for problems with large-scale datasets. In contrast, sample reduction techniques are considered to be the most straightforward and effective ways to use SVM for big data.

There exist five main groups of approaches for sample reduction, including data geometry analysis methods, neighborhood analysis methods, evolutionary methods, active learning methods and random sampling methods. These algorithms seem to have certain shortcomings and limitations. For example, data geometry analysis methods usually require clustering [35], which can be time-consuming and need to determine the number of clusters manually. Neighborhood analysis methods [36] usually need to calculate all distances between samples in training set, which have high time complexity and are sensitive to noise. Genetic algorithm [37] is one of the important evolutionary methods for SVM sample reduction. However, it is a challenging to effectively determine the search space and ensure appropriate convergence speed, which is also the biggest disadvantage of genetic algorithm. Active learning algorithms can also be used to select training samples [38, 39], but these methods usually require multiple iterations and take a long time for large-scale datasets. Random sampling methods are fast and effective, but that ignore the relationship between samples, which may lead to a lot of useful information not being extracted.

In this paper, a new parallel data geometry analysis (PDGA) algorithm was proposed to select training data for SVM, thus solve the problems of long execution time and insufficient ability to maintain classification accuracy of existing algorithms. This method extracts redundant samples from training set that cannot be SVs, and ensures that potential SVs are not eliminated, to reduce the size of training set as much as possible. The properties of this algorithm are summarized as follows.

1. The non-SVs removal is more thorough, will not be affected by the correlation and dimensional differences of data different attributes, and will not increase the risk of SVs being removed by mistake. It is more stable and reliable;
2. For the case that data distribution shape is not convex or the centroid is not in geometric center, this algorithm can also well extract the potential SVs located on the boundary;
3. This algorithm can be carried out in parallel, and the time complexity of proposed algorithm increases approximately linearly with sample size, which improves the efficiency of training SVM on large-scale datasets.

The rest of this paper is organized as follows. In Section 2, we provide a brief review of SVM and sample reduction algorithms. We introduce the basic theory of the proposed algorithm and analyze time complexity in Section 3. In Section 4, we experiment the proposed algorithm on artificial dataset and 6 real datasets, and compare it with the most competitive algorithms. Finally, we summarize the paper and discuss future work in Section 5.

2. A Brief Review of SVM and Sample Reduction Algorithms

2.1. Support Vector Machine

SVM is a classical binary classification algorithm with excellent generalization ability. Take a linear binary classification problem as an example, for dataset $\{(x_i, y_i) \mid x_i \in R^n, y_i \in \{-1, +1\}, i = 1, 2, \dots, m\}$, where x_i is a n -dimensional vector, y_i is the category label of x_i . The essence of SVM is to find a hyperplane $\omega \cdot x + b = 0$ in a n -dimensional space that can separate different types of data as much as possible. According to the theory of SVM, the problem of finding the optimal hyperplane is transformed into a QP problem as follows:

$$\begin{aligned} \min \varphi(\omega) &= \frac{1}{2} \|\omega\|^2 \\ \text{s.t. } y_i (\omega^T x_i + b) &\geq 1, i = 1, 2, \dots, m. \end{aligned} \quad (1)$$

If the training data is not linear separable, in order to tolerate some samples that do not meet the constraint, relaxation variables ξ need to be introduced, and the problem to be optimized is accordingly modified as

$$\begin{aligned} \min \varphi(\omega, \xi) &= \frac{1}{2} \|\omega\|^2 + C \sum_{i=1}^m \xi_i \\ \text{s.t. } y_i (\omega^T x_i + b) &\geq 1 - \xi_i, i = 1, 2, \dots, m, \\ \xi_i &\geq 0, i = 1, 2, \dots, m. \end{aligned} \quad (2)$$

where $C > 0$, is used to control the cost of violating the constraint. By introducing Lagrange multipliers, the dual form of Eq. 2 minimization problem is equivalent to

$$\begin{aligned} \max L(\alpha) &= \sum_{i=1}^m \alpha_i - \frac{1}{2} \sum_{i=1}^m \sum_{j=1}^m \alpha_i \alpha_j y_i y_j (x_i \cdot x_j) \\ \text{s.t. } 0 &\leq \alpha_i \leq C, i = 1, 2, \dots, m, \\ \sum_{i=1}^m y_i \alpha_i &= 0. \end{aligned} \quad (3)$$

where α is a m -dimensional vector, and α_i is its component, which is the i -th Lagrange multiplier. By solving quadratic convex programming problem of Eq. 3, the optimal decision function of SVM can be obtained as

$$f(x) = \text{sgn} \left(\sum_{i=1}^m \alpha_i y_i (x_i \cdot x) + b \right), \quad (4)$$

where α_i corresponds to the solution of the QP problem, and b is the optimal offset of the hyperplane.

For nonlinear data in original space, SVM implicitly maps vectors to a high-dimensional feature space with the help of kernel function [9], and searches for the optimal hyperplane in this space, so as to solve the linearly unfractionable problem in original space. In the high-dimensional feature space mapped by kernel function, the optimal SVM hyperplane becomes

$$f(x) = \text{sgn} \left(\sum_{i=1}^m \alpha_i y_i K(x_i, x) + b \right), \quad (5)$$

where $K(x_i, x_j) = \phi(x_i) \cdot \phi(x_j)$, $\phi: R^n \rightarrow F$ is a mapping of vector x from the input space R^n to the eigenspace F . With the help of kernel function, we do not need to find specific mapping expression, only need to calculate inner product of two eigenvectors $\phi(x_i)$ and $\phi(x_j)$ by kernel function.

The key to training SVM is to solve the QP problem of Eq. 3, whose time complexity is $O(m^3)$ and memory complexity is $O(m^2)$ [20], where m is the number of samples in training set. In fact, there is only a small part

of $\alpha_i > 0$ by solving Eq. 3, and the corresponding samples x_i are called support vectors (SVs), which can actually define the optimal hyperplane. Most of $\alpha_i = 0$, and the corresponding samples are non-support vectors (non-SVs), which make no contribution to the optimal hyperplane of SVM training, but waste memory space and increase training time. That is to say, the hyperplane trained by SVs alone is completely equivalent to the hyperplane trained by all data. In recent years, in order to accelerate the training speed of SVM, there are many methods devoted to extracting SVs or removing non-SVs from training set. Such methods are called sample reduction [40] or instance selection [41].

2.2. Sample Reduction for SVM

Simple random sampling methods (SRS) [42] are the most direct ways to reduce training samples. These methods randomly extract a certain proportion of samples from training set to train SVM, so as to improve the training efficiency of SVM. Among them, uniform random sampling [43] may be the most robust method to reduce training samples. However, these methods are not reliable, because they ignore the structural information between samples, which lead to uncertain results [44].

The neighborhood analysis methods are based on the fact that only samples belonging to different categories and located near the hyperplane are likely to become SVs [45]. The purpose of neighborhood analysis is to extract data point pairs that are close to each other and belong to different categories. Compressed Nearest Neighbor algorithm (CNN) [46] is one of the neighborhood analysis methods. However, CNN needs to calculate all distances between samples, the time complexity is $O(m^2)$, and CNN is not noise tolerant. Shin et al. [47] proposed a method to select samples of overlapping regions near decision boundaries, because SVs generally exist in different categories of overlapping regions. Jiantao et al. [48] proposed a fuzzy clustering method to select samples existing in overlapping regions. However, these two approaches do not perform well when the class distributions do not overlap.

There are few evolutionary methods to reduce training set for SVM, and genetic algorithm [37] is one of the most important methods. It evolves a group of solutions (individuals), and the fitness of individuals corresponds to the classification accuracy of SVM. However, it is a challenging problem of how to effectively determine search interval of genetic algorithm and speed up convergence rate. Pighetti et al. [49] enhanced genetic evolution by using a locality sensitive hashing method, which found the nearest vector in training set for all generated vectors in the optimization process and applied it to solve the multi-classification problem. Although this method is novel, it did not provide a stop criterion for the optimization of multi-classification tasks.

Active learning methods [50] are a group of semi supervised learning approaches, which can also be used to reduce the size of SVM training set. They initially divide the training set into a labeled dataset and an unlabeled dataset (corresponding to SVs and non-SVs). Firstly, the classifier is trained with the tagged set. According to the pre-set criteria, the selected unlabeled samples are labeled and added to the tagged dataset. Then, the classifier is retrained with the updated tagged set, and the iteration continues until the predefined conditions are satisfied. Wang et al. [51] proposed an active learning algorithm for SVM sample selection based on neighborhood entropy measure. On the basis of tagged samples, they introduced the concept of neighborhood entropy to analyze the uncertainty of untagged samples. This model improves the classification ability of SVM. However, active learning algorithms have shortcoming in the whole dataset could be scanned for many times, thus the convergence time may be very long.

Data geometry analysis methods [52] exploit the information about training set structure to extract potential SVs. These methods generally calculate the distance from each vector to its centroid of the category to judge whether the vector is located on the interior or the boundary. Samples on the boundary have a high probability of becoming SVs. These methods can be roughly classified into two groups, the first encompasses clustering-based techniques, whereas the second contains non-clustering algorithms.

