

# A New K-Nearest Neighbor Classification Method Based on Belief Functions in Wireless Sensor Networks



Yang Zhang<sup>1\*</sup>, Deyang Yuan<sup>2</sup>

<sup>1</sup> Key Laboratory of Communication and Information Systems, Beijing Municipal Commission of Education, School of Electronic and Information Engineering, Beijing Jiaotong University, Beijing 100044, China  
zhang.yang@bjtu.edu.cn

<sup>2</sup> Nanyang Fire Rescue Detachment, Nanyang, Henan 473000, China  
deyangyuan@126.com

Received 17 April 2021; Revised 1 May 2021; Accepted 17 May 2021

**Abstract.** In wireless sensor networks, the classification of uncertain data reported by sensor nodes is an open issue because the given attribute information can be insufficient for making a correct specific classification of the objects. Although the traditional Evidential  $k$ -Nearest Neighbor ( $E_kNN$ ) algorithm can effectively model the uncertainty, it is easy to misjudge the target data to the incorrect class when the observed sample data is located in the feature overlapping region of training samples of different classes. In this paper, a novel Evidential  $k$ -Nearest Neighbor ( $NE_kNN$ ) algorithm is proposed based on the evidential editing method. The main idea of  $NE_kNN$  algorithm is to consider the expected value and standard deviation of various training sample data sets, and use normalized Euclidean distance to assign class labels with basic belief assignment (BBA) structure to each training sample, so that training samples in overlapping region can offer more abundant and diverse class information. Further,  $E_kNN$  classification of the observation sample data is carried out in the training sample sets of various classes, and mass functions of the target to be tested under this class are obtained, and Redistribute Conflicting Mass Proportionally Rule 5 (PCR5) combination rule is used to conduct global fusion, thus obtaining the global fusion results of the targets. The experimental results show that this algorithm has better performance than other classification methods based on  $k$ -nearest neighbor. Several experiments using both simulation and real data sets are presented at the end of this paper. The results indicate that the  $NE_kNN$  algorithm can effectively improve the classification accuracy.

**Keywords:** evidence theory, uncertain data, target classification, combination rule

## 1 Introduction

When using the sensor's observation data to carry out the local classification of targets, the sensor's observation data contains a lot of imprecise information due to various interferences [1]. For example, some sample data comes from different categories of targets but they are very similar, that is, the sample data of different categories may partially overlap, which brings great challenges to traditional target classification tasks [2]. In the classification task with supervision, the sensor's observation data may be in the overlapping area of different categories of training samples, it is difficult for the traditional voting  $kNN$  classifier to accurately classify the target at this time. For this reason, many scholars have fully considered the distance relationship between the target and its neighbors, and proposed the fuzzy  $kNN$  (Fuzzy  $kNN$ ,  $FkNN$ ) classification algorithm [3]. This algorithm allows the target to belong to different categories with different fuzzy membership degrees, which obtains the better classification effect than voting  $kNN$  [4].

---

\* Corresponding Author

Dempster-Shafer evidence theory, also referred as evidential reasoning or belief functions theory, has been proved to be valuable as a solution for dealing with uncertain and inaccurate data [5], and it has been widely applied in sorts of applications, for example, state estimation [6], target recognition [7], data classification [2], and information fusion [8], and etc. As the extension of probability theory, evidence theory provides a series of functions and operations defined on the power set of the identification framework, which can effectively reason and model the uncertainty, and can provide more abundant category information than fuzzy membership [9]. Therefore, many scholars have combined evidence theory with traditional classification algorithms with supervision and developed a series of evidence classification algorithms. Among them, the most representative is the evidential  $k$ NN (Evidential  $k$ -Nearest Neighbor,  $E_k$ NN) algorithm proposed by Scholars such as Denoeux, et al. [10-11]. This algorithm is simple and direct with low error rate, so it is very suitable for the target classification task of sensor nodes. However, the  $E_k$ NN algorithm only considers the factor of the distance between the target and the training samples, and does not treat all training samples differently [12]. It is assumed that the target sample is in the overlapping area of the training set, and the target data is far from the sample points of the same category, and closer to the sample points of other categories, if the  $E_k$ NN algorithm is used at this time, the evidence formed by the sample points that are closer to the target will be given a large mass value, then when making the decision after the evidence is fused, it is easy to misjudge the target data into other categories. In [13], it is further pointed out that since the  $E_k$ NN algorithm treats imprecise training samples from overlapping regions as training samples that truly represent the distribution of the target category, it will have a greater negative impact on the final classification effect. In order to solve this problem, it is necessary to preprocess the original training samples with the evidence editing method based on  $E_k$ NN, and replace the category labels of the original training samples with basic belief assignment, it can better characterize the inaccuracy of the overlapping regions of categories. However, in [13], it is proposed that the evidence editing method will make the edited evidence have a higher correlation. When subsequently fusing the evidence constructed by the target's neighbors, it is necessary to evaluate the correlation between the evidences, and to search for the corresponding fusion rules according to the degree of correlation between the evidences. Therefore, this method has the problems of high algorithm complexity and excessive calculation, which is not suitable for sensor nodes with limited energy. In addition to the evidence editing method, in [14], it is pointed out that if a target falls in the overlapping area of the training set, to first consider using  $E_k$ NN to classify the target in the training set of each category and then perform the evidence fusion of the classification results can also suppress the influence of other categories of training samples in the overlapping area on the fusion result. The improved  $E_k$ NN algorithm (Improved Evidential  $k$ -Nearest Neighbor,  $IE_k$ NN) was also proposed to improve the performance of the  $E_k$ NN algorithm [15], however,  $IE_k$ NN does not edit the samples, while directly uses the original training samples for classification [16].

