



# Efficient $k$ NN classification algorithm for big data



Zhenyun Deng, Xiaoshu Zhu\*, Debo Cheng, Ming Zong, Shichao Zhang\*

Guangxi Key Lab of Multi-Source Information Mining & Security, Guangxi Normal University, Guilin, Guangxi 541004, China

## ARTICLE INFO

### Article history:

Received 17 March 2015

Received in revised form

1 July 2015

Accepted 7 August 2015

Available online 5 February 2016

### Keywords:

Classification

$k$ NN

Big data

Data cluster

Cluster center

## ABSTRACT

$K$  nearest neighbors ( $k$ NN) is an efficient lazy learning algorithm and has successfully been developed in real applications. It is natural to scale the  $k$ NN method to the large scale datasets. In this paper, we propose to first conduct a  $k$ -means clustering to separate the whole dataset into several parts, each of which is then conducted  $k$ NN classification. We conduct sets of experiments on big data and medical imaging data. The experimental results show that the proposed  $k$ NN classification works well in terms of accuracy and efficiency.

© 2016 Elsevier B.V. All rights reserved.

## 1. Introduction

In the era of big data, it is the most important to efficiently learn from large scale in all kinds of real applications, such as classification and clustering. Therefore, it is very obvious to scale traditional classification algorithms, such as decision trees, support vector machine, Naive Bayes, neural network, and  $k$  nearest neighbors ( $k$ NN), so that these methods can be easily used in big data. Due to the simplicity, easy-understand and relatively high performance of  $k$ NN, this paper focuses on scaling the  $k$ NN classification into the application of big data [13].

The previous  $k$ NN method first selects  $k$  nearest training samples for a test sample, and then predicts the test sample with the major class among  $k$  nearest training samples. However,  $k$ NN needs to compute the distance (or similarity) of all training samples for each test sample in the process of selecting  $k$  nearest neighbors [24,25]. The high cost (i.e., linear time complexity over the sample size) prohibits the traditional  $k$ NN method to be used in big data [21]. Obviously, conducting  $k$ NN algorithm in the memory of a PC should be an interesting issue.

Inspired by the recent progress on big data, this paper devised a new  $k$ NN method for dealing with big data. Specifically, we propose to first conduct a  $k$ -means clustering to separate the whole dataset into several parts [26]. And then we select the nearest cluster as the training samples and conduct the  $k$ NN classification. The time complexity of the proposed algorithm is linear to the sample size. The experimental results on real datasets including

medical imaging datasets indicated that the proposed method achieved better performance than conventional  $k$ NN method in terms of both classification performance and time cost [15].

The rest of paper is organized as follows: in Section 2, we provide a brief review the previous fast  $k$ NN methods. We then propose the new  $k$ NN classification in Section 3. The experiment results are presented in Section 4. Finally, we give the conclusions and our future work in Section 5.

## 2. Related work

The research of  $k$ NN method has been becoming a hot research topic in data mining and machine learning since the algorithm was proposed in 1967. To apply for the traditional  $k$ NN method in big data, the previous literatures can be often categorized into two parts, i.e., fast finding the nearest samples [21] and selecting representatives samples (or removing some samples) to reduce the calculation of  $k$ NN [22]. For instance, Zhang proposed a Certainly Factor (CF) measure to deal with the unsuitability of skewed class distribution in  $k$ NN methods [14]. Li et al. proposed a density-based method for reducing the amount of training data [9]. Zhao et al. proposed a new algorithm based on the use of labeled samples and add the screening process condition, it making the new algorithm in time complexity have significantly decreased, and no significant effect on algorithm result [16]. These methods were mainly applied for fast search [17,18], dimension reduction [19,20], and improving the efficiency of the algorithms.

$k$ NN algorithm computes the distance between each training sample and test samples in the dataset and then returns  $k$  closest samples. Its time complexity is linearly and is guaranteed to find

\* Corresponding authors.

E-mail address: [zhangsc@mailbox.gxnu.edu.cn](mailto:zhangsc@mailbox.gxnu.edu.cn) (S. Zhang).

exact  $k$  nearest neighbors. However, the computational complexity of the linear search method is proportional to the size of the training dataset for each test sample, it is  $O(nd)$ , where  $n$  is the size of the training dataset and  $d$  is the dimensionality [11,12]. This complexity is expensive for big data. Since the  $k$ NN method is not training process, we propose to introduce a new training process for  $k$ NN, which blocks training dataset by a clustering algorithm with linear complexity. During the testing process, for each test sample, we find the  $k$  nearest cluster centers and conduct a clustering for each test sample, and then construct a new classification model base on each cluster [21]. In particular, the samples within a cluster has high similarity. Thus, comparing to the traditional  $k$ NN method, the proposed algorithm not only reduces the time complexity of  $k$ NN, but also does not add significantly effect on classification accuracy [22].

### 3. Method

In this section, we introduce two processes in our algorithm, namely training process and testing process. The training process is designed to select a nearest cluster for each test sample as its new training dataset, and the testing process is used to classify each test sample by  $k$ NN algorithm within its nearest cluster.

#### 3.1. Training process

Clustering is one of the fundamental technique in data mining because it can be used for database segmentation and data compression and can also be employed for data preprocessing of data mining. Clustering is designed to group a set of samples in such a way that samples in the same group (cluster) are more similar to each other than to those in other groups. That is, samples with high similarity within a cluster and low similarity between clusters.