Clustering-based algorithms have been intensively studied for selecting refined training sets, Wang and Shi [53] proposed an SR-DSA algorithm, which first uses agglomerative hierarchical clustering (AHC) to cluster each category separately. Then calculate the Mahalanobis distance from the vector of each cluster to its centroid, and delete the "inner points" which are very close to the centroid. Finally, remove "exterior points" that are distant

130 from the opposite class. This method can maintain the classification accuracy well while reducing training samples, but it has two disadvantages, one is the time complexity of clustering is relatively high, the other is both the number of clustering and the proportion of retained samples need to be determined artificially, which depends on experience heavily. These are also common disadvantages of clustering based methods.

In the non-clustering-based methods, the SE algorithm proposed by Liu et al.[54] is a very competitive technique, which can adapt to the case where the centroid is not located in the geometric center of this class due to uneven distribution of instances in vector space. However, too many parameters introduced by this method artificially may lead to too many iterations of sample reduction and increase the time complexity. In addition, this algorithm uses Euclidean distance to measure the distance from each sample to its centroid, which cannot reflect the structural information of the sample distribution well, and is not applicable to samples with
140 large differences in the distribution range of each attribute.

Data geometry analysis methods are a fast and direct sample reduction techniques, which have attracted the attention of many researchers. In general, it is very difficult to extract the SVs directly from training set, mainstream geometric analysis methods usually extract redundant samples that are unlikely to become SVs. However, most existing geometric analysis methods have the limitations of high time complexity or insufficient removal ability of non-SVs. We proposed a new approach to reduce the SVM samples. Firstly, the non-SVs gathered in the interior are removed by calculating the spatial density of the hyperellipsoid, and then the non-SVs near the boundary are further removed by cosine angle distance judgment, to eliminate redundant data to the greatest extent. Since this method can be implemented in parallel, the efficiency of dataset processing was greatly improved.

150 3. Sample Reduction of Parallel Data Geometric Analysis for Support Vector Machine

Some algorithms compare the Euclidean distances from different samples to the centroid to determine whether the samples are located in the interior or the boundary of this class, which is not stable. Because the Euclidean distance is easily affected by the dimensionality of different attributes, and the distance only contains the information of a single sample point, without considering the overall distribution of training set. If the centroid is not located in the geometric center of this class, that may cause some SVs at the boundary are discarded. As shown in Fig. 1, points A , B and C are all located on the boundary of class distribution and are likely to become SVs. They should be retained with equal probability in the sample reduction process. If compare their Euclidean distance to the centroid O , the probability of discarding point A is greater than that of point B when the radius of the reducing sphere increases, because the data is more distributed in the direction of e_1 . Point B has a higher probability of being discarded than C , because the centroid is not located in the
160 geometric center of this class, resulting in point C being farther away from the centroid, which is obviously unreasonable.

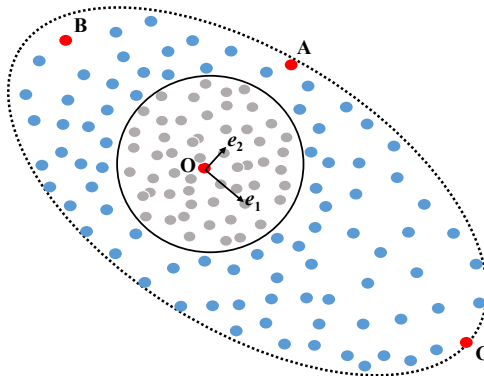


Figure 1: Redundant samples are deleted based on the data geometry analysis method. Point O represents the centroid of this class, and the samples in the spherical region are discarded.

Compared with Euclidean distance, Mahalanobis distance is a more reasonable distance measure, because it is calculated by the first and second order statistical information of samples, which implies the structural

information of samples. Mahalanobis distance can eliminate the influence of dimensional difference and linear correlation of different attributes of samples.

In our algorithm, Mahalanobis distance is introduced to measure the distance from each sample to its centroid. On this basis, the spatial density of hyperellipsoid is defined to judge whether the sample is on the boundary or the inside. This section first introduces the theory of Mahalanobis distance, gives the calculation details of hyperellipsoid spatial density, then explains how to calculate the cosine angle distance between samples, finally summarizes the implementation steps of our algorithm, and analyzes the time complexity.

3.1. Maharanobis Distance

Let X be a $m \times n$ sample matrix, containing m random observations $x_i, i = 1, 2, \dots, m$. n represents the number of features of each sample. The definition of Mahalanobis distance is

$$d^2(x_i, X) = (x_i - \mu) \Sigma^{-1} (x_i - \mu)^T, \quad (6)$$

where $d^2(x_i, X)$ is the square of Mahalanobis distance from sample x_i to population X , μ is the mean vector of sample matrix X , $\mu = \frac{1}{m} \vec{1}^T X$, $\vec{1}$ denotes a m -dimensional all one vector. Σ is the covariance matrix of sample matrix X , $\Sigma = \frac{1}{m} X^T X - \frac{1}{m^2} X^T \vec{1} \vec{1}^T X$. Since Σ is a real symmetric and positive semidefinite matrix, it can be orthogonally similar diagonalized, that is, there exists an orthogonal matrix P , which makes the

$$\Sigma = P \Lambda P^T, \quad (7)$$

where $P = [e_1, e_2, \dots, e_n]$ satisfy $e_i \cdot e_k = 0, i \neq k, i, k = 1, 2, \dots, n$. Λ is a diagonal matrix, the elements on the diagonal are the eigenvalues of Σ , and $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_n \geq 0$. The Mahalanobis distance in Eq. 6 is rewritten as

$$d^2(x_i, X) = (x_i - \mu) P \Lambda^{-1} P^T (x_i - \mu)^T. \quad (8)$$

If the first d eigenvalues greater than 0 are selected, accordingly, the orthogonal matrix $P = [e_1, e_2, \dots, e_d]$, let $z_i = (x_i - \mu) P$, then Eq. 8 becomes

$$d^2(x_i, X) = z_i \Lambda^{-1} z_i^T = \sum_{j=1}^d \frac{z_{ij}^2}{\lambda_j}. \quad (9)$$

It can be seen that the orthogonal matrix P projects x_i to another d -dimensional space, and the attributes of the projected samples are orthogonal to each other, eliminate the influence of dimensional difference and linear correlation of different attributes of X . We may as well abbreviate the Mahalanobis distance from x_i to the X as d_i and further transform Eq. 9 to obtain

$$\frac{z_{i1}^2}{(d_i \sigma_1)^2} + \frac{z_{i2}^2}{(d_i \sigma_2)^2} + \dots + \frac{z_{id}^2}{(d_i \sigma_d)^2} = 1, \quad (10)$$

where $\lambda_j = \sigma_j^2$, it also represents the sample variance of $(X - \mu)P$ on the j -th dimension. Under the measure of Mahalanobis distance, the projection of sample x_i to $x_i P$ must be located on the hyperellipsoid with centroid μP and axis $d_i \sigma_1, d_i \sigma_2, \dots, d_i \sigma_d$. For different samples in X , $\sigma_1, \sigma_2, \dots, \sigma_d$ is fixed, so d_i can also be regarded as the generalized radius of the hyperellipsoid. As shown in Fig. 2, where Fig. 2(a) is the distribution of original dataset, and Fig. 2(b) is the distribution of the projected dataset. The data are distributed as hyperellipsoids, and Mahalanobis distances from data points on the same hyperellipsoid to the centroid are equal. Mahalanobis distance can better adapt to the shape of data distribution.

3.2. Spatial Density of Hyperellipsoid

Shen et al. [55] pointed out that samples belonging to the same category are generally denser as they are closer to the centroid, and sparser as they are farther away from the centroid. That is to say, SVs are generally located in a relatively sparse boundary region. We define the spatial density of hyperellipsoid based on Mahalanobis distance, and train SVM by removing dense interior points and leaving only sparse boundary

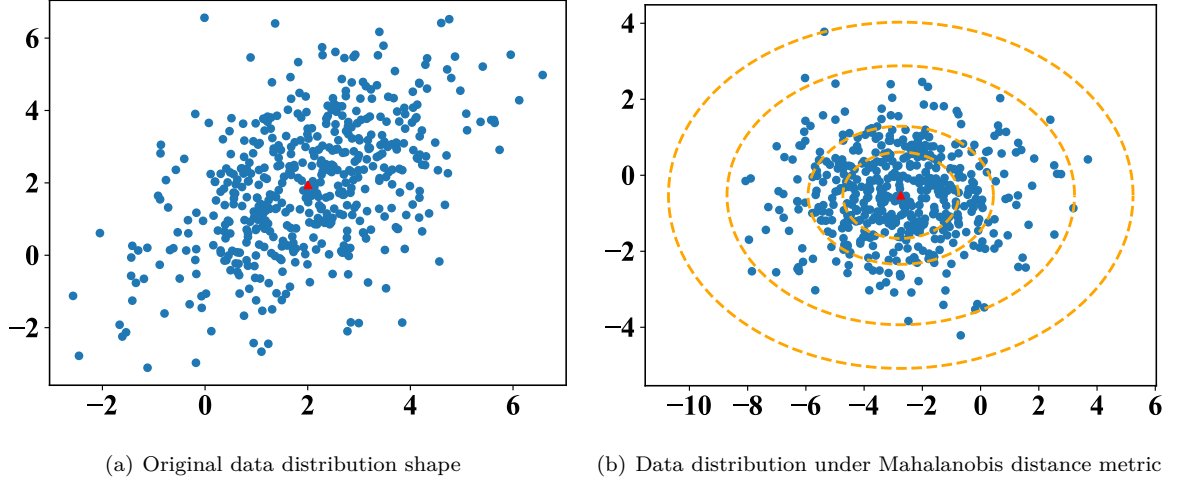


Figure 2: The effect of Mahalanobis distance on data conversion

points. The spatial density reflects the overall information of the sample distribution, which is more stable than the distance measure based on a single sample.