In order to effectively model and reason about imprecise data, this paper proposes a  $NE_k$ NN (New Evidential  $k$ -Nearest Neighbor,  $NE_k$ NN) algorithm. The  $NE_k$ NN algorithm proposes a simple evidence preprocessing method under the framework of evidence theory. This method only considers the expected value and standard deviation of the training sample sets of each category, thereby avoiding the evidence correlation that may be caused by the original evidence editing method. On this basis, by fusing the classification results of the target to be tested in the training sample set of each category, the  $NE_k$ NN obtain a more accurate identification and judgment of the target.

The other parts of this paper are arranged as follows. Section 2 fundamentally introduces the basis of evidence theory. Section 3 focuses on the original training data preprocessing method, and design the  $NE_k$ NN classification algorithm after preprocessing. Section 4 comprehensively evaluates and analyzes the classification performance of the proposed  $NE_k$ NN algorithm based on simulation and real data sets, and finally summarize the work of this paper in Section 5.

## 2 Basics of Belief Functions Theory

The Dempster-Shafer evidence theory introduced by Shafer is also known as belief functions theory [17]. In this theory, the frame of discernment  $\Omega$  is a finite set, whose elements are exhaustive and mutually exclusive, and it is denoted as  $\Omega = \{w_1, w_2, \dots, w_i, w_o\}$ .  $2^\Omega$  is the power set of the frame of discernment, which represents the set of all possible subsets of  $\Omega$ , indicated by  $2^\Omega = \{\emptyset, \{w_1\}, \dots, \{w_n\}, \{w_1, w_2\},$

..., {w<sub>1</sub>, w<sub>2</sub>, ..., w<sub>i</sub>}, ..., Ω} Given an object *X*, it can be classified as any singleton element and any sets of elements in 2<sup>Ω</sup> with a basic belief assignment (BBA). The BBA is also known as the mass function, which is a mapping  $m:2^{\Omega} \rightarrow [0,1]$  satisfying  $\sum_{A \in 2^{\Omega}} m(A) = 1, m(\phi) = 0$ . The function  $m(A)$  is used to quantify the degree of belief that is exactly assigned to the subsets *A* of Ω. If  $m(A) > 0$ , the subset *A* can be called the focal elements of the mass function  $m(\cdot)$ . The mass values assigned to compound elements can reflect the imprecise observation of object *X*.

The mass function  $m(\cdot)$  is always associated with three main functions, including the belief function  $Bel(\cdot)$ , the plausibility function  $Pl(\cdot)$  and the pignistic probability function  $BetP(\cdot)$ , which are defined as follows, respectively:

$$Bel(B) = \sum_{A \subseteq B} m(A) \tag{1}$$

$$Pl(B) = \sum_{A \cap B \neq \phi} m(A) \tag{2}$$

$$BetP(w) = \sum_{w \in A, A \subseteq \Omega} \frac{1}{|A|} m(A) \tag{3}$$

where  $m(\cdot)$  is the focal elements on Ω, and  $|A|$  denotes the cardinality of focal elements *A*. All three functions can be employed to make a decision on an unknown object according to a few rules, such as selecting the class with maximum  $BetP$ .

Assuming that there are two pieces of evidence denoted by  $m_1$  and  $m_2$ , the popular Dempster’s combination rule can be used to combine them as follows:

$$m_{\oplus}(A) = m_1(B) \oplus m_2(C) = \begin{cases} 0, & B \cap C = \phi \\ \frac{\sum_{B \cap C = A, \forall B, C \subseteq \Omega} m_1(B) \times m_2(C)}{1 - \sum_{B \cap C = \phi, \forall B, C \subseteq \Omega} m_1(B) \times m_2(C)} & B \cap C \neq \phi \end{cases} \tag{4}$$

where  $\sum_{B \cap C = \phi, \forall B, C \subseteq \Omega} m_1(B) \times m_2(C)$  represents the conflict between  $m_1$  and  $m_2$ , which is used to redistribute the conflicting mass values. Dempster’s combination rule is commutative and associative. It provides a simple and flexible solution for data fusion problems.