The recent clustering methods can be parted into the following categories: density-based clustering, grid-based clustering, partitioning clustering and hierarchical clustering, respectively. Even though the previous clustering methods showed good performance, but they are limited in its applicability to big data due to their high computational complexity. To address this, this paper considers to employ a clustering algorithm satisfying the following two advantages: low complexity; scales linearly [1,5], respectively. To this end, we used the Landmark-based Spectral Clustering (LSC) [4] in this paper. The rationale of LSC is to select  $p$  ( $\ll n$ ) representative sample as landmarks and represents the original samples as the linear combinations of these landmarks. Different from that traditional spectral clustering method use the entire samples to represent each sample, the LSC significantly reduces the complexity of affinity matrix. At the same time, the complexity of solving the eigenvalue problem is reduced to scales linearly [24].

LSC tries to compress the original samples by finding a set of basis vectors and the representation with respect to the bases for each sample, i.e., searching for  $p$  representative samples. In this way, we have two simple and effective methods to select landmark sample from original sample, such as random sampling and  $k$ -means-based method. Random sampling randomly selects samples as landmark samples while the  $k$ -means-based method first conducts clustering on all samples several times (no need to convergence) and then uses the cluster centers as the landmark samples. In this paper, we repeat  $k$ -means algorithm 10 times and then use the cluster centers as the landmark samples.

First, we treat every sample as a basis vector to construct landmark matrix  $\mathbf{Z}$ . Note that LSC uses  $p$  landmarks to represent each sample  $\mathbf{X} = [\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n] \in R^{m \times n}$ . Thus, we need to find the matrix  $\mathbf{W}$  which is projection matrix of  $\mathbf{X}$  at the landmark matrix  $\mathbf{Z}$  [10]. The projection function can be defined as follows:

$$w_{ji} = \frac{K_h(\mathbf{x}_i, \mathbf{z}_j)}{\sum_{j' \in Z_{<i>}} K_h(\mathbf{x}_i, \mathbf{z}_{j'})}, \quad j \in Z_{<i>}. \quad (1)$$

Where  $\mathbf{z}_i$  is  $i$ -th column vector of  $\mathbf{Z}$ , and  $\mathbf{Z}_{<i>}$  denote a sub-matrix of  $\mathbf{Z}$  composed of  $r$  nearest landmarks of  $\mathbf{x}_i$ . Here we need  $O(pmn)$  to construct  $\mathbf{W}$ .  $K_h(\cdot)$  is a kernel function with a bandwidth  $h$ . The Gaussian kernel  $K_h(\mathbf{x}_i, \mathbf{z}_j) = \exp(-\|\mathbf{x}_i - \mathbf{z}_j\|^2 / 2h^2)$  is one of the most commonly used. Then we conduct spectral analysis on landmark-based graph and compute the graph matrix as:

$$\mathbf{G} = \hat{\mathbf{W}}^T \hat{\mathbf{W}}. \quad (2)$$

which has a very efficient eigen-decomposition. In this method, we choose  $\hat{\mathbf{W}} = \mathbf{D}^{-1/2} \mathbf{W}$  where  $\mathbf{D}$  is the row sum of  $\mathbf{W}$ . Note that each column of  $\mathbf{W}$  sums up to 1 and thus the degree matrix of  $\mathbf{G}$  is  $\mathbf{I}$ .

Let the Singular value Decomposition (SVD) of  $\hat{\mathbf{W}}$  is as follows:

$$\hat{\mathbf{W}} = \mathbf{U} \mathbf{\Sigma} \mathbf{V}^T. \quad (3)$$

where  $\mathbf{U} = [\mathbf{u}_1, \mathbf{u}_2, \dots, \mathbf{u}_k] \in R^{p \times p}$  is called as the left singular vectors of the first  $k$  eigenvectors of  $\hat{\mathbf{Z}} \hat{\mathbf{Z}}^T$ ,  $\mathbf{V} = [\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_k] \in R^{n \times p}$  is called as right singular vectors of the  $k$  eigenvectors of  $\hat{\mathbf{Z}}^T \hat{\mathbf{Z}}$ . We compute  $\mathbf{U}$  within  $O(p^3)$ , linear to the sample size.  $\mathbf{V}$  can be computed as:

$$\mathbf{V}^T = \mathbf{\Sigma}^{-1} \mathbf{U}^T \hat{\mathbf{W}}. \quad (4)$$

The overall time complexity of  $\mathbf{V}$  is  $O(p^3 + p^2n)$ , which is a significant reduction from  $O(n^3)$  by considering  $p \ll n$ . Each row of  $\mathbf{V}$  is a sample and apply  $k$ -means to get the clusters. Due to the time complexity from  $O(n^3)$  to  $O(n)$ , the LSC algorithm substantially reduces computational time. Thus, the proposed algorithm with low-complexity is very suitable to be applied in the domain of big data.

Finally, the pseudo of LSC is presented in Algorithm 1.

**Algorithm 1.** The pseudo of LSC.

- Input:**  $n$  data points  $x_1, x_2, \dots, x_n \in R^m$ ; Cluster number  $k$ ;  
**Output:**  $k$  clusters;
- 1 Produce  $p$  landmark points using the  $k$ -means method;
  - 2 Construct a landmark matrix  $\mathbf{Z}$  between data points and landmark samples, with the affinity calculated according to Eq. (1);
  - 3 Compute the first  $k$  eigenvectors of  $\mathbf{W}\mathbf{W}^T$ , denoted by  $\mathbf{U} = [\mathbf{u}_1, \mathbf{u}_2, \dots, \mathbf{u}_k]$ ;
  - 4 Compute  $\mathbf{V} = [\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_k]$  according to Eq. (4);
  - 5 Each row of  $\mathbf{V}$  is a data point and apply  $k$ -means to get the clusters.