Given a training dataset $D = \{(x_i, y_i) \mid x_i \in R^n, y_i \in \{1, 2, \dots, C\}, i = 1, 2, \dots, m.\}$, x_i is a sample of D , y_i is the class label of x_i . Assuming that each sample can only belong to one category, then D can be expressed as $D = D_1 \cup D_2 \cup \dots \cup D_C$ according to different labels, where D_1, D_2, \dots, D_C represent C subsets of different categories. Correspondingly, X_1, X_2, \dots, X_C are used to represent the sample matrix, where the size of X_c is $m_c \times n$, m_c represents the number of samples in D_c , and n represents the number of features of each sample. Let μ_c be the centroid of X_c , and Σ_c be the covariance matrix of X_c . For a data point $(x_i, y_i) \in D_c, x_i \in R^n, y_i = c$, the Mahalanobis distance between x_i and the population X_c can be calculated according to Eq. 7 and Eq. 8. After sorting by Mahalanobis distance, the dataset is represented by $\{sx_1, sx_2, \dots, sx_{m_c}\}$, and satisfies $d_1 \leq d_2 \leq \dots \leq d_{m_c}$. To simplify notation, denote the distance from sx_i to X_c by d_i . According to the definition of Mahalanobis distance, sx_i must be on the hypersphere with μ_c as the center and d_i as the generalized radius, then the spatial density of the hyperellipsoid where sx_i is located can be defined as

$$SD_i = \frac{i}{V_i}, \quad (11)$$

where SD_i represents the spatial density of the hyperellipsoid where point sx_i is, the number of samples in the envelope of hyperellipsoid is i , V_i represents the volume of the hyperellipsoid, and the calculation formula is (see the proof in Reference [56])

$$V_i = \frac{\pi^{\frac{d}{2}}}{\Gamma(\frac{d}{2} + 1)} \prod_{j=1}^d (\sigma_j d_i), \quad (12)$$

where d is the number of non-zero eigenvalues in matrix Σ_c , and Γ is the gamma function, $\sigma_j = \sqrt{\lambda_j}$. As shown in Fig. 3, the solid line represents the interface between sparse and dense samples. SVs are located in the sparse region outside the hyperellipsoid interface.

In order to determine the boundary between sparse and dense samples, we normalize the density according to Eq. 13

$$SD_i^{\text{norm}} = \frac{SD_i - \min_i SD_i}{\max_i SD_i - \min_i SD_i}. \quad (13)$$

¹⁹⁰ SD_i^{norm} is the sample density after normalization. When a certain SD_i^{norm} drops close to the artificial threshold ξ , then take d_i as the generalized radius and delete all points within the hyperellipsoid.

3.3. Cosine Angle Distance Analysis

For the class distribution is non-convex, if the redundant samples are removed by constructing a reduction sphere, some useful samples may be deleted [54]. In our algorithm, samples are removed based on the spatial

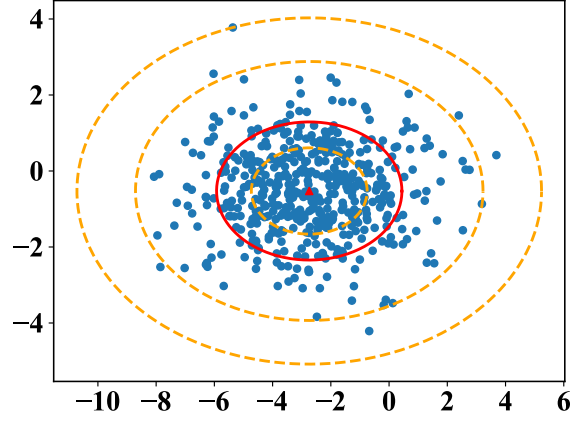


Figure 3: The clustered data points inside the hyperellipsoid are discarded, leaving the sparse boundary points as training samples.

density of hyperellipsoid, and the selection of density threshold should also be cautious, but this may lead to incomplete extraction of redundant samples. In order to solve this problem, we propose a method of cosine angle distance analysis between samples to determine whether the remaining samples are on the boundary. The basic idea is to judge whether there are samples on the line between the outermost sample and the centroid. If there are samples, it means that these samples cannot be located on the boundary of class distribution and are unlikely to become SVs, as shown in Fig. 4.

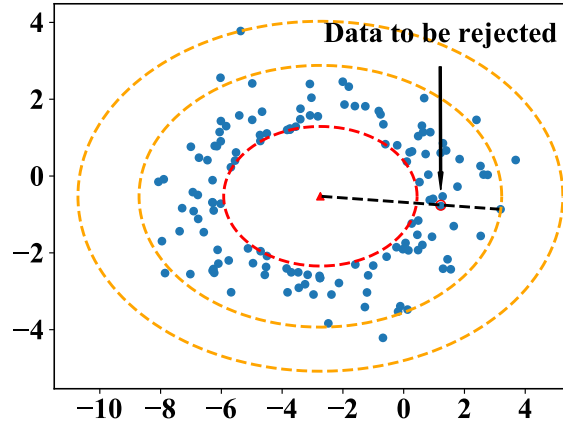


Figure 4: The samples on the line between the centroid and the outermost point should be removed as redundant samples.

200

After the dense internal points are deleted based on the spatial density of hyperellipsoid, it is possible to assume that the training set D_c labeled with class c is reduced to D'_c . Then, the formula for calculating the cosine angle distance between samples in D'_c is

$$\theta_{ij} = \cos^{-1} \left(\frac{(x_i - \mu_c) \cdot (x_j - \mu_c)}{\|x_i - \mu_c\| \|x_j - \mu_c\|} \right), \quad (14)$$

where x_i, x_j are any two different training samples in D'_c , μ_c is the centroid of the training set D_c , θ_{ij} represents the angle between vector $x_i - \mu_c$ and $x_j - \mu_c$ and $\cos^{-1}(\cdot)$ represents the inverse function of cosine value. If $\theta_{ij} = 0$ and $d_i > d_j$ (d_i, d_j are the Mahalanobis distance from x_i, x_j to the centroid), it means that x_j is located between the line connected by x_i and μ_c , thus cannot be located on the boundary of class distribution. Therefore, it should be removed from the training set. To simplify the calculation, let's directly calculate the cosine distance between the vector $x_i - \mu_c$ and $x_j - \mu_c$, where $\theta_{ij} = 0$ is equivalent to $\cos \theta_{ij} = 1$,

$$\cos \theta_{ij} = \frac{(x_i - \mu_c) \cdot (x_j - \mu_c)}{\sqrt{(x_i - \mu_c) \cdot (x_i - \mu_c)} \sqrt{(x_j - \mu_c) \cdot (x_j - \mu_c)}}. \quad (15)$$

After cosine angle analysis, the sample points left are shown in Fig. 5. The cosine angle analysis can not only minimize the training samples, but also ensure that SVs located on the boundary will not be deleted.

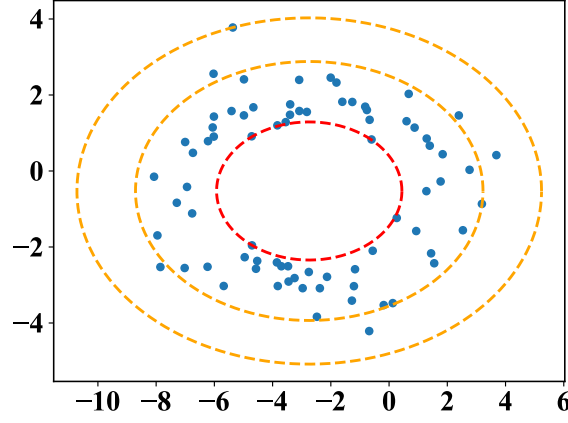


Figure 5: After cosine angle analysis, the boundary points which are not collinear with each other are left as training samples.

3.4. Sample Reduction Algorithm

The sample reduction of parallel data geometric analysis algorithm (PDGA) is summarized as follows and gives the pseudo code of the algorithm.

Step1 Computing Mahalanobis distance matrix.

For training dataset D_c with class label c , the Mahalanobis distances from all the samples to the population are calculated by using Eq. 7 and Eq. 8.

Step2 Calculate spatial density and delete dense interior points.

210 Sorting samples in D_c based on the Mahalanobis distances from small to large. Then, use Eq. 11 and Eq. 12 to calculate the spatial density of each point on its hyperellipsoid, and normalized by Eq. 13. Given the density threshold of hyperellipsoid ξ , for $\forall i \in 1, 2, \dots, m_c$, if $SD_i > \xi$, then take $\max_i d_i$ as the generalized radius to remove all the data inside the hyperellipsoid and update the sample set to D'_c .

Step3 Cosine angle distance analysis, delete the internal non-SVs which close to the boundary.

Calculating the cosine distance between any two samples in D'_c based on Eq. 15. If there is $\cos \theta_{ij} = 1$, and $d_i > d_j$, then remove x_j from D'_c .

Step4 Repeat step 3 until all points does not meet the condition, and the training set is updated to D''_c .

Step5 Repeat step 1 to step 4, until all categories of training set are processed.

