

Cone Cluster Labeling for Support Vector Clustering

Sei-Hyung Lee

Department of Computer Science
University of Massachusetts Lowell
MA 01854, U.S.A.
slee@cs.uml.edu

Karen M. Daniels

Department of Computer Science
University of Massachusetts Lowell
MA 01854, U.S.A.
kdaniels@cs.uml.edu

Abstract

Clustering forms natural groupings of data points that maximize intra-cluster similarity and minimize inter-cluster similarity. Support Vector Clustering (SVC) is a clustering algorithm that can handle arbitrary cluster shapes. One of the major SVC challenges is a cluster labeling performance bottleneck. We propose a novel cluster labeling algorithm that relies on approximate coverings both in feature space and data space. Comparison with existing cluster labeling approaches suggests that our strategy improves efficiency without sacrificing clustering quality.

1 Introduction

1.1 Clustering Overview Clustering is a natural grouping or unsupervised classification of data into groups [6]. Most clustering algorithms use one or a combination of the following techniques: graph-based [5, 9, 10, 7], density-based [8], model-based methods using either a statistical approach or a neural network approach, or optimization of a clustering criterion function. Constructing cluster boundaries is another popular technique [10, 11]. SVC is similar to a boundary-finding clustering method except that it only finds certain points (called *support vectors*) on the boundary of each cluster.

1.2 Support Vector Clustering As with support vector machines [4], SVC uses a nonlinear mapping of the data from the data space into a high-dimensional feature space. Whereas support vector machines use a linear separator in feature space in order to separate and classify data points, SVC uses a minimal hypersphere encompassing feature space images of data points [2, 3]. The minimal hypersphere corresponds to several disjoint bounded regions in the data space that are contours and are interpreted as clusters. The cluster labeling challenge of SVC is to associate each data point with a cluster, using only the available operations, without explicitly constructing the contours in data space.

Given a finite set $\mathcal{X} \subseteq \mathcal{R}^d$ of N distinct data points, the minimal hypersphere of radius R enclosing all data points' images in the feature space can be described by the following, as in [2]:

$$(1.1) \quad \|\Phi(x) - a\|^2 \leq R^2 \quad \forall x \in \mathcal{X},$$

where Φ is a nonlinear mapping from data space to feature space, $\Phi(x)$ is the feature space image of data point x , $\|\cdot\|$ is the Euclidean norm, and a is the center of the sphere. Therefore, data points can be categorized into three groups based on the location of their images in feature space: (1) points whose images are on the surface of the minimal hypersphere are *Support Vectors* (SVs), (2) points whose images are outside of the minimal hypersphere are *Bounded Support Vectors* (BSVs), and (3) points whose images are inside the minimal hypersphere.

The mapping from data space to feature space is determined by a kernel function, $K : \mathcal{X} \times \mathcal{X} \rightarrow \mathcal{R}$, that defines the inner product of image points. As in other SVC literature [1, 2], we use a Gaussian kernel given by Eq. 1.2 below:

$$(1.2) \quad K(x_i, x_j) = e^{-q\|x_i - x_j\|^2} = \langle \Phi(x_i) \cdot \Phi(x_j) \rangle,$$

where q is the width of the Gaussian kernel. From Eq. 1.2, it immediately follows that,

$$(1.3) \quad K(x_i, x_i) = 1.$$

This implies that all data point images are on the surface of the unit ball in feature space.

Cluster labeling, the focus of this paper, has a goal of grouping together data points that are in the same contour obtained from the minimal hypersphere. Prior work is discussed in Section 1.3. In that work, cluster labeling has been an SVC performance bottleneck.

1.3 Cluster Labeling Traditional SVC cluster labeling algorithms are Complete Graph (CG) [2], Support Vector Graph (SVG) [2], Proximity Graph (PG)

[7], and Gradient Descent (GD) [1]. All these algorithms group together data points by representing pairs of data points using an adjacency structure (typically a matrix). Each element records whether there is sufficient evidence to conclude that the associated pair of data points is in the same contour and therefore the same cluster. Each connected component in the adjacency matrix is regarded as a cluster. CG requires a $O(N^2)$ sized adjacency structure because it represents all data point pairs. SVG represents pairs in which one point is a support vector, so its adjacency structure only uses $O(N_{sv}N)$ space, where N_{sv} is the number of support vectors. PG forms an adjacency structure from a proximity graph that has only $O(N)$ edges. The GD method finds Stable Equilibrium Points (SEPs) that are the nearest minimum point (of Eq. 1.4 below) for each data point and then tests pairs of SEPs. GD's adjacency structure therefore uses $O(N_{sep}^2)$ space, where N_{sep} is the number of SEPs.

$$\begin{aligned}
 R^2(x) &= \|\Phi(x) - a\|^2 \\
 (1.4) \quad &= 1 - 2 \sum_j \beta_j K(x_j, x) + \sum_{i,j} \beta_i \beta_j K(x_i, x_j),
 \end{aligned}$$

where $x \in \mathcal{R}^d$, β_i is the Lagrange multiplier for x_i , and $a = \sum_i \beta_i \Phi(x_i)$ is the center of the minimal hypersphere [2].