#### 3.2. Testing process

Suppose that we already produce  $k$  clusters and cluster centers by LSC algorithm, we then find out the nearest cluster center for each test sample and use the corresponding cluster as the new training dataset for each test sample. Due to conducting the clusters with high similarity within a cluster, we apply  $k$ NN to classify test samples in the new training dataset. In this way, the proposed algorithm still guarantees relatively high classification accuracy. We list the pseudo of our proposed method in Algorithm 2.

**Algorithm 2.** The pseudo of LC-kNN algorithm.

**Input:** Training dataset, test samples  $Y$ ;

**Output:** Class label;

- 1 Produce  $m$  cluster centers using LSC algorithm, denoted by  $C_1, C_2, C_3, \dots, C_m$ ;
- 2 Compute the distance  $D(y, C_i)$  between test sample  $y$  and each cluster center, denoted by  $D(y, C_i)$ ,  $i = 1, 2, \dots, m$ ;
- 3 Compute the nearest cluster center  $C_i$  to  $y$ ,  $C_i = \min\{D(y, C_i)\}$ ,  $i = 1, 2, \dots, m$ ;
- 4 Using the corresponding cluster of  $C_i$  as the training dataset, denoted by  $NewX_i$ ;
- 5 Apply to kNN algorithm to predict  $y$  in the training dataset.

From Algorithm 2, the number of  $NewX_i$  is much smaller than the size of the training dataset. When the size of  $m$  is large, LC-kNN is easy to reduce the calculation of kNN and improves classification quality. However, with the increase of  $m$ , the overhead of clustering will increase and the classification accuracy will be lower at the same time. Therefore, in order to avoid this situation, the number of cluster  $m$  needs to be set in a reasonable value.

In general, the larger the value of  $m$ , the higher the classification efficiency, which lead to a higher classification accuracy. However, if the training dataset distribution is relatively concentrated, which will lead to low accuracy. In addition, if  $m$  is small, i.e.,  $m = 1$ , that is standard kNN algorithm, which prohibits the kNN to be used in big data. Thus we set  $m$  in a suitable range. Assuming that the proposed algorithm need  $M$  memory space, the memory of PC is  $M_1$  and the smallest class of training sample is  $n_0$ . To this end, the value range of  $m$  is expressed as follows:

$$\frac{M}{M_1} < m < n_0. \quad (5)$$

Note that we conduct a  $k$ -means clustering to separate the whole dataset into several parts, each of which we conduct kNN classification, the selection of  $k$  value is important in the later process. For example, Lall mentioned that the suitable setting of  $k$  value should satisfy  $k = \sqrt{n}$  for the datasets with sample size larger than 100 [8]. However, such a setting has been proved to not suitable for all cases of datasets. Considering that the samples has a strong correlation in sub-cluster, so the  $k$  value not be setting too large. Thus, the selection of  $k$  value should be set as small as possible in the case of high classification accuracy.

## 4. Experiments

As mentioned in the previous sections, our proposed algorithm is an extension of kNN. Thus, in order to show the effectiveness of LC-kNN algorithm, we took the kNN as the baseline and made a comparison between kNN, LC-kNN and RC-kNN (random clustering kNN), and several real datasets from UCI [2], medical imaging datasets [6,7], and LIBSVM [3] datasets.

### 4.1. Experimental results of LC-kNN and RC-kNN with different value of $m$

The value of  $m$  is very important to the LC-kNN algorithm, because it directly affects the final performance and real applications. Thus, in order to select appropriate  $m$  for bringing great effect to LC-kNN, a group of experiments were conducted on datasets by choosing different values of  $m$  for LC-kNN and RC-kNN. Specifically, the LC-kNN and RC-kNN in this group experiments were carried out on the datasets with  $m = 10, 15, 20, 25$  and  $30$ , respectively. The comparison of classification performance and

**Table 1**

Classification accuracy (mean  $\pm$  standard deviation) and Time cost (seconds: mean  $\pm$  standard deviation) on USPS dataset at different values of  $m$ .

$m$	Criterion	RC-kNN	LC-kNN
10	Accuracy	0.9027 $\pm$ 1.6498e-005	<b>0.9355 <math>\pm</math> 7.1306e-006</b>
	Time	3.5589 $\pm$ 0.0107	<b>3.7605 <math>\pm</math> 0.0242</b>
15	Accuracy	0.8964 $\pm$ 5.1803e-005	<b>0.9338 <math>\pm</math> 4.1625e-006</b>
	Time	2.4857 $\pm$ 0.0032	<b>2.7260 <math>\pm</math> 0.0077</b>
20	Accuracy	0.8770 $\pm$ 7.4889e-005	<b>0.9300 <math>\pm</math> 4.9238e-006</b>
	Time	2.3202 $\pm$ 0.0010	<b>2.5157 <math>\pm</math> 0.01928</b>
25	Accuracy	0.8793 $\pm$ 4.9917e-005	<b>0.9284 <math>\pm</math> 1.0637e-005</b>
	Time	1.8586 $\pm$ 0.0008	<b>1.9971 <math>\pm</math> 0.0042</b>
30	Accuracy	0.8607 $\pm$ 4.6629e-005	<b>0.9275 <math>\pm</math> 1.1596e-005</b>
	Time	1.6441 $\pm$ 0.0002	<b>1.9249 <math>\pm</math> 0.0023</b>

**Table 2**

Classification accuracy (mean  $\pm$  standard deviation) and Time cost (seconds: mean  $\pm$  standard deviation) on MNIST dataset at different values of  $m$ .