Step6 The reduced training dataset $D'' = D''_1 \cup D''_2 \cup \dots \cup D''_C$ is obtained and used to train SVM.

220 It is noted that this algorithm processes training set of each category separately, which has two significant advantages. On the one hand, it is easy to apply the algorithm to multi-classification datasets. On the other hand, the whole training set can be processed in parallel according to different class label, which can significantly accelerate the training speed on large-scale datasets. The parallel processing pattern for this algorithm is shown in Fig. 6.

The pseudo code of the algorithm is shown in algorithm 1. In the parallel processing mode, the input dataset is allocated to different CPUs according to categories for processing at the same time. This mode is more efficient than previous data geometric structure analysis methods.

3.5. Complexity Analysis

The time consumed by this algorithm is mainly in the following two stages.

- 230 (1) The PDGA sample reduction algorithm is used for each subset of different categories, and the time spent is mainly concentrated on calculating the spatial density and cosine distance of the subset. Taking subset D_c as an example, we need to decompose the covariance matrix Σ_c with time complexity of $O(n^3)$, and the time complexity of calculating the Mahalanobis distance and spatial density of each point is $O(m_c)$. When further extracting redundant data, the cosine distances between samples need to be calculated with the time complexity of $O(m'_c)^2$.

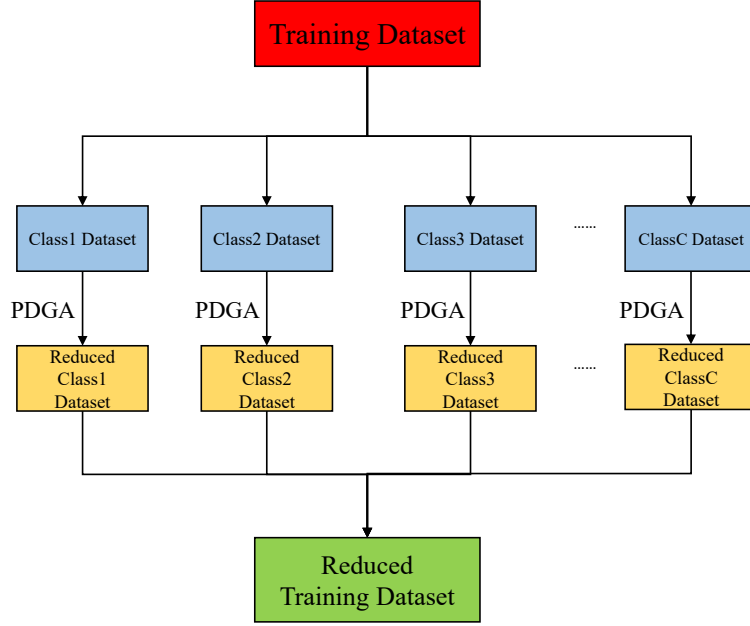


Figure 6: Parallel processing mode of PDGA algorithm.

- (2) Training SVM with reduced dataset, the time complexity is $O\left(\sum_{c=1}^C m_c''\right)^2$ (Time complexity of SVM implementation based on LIBSVM).

The total time complexity of the algorithm is $\max_c \left\{ O\left(n^3 + m_c + (m_c')^2\right) \right\} + O\left(\sum_{c=1}^C m_c''\right)^2$. For large-scale datasets, since $n \ll m$, $m_c'' \ll m_c' \ll m_c \ll m$, the execution time of this algorithm is far less than the time $O(m^2)$ needed to directly train SVM with all samples.

4. Experiments

In order to verify the rationality and performance of the algorithm, we have carried out experiments on a artificial dataset and 6 real datasets, and compared with several other competitive algorithms for data geometry analysis. The abbreviations of the mainstream algorithms used for comparison are shown in Table 1.

Table 1: Abbreviations of the competitive comparison algorithms.

Abbreviation	Describe	Reference
LIBSVM	Directly use all the data of the original training set to train SVM	[26]
SR-DSA	The sample reduction by data structure analysis algorithm	[53]
SE	Shell extraction	[54]
PSCC	Pre-Selection sample based on class centroid	[57]

The contents to be explored in our experiments are (i) the ability of reduction algorithms to maintain classification accuracy; (ii) The total time of algorithms execution (including processing datasets and training SVM with reduced datasets). The experiments are performed on a PC with Intel(R) Pentium(R) CPU G2030 at 3.0 GHz 8 GB RAM, Windows 10, 64 bit Operating System.

4.1. Datasets

The experiments were carried out on a artificial dataset, 4 low-dimensional UCI real datasets and 2 high-dimensional text datasets. The following is a brief introduction to those datasets.

Artificial dataset includes 3000 samples subject to two-dimensional extreme value distribution, which is used to simulate the data that the centroid is not in the geometric center and the distribution range of each attribute

Algorithm 1 A new sample reduction of parallel data geometric analysis algorithm (PDGA)

Input: Training dataset D_c , density threshold ξ ;

Output: Reduced training dataset D_c'' ;

```
1: Compute  $\mu_c = \frac{1}{m_c} \vec{1}_c^T X_c$ ,  $\Sigma_c = \frac{1}{m_c} X_c^T X_c - \frac{1}{m_c^2} X_c^T \vec{1}_c \vec{1}_c^T X_c$ ,  $\Sigma_c = P_c \Lambda_c P_c^T$ ;  
2: for  $i = 1$  to  $m_c$  do  
3:    $d_i = \sqrt{(x_i - \mu_c)^T P_c \Lambda_c^{-1} P_c^T (x_i - \mu_c)}$ ;  
4:    $SD_i = \frac{i}{V_i}$ ;  
5: Sorting data according to Mahalanobis distance:  $\{sx_1, sx_2, \dots, sx_{m_c}\}$  satisfy  $d_1 \leq d_2 \leq \dots \leq d_{m_c}$ ;  
6:  $SD_i \leftarrow \frac{SD_i - \min_i SD_i}{\max_i SD_i - \min_i SD_i}$   
7: for  $i = 1$  to  $m_c$  do  
8:   if  $SD_i > \xi$  then  
9:     Delete  $sx_1, sx_2, \dots, sx_i$ , and  $D'_c = \{sx_{i+1}, sx_{i+2}, \dots, sx_{m_c}\}$  is reserved;  
10: for  $k = i + 1$  to  $m_c$  do  
11:   for  $j = i + 1$  to  $m_c$  do  
12:     compute  $\cos\theta_{kj}$  according to Eq. 15;  
13:     if  $\cos\theta_{kj} = 1$  then  
14:       if  $d_k > d_j$  then  
15:         delete  $sx_j$   
16:       else  
17:         delete  $sx_k$   
18: Output  $D_c''$ , the data in  $D_c''$  are the remaining samples in  $D'_c$ .
```

is quite different. It contains 3 categories, each category contains 1000 samples. The number and distribution of datasets used for testing remain the same as that of training sets.

The 4 low-dimensional UCI datasets are described below. (i) Dermatology is used to determine the type of Erythematous-Squamous disease, which contains 6 classes and 34 attributes. (ii) Letter is a database of character image features, which is used to identify the letter. (iii) Mushrooms includes descriptions of hypothetical samples corresponding to 23 species of gilled mushrooms in the Agaricus and Lepiota Family. Each species is identified as definitely edible, definitely poisonous, or of unknown edibility and not recommended. This latter class was combined with the poisonous one. (iv) Adult dataset is from the UCI machine learning database repository. The goal of this data is to predict whether a household has an income greater than \$50,000. The dataset contains 48,842 samples. Each datum has 14 attributes.

The 2 high-dimensional datasets are 20-Newsgroups and Reuters-21578. 20-Newsgroups contains 19,997 messages from 20 kinds of news, including 4% reprint. Reuters-21578 is a group of 1987 Reuters news. For text datasets, first remove the texts with multiple labels, delete the empty documents, then remove the category of less than 30 texts, and finally use the Word2vec method to convert each text into a 300-dimensional word vector.

In this paper, SVM was implemented based on LIBSVM. For each dataset, parameter tuning is carried out by using grid search first, so as to find the optimal parameters on each dataset. Before the experiment, the kernel type (-t), loss function (-c) and gamma function (-g) of LIBSVM were adjusted, where -t traverses 0, 1, 2; -c traversal $2^{-7}, 2^{-6}, \dots, 2^6, 2^7$, -g traversal $10^{-1}, 10^{-2}, \dots, 10^{-5}$. In order to obtain the parameters corresponding to the best classification performance, 225 parameter optimization experiments were carried out on each dataset. The information and parameter setting details of each dataset are shown in Table 2.

We use F_1 score to measure the accuracy of SVM. F_1 score is the harmonic average of precision and recall, and is defined as

$$F_1 = 2 \frac{\text{precision} \cdot \text{recall}}{\text{precision} + \text{recall}}.$$

In order to test the stability of each algorithm, the five-fold cross validation was adopted on real datasets of

Table 2: Dataset information and parameter setting.