Deciding if a pair of data points x_i and x_j is in the same contour is problematic because data space contours cannot be explicitly constructed due to the nature of the feature space mapping Φ . The aforementioned algorithms therefore use an indirect strategy that relies on the fact that every path connecting two data points in different data space contours exits the contours in data space and the image of the path in feature space exits the minimal hypersphere. The algorithms use, as a path in data space, the straight line segment $\overline{x_i x_j}$ connecting x_i and x_j . The line segment $\overline{x_i x_j}$ is sampled. For sample point x' in data space, $\Phi(x')$ is outside the minimal hypersphere if $R^2(x') > R^2$, where $R^2(x')$ is defined by Eq. 1.4. If the image of every sample point along $\overline{x_i x_j}$ is inside the minimal hypersphere, then x_i and x_j are determined to be in the same contour (hence cluster).

Unfortunately, the sampling approach creates a running time versus accuracy tradeoff. If m is the number of sample points along a line segment, then solving Eq. 1.4 for each sample point introduces a multiplicative factor of mN_{sv} beyond the time proportional to the size of the adjacency structure. The algorithms limit m to be a small constant (typically 10 to 20) in order to limit running time. However, small values of m can also cause two types of errors, false positive and false negative, as shown in Figure 1.

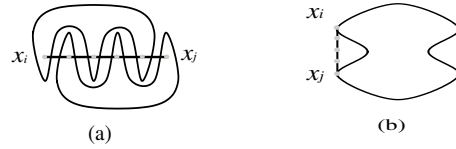


Figure 1: Problems in using line segment $\overline{x_i x_j}$, depicted in data space: (a) sample points on $\overline{x_i x_j}$ are all inside the minimal hypersphere although x_i and x_j are in different contours; (b) all sample points are outside the minimal hypersphere although x_i and x_j are in the same contour.

1.4 Overview This paper avoids the cluster labeling problems listed in Section 1.3 by using a novel approach that decides if two data points are in the same cluster without sampling a path between them. The main idea of this paper is to try to cover a key portion of the minimal hypersphere in feature space using cones that are anchored at each support vector in the feature space and also correspond to hyperspheres in data space. The union of the hyperspheres forms an approximate covering of the data space contours. The union need not be constructed; pairs of support vectors can be quickly tested during the cluster labeling process and then the remaining data points can be easily clustered. The algorithm¹, presented in Section 2, does not use sample points. It works for data sets of arbitrary dimension, and has been tested for up to 200 dimensions. The results in Section 3 show that our new algorithm is faster than traditional cluster labeling algorithms.

2 Cone Cluster Labeling

2.1 High-Level Approach We provide a new cluster labeling algorithm that is quite different from traditional methods. For a given kernel width value, our method does not sample a path between two data points in order to judge if they belong in the same cluster. Instead, we leverage the geometry of the feature space to help perform cluster labeling in the data space. First we find an approximate covering for the minimal hypersphere in feature space. This is described in Section 2.2. Strictly speaking, the covering is not for the minimal hypersphere but for the intersection P between the surface of the unit ball and the minimal hypersphere.

The approximate covering consists of a union of cone-shaped regions. One region is associated with each support vector's feature space image. Let $V = \{v_i | v_i \text{ is a support vector, } 1 \leq i \leq N_{sv}\}$ be the set of SVs for a given q value. The region for support vector v_i is called a support vector cone and is denoted by \mathcal{E}_{v_i} .

¹An earlier version of this work appears in [14].

We call our algorithm ‘Cone Cluster Labeling’ (CCL) because of its reliance on cones.

Let B be the surface of the feature space unit ball. A hypersphere S_{v_i} centered on v_i in the data space maps to a subset of $\mathcal{E}_{v_i} \cap B$. Section 2.3 derives the radius of \mathcal{S}_{v_i} , which is shown to be the same for all support vectors. Having only one radius contributes to the speed of CCL. The union $\cup_i(\mathcal{S}_{v_i})$ is an approximate covering of the data space contours P' , where $\Phi(P') \subseteq P$. However, the union is not explicitly constructed. Rather, cluster labeling is done in two phases, as described in Section 2.4. The first phase clusters support vectors while the second clusters the remaining data points. We regard two support vectors v_i and v_j as connected if their hyperspheres overlap: $(\mathcal{S}_{v_i} \cap \mathcal{S}_{v_j}) \neq \emptyset$. Forming the transitive closure of the connected relation yields a set of support vector clusters. The final step is a clustering for data points that are not support vectors. For each such data point, we assign the closest support vector’s cluster label. Proofs are in [14].