$m$	Criterion	RC-kNN	LC-kNN
10	Accuracy	0.7221 $\pm$ 4.8878e-005	<b>0.8389 <math>\pm</math> 3.1656e-005</b>
	Time	2.9369 $\pm$ 0.0508	<b>3.5504 <math>\pm</math> 0.0927</b>
15	Accuracy	0.6840 $\pm$ 2.3333e-004	<b>0.8364 <math>\pm</math> 2.3136e-005</b>
	Time	2.8905 $\pm$ 0.0456	<b>3.1222 <math>\pm</math> 0.01397</b>
20	Accuracy	0.6657 $\pm$ 2.4739e-004	<b>0.8353 <math>\pm</math> 3.3233e-005</b>
	Time	2.0564 $\pm$ 0.0011	<b>2.1490 <math>\pm</math> 0.0065</b>
25	Accuracy	0.6478 $\pm$ 2.2689e-004	<b>0.8338 <math>\pm</math> 8.7844e-005</b>
	Time	1.8240 $\pm$ 0.0020	<b>2.1148 <math>\pm</math> 0.0094</b>
30	Accuracy	0.6396 $\pm$ 6.9156e-005	<b>0.8313 <math>\pm</math> 3.8678e-005</b>
	Time	1.5457 $\pm$ 0.0002	<b>1.7274 <math>\pm</math> 0.0011</b>

**Table 3**

Classification accuracy (mean  $\pm$  standard deviation) and Time cost (seconds: mean  $\pm$  standard deviation) on GISETTE dataset at different values of  $m$ .

$m$	Criterion	RC-kNN	LC-kNN
10	Accuracy	0.9311 $\pm$ 5.0989e-005	<b>0.9526 <math>\pm</math> 1.4511e-005</b>
	Time	23.3933 $\pm$ 0.9677	<b>28.5940 <math>\pm</math> 3.2405</b>
15	Accuracy	0.9252 $\pm$ 1.0573e-004	<b>0.9494 <math>\pm</math> 1.3378e-005</b>
	Time	18.0106 $\pm$ 0.2434	<b>23.1904 <math>\pm</math> 1.0894</b>
20	Accuracy	0.9166 $\pm$ 2.8267e-005	<b>0.9411 <math>\pm</math> 5.4699e-004</b>
	Time	12.7685 $\pm$ 0.0966	<b>16.2759 <math>\pm</math> 0.8880</b>
25	Accuracy	0.9150 $\pm$ 7.0000e-005	<b>0.9321 <math>\pm</math> 6.4810e-004</b>
	Time	9.9201 $\pm$ 0.3696	<b>13.8645 <math>\pm</math> 1.5093</b>
30	Accuracy	0.9079 $\pm$ 1.0366e-004	<b>0.9192 <math>\pm</math> 5.3796e-004</b>
	Time	8.4064 $\pm$ 0.0784	<b>11.3922 <math>\pm</math> 0.0658</b>

**Table 4**

Classification accuracy (mean  $\pm$  standard deviation) and Time cost (seconds: mean  $\pm$  standard deviation) on LETTER dataset at different values of  $m$ .

$m$	Criterion	RC-kNN	LC-kNN
10	Accuracy	0.7892 $\pm$ 3.8822e-005	<b>0.9495 <math>\pm</math> 1.0760e-006</b>
	Time	3.2391 $\pm$ 0.0015	<b>3.2994 <math>\pm</math> 0.0010</b>
15	Accuracy	0.7932 $\pm$ 3.7106e-005	<b>0.9469 <math>\pm</math> 5.5751e-006</b>
	Time	3.3808 $\pm$ 0.0435	<b>3.4334 <math>\pm</math> 0.0585</b>
20	Accuracy	0.6815 $\pm$ 1.3812e-004	<b>0.9451 <math>\pm</math> 1.9756e-006</b>
	Time	3.0938 $\pm$ 5.8392e-004	<b>3.1285 <math>\pm</math> 3.1243e-004</b>
25	Accuracy	0.7279 $\pm$ 5.6480e-005	<b>0.9423 <math>\pm</math> 5.2818e-006</b>
	Time	3.3950 $\pm$ 0.0018	<b>3.4813 <math>\pm</math> 0.0054</b>
30	Accuracy	0.6214 $\pm$ 9.8480e-005	<b>0.9403 <math>\pm</math> 3.9204e-006</b>
	Time	3.0889 $\pm$ 1.3000e-003	<b>3.1168 <math>\pm</math> 3.8514e-004</b>

time cost of two algorithms (i.e., RC-kNN and LC-kNN, respectively) with different values of  $m$  are shown in Tables 1–9.

From Tables 1–9, we found that the proposed two algorithms needed less time with the larger number of clusters  $m$ , while more time cost for the smaller number of  $m$ . For example, when  $m$  value

**Table 5**

Classification accuracy (mean  $\pm$  standard deviation) and Time cost (seconds: mean  $\pm$  standard deviation) on PENDIGITS dataset at different values of  $m$ .