Dataset	# samples	# features	# classes	-t	-c	-g
Artificial dataset	3000	2	3	2	1	10^{-1}
Dermatology	358	34	6	0	2^7	10^{-4}
Letter	20000	16	26	1	2^1	10^{-2}
Mushroom	8124	22	2	2	2^4	10^{-1}
Adult	48842	14	2	2	2^4	10^{-1}
20-NewsGroup	18828	300	20	2	2^6	10^{-1}
Reuters	9931	300	30	2	2^6	10^{-1}

UCI and 2 high-dimensional text datasets, it means that, the training set was divided into five parts randomly, one of which was selected for testing and the remaining four subsets were used for training each time, to avoid the impact of classification performance due to random division of the training set. It is noted that our algorithm and several other algorithms can control the parameters to determine the number of samples to be retained. In order to explore the influence of parameter setting on the experimental results, the experiment studies the classification accuracy and the change of training time under different sample retention rates. The sample retention rate R is defined as the ratio of reduced sample subset size to the training set size, that is

$$R = \frac{\# \text{ samples in subset}}{\# \text{ samples in original training set}}.$$

It should be noted that in order to keep the same samples with other algorithms, PDGA algorithm does not carry out cosine distance analysis when R is greater than 0.3, only when R is less than 0.3.

4.2. Experiments on Artificial Dataset

The experiment is first carried out on artificial dataset, which is used to visually demonstrate the process of reducing the size of training set by the proposed algorithm. The experiment is divided as (i) reduce the training set to get the sample subset; (ii) train SVM with sample subset to get decision plane.

In these five algorithms, our algorithm and SR-DNA, SE are to retain the samples on the boundary of data distribution. In order to show the process of extracting redundant data intuitively, we set the sample retention rate R to 0.7, 0.4, and 0.1 respectively to observe the ability of these algorithm to retain boundary samples. The results are shown in Fig. 7.

It can be seen from Fig. 7 that our algorithm can accurately retain the samples on the boundary under different retention rates R . For SR-DNA algorithm, the internal data deletion is not complete. When the sample retention rate is very small, SE algorithm has the risk of mistakenly deleting SVs.

In order to explore the accuracy of training SVM with sample subsets obtained by different reduction algorithms, we conducted another group of experiments on artificial dataset. In this group of experiments, our algorithm controls R at about 0.3 by adjusting the density threshold ξ , and further reduces the size of training set by cosine angle distance analysis. Finally, the sample subset only accounts for 0.09 of the original dataset. Other algorithms also try to keep the same dataset size as our algorithm by controlling the parameters. Finally, we use the obtained sample subset to train SVM, and draw the classification hyperplane and SVs. The results are shown in Fig. 8.

Fig. 8(a) is the decision plane trained by all training samples. It can be seen from Fig. 8(b) that the PDGA algorithm well preserves the boundary points of each class. The SVs in Fig. 8(b) is highly consistent with the SVs in Fig. 8(a), and the hyperplanes in Fig. 8(a) and Fig. 8(b) are almost identical. Fig. 8(c) is the result of SR-DNA algorithm, it does not perform well for the samples whose centroid is not in the geometric center, and the classification plane is quite different from the real plane. Fig. 8(d) is the result of PSCC algorithm. It is worth noting that this algorithm is designed based on the cosine distance from each class of samples to the centroid of different classes of samples, and almost loses most of SVs in original training set. The processing result of SE algorithm is shown in Fig. 8(e). The algorithm is based on Euclidean distance to delete redundant

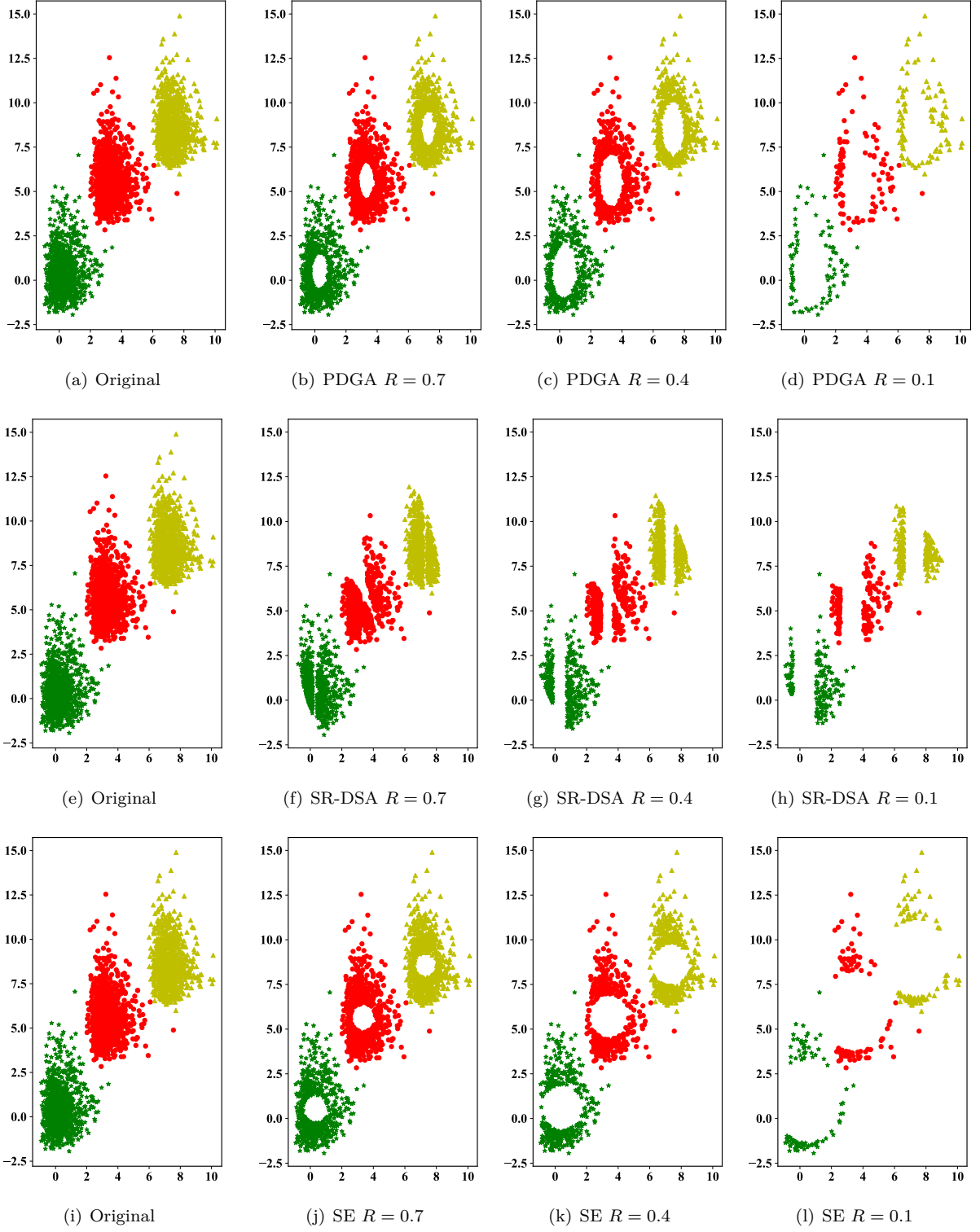


Figure 7: Sample subsets with different retention rates.

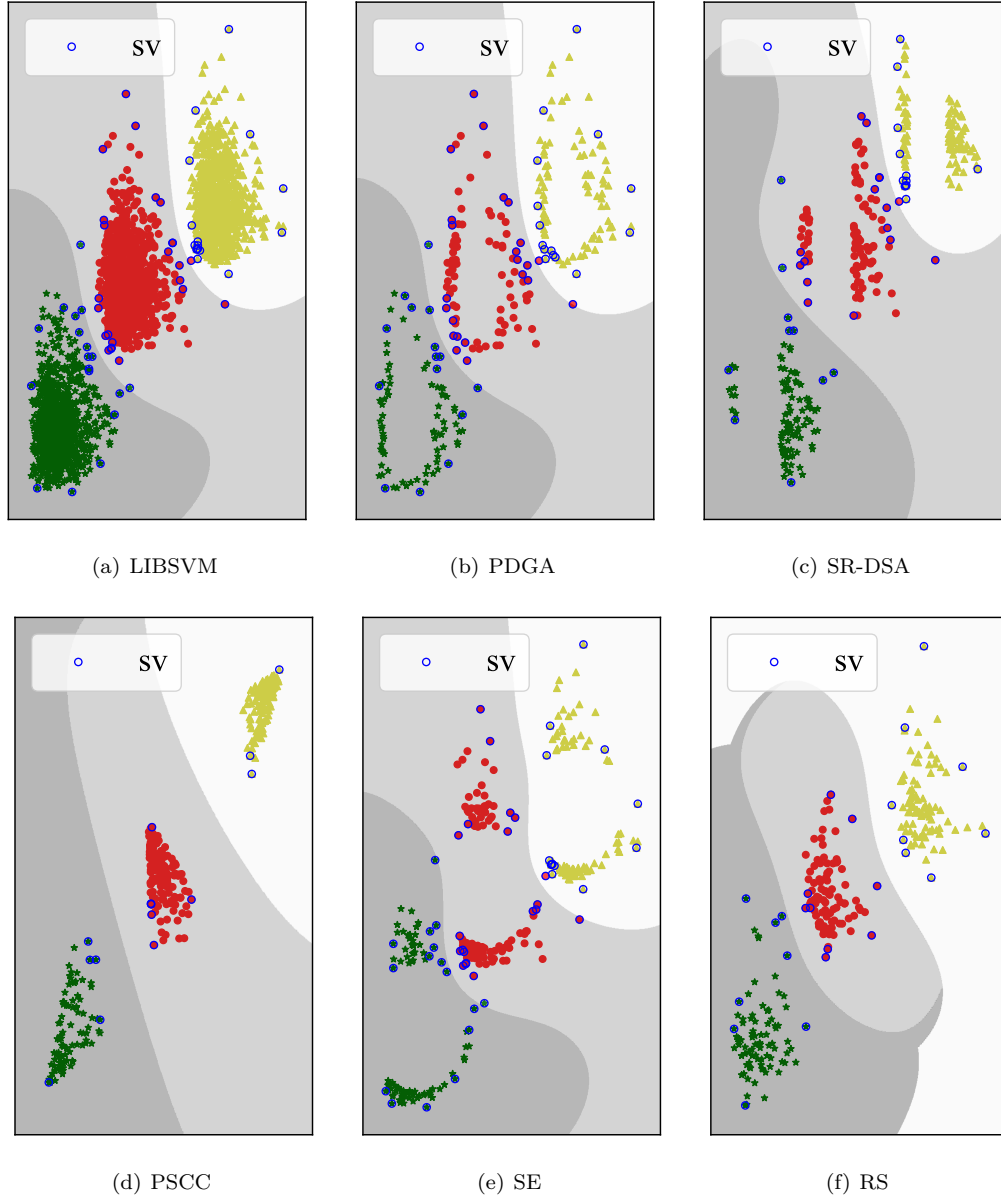


Figure 8: Performance comparison of several reduction algorithms on 2-dimensional artificial datasets.