### 3 The New Evidential k-Nearest Neighbor Algorithm

In order to overcome the limitations of EkNN, a new EkNN classification algorithm is proposed in this section. The algorithm uses the method of preprocessing training samples to replace the category labels of the original samples with the basic belief assignment, so as to better describe the uncertainty of the training samples in the overlapping regions of the categories. In order to avoid the pre-processed evidence from generating greater correlation, the newly obtained category labels with the basic belief assignment structure for each sample are constructed based on the Mahalanobis distance from the evidence to the center of the corresponding category. In the subsequent classification of the target, it is first to find the *k* nearest neighbors of the input sample in each category of training sample set, construct *k* nearest neighbor evidence describing the respective classification information, and perform fusion to obtain the mass function under this category of condition, then the global fusion of evidence between categories is performed based on the mass function generated by each category, and the final classification result is obtained.

#### 3.1 Preprocessing of Training Samples

In the evidence editing method proposed in [13], every training sample will regenerate a category label with a basic belief assignment structure based on the *k*-nearest neighbor data in the training sample set. In this process, the same training sample may participate in the calculation of multiple sample category labels, which will cause the edited training samples to be not independent anymore, and make the

evidence provided by these samples have a certain correlation, and then affect the subsequent fusion process. In order to avoid evidence-related problems, this section focuses on the preprocessing method of training samples based on Mahalanobis distance. The concept of Mahalanobis distance belongs to the theory of multivariate statistical analysis [18]. It is a discriminant method that uses the distance between the sample to be judged and each population as the measurement scale to judge the attribution of the sample. When processing numerical data in wireless sensor networks, Mahalanobis distance comprehensively considers the two statistical characteristics of the expected value and standard deviation of each category in the true distribution. It avoids discussing the correlation caused by the specific distribution of sample data. At the same time, compared with Euclidean distance, Mahalanobis distance can also eliminate the interference of the correlation between attribute variables, which is more reasonable. The Mahalanobis distance used in this section can also be called the normalized Euclidean distance.

To consider a  $M$ -class problem, where the object may belong to  $M$  different classes, and  $\Omega = \{w_1, \dots, w_M\}$  is the set of all classes. It is supposed that the training sample set is  $Y = \{y_1, \dots, y_g\}$ . First, the attribute information of each category of training sample can be used to calculate the center vector of the category. The center of  $c_i (i = 1, \dots, M)$  can be expressed as:

$$c_i = \frac{1}{s_i} \sum_{y_j \in w_i} y_j \quad (5)$$

where  $s_i$  is the number of training samples of class  $w_i$ .

For each training sample  $y_h (h = 1, \dots, g)$ , sample preprocessing is performed according to the distance from it to the center. The distance requires to fully consider the degree of dispersion of each category of sample distribution, that is, the size of the standard deviation. Therefore, the Mahalanobis distance is used as the measurement scale of distance here, which is:

$$d_h^{w_i} = \sqrt{\sum_{k=1}^p \left( \frac{y_h(k) - c_i(k)}{\delta_i(k)} \right)^2} \quad (6)$$

where  $\delta_i(k)$  is the standard deviation of the training data set of class  $w_i$ ,  $y_h(k)$  and  $c_i(k)$  are the values of the attribute vector  $y_h$  and center  $c_i$  on the  $k$ -th dimension respectively, and  $p$  is the number of dimensions.

The smaller the distance  $d_h^{w_i}$  is, the more likely the training sample  $y_h$  belongs to the category  $w_i$ . If  $y_h$  is farther from the center  $c_i$ , the less likely it is that  $y_h$  belongs to category  $w_i$ . Therefore, the support of  $y_h$  belonging to category  $w_i$  is:

$$s_h(w_i) = e^{-d_h^{w_i}} \quad (7)$$

The BBA  $m_h$  should correspond to the normalized  $s_h(w_i)$ , formally defined by:

$$m_h = \frac{s_h(w_i)}{\sum_{l=1}^M s_h(w_l)} \quad (8)$$

The above mass function  $m_h$  provides more powerful information to characterize the uncertainty for training sample  $y_h$  than the original class label  $w_i \in \Omega$ , and it can be considered as a new soft class label of the sample  $y_h$ . As a consequence, the new training sample set with soft class labels  $Y' = \{y_1, \dots, y_g\}$  is adopted for the target classification task in this paper.

