



A Proximity Weighted Evidential k Nearest Neighbor Classifier for Imbalanced Data

Md. Eusha Kadir¹(✉), Pritom Saha Akash¹, Sadia Sharmin², Amin Ahsan Ali³, and Mohammad Shoyaib¹

¹ Institute of Information Technology, University of Dhaka, Dhaka, Bangladesh
{bsse0708,bsse0604}@iit.du.ac.bd, shoyaib@du.ac.bd

² Islamic University of Technology, Gazipur, Bangladesh
sharmin@iut-dhaka.edu

³ Independent University, Dhaka, Bangladesh
aminali@iub.edu.bd

Abstract. In k Nearest Neighbor (k NN) classifier, a query instance is classified based on the most frequent class of its nearest neighbors among the training instances. In imbalanced datasets, k NN becomes biased towards the majority instances of the training space. To solve this problem, we propose a method called Proximity weighted Evidential k NN classifier. In this method, each neighbor of a query instance is considered as a piece of evidence from which we calculate the probability of class label given feature values to provide more preference to the minority instances. This is then discounted by the proximity of the neighbor to prioritize the closer instances in the local neighborhood. These evidences are then combined using Dempster-Shafer theory of evidence. A rigorous experiment over 30 benchmark imbalanced datasets shows that our method performs better compared to 12 popular methods. In pairwise comparison of these 12 methods with our method, in the best case, our method wins in 29 datasets, and in the worst case it wins in least 19 datasets. More importantly, according to Friedman test the proposed method ranks higher than all other methods in terms of AUC at 5% level of significance.

Keywords: Classifier · Imbalanced learning · k NN · Evidence theory

1 Introduction

Classification is one of the most important tasks in machine learning. Numerous classification approaches, such as k Nearest Neighbor (k NN) [9], Decision Tree (DT), Naïve Bayes (NB), and Support Vector Machine, have been well developed and applied in many applications. However, most of the classifiers face serious trouble for imbalanced class distribution and thus learning from the imbalanced dataset is one of the top ten challenging problems in data mining research [20].

To solve class imbalance problem, various strategies have already been proposed which can be grouped into two broad categories namely data oriented and algorithm oriented approaches. Data oriented approaches use sampling techniques. In order to make dataset balanced, the sampling techniques either over-sample the minority instances or select instances (under-sample) from the majority class. A sampling technique namely Synthetic Minority Over-sampling TEchnique (SMOTE) has been proposed that increases the number of minority class instances by creating artificial and non-repeated samples [4].

In contrast, algorithm oriented approaches are the modifications of traditional algorithms such as DT and k NN. The modified DTs for imbalanced classification are Hellinger Distance DT (HDDT) [5], Class Confidence Proportion DT (CCPDT) [13] and Weighted Inter-node Hellinger Distance DT (iHDwDT) [1]. These DTs use different splitting criteria while selecting a feature in split point.

k NN is one of the simplest classifiers. Despite its simplicity, k NN is considered as one of the top most influential data mining algorithms [19]. Traditional k NN finds the k closest instances from the training data to a query instance and treats all neighbors equally. Dudani has proposed a distance based weighted k NN which provides more weights to closer neighbors [8]. Another variant of k NN approach, Generalized Mean Distance based k NN (GMDKNN) [10], has been presented by introducing multi-generalized mean distance and the nested generalized mean distance. All these variants of k NN are sensitive to the majority instances and thus perform poorly for imbalanced datasets.

Considering this imbalance problem, several researchers extended k NN for imbalanced datasets [7, 11, 12]. In Exemplar-based k NN (k ENN) [11], Li and Zhang expand the decision boundary for the minority class by identifying the exemplar minority instances. A weighting algorithm namely Class Confidence Weighted k NN (CCWKNN) has been presented in [12] where the probability of feature values given the class labels is considered as weight. Dubey and Pudi have proposed a weighted k NN (WKNN) [7] which considers the class distribution in a wider region around a query instance. The class weight for each training instance is estimated by taking the local class distributions into account.