samples, which is not effective for datasets with large differences in the distribution range of each dimension. It removes some SVs mistakenly on the boundary. We also use random sampling algorithm to extract a part of samples to train SVM, and the results are shown in Fig. 8(f). This method loses most of the SVs in original dataset.

We also counted the number of mis-classified samples of each geometric analysis algorithm on the testing set, and the experimental results are shown in Table 3. It is noted that the classification accuracy of PDGA algorithm is consistent with LIBSVM when R is only 0.09, while the classification accuracy of PSCC algorithm is seriously reduced.

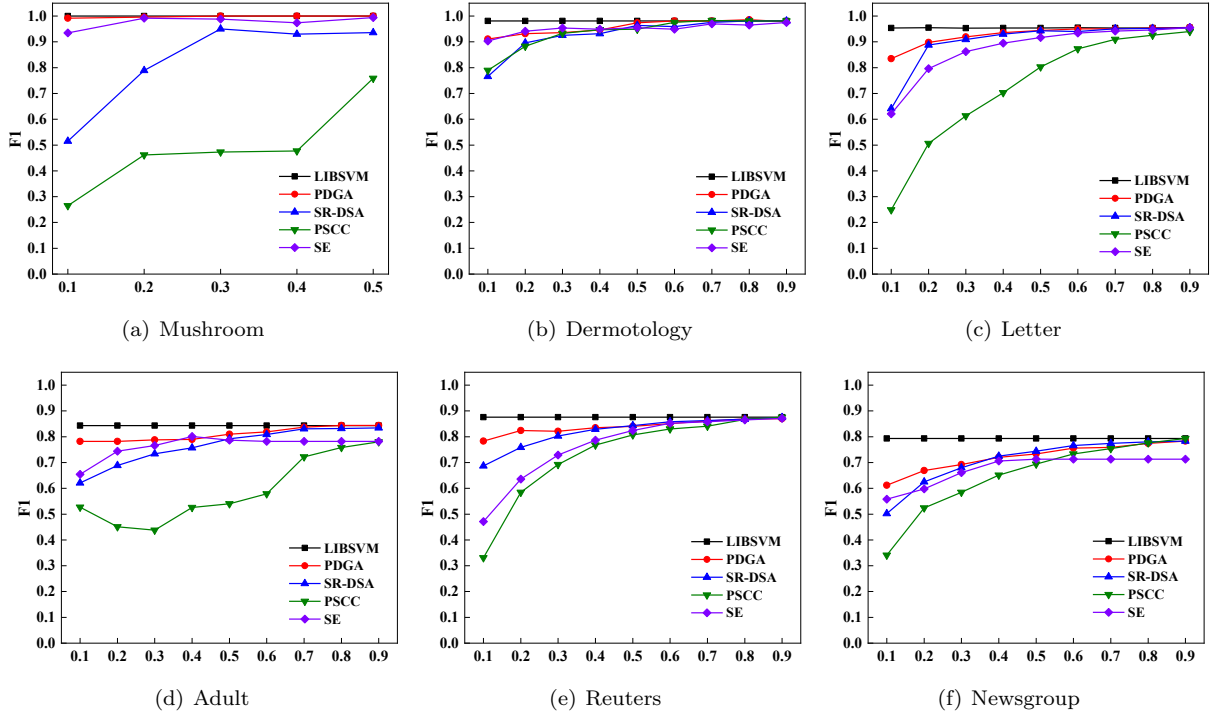
4.3. Experiments on Real Datasets

On real datasets, we will study the ability of the algorithm to maintain classification accuracy and the total time of each algorithm running. In order to facilitate comparison, we adjust the parameters of several algorithms to observe the classification accuracy and running time of several algorithms at the same level of sample retention R . Since the SR-DSA algorithm based on clustering can reduce the amount of data by at least half in mushroom dataset, thus R of all algorithms traverse 0.5, 0.4, \dots , 0.1. On the other datasets, let R traverse 0.9, 0.8, \dots , 0.1.

Table 3: Performance comparison of several algorithms on artificial datasets.

Algorithms	# reserved samples	# mis-classified
LIBSVM	3000	2
PDGA	264	2
SR-DSA	252	3
PSCC	282	48
SE	255	5

Firstly, the ability of each algorithm to maintain the classification accuracy in the process of reducing samples is compared. Experimental results on several real datasets are shown in Fig. 9. The horizontal axis represents the sample retention rate R , and the vertical axis represents the change of F_1 score on the testing set. It should be noted that the LIBSVM shown in the figure is always the result of training with all data.

Figure 9: Comparison of F_1 score of several SVM sample reduction algorithms on real datasets.

It is easy to see from Fig. 9(a) that the PDGA algorithm has outstanding performance on Mushroom dataset. When R drops to 0.1, the classification accuracy of PDGA algorithm is consistent with that of the LIBSVM trained with all the samples, while the accuracy of SE algorithm decreases by 0.1, the SR-DSA by 0.5, and the PSCC algorithm by about 0.75.

Fig. 9(b) shows that when R is higher than 0.5, the classification accuracy of several algorithms on Dermatology dataset is almost unaffected. However, when $R < 0.5$, the classification accuracy of several algorithms has a downward trend, but PDGA algorithm and SE algorithm can still maintain the classification accuracy above 0.9 when $R = 0.1$. SR-DSA drops to about 0.8 and PSCC to about 0.75.

Fig. 9(c) shows the experimental results on Letter dataset. When $R \geq 0.5$, the classification performance of PDGA, SR-DSA and SE algorithms are very close to that of LIBSVM. When $R \geq 0.2$, the F_1 score of PDGA and SR-DSA algorithms are not significantly different, and the F_1 score is maintained at about 0.9. When $R = 0.1$, the F_1 score of several algorithms all showed a certain decrease, but PDGA still ranked first, with its F_1 score above 0.83, while the other algorithms were reduced to below 0.7.

On Adult dataset (see Fig. 9(d)), the change trend of PDGA was relatively slow. When $R \geq 0.5$, the F_1 score of PDGA and SR-DSA was not affected much, while when R decreased to 0.4, the F_1 score of PDGA

algorithm decreased slightly. When R continues to decrease, the F_1 score of PDGA does not change much. However, the F_1 score of other algorithms continue to decline, finally dropped below 0.7.

Fig. 9(e) and 9(f) are the results of all algorithms on 2 high-dimensional text datasets. On Reuters dataset, when $R \geq 0.5$, the F_1 score of PDGA and SR-DSA are almost unaffected. On 20-Newsgroups dataset, F_1 score of PDGA and SR-DSA are almost unaffected when $R \geq 0.6$. When $R < 0.5$, the F_1 score of several algorithms decrease rapidly, but PDGA algorithm decrease the least. On Reuters dataset, the F_1 score of PDGA only decreases by about 0.1 when $R = 0.1$, but other algorithms drop at least 0.2. On 20-NewGroups dataset, the PDGA algorithm decreases by about 0.2 when $R = 0.1$, while the other algorithms decrease by at least 0.25.

From the experimental results on several datasets, PDGA algorithm maintain ability of classification accuracy is better than that of several other algorithms. In most datasets, PDGA algorithm can reduce the sample size by about 50%, while maintain almost no impact on the classification accuracy. Even as R dropped to 0.1, PDGA algorithm can still obtain a better classification accuracy than several other algorithms. However, PSCC algorithm performs poorly on these datasets.

We also compared the running time of several algorithms, including the time spent on the sample reduction process and on training SVM. The experimental results are shown in Fig. 10.