### 3.2 Classification with Preprocessed Training Samples

After preprocessing, the next problem to be solved is how to classify the newly observed unknown target  $x \in R^P$  based on the preprocessed training samples. Different from the general classification problem, the

category of training sample used here is represented by the structure of basic belief distribution, so it is necessary to improve the original evidence  $k$ NN classification algorithm accordingly to enable it to classify  $x$  reasonably using the category label of this structure. For the target  $\Omega = \{w_1, \dots, w_M\}$  of categories  $M$ , it is to first establish training sample sets for each category based on the total training samples, and then refer to the  $k$ NN classification algorithm to generate evidence that can be combined in various training sets based on the training samples and the feature data of the target to be tested. The entire classification process can be divided into the construction and fusion of mass functions under each category, and the global fusion of the fusion results between categories, which will be introduced separately below.

### 3.2.1 The Construction and Fusion of Mass Functions

To consider the  $k$  nearest neighbor samples of the target  $x$  to be tested in the training samples of category  $w_i (i=1, \dots, M)$ , if one of the training samples is very close to the sample  $x$  to be tested, the training sample provides a more reliable evidence for the classification of the sample to be tested. Conversely, if the distance is far, the reliability of the evidence provided by the training sample is relatively small. According to the evidence  $k$ NN algorithm, it is to choose Euclidean distance as the measurement scale to calculate the distance between the target and the training sample. It is assumed that the set of  $k$  nearest neighbor samples of target  $x$  in category  $w_i$  is,  $\Gamma_i = \{(y_1, d_1), \dots, (y_k, d_k)\}$ ,  $d_j (j=1, \dots, k)$  is the Euclidean distance between neighbor  $y_j$  and target  $x$ ,  $m_j$  is the category label of  $y_j$ ,  $\beta_j$  is the reliability of classifying  $x$  based on sample  $y_j$ , then the evidential mass function  $m'_j$  provided by  $y_j$  for the classification of target  $x$  can be expressed as:

$$\begin{cases} m'_j(w_i) = \beta_j m_j(w_i), i = 1, \dots, M \\ m'_j(\Omega) = \beta_j m_j(\Omega) + (1 - \beta_j) \end{cases} \quad (9)$$

where the reliability  $\beta_j$  is determined by the Euclidean distance  $d_j$  between  $y_j$  and the target  $x$ . The greater the distance between the two, the lower the corresponding reliability, that is, the reliability  $\beta_j$  and  $d_j$  show a decreasing relationship, which can be expressed as:

$$\beta_j = e^{(-d_j/\bar{d}^i)} \quad (10)$$

where  $\bar{d}^i$  is the average distance between all training samples in category  $w_i$ .

In order to classify the unknown target  $x$ , the  $k$  mass functions constructed by the  $k$  nearest neighbor samples  $y_j (j=1, \dots, k)$  in the category  $w_i$  need to be fused to obtain the classification result of the target by the training samples of the category  $w_i$ . In the fusion process, considering that the mass functions provided by the same category of training samples have high consistency, the Dempster combination rule can be used directly for the fusion operation, it can be expressed as:

$$m_i = m'_1 \oplus m'_2 \oplus \dots \oplus m'_k \quad (11)$$

where  $\oplus$  is Dempster combination operation.

Considering that there are a total of  $M$  categories of targets, a set of  $M$  nearest neighbor samples of the target can be generated, namely. According to equation (11), the mass function set  $\Gamma = m_1, \dots, m_M$  under categories  $M$  can be obtained.

### 3.2.2 The Global Fusion of the Fusion Results Between Classes

For the mass function set of categories  $M$  targets, the mass value constructed by the samples of category  $w_i$  is mainly assigned to the corresponding focal element, that is  $m_i(w_i)$ . Therefore, it can be considered that the distribution of mass values of different categories is different, and there will be certain conflicts between the mass functions obtained by equation (11). At this time, if the Dempster combination rule is

used for global fusion, a fusion result that contradicts the facts may be obtained. Therefore, when fusing between categories, this paper uses PCR5 (Redistribute Conflicting Mass Proportionally Rule 5, PCR5) combination rules to accurately and reasonably allocate conflict information. The PCR5 rule is an evidence fusion rule proposed by Desert and Smarandache for conflicting data. This rule can accurately distribute the conflict information proportionally according to the mass values of the two parties in the conflict, which is very suitable for combining high conflict evidence. While compared with the Dempster rule, it is a more conservative combination method, and the convergence speed of the fusion result is relatively slow. Assuming that B and C are two independent evidences to be combined, the corresponding focal elements are  $B_j$  and  $C_j$ , and the mass functions are  $m_1$  and  $m_2$  respectively, then the PCR5 rule can be expressed as [19]:

$$m(x) = \sum_{\substack{B_i, C_j \in 2^\Omega \\ B_i, C_j = x}} m_1(B_i) m_2(C_j) + \sum_{\substack{y \in 2^\Omega \\ x \cap y = \emptyset}} \left[ \frac{m_1(x)^2 m_2(y)}{m_1(x) + m_2(y)} + \frac{m_2(x)^2 m_1(y)}{m_2(x) + m_1(y)} \right] \quad (12)$$

where  $x$  and  $y$  are two focal elements of evidence body  $B$  and  $C$  with conflicting information.