2.2 Approximate Covering in Feature Space

This section forms a collection of support vector cones that approximately cover P . Let $\Theta_i = \angle(\Phi(v_i)Oa)$, where O is the feature space origin and a is the center of the minimal hypersphere (see Figure 2(a)). In feature space, each support vector has its own cone \mathcal{E}_{v_i} that covers a part of P . We define the support vector cone \mathcal{E}_{v_i} to be the infinite cone with axis $\overrightarrow{\Phi(v_i)}$ and base angle Θ_i . Lemma 2.1 below shows that $\Theta_i = \Theta_j = \Theta$ for all $v_i, v_j \in V$.

LEMMA 2.1. $\angle(\Phi(v_i)Oa) = \angle(\Phi(v_j)Oa), \forall v_i, v_j \in V$.

To approximately cover P , we denote by a' the intersection of \vec{a} with the surface of the unit ball (see Figure 2(b)). The point a' is a common point of intersection for all the support vector cones. Thus, $((\cup_i(\mathcal{E}_{v_i})) \cap P) \approx P$.

2.3 Approximate Covering in Data Space

To approximately cover P' in data space using support vector cones, we find a hypersphere \mathcal{S}_{v_i} in data space associated with each \mathcal{E}_{v_i} in feature space. Since all support vector cones have the same base angle Θ , all \mathcal{S}_{v_i} have the same radius and each is centered at v_i .

LEMMA 2.2. *Each data point whose image is inside $(\mathcal{E}_{v_i} \cap P)$ is at distance $\leq \|v_i - g_i\|$ from v_i , where $g_i \in \mathcal{R}^d$ is such that $\angle(\Phi(v_i)O\Phi(g_i)) = \Theta$.*

The claim implies that $(\mathcal{E}_{v_i} \cap P)$ corresponds to a hypersphere \mathcal{S}_{v_i} in the data space centered at v_i with radius $\|v_i - g_i\|$. Since $(\cup_i(\mathcal{E}_{v_i})) \cap P$ approximately covers P , $\cup_i(\mathcal{S}_{v_i})$ approximately covers P' .

The next task is to obtain $\|v_i - g_i\|$ in data space. Because $\|\Phi(v_i)\| = 1 = \|a'\|$, $\cos(\Theta) = \cos(\angle(\Phi(v_i)O\Phi(g_i))) = \cos(\angle(\Phi(v_i)Oa')) = \langle \Phi(v_i) \cdot a' \rangle$. Thus, we can solve for $\|v_i - g_i\|$ as follows:

$$(2.5) \quad \|v_i - g_i\| = \sqrt{-\frac{\ln(\cos(\Theta))}{q}}.$$

Note that because $\Theta = \angle(\Phi(v_i)Oa')$ is the same for all $v_i \in V$, all \mathcal{S}_{v_i} have the same radii. We therefore denote $\|v_i - g_i\|$ by Z . We now need to obtain $\langle \Phi(v_i) \cdot a' \rangle$ in order to calculate $\cos(\Theta)$. To do this, we first show in Lemma 2.3 that $\langle \Phi(v_i) \cdot a' \rangle = \|a\|$, $\forall v_i \in V$ (see Figure 2 (c)). We then show that $\|a\| = \sqrt{1 - R^2}$ (see Figure 2 (d)).

LEMMA 2.3. $\langle \Phi(v_i) \cdot a' \rangle = \|a\|, \forall v_i \in V$.

COROLLARY 2.1. $\overrightarrow{\Phi(v_i)a}$ is orthogonal to \vec{a} .

LEMMA 2.4. $\langle \Phi(v_i) \cdot a \rangle = 1 - R^2 = \|a\|^2, \forall v_i \in V$.

Consequently, we have:

$$(2.6) \quad Z = \sqrt{-\frac{\ln(\sqrt{1 - R^2})}{q}}.$$

2.4 Assign Cluster Labels Table 1 below shows the CCL algorithm. For the given q value, it first computes Z using Eq. 2.6. Next, support vectors are clustered by finding connected components in the resulting adjacency structure. Connected components can be efficiently found using a standard algorithm such as Depth First Search. Each connected component corresponds to a cluster. Therefore, the output of FindConnComponents() is an array of size N that contains cluster labels. Finally, the remaining data points are clustered. Each is assigned the cluster label of the nearest support vector.