$m$	Criterion	RC-kNN	LC-kNN
10	Accuracy	0.9452 $\pm$ 3.5382e-005	<b>0.9721 <math>\pm</math> 4.7991e-006</b>
	Time	2.3380 $\pm$ 0.0041	<b>2.4056 <math>\pm</math> 0.0101</b>
15	Accuracy	0.9316 $\pm$ 1.0341e-004	<b>0.9711 <math>\pm</math> 6.0196e-006</b>
	Time	2.5451 $\pm$ 0.0011	<b>2.5709 <math>\pm</math> 0.0089</b>
20	Accuracy	0.9163 $\pm$ 1.5515e-004	<b>0.9700 <math>\pm</math> 2.5390e-006</b>
	Time	2.2233 $\pm$ 6.4795e-005	<b>2.2554 <math>\pm</math> 2.1569e-004</b>
25	Accuracy	0.9216 $\pm$ 1.5677e-004	<b>0.9687 <math>\pm</math> 3.5642e-006</b>
	Time	2.5270 $\pm$ 0.0056	<b>2.5468 <math>\pm</math> 0.0083</b>
30	Accuracy	0.9088 $\pm$ 1.8409e-004	<b>0.9683 <math>\pm</math> 1.5809e-006</b>
	Time	2.1805 $\pm$ 7.4785e-005	<b>2.2022 <math>\pm</math> 8.9611e-005</b>

**Table 6**

Classification accuracy (mean  $\pm$  standard deviation) and Time (seconds: mean  $\pm$  standard deviation) on SATIMAGE dataset at different values of  $m$ .

$m$	Criterion	RC-kNN	LC-kNN
10	Accuracy	0.8603 $\pm$ 8.9122e-005	<b>0.8883 <math>\pm</math> 8.1139e-006</b>
	Time	1.2868 $\pm$ 6.5495e-005	<b>1.3027 <math>\pm</math> 1.1429e-004</b>
15	Accuracy	0.7917 $\pm$ 3.1680e-005	<b>0.9468 <math>\pm</math> 3.7244e-006</b>
	Time	3.8583 $\pm$ 0.0332	<b>3.9337 <math>\pm</math> 0.0152</b>
20	Accuracy	0.8418 $\pm$ 8.8847e-005	<b>0.8884 <math>\pm</math> 6.8028e-006</b>
	Time	1.2292 $\pm$ 3.2277e-005	<b>1.2463 <math>\pm</math> 4.3126e-005</b>
25	Accuracy	0.7283 $\pm$ 1.1039e-004	<b>0.9421 <math>\pm</math> 8.8449e-006</b>
	Time	3.5062 $\pm$ 0.0061	<b>3.6287 <math>\pm</math> 0.0052</b>
30	Accuracy	0.8312 $\pm$ 3.5146e-004	<b>0.8878 <math>\pm</math> 4.9556e-006</b>
	Time	1.2225 $\pm$ 1.8176e-005	<b>1.2396 <math>\pm</math> 1.7711e-005</b>

**Table 7**

Classification accuracy (mean  $\pm$  standard deviation) and Time cost (seconds: mean  $\pm$  standard deviation) on ADNC dataset at different values of  $m$ .

$m$	Criterion	RC-kNN	LC-kNN
10	Accuracy	0.7024 $\pm$ 0.0010	<b>0.7667 <math>\pm</math> 0.0019</b>
	Time	0.0310 $\pm$ 2.7390e-004	<b>0.0333 <math>\pm</math> 1.9201e-005</b>
15	Accuracy	0.6857 $\pm$ 0.0014	<b>0.7500 <math>\pm</math> 0.0013</b>
	Time	0.0308 $\pm$ 7.9360e-006	<b>0.0325 <math>\pm</math> 3.3635e-005</b>
20	Accuracy	0.6619 $\pm$ 0.0029	<b>0.7143 <math>\pm</math> 0.0028</b>
	Time	0.0295 $\pm$ 3.0063e-006	<b>0.0313 <math>\pm</math> 6.8908e-006</b>
25	Accuracy	0.6548 $\pm$ 0.0039	<b>0.7071 <math>\pm</math> 0.0038</b>
	Time	0.0281 $\pm$ 4.8219e-007	<b>0.0286 <math>\pm</math> 7.0881e-007</b>
30	Accuracy	0.6619 $\pm$ 0.0015	<b>0.7190 <math>\pm</math> 0.0031</b>
	Time	0.0306 $\pm$ 1.2641e-005	<b>0.0326 <math>\pm</math> 5.3013e-006</b>

**Table 8**

Classification accuracy (mean  $\pm$  standard deviation) and Time cost (seconds: mean  $\pm$  standard deviation) on psMCI dataset at different values of  $m$ .