**Example 1:** The evidence body  $m_1$  confirms that the mass value of that the target belongs to category  $w_1$  is 0.9, and the mass value of that the target belongs to category  $w_3$  is 0.1. The evidence body  $m_2$  confirms that the mass value of that the target belongs to category  $w_2$  is 0.9, and the mass value of that the target belongs to category  $w_3$  is 0.1.

After combining  $m_1$  and  $m_2$  by the Dempster combination rule, it can be obtained that:

$$m_{Dempster}(w_1) = 0, m_{Dempster}(w_2) = 0, m_{Dempster}(w_3) = 1.$$

After combining  $m_1$  and  $m_2$  by the PCR5 combination rule, it can be obtained that:

$$m_{PCR5}(w_1) = 0.486, m_{PCR5}(w_2) = 0.486, m_{PCR5}(w_3) = 0.028.$$

It can be seen that the original two evidences respectively believe that the target belongs to  $w_1$  and  $w_2$ , and the reliability values provided are all 0.9. The Dempster rule offers a fusion result contrary to  $m_1$  and  $m_2$ , which is obviously not reasonable. While PCR5 believes that the mass values of that the target belongs to  $m_1$  and  $m_2$  are still the same and are much higher than the mass value of that the target belongs to  $w_3$ . The fusion result is more reasonable and credible than the result of Dempster's rule.

Therefore, considering the inconsistency of evidence between categories, for the mass function set of the categories  $M$  observation target, it is necessary to use the PCR5 combination rule for fusion, and it can be obtained that:

$$m = m_1 \overset{PCR5}{\oplus} m_2 \overset{PCR5}{\oplus} \cdots \overset{PCR5}{\oplus} m_M \quad (13)$$

where  $\overset{PCR5}{\oplus}$  represents PCR5 combination operation.

The final fusion result  $m$  can be calculated according to equation (13). According to the mass value of each category assigned to it, the final recognition result can be made on the target  $x$ , that is, the unknown target  $x$  is assigned to the category with the maximum mass value.

The entire process of the NEkNN classification algorithm proposed in this paper is shown in Fig. 1.

## 4 Experimental Results

This section is to use two experiments to compare the NEkNN classification algorithm mentioned in this paper with voting kNN, EkNN, IEkNN, and to illustrate the effectiveness of NEkNN. In the experiments, the parameters in EkNN are optimized according to the existing method [11]. The two experiments use simulation data and the standard test data sets of UCI database to comprehensively compare and analyze the misclassification rate of the proposed method and other classification methods based on  $k$ -nearest neighbors.