CCL(\mathcal{X}, q, V)	
compute Z for q using Eq. 2.6	
AdjacencyMatrix \leftarrow Construct connectivity matrix	
// assign cluster labels to the support vectors	
Labels \leftarrow FindConnComponents(AdjacencyMatrix)	
// assign cluster label to the rest of data	
for each $x \in \mathcal{X}$, where $x \notin V$	
$idx \leftarrow$ find the nearest SV from x	
Labels[x] \leftarrow Labels[x_{idx}]	
end for	
return Labels	

Table 1: Main algorithm of Cone Cluster Labeling

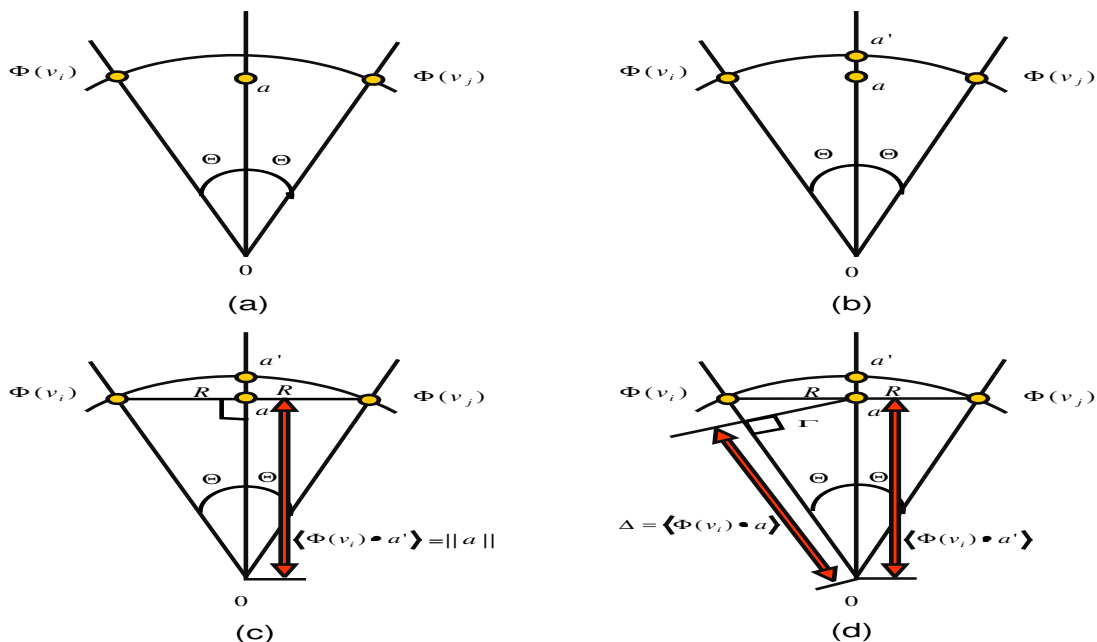


Figure 2: Developing an approximate cover of part of P , where v_i and v_j are support vectors and $\Phi(v_i)$ and $\Phi(v_j)$ are their feature space images, respectively. (a) $\Theta = \angle(\Phi(v_i)Oa) = \angle(\Phi(v_j)Oa)$. (b) a' is intersection of \vec{a} with the surface of the unit ball. (c) the length $\|a\|$ is $\langle \Phi(v_i) \cdot a' \rangle$. (d) $\Delta = \langle \Phi(v_i) \cdot a \rangle = 1 - R^2 = \|a\|^2$. Note this is only a two-dimensional illustration; the actual feature space is high-dimensional.

The worst-case asymptotic running time complexity for `FindConnComponents()` is in $O(N_{sv}^2)$. The time complexity of the CCL **for** loop is in $O((N - N_{sv})N_{sv})$. Therefore, this algorithm uses time in $O(NN_{sv})$ for each q value. Note that, unlike previous cluster labeling algorithms, this time does not depend on a number of sample points. Detailed execution time comparisons are given in Section 3.

3 Results

CCL provides high-quality clustering using less computation time than existing cluster labeling algorithms.

We compare three different types of execution times as well as total time of CCL with respect to 4 cluster labeling approaches introduced in Section 1.3: CG, SVG, PG, and GD. Times are for 1) preprocessing and constructing the adjacency matrix, 2) finding connected components from the adjacency matrix (this clusters some data points), and 3) clustering remaining data points. Worst-case asymptotic time complexity and actual running time comparisons are given in Table 2 with the data set of Figure 3.

The total of the three actual execution times is mea-



Figure 3: Left: dataset ($N=98$), right: cluster result (2 clusters). Data is a subset of data from the authors of [10].

sured and divided by the number of q values² to compute average times. Since a cluster labeling algorithm receives q as an input and produces a cluster result, the average time is appropriate as a measurement for execution time comparison. To test high-dimensional datasets without outliers or strongly overlapping clusters, we created datasets with different dimensions. CCL worked well with these high-dimensional data. Details of these and all of our two and three-dimensional experiments can be found at [14].

²A list of q values is generated by our method in [13].