$m$	Criterion	RC-kNN	LC-kNN
10	Accuracy	0.4792 $\pm$ 0.0082	<b>0.5833 <math>\pm</math> 0.0089</b>
	Time	0.0168 $\pm$ 4.5147e-007	<b>0.0171 <math>\pm</math> 1.4130e-006</b>
15	Accuracy	0.5125 $\pm$ 0.0019	<b>0.6042 <math>\pm</math> 0.0040</b>
	Time	0.0167 $\pm$ 9.4536e-007	<b>0.0173 <math>\pm</math> 5.6957e-007</b>
20	Accuracy	0.5458 $\pm$ 0.0060	<b>0.6500 <math>\pm</math> 0.0035</b>
	Time	0.0170 $\pm$ 1.2416e-006	<b>0.0188 <math>\pm</math> 1.1793e-005</b>
25	Accuracy	0.5333 $\pm$ 0.0053	<b>0.6417 <math>\pm</math> 0.0035</b>
	Time	0.0163 $\pm$ 4.9822e-007	<b>0.0286 <math>\pm</math> 7.0881e-007</b>
30	Accuracy	0.5000 $\pm$ 0.0042	<b>0.6125 <math>\pm</math> 0.0019</b>
	Time	0.0166 $\pm$ 8.6385e-007	<b>0.0250 <math>\pm</math> 2.2840e-004</b>

is 10, i.e., we separated the whole dataset into 10 parts. We then conducted  $k$ NN classification for each part, and obtained about one-tenth time cost of the method conducting in the whole dataset.

Besides, considering that the proposed two algorithms are the extension of  $k$ NN, their classification accuracies should as close as

**Table 9**

Classification accuracy (mean  $\pm$  standard deviation) and Time cost (seconds: mean  $\pm$  standard deviation) on MCINC dataset at different values of  $m$ .

$m$	Criterion	RC-kNN	LC-kNN
10	Accuracy	0.5476 $\pm$ 0.0049	<b>0.6159 <math>\pm</math> 0.0045</b>
	Time	0.0468 $\pm$ 2.8577e-005	<b>0.0506 <math>\pm</math> 5.9543e-005</b>
15	Accuracy	0.5397 $\pm$ 0.0032	<b>0.5633 <math>\pm</math> 0.0036</b>
	Time	0.0505 $\pm$ 2.7585e-005	<b>0.0538 <math>\pm</math> 4.6465e-005</b>
20	Accuracy	0.5458 $\pm$ 0.0060	<b>0.6500 <math>\pm</math> 0.0035</b>
	Time	0.0452 $\pm$ 4.9222e-006	<b>0.0480 <math>\pm</math> 1.1608e-005</b>
25	Accuracy	0.5222 $\pm$ 0.0015	<b>0.5746 <math>\pm</math> 0.0008</b>
	Time	0.0489 $\pm$ 1.5349e-005	<b>0.0540 <math>\pm</math> 4.5819e-005</b>
30	Accuracy	0.5016 $\pm$ 0.0245	<b>0.5984 <math>\pm</math> 0.0021</b>
	Time	0.0490 $\pm$ 2.6160e-005	<b>0.0536 <math>\pm</math> 4.5625e-005</b>

**Table 10**

Classification Accuracy and Time Cost of three algorithm on nine dataset.

Dataset	RC-kNN		LC-kNN		kNN	
	Accuracy	Time	Accuracy	Time	Accuracy	Time
USPS	0.9027	3.5589	<b>0.9355</b>	<b>3.7605</b>	0.9482	32.8764
MNIST	0.7221	2.9369	<b>0.8389</b>	<b>3.5504</b>	0.8635	24.1575
GISETTE	0.9311	23.3933	<b>0.9526</b>	<b>28.594</b>	0.9660	217.3327
LETTER	0.7892	3.2391	<b>0.9495</b>	<b>3.2994</b>	0.9518	19.8246
PENDIGITS	0.9452	2.3380	<b>0.9721</b>	<b>2.4056</b>	0.9780	7.2982
SATIMAGE	0.8603	1.2868	<b>0.8883</b>	<b>1.3027</b>	0.9065	3.5525
ADNC	0.7024	0.0310	<b>0.7667</b>	<b>0.0333</b>	0.7857	0.0355
psMCI	0.4792	0.0168	<b>0.5833</b>	<b>0.0171</b>	0.6042	0.0176
MCINC	0.5476	0.0468	<b>0.6159</b>	<b>0.0506</b>	0.6413	0.0575

possible  $k$ NN. We found that two algorithms have high classification accuracy with the larger number of clusters  $m$ , while these two algorithms have low classification accuracy for the small number of  $m$ . In addition, from Table 10, we found that two algorithms were closer to  $k$ NN classification accuracy with  $m=10$ . Thus, we set  $m=10$  in the next section.

#### 4.2. To determine parameter $k$

In this section, in order to select the appropriate  $k$  value, a group of experiments were conducted on nine datasets with  $m=1$  to 20, each point represent the mean of 10 results in the figures.

From Fig. 1, with the increase of  $k$  value, the overall of classification accuracy decreases. That is because that the smaller the training dataset is, the more the classification accuracy is. Hence, the difference between the samples is significant and the classification accuracy will reduce. We can make a conclusion that we should choose a suitable  $k$  value in this case. According to the previous analysis,  $k$  value should be set as small as possible in the case of the higher accuracy of the proposed algorithm. From Fig. 1, the classification accuracy was 0.8986 with  $k=3$ , which was higher than  $k=1$  on satimage dataset. But the LC-kNN performance has higher than  $k$ NN with  $k=1$  in Table 10. Thus we set  $k=1$  in the next experiment.

#### 4.3. Performance comparison of $k$ NN, RC-kNN and LC-kNN

In this experiment, according to previous analysis, we set  $m=10$ ,  $k=1$ . Then, we use classification accuracy and running time as the evaluations for the classification task [23]. The shorter time and higher accuracy the algorithm is, the better the performance is. We report our experimental results in Table 10.