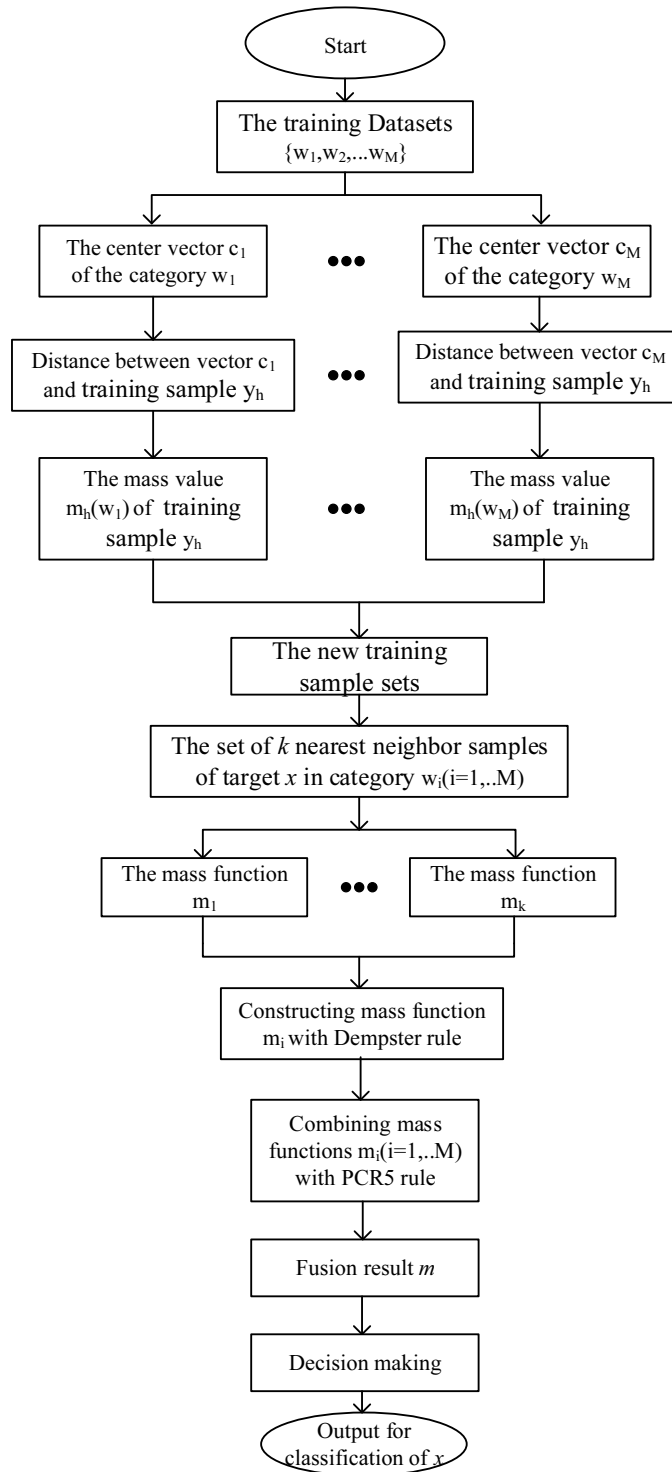


Fig. 1. Flowchart of the NEkNN algorithm

#### 4.1 Experiment 1 (simulation Data Sets)

In this target recognition simulation experiment, a classification problem of 3-class target  $\Omega = \{w_1, w_2, \dots, w_3\}$  is considered. After the sensor's observation data is preprocessed by data association and feature extraction, the training database is used to classify the feature data containing these three categories of targets. Assuming that the feature data is a three-dimensional vector, and the observation data and training data are generated from three three-dimensional data sets that obey the Gaussian distribution, then its mean and standard deviation have the following characteristics:

In Table 1, the three characteristics of each category of data have the same standard deviation. For example, the probability density functions of the three attribute data of category  $w_2$  are:  $x_1 | w_2 \sim N(-1, 1)$ ,  $x_2 | w_2 \sim N(1, 1)$ ,  $x_3 | w_2 \sim N(0, 1)$ , their standard deviation is the same as 1, and it randomly generates  $3 \times 100$  test samples and  $3 \times 200$  training samples. It is to select  $k$ NN,  $E_k$ NN,  $IE_k$ NN to compare and analyze with  $NE_k$ NN proposed in this paper, take the value of the adjacent number  $k$  from 5 to 15, and take the average of 10 simulation results as the error rate of the test data set. The classification results are shown in Table 2.

**Table 1.** 3-class data set with 3D Gaussian distributions

Label	$\mu_1$	$\mu_2$	$\mu_3$	Standard deviation
$w_1$	1	1	1	1
$w_2$	-1	1	0	1
$w_3$	0	-1	1	2

**Table 2.** 3-class data set with 3D Gaussian distributions

$k$	$k$ NN	$E_k$ NN	$IE_k$ NN	$NE_k$ NN
$k = 5$	32.73	29.12	25.12	23.60
$k = 6$	32.68	28.33	25.19	23.97
$k = 7$	31.19	27.87	25.38	24.06
$k = 8$	33.62	27.61	24.67	24.11
$k = 9$	31.44	27.59	24.71	23.86
$k = 10$	32.25	27.55	25.09	24.03
$k = 11$	31.28	27.43	25.02	23.79
$k = 12$	31.42	26.96	24.91	23.45
$k = 13$	30.99	26.93	24.47	23.52
$k = 14$	32.94	26.86	24.52	23.35
$k = 15$	32.17	26.98	24.59	23.35

It can be seen from Table 2 that the  $E_k$ NN and  $IE_k$ NN methods are better than the traditional voting  $k$ NN and can effectively improve the classification accuracy. The  $NE_k$ NN algorithm proposed in this paper can better characterize the imprecision of the sample data in the overlapping area of the category and improve the classification accuracy of the data by using the basic belief assignment to replace the original category label. Therefore, compared with the  $E_k$ NN and  $IE_k$ NN algorithms, it has a smaller classification error rate, especially when the number of neighbors is small, the performance improvement is more significant. In addition, it can be found that compared with other classification methods based on  $k$ -nearest neighbors,  $NE_k$ NN method is less sensitive to the value of  $k$ -nearest neighbor.

#### 4.2 Experiment 2 (Real Data Sets)

In this experiment, it is to use three commonly used data sets from the UCI database (such as Iris, Ecoli and Wine) to verify and analyze the classification performance of  $NE_k$ NN. In the Ecoli data set, three categories of relatively similar data were selected, that is, cp, im and imU. The specifications of the selected data sets are given in Table 3.