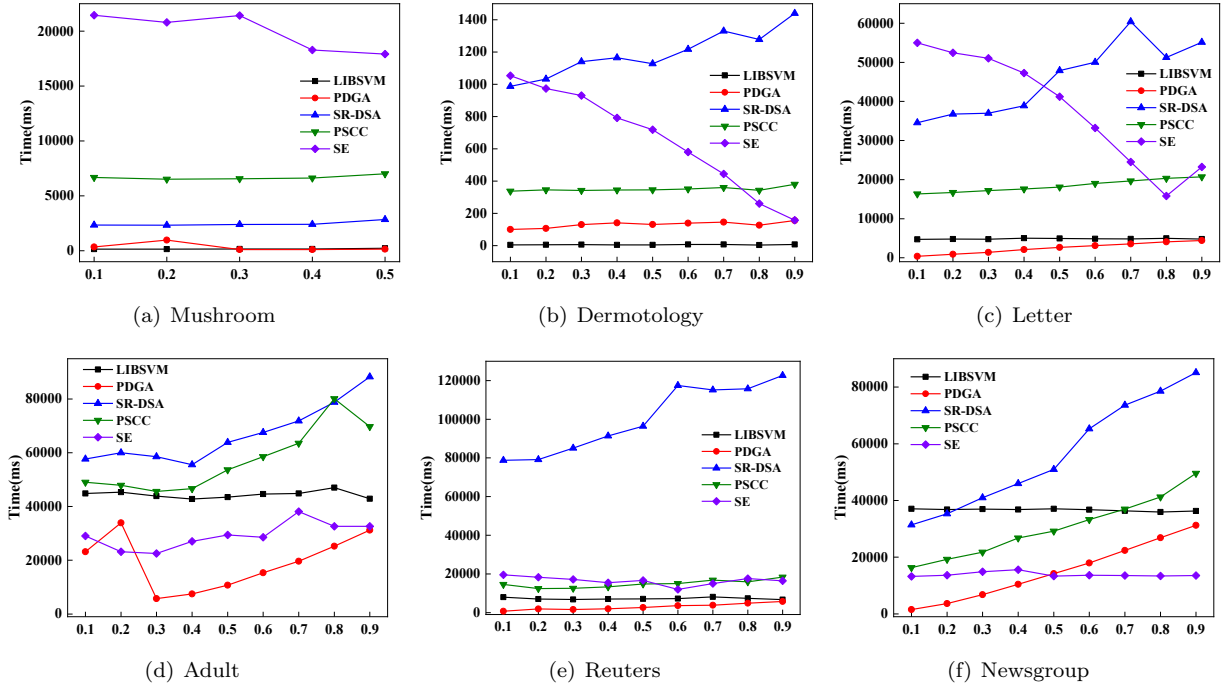


Figure 10: The total running time of several algorithms.

As can be seen from Fig. 10(b), for Dermatology dataset with the minimum sample size (only about 286 training samples), it is faster to train SVM directly with all the data. Figure 10(a) shows that the running time of PDGA and LIBSVM is almost the same. However, with the increase of sample size and sample dimension, the advantage of sample reduction algorithm becomes prominent. On several datasets with a sample size of over 8000, the total time consumed by PDGA was less than that of LIBSVM. In particular, almost all sample reduction algorithms consume less time than LIBSVM in the high-dimensional and large-scale 20-Newsgroups dataset as the R decreases. This also reflects the advantage of SVM with sample reduction on large-scale datasets.

It is noted that the time consumed by SE algorithm in Fig. 10(a), 10(b) and 10(c) increases with the decrease of R . The reason for this phenomenon is that SE algorithm needs several iterations to reduce samples. The smaller R is, the more iterations are, which increases the running time. In Fig. 10(d), 10(e) and 10(f), the SE algorithm consumes less time to reduce the sample size than training SVM with the reduced dataset, so there is no increase in the consumption time. Fig. 10 shows that the clustering-based SR-DSA algorithm consumes much time on most datasets, which is also the biggest disadvantage of the clustering-based sample reduction

algorithm, especially when it is applied to the clustering problem of large-scale datasets.

In order to quantify the efficiency of PDGA, we count the maximum sample size that PDGA can reduce when the decrease degree of F_1 score is less than 0.02, and calculate the proportion of PDGA running time to LIBSVM running time. The results are shown in Table 4.

Table 4: The time saving of PDGA algorithm without affecting the classification accuracy.

Datasets	LIBSVM			PDGA		
	F_1	Time(ms)	R	F_1	Time(ms)	Rate of running time
Dermatology	0.981	3.8	0.5	0.974	131	—
Letter	0.954	4937	0.5	0.944	2643.6	53.55%
Mushroom	1	226.8	0.1	1	150	66.14%
Adult	0.843	43472	0.7	0.837	19668	45.24%
Reuters	0.876	7083	0.7	0.869	3813	53.83%
20-NewsGroups	0.793	37052	0.8	0.774	26869	72.52%

It can be seen from Table 4 that it takes less time to train SVM directly with all samples in very small sample size of Dermatology dataset. On mushroom dataset, PDGA algorithm can reduce 90% of the samples, while the classification accuracy does not changed, and the running time of the algorithm only accounts for 66.14% of LIBSVM. On Adult dataset, the running time of PDGA only accounts for 45% of LIBSVM, while the F_1 score only decreases by 0.006, which almost has no impact on the classification accuracy. On Letter dataset, PDGA can remove 50% of the training samples, which takes only half of LIBSVM time, and the classification accuracy is only reduced by 0.01. The experimental results on 2 high-dimensional text datasets also show that PDGA algorithm can save nearly 30%-50% of LIBSVM time without significantly reducing classification accuracy.

5. Conclusion and Future Work

This paper proposed a parallel data geometry analysis method (PDGA) for SVM sample reduction to extract the data geometry structure of SVM training set. It can minimize the number of samples without reducing the classification accuracy significantly, thus improving the training speed of SVM. The experimental results on artificial dataset shows that PDGA can retain the potential SVs well on all kinds of boundaries in training set, even for datasets with irregular geometric distribution. The experimental results on 6 real datasets show that PDGA is always the best compared with the same kind of algorithms in maintaining classification accuracy, and the time consumed is also very small. In the best case, more than 90% samples can be reduced without affecting the classification accuracy.

In the future work, we will try to extend the algorithm on high-dimensional and large-scale datasets, to simultaneously reduce redundant samples and invalid features, thus achieve effective expansion.

References

- [1] J. Cervantes, F. Garcia-Lamont, L. Rodríguez-Mazahua, A. Lopez, A comprehensive survey on support vector machine classification: Applications, challenges and trends, *Neurocomputing* 408 (2020) 189–215. doi:<https://doi.org/10.1016/j.neucom.2019.10.118>.
- [2] Y. Wang, Z. Wang, Q. Hu, Y. Zhou, H. Su, Hierarchical semantic risk minimization for large-scale classification, *IEEE Transactions on Cybernetics PP* (99) (2021) 1–13.
- [3] Multi independent latent component extension of naive bayes classifier, *Knowledge-Based Systems* 213 (2021) 106646. doi:<https://doi.org/10.1016/j.knosys.2020.106646>.
- [4] L. Jiang, C. Li, S. Wang, L. Zhang, Deep feature weighting for naive bayes and its application to text classification, *Engineering Applications of Artificial Intelligence* 52 (Jun.) (2016) 26–39. doi:[10.1016/j.engappai.2016.02.002](https://doi.org/10.1016/j.engappai.2016.02.002).

- 400 [5] R. J. Prokop, A. P. Reeves, A survey of moment-based techniques for unoccluded object representation and recognition, *Cvgip Graphical Models & Image Processing* 54 (5) (1992) 438–460. doi:10.1016/1049-9652(92)90027-U.
- [6] A. Trabelsi, Z. Elouedi, E. Lefevre, Decision tree classifiers for evidential attribute values and class labels, *Fuzzy Sets and Systems* 366 (11 2018). doi:10.1016/j.fss.2018.11.006.
- [7] M. Fratello, R. Tagliaferri, Decision trees and random forests, *Encyclopedia of Bioinformatics and Computational Biology* 1 (2019) 374–383. doi:https://doi.org/10.1016/B978-0-12-809633-8.20337-3.
- [8] T. Leonard, Logistic Regression, 2020, pp. 139–152. doi:10.1201/9781003073109-8.
- [9] P. Skryjomski, B. Krawczyk, A. Cano, Speeding up k-nearest neighbors classifier for large-scale multi-label learning on gpus, *Neurocomputing* 354 (2019) 10–19. doi:10.1016/j.neucom.2018.06.095.
- 410 [10] V. Vapnik, R. Izmailov, Reinforced svm method and memorization mechanisms, *Pattern Recognition* (2021) 108018doi:https://doi.org/10.1016/j.patcog.2021.108018.
- [11] V. N. Vapnik, Statistical learning theory, *Encyclopedia of the ences of Learning* 41 (4) (1998) 3185–3185. doi:10.1007/978-1-4419-1428-6_5864.
- [12] C. J. C. Burges, A tutorial on support vector machines for pattern recognition, *Data Mining and Knowledge Discovery* 2 (2) (1998) 121–167. doi:10.1023/A:1009715923555.
- [13] N. Cristianini, J. Shawe-Taylor, An introduction to support vector machines and other kernel-based learning methods: Preface (2000). doi:10.1017/CB09780511801389.001.
- [14] T. K. Bhowmik, P. Ghanty, A. Roy, S. K. Parui, Svm-based hierarchical architectures for handwritten bangla character recognition, *International Journal on Document Analysis & Recognition* 12 (2) (2009) 97–108. doi:10.1007/s10032-009-0084-x.
- 420 [15] V. Naik, A. Desai, Online handwritten gujarati character recognition using svm, mlp, and k-nn, 2017, pp. 1–6. doi:10.1109/ICCCNT.2017.8203926.
- [16] X. Liang, L. Zhu, D.-S. Huang, Multi-task ranking svm for image cosegmentation, *Neurocomputing* 247 (03 2017). doi:10.1016/j.neucom.2017.03.060.
- [17] Y. Tang, Deep learning using support vector machines, *CoRR abs/1306.0239* (2013). arXiv:1306.0239. URL <http://arxiv.org/abs/1306.0239>
- [18] Y. Chen, Z. Lin, X. Zhao, G. Wang, Y. Gu, Deep learning-based classification of hyperspectral data, *Selected Topics in Applied Earth Observations and Remote Sensing, IEEE Journal of* 7 (2014) 2094–2107. doi:10.1109/JSTARS.2014.2329330.
- 430 [19] P. Liu, K.-K. R. Choo, L. Wang, F. Huang, Svm or deep learning? a comparative study on remote sensing image classification, *Soft Computing* 21 (12 2017). doi:10.1007/s00500-016-2247-2.
- [20] J. Nalepa, M. Kawulok, Adaptive memetic algorithm enhanced with data geometry analysis to select training data for svms, *Neurocomputing* 185 (12 2015). doi:10.1016/j.neucom.2015.12.046.
- [21] J. Qiu, Q. Wu, G. Ding, Y. Xu, S. Feng, A survey of machine learning for big data processing, *EURASIP Journal on Advances in Signal Processing* 2016 (05 2016). doi:10.1186/s13634-016-0355-x.
- [22] T. Joachims, Making large-scale svm learning practical, *Advances in Kernel methods - support vector learning*: 169 - 184 (11 1998). doi:10.17877/DE290R-14262.
- [23] Y. Ma, X. Liang, G. Sheng, J. Kwok, M. Wang, G. Li, Noniterative sparse ls-svm based on globally representative point selection, *IEEE Transactions on Neural Networks and Learning Systems* PP (2020) 1–11. doi:10.1109/TNNLS.2020.2979466.
- 440