From Table 10, we can observe that the proposed RC-kNN and LC-kNN improved by 7–9 times than the  $k$ NN (close to the number of cluster), in terms of time cost. In the evaluation of classification



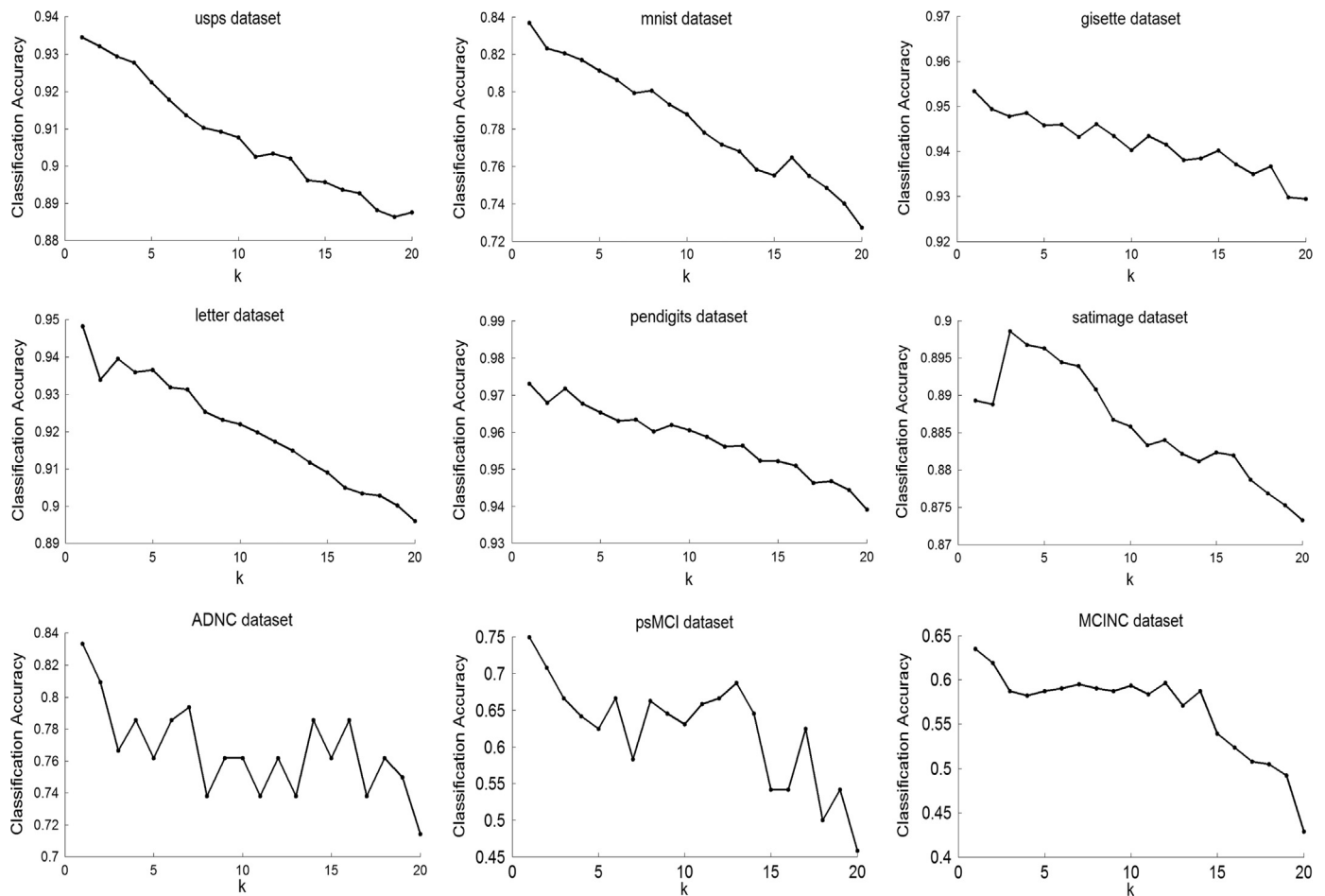


Fig. 1. Classification accuracy of LC-kNN on nine dataset with different  $k$ .

accuracy, the LC-kNN and RC-kNN is lower by 1–2.6% and 3.3–14% than kNN. Therefore, according to the experimental results, we may the conclusion that the LC-kNN works well in terms of classification accuracy and time.

## 5. Conclusions and future work

In this work, we have proposed an efficient kNN classification to conduct a  $k$ -means clustering to separate the whole dataset into several parts. We then conducted kNN classification for each part. To do this, we parted the conventional kNN method into two processes, namely training process and testing process. Moreover, we analyzed the suitable value for the parameters, such as  $m$  and  $k$ . Furthermore, we took the kNN as the baseline and conducted a groups of compared experiments among kNN, LC-kNN and RC-kNN. The experimental results showed that the proposed kNN classification worked well in terms of accuracy and efficiency, and it is appropriate to deal with big data.

## Acknowledgments

This work was supported in part by the China 973 Program under Grant 2013CB329404, in part by the National Natural Science Foundation of China under Grant 61450001, Grant 61263035, and Grant 61573270, in part by the Guangxi Natural Science Foundation under Grant 2012GXNSFGA060004 and Grant

2015GXNSFCB139011, in part by the China Postdoctoral Science Foundation under grant 2015M57570837, in part by the Guangxi 100 Plan, in part by the Guangxi Collaborative Innovation Center of Multi-Source Information Integration and Intelligent Processing, and in part by the Guangxi Bagui Scholar Teams for Innovation and Research Project.