**Table 3.** Specifications of the selected data sets

Data Set	Categories	Attributes	Samples
Iris	3	4	150
Ecoli	3	7	255
Wine	3	13	178

The  $k$ -fold cross validation method is used to verify the classification effect of various classification algorithms on the data sets Iris, Ecoli and Wine. Generally speaking, the value of  $k$  is an uncertain



parameter. The 5-fold cross validation method selected in this experiment is to divide each category of sample data in all data sets into 5 parts of the same capacity, and take one of them as the test sample, and the remaining 4 as the training sample. In this way, each sample can be used as training data or as a test sample. The experimental result is the average misclassification rate of 5 rounds of testing. The misclassification rates of  $k$ NN,  $E$  $k$ NN,  $IE$  $k$ NN and  $NE$  $k$ NN of different values of  $k$  on different data sets are shown in Table 4 to Table 6.

**Table 4.** Classification results of different methods for Iris data (%)

$k$	$k$ NN	$E$ $k$ NN	$IE$ $k$ NN	$NE$ $k$ NN
$k = 5$	6.79	6.04	3.84	2.71
$k = 6$	6.61	5.68	4.63	3.08
$k = 7$	6.45	5.62	3.90	3.10
$k = 8$	5.59	4.92	4.77	2.69
$k = 9$	5.59	4.87	3.41	1.94
$k = 10$	4.83	4.95	2.72	1.94
$k = 11$	5.72	4.52	3.67	3.15
$k = 12$	5.76	4.52	3.67	3.12
$k = 13$	5.69	3.76	3.23	2.77
$k = 14$	5.72	3.79	3.23	2.77
$k = 15$	6.24	3.79	3.23	2.89

**Table 5.** Classification results of different methods for Ecoli data (%)

$k$	$k$ NN	$E$ $k$ NN	$IE$ $k$ NN	$NE$ $k$ NN
$k = 5$	10.98	10.56	8.41	6.28
$k = 6$	11.72	10.58	8.59	6.67
$k = 7$	11.63	9.97	8.81	5.92
$k = 8$	12.07	10.88	7.98	6.50
$k = 9$	11.93	11.53	8.76	5.27
$k = 10$	12.88	11.81	8.16	4.35
$k = 11$	12.62	12.25	7.94	4.40
$k = 12$	13.21	11.80	8.76	4.37
$k = 13$	13.43	11.55	7.15	5.14
$k = 14$	13.43	12.01	8.16	4.73
$k = 15$	14.10	12.42	9.16	5.39

**Table 6.** Classification results of different methods for Wine data (%)

$k$	$k$ NN	$E$ $k$ NN	$IE$ $k$ NN	$NE$ $k$ NN
$k = 5$	28.50	28.47	19.71	17.24
$k = 6$	28.32	25.69	23.59	16.73
$k = 7$	28.67	28.39	20.36	17.24
$k = 8$	30.49	29.65	21.14	18.09
$k = 9$	27.42	27.21	20.47	16.54
$k = 10$	28.71	27.06	19.32	16.52
$k = 11$	26.48	25.80	20.53	17.30
$k = 12$	27.11	26.91	19.38	18.01
$k = 13$	27.15	26.65	19.33	17.31
$k = 14$	26.83	25.53	19.26	16.13
$k = 15$	25.95	26.68	18.73	16.87

As can be seen from Table 4 to Table 6, the  $NE$  $k$ NN proposed in this paper always has the lowest misclassification rate, because many samples that are difficult to accurately classify can obtain richer category information through the replaced category labels, which makes the pre-processed sample data

provide better classification performance than other classification methods based on  $k$  nearest neighbors, especially when the nearest neighbor data is less, the good classification effect can also be obtained. At the same time, the NE $k$ NN algorithm is not very sensitive to the value of the  $k$ -nearest neighbor in the real data set, and can achieve a low misclassification rate when only considering a small number of neighbor samples. The reason is that the involved imprecise information of each training sample can be well characterized with the evidential editing procedure, and thus a reasonable classification result can be obtained even with a small value of  $k$ .

## 5 Conclusion

In order to effectively express and reason about imprecise data, this paper proposes an evidence editing classification method based on E $k$ NN. Before identifying and classifying the sample to be tested, this method uses a class label with a basic belief assignment structure to replace the original numerical category of the training sample, so that the training sample in the class overlapping area can provide more abundant and more diverse category information. And lay a better foundation for the follow-up  $k$ NN classification process. From the comparative analysis of related experiments, it can be found that for imprecise sample data, the NE $k$ NN algorithm can obtain better classification performance than other classification algorithms based on  $k$ -nearest neighbors. At the same time, the algorithm is not very sensitive to the value of the neighbor  $k$  of the observation data, which enables the sensor node to obtain good target recognition results even with less neighbor data, it has extremely important application value for wireless sensor networks with limited computing power and energy supply.