- [24] J. Platt, Sequential minimal optimization: A fast algorithm for training support vector machines, *Advances in Kernel Methods-Support Vector Learning* 208 (07 1998).
- [25] G. Galvan, M. Lapucci, C.-J. Lin, M. Sciandrone, A two-level decomposition framework exploiting first and second order information for svm training problems, *Journal of Machine Learning Research* 22 (23) (2021) 1–38.
URL <http://jmlr.org/papers/v22/19-632.html>
- [26] Libsvm, *Acm Transactions on Intelligent Systems & Technology* (2012). doi:10.1145/1961189.1961199.
- [27] H. Graf, E. Cosatto, L. Bottou, I. Durdanovic, V. Vapnik, *Parallel support vector machines: The cascade svm.*, 2004.
- 450 [28] B.-L. Lu, K.-A. Wang, Y. Wen, Comparison of parallel and cascade methods for training support vector machines on large-scale problems, 2004, pp. 3056 – 3061 vol.5. doi:10.1109/ICMLC.2004.1378557.
- [29] B. Scholkopf, A. J. Smola, *Learning with Kernels: Support Vector Machines, Regularization, Optimization, and Beyond*, MIT Press, Cambridge, MA, USA, 2001.
- [30] Y.-J. Lee, S.-Y. Huang, Reduced support vector machines: A statistical theory, *IEEE transactions on neural networks / a publication of the IEEE Neural Networks Council* 18 (2007) 1–13. doi:10.1109/TNN.2006.883722.
- [31] F. Cheng, J. Chen, J. Qiu, L. Zhang, A subregion division based multi-objective evolutionary algorithm for svm training set selection, *Neurocomputing* 394 (02 2020). doi:10.1016/j.neucom.2020.02.028.
- [32] J. Nalepa, M. Kawulok, Selecting training sets for support vector machines: a review, *Artificial Intelligence Review* 52 (08 2019). doi:10.1007/s10462-017-9611-1.
- 460 [33] L. Guo, S. Boukir, Fast data selection for svm training using ensemble margin, *Pattern Recognition Letters* 51 (2015) 112–119. doi:10.1016/j.patrec.2014.08.003.
- [34] Y. Lin, F. Lv, S. Zhu, M. Yang, T. Cour, K. Yu, L. Cao, Large-scale image classification: Fast feature extraction and svm training, Vol. 1, 2011, pp. 1689–1696. doi:10.1109/CVPR.2011.5995477.
- [35] A. Lyhyaoui, M. Martinez-Ramon, I. Jiménez, M. Vazquez, J. L. Sancho-Gómez, A. Figueiras-Vidal, Sample selection via clustering to construct support vector-like classifiers, *Neural Networks, IEEE Transactions on* 10 (1999) 1474 – 1481. doi:10.1109/72.809092.
- [36] G. GATES, The reduced nearest neighbor rule, *IEEE Transactions on Information Theory - TIT* (05 1972).
- [37] M. Kawulok, J. Nalepa, Support vector machines training data selection using a genetic algorithm, in: G. Gimel'farb, E. Hancock, A. Imiya, A. Kuijper, M. Kudo, S. Omachi, T. Windeatt, K. Yamada (Eds.), *Structural, Syntactic, and Statistical Pattern Recognition*, Springer Berlin Heidelberg, Berlin, Heidelberg, 2012, pp. 557–565. doi:10.1007/978-3-642-34166-3_61.
- 470 [38] D. Musicant, A. Feinberg, Active set support vector regression, *Neural Networks, IEEE Transactions on* 15 (2004) 268 – 275. doi:10.1109/TNN.2004.824259.
- [39] F. Alamdar, S. Ghane, A. Amiri, On-line twin independent support vector machines, *Neurocomputing* 186 (01 2016). doi:10.1016/j.neucom.2015.12.062.
- [40] D. Wilson, T. Martinez, Reduction techniques for instance-based learning algorithms, *Machine Learning* 38 (2000) 257–286. doi:10.1023/A:1007626913721.

- 480 [41] N. Jankowski, M. Grochowski, Comparison of instances selection algorithms i. algorithms survey, in: L. Rutkowski, J. H. Siekmann, R. Tadeusiewicz, L. A. Zadeh (Eds.), *Artificial Intelligence and Soft Computing - ICAISC 2004*, Springer Berlin Heidelberg, Berlin, Heidelberg, 2004, pp. 598–603. doi:10.1007/978-3-540-24844-6_90.
- [42] J. Balcázar, Y. Dai, O. Watanabe, A random sampling technique for training support vector machines, 2001, pp. 119–134. doi:10.1007/3-540-45583-3_11.
- [43] Zhu, Yang, Gao, GB, Yin, TM, Neighbors’ distribution property and sample reduction for support vector machines, *APPL SOFT COMPUT* (2014). doi:10.1016/j.asoc.2013.12.009.
- [44] X. Li, J. Cervantes, W. Yu, Fast classification for large data sets via random selection clustering and support vector machines, *Intell. Data Anal.* 16 (6) (2012) 897–914.
- 490 [45] S. Abe, T. Inoue, Fast training of support vector machines by extracting boundary data, 2001. doi:10.1007/3-540-44668-0_44.
- [46] Hart, P., The condensed nearest neighbor rule (corresp.) 14 (3) (1968) 515–516. doi:10.1109/tit.1968.1054155.
- [47] H. Shin, S. Cho, Neighborhood property-based pattern selection for support vector machines, *Neural computation* 19 (2007) 816–55. doi:10.1162/neco.2007.19.3.816.
- [48] X. Jiantao, M. He, W. Yuying, F. Yan, A fast training algorithm for support vector machine via boundary sample selection, 2004, pp. 20 – 22 Vol.1. doi:10.1109/ICNNSP.2003.1279203.
- [49] R. Pighetti, D. Pallez, F. Precioso, Improving svm training sample selection using multi-objective evolutionary algorithm and lsh, 2015. doi:10.1109/SSCI.2015.197.
- [50] A. Vlachos, Active learning with support vector machines (2021) 313–326doi:10.1002/widm.1132.
- 500 [51] R. Wang, S. Kwong, Sample selection based on maximum entropy for support vector machines, Vol. 3, 2010, pp. 1390–1395. doi:10.1109/ICMLC.2010.5580848.
- [52] W. Wang, Z. Xu, A heuristic training for support vector regression, *Neurocomputing* 61 (2004) 259–275. doi:10.1016/j.neucom.2003.11.012.
- [53] D. Wang, L. Shi, Selecting valuable training samples for svms via data structure analysis, *Neurocomputing* 71 (2008) 2772–2781. doi:10.1016/j.neucom.2007.09.008.
- [54] C. Liu, W. Wang, M. Wang, F. Lv, M. Konan, An efficient instance selection algorithm to reconstruct training set for support vector machine, *Knowledge-Based Systems* 116 (2017) 58–73. doi:https://doi.org/10.1016/j.knosys.2016.10.031.
- 510 [55] X.-J. Shen, L. Mu, Z. Li, H.-X. Wu, J. Gou, X. Chen, Large-scale support vector machine classification with redundant data reduction, *Neurocomputing* 172 (05 2015). doi:10.1016/j.neucom.2014.10.102.
- [56] K. Teeyapan, N. Theera-Umpon, S. Auephanwiriyakul, Ellipsoidal support vector data description, *Neural Computing and Applications* 28 (12 2017). doi:10.1007/s00521-016-2343-3.
- [57] L. Yu, W. Yi, D. He, y. Lin, Fast reduction for large-scale training data set, *Journal of Southwest Jiaotong University* 42 (08 2007). doi:10.1016/S1874-8651(08)60023-X.