## References

- [1] A. Andrew, M. Jordan, Y. Weiss, On spectral clustering: analysis and algorithm, in: *Proceedings of Advance in Neural Information Processing Systems* 14, MIT Press.
- [2] K. Bache, M. Lichman, *UCI Mach. Learn. Repos.* (2013).
- [3] C.-C. Chang, C.-J. Lin, LIBSVM: a library for support vector machines, *ACM Trans. Intell. Syst. Technol.* (2011) 2:27:1–27:27, Software available at (<http://www.csie.ntu.edu.tw/~cjlin/libsvm>).
- [4] X. Chen, D. Cai, Large scale spectral clustering with landmark-based representation, *AAAI* (2011) 313–318.
- [5] M. Filippone, F. Camestra, F. Masulli, et al., A survey of kernel and spectral methods for clustering, *Pattern Recognit.* 41 (2007) 176–190.
- [6] Y. Gao, M. Wang, Z. Zha, J. Shen, X. Li, X. Wu, Visual-textual joint relevance learning for tag-based social image search, *IEEE Trans. Image Process.* 22 (1) (2013) 363–376.
- [7] Y. Gao, M. Wang, D. Tao, R. Ji, Q. Dai, 3D object retrieval and recognition with hypergraph analysis, *IEEE Trans. Image Process.* 21 (9) (2012) 4290–4303.
- [8] U. Lall, A. Sharma, A nearest neighbor bootstrap for resampling hydrologic time series, *Water Resour. Res.* 32 (3) (1996) 679–693.
- [9] R. Li, Y. Hu, A density-based method for reducing the amount of training data in kNN text classification, *J. Comput. Res. Dev.* 41 (4) (2004) 539–545.
- [10] W. Liu, J. He, S. Chang, Large graph construction for scalable semi-supervised learning, in: *Proceedings of the 27th International Conference on Machine Learning*, 2010.

- [11] X. Wu, C. Zhang, S. Zhang, Database classification for multi-database mining, *Inform. Syst.* 30 (1) (2005) 71–88.
- [12] X. Wu, S. Zhang, Synthesizing high-frequency rules from different data sources, *IEEE Trans. Knowl. Data Eng.* 15 (2) (2003) 353–367.
- [13] X. Wu, C. Zhang, S. Zhang, Efficient mining of both positive and negative association rules, *ACM Trans. Inform. Syst.* 22 (3) (2004) 381–405.
- [14] S. Zhang, KNN-CF approach: incorporating certainty factor to knn classification, *IEEE Intell. Inform. Bull.* 11 (1) (2010) 24–33.
- [15] Y. Zhao, S. Zhang, Generalized dimension-reduction framework for recent-biased time series analysis, *IEEE Trans. Knowl. Data Eng.* 18 (2) (2006) 231–244.
- [16] D. Zhao, Wei Zou, G. Sun, A fast image classification algorithm using support vector machine, *ICCTD* (2010).
- [17] X. Zhu, Z. Huang, H. Cheng, J. Cui, H. Shen, Sparse hashing for fast multimedia search, *ACM Trans. Inform. Syst.* 31 (2) (2013) 9.
- [18] X. Zhu, Z. Huang, H. Shen, X. Zhao, Linear cross-modal hashing for efficient multimedia search, *ACM Multimed.* (2013) 143–152.
- [19] X. Zhu, Z. Huang, H. Shen, J. Cheng, C. Xu, Dimensionality reduction by mixed kernel canonical correlation analysis, *Pattern Recognit.* 45 (8) (2012) 3003–3016.
- [20] X. Zhu, Z. Huang, Y. Yang, H. Shen, C. Xu, J. Luo, Self-taught dimensionality reduction on the high-dimensional small sized data, *Pattern Recognit.* 46 (1) (2013) 215–229.
- [21] X. Zhu, L. Zhang, Z. Huang, A sparse embedding and least variance encoding approach to hashing, *IEEE Trans. Image Process.* (2014).
- [22] X. Zhu, S. Zhang, Z. Jin, Z. Zhang, Z. Xu, Missing value estimation for mixed-attribute datasets, *IEEE Trans. Knowl. Data Eng.* 23 (1) (2011) 110–121.
- [23] X. Zhu, Z. Huang, J. Cui, H.T. Shen, Video-to-shot tag propagation by graph sparse group lasso, *IEEE Trans. Multimed.* 15 (3) (2013) 633–646.
- [24] X. Zhu, H.-I. Suk, D. Shen, A novel matrix-similarity based loss function for joint regression and classification in ad diagnosis, *NeuroImage* 100 (2014) 91–105.
- [25] X. Zhu, H.-I. Suk, D. Shen, Matrix-similarity based loss function and feature selection for alzheimer's disease diagnosis, *CVPR* (2014) 3089–3096.
- [26] X. Zhu, X. Li, S. Zhang, Block-row sparse multiview multilabel learning for image classification, *IEEE Trans. Cybern.* 46 (2) (2016) 450–461.



**Shichao Zhang** is a Distinguished Professor and the director of Institute of School of Computer Science and Information Technology at the Guangxi Normal University, Guilin, China. He holds a Ph.D. degree in Computer Science from Deakin University, Australia. His research interests include data analysis and smart pattern discovery. He has published over 50 international journal papers and over 60 international conference papers.