In our research, the training sample preprocessing is performed according to the distance from it to the center of the corresponding category. However, when dealing with some special irregular sample data sets, their categories may not be represented by class center vectors. As a result, the future work mainly involves the following two aspects: (1) finding a more efficient strategy to preprocess the training samples to improve the classification accuracy; (2) designing more credible combination rules to deal with the uncertain data in IOT environment.

## Acknowledgements

This work was supported by the Fundamental Research Funds for the Central Universities 2019RC044.

## References

- [1] O. Bamgboye, X. Liu, P. Cruickshank, Towards modelling and reasoning about uncertain data of sensor measurements for decision support in smart spaces, in: Proc. 2018 IEEE 42nd annual computer software and applications conference (COMPSAC), 2018.
- [2] Z.-G. Liu, Q. Pan, J. Dezert, G. Mercier, Hybrid classification system for uncertain data, IEEE Transactions on Systems, Man, and Cybernetics: Systems 47(10)(2017) 2783-2790.
- [3] A.M. Andrew, A.Y.M. Shakaff, A. Zakaria, R. Gunasagaran, E. Kanagaraj, S.M. Saad, Fuzzy K-Nearest Neighbour (FkNN) Based Early Stage Fire Source Classification in Building, in: Proc. 2018 IEEE Conference on Systems, Process and Control, 2018.
- [4] C. Viji, J.B. Raja, R.S. Ponmagal, S.T. Suganthi, P. Parthasarathi, S. Pandiyan, Efficient Fuzzy based K-Nearest Neighbour Technique for Web Services Classification, Microprocessors and Microsystems 76(2020) 103097.
- [5] Y. Deng, Uncertainty measure in evidence theory, Science China (Information Sciences) 63(11)(2020) 5-23.
- [6] Z. Zhang, W. Zhang, H.-C. Chao, C.-F. Lai, Toward Belief Function-Based Cooperative Sensing for Interference Resistant Industrial Wireless Sensor Networks, IEEE Transactions on Industrial Informatics 12(6)(2016) 2115-2126.

- [7] Y. Han, Y. Deng, An Evidential Fractal Analytic Hierarchy Process Target Recognition Method, *Defence Science Journal* 68(4)(2018) 367-373.
- [8] Y. Li, S. Yao, R. Zhang, C. Yang, Analyzing host security using D- S evidence theory and multisource information fusion, *International Journal of Intelligent Systems* 36(2)(2021) 1053-1068.
- [9] J. Luo, L. Shi, Y. Ni, Uncertain Power Flow Analysis Based on Evidence Theory and Affine Arithmetic, *IEEE Transactions on Power Systems* 33(1)(2018) 1113-1115.
- [10] Z.-G. Su, T. Denoeux, Y.-S. Hao, M. Zhao, Evidential K-NN classification with enhanced performance via optimizing a class of parametric conjunctive t-rules, *Knowledge-Based Systems* 142(2018) 7-16.
- [11] T. Denoeux, O. Kanjanatarakul, S. Sriboonchitta, A new evidential K-nearest neighbor rule based on contextual discounting with partially supervised learning, *International Journal of Approximate Reasoning* 113(2019) 287-302.
- [12] S. Faziludeen, P. Sankaran, ECG beat classification using evidential K-nearest neighbours, *Procedia Computer Science* 89(2016) 499-505.
- [13] L. Jiao, X. Geng, Q. Pan, EEkNN: k-Nearest Neighbor Classifier with an Evidential Editing Procedure for Training Samples, *Electronics* 8(5)(2019) 592.
- [14] Z.-G. Liu, Y. Liu, J. Dezert, Q. Pan, Classification of incomplete data based on belief functions and K-nearest neighbors, *Knowledge-Based Systems* 89(2015) 113-125.
- [15] Y. Zhang, J. Hou, Z.-G. Liu, Q. Pan, A New Evidential K-Nearest Neighbors Data Classification Method, *Fire Control and Command Control* 38(9)(2013) 58-61.
- [16] X.-L. Chen, P.-H. Wang, Y.-S. Hao, M. Zhao, Evidential KNN-based condition monitoring and early warning method with applications in power plant, *Neurocomputing* 315(2018) 18-32.
- [17] X. Xu, J. Zheng, J.-B. Yang, D.-L. Xu, Y.-W. Chen, Data classification using evidence reasoning rule, *Knowledge-Based Systems* 116(2017) 144-151.
- [18] H. Sarmadi, A. Entezami, B. Saeedi Razavi, K.-V. Yuen, Ensemble learning- based structural health monitoring by Mahalanobis distance metrics, *Structural Control and Health Monitoring* 28(2)(2021) e2663.
- [19] J. Dezert, F. Smarandache, Canonical decomposition of dichotomous basic belief assignment, *International Journal of Intelligent Systems* 35(7)(2020) 1105-1125.