The purpose of these existing studies is to improve the overall performance for imbalanced data. However, these methods overlook the problem of uncertainty which is prevalent in almost all datasets [18]. The reason behind this uncertainty is that the complete statistical knowledge associated with the conditional density function of each class is hardly available [6]. To address this problem, k NN has been extended using Dempster-Shafer Theory of evidence (DST) to better model uncertain data named Evidential k NN (EKNN) [6]. In EKNN, each neighbor assigns basic belief on classes based on a distance measure. Nevertheless, this approach again does not take consideration of the class imbalance problem.

To address these aforementioned problems, we propose a Proximity weighted Evidential k NN (PE k NN) classifier and make the following contributions. Firstly, we have proposed a confidence (posterior) assignment procedure on each neighbor of a query instance. Secondly, we have also proposed to use proximity of a

neighbor as a weight to discount the confidence of a neighbor. It is shown that, this weighted confidence increases the likelihood of classifying a minority class. Thirdly, DST framework is used to combine decisions from different neighbors.

2 Dempster-Shafer Theory of Evidence

Dempster-Shafer theory of evidence is a generalized form of Bayesian theory. It assigns degree of belief for all possible subsets of the hypothesis set. Let, $C = \{C_1, \dots, C_M\}$ be a finite hypothesis set of mutually exclusive and exhaustive hypotheses. The belief in a hypothesis assigned based on a piece of evidence is ranged numerically as $[0, 1]$. A Basic Belief Assignment (BBA) is a function $m : 2^C \rightarrow [0, 1]$ which satisfies the following properties:

$$m(\emptyset) = 0 \quad \text{and} \quad \sum_{A \subseteq C} m(A) = 1 \quad (1)$$

where $m(A)$ is a degree of belief (referred as mass) which reflects how strongly A is supported by the piece of evidence. $m(C)$ represents the degree of ignorance.

Several pieces of evidence characterized by their BBAs can be fused using Dempster's rule of combination [16]. For two BBAs $m_1(\cdot)$ and $m_2(\cdot)$ which are not totally conflicting, the combination rule can be expressed using Eq. (2).

$$m(A) = \frac{\sum_{B \cap C = A} m_1(B)m_2(C)}{1 - \sum_{B \cap C = \emptyset} m_1(B)m_2(C)} \quad A \neq \emptyset \quad (2)$$

where $A, B, C \in 2^C$ and $\sum_{B \cap C = \emptyset} m_1(B)m_2(C) < 1$.

For decision making, Belief, Plausibility and betting Probability (P_{bet}) are usually used. For a singleton class A , $P_{bet}(A)$ is derived in Eq. (3) where $|B|$ represents the cardinality of the element B .

$$P_{bet}(A) = \sum_{A \subseteq B} \frac{|A \cap B|}{|B|} \times m(B) \quad (3)$$

3 Proximity Weighted Evidential k NN (PE k NN)

k NN faces difficulty in imbalanced datasets as it treats all neighbors of the query instance equally and most of the neighbors will be of the majority class. To deal with this issue, the proposed algorithm attempts to provide more importance to neighbors with a higher proximity weighted confidence. Here, confidence of an instance indicates a conditional probability of that instance based on training data. Algorithms such as NB also uses conditional probability while classifying a query instance. However, the performance of NB degrades due to the poor estimation of the conditional density of the query instance associated with each class. In contrast, PE k NN computes conditional probability of neighborhood instances rather than query instance. Furthermore, as uncertainty is prevalent in

almost all datasets [18]. This is more significant for imbalanced datasets where little information is available for the minority class. To deal with this issue, PE k NN uses DST to combine the evidences provided by each neighbor.

For a new query instance (x_t), PE k NN first finds k closest neighbors according to some distance measurement (e.g. Euclidean distance). Let, $S(x_t, k)$ be the set of k closest neighbors of x_t and each member of $S(x_t, k)$ is considered as a piece of evidence which assigns mass values for each subset of C known as BBA.

Now, consider x_i as the i -th neighbor of x_t belonging to class C_q . As x_i is a piece of evidence belonging to C_q , some part of its belief will be committed to C_q . The rest of the belief can not be distributed to any other subset of C except itself. The BBA provided by x_i can be represented by Eq. (4), (5) and (6) where $0 < \beta_0 < 1$.

$$m_i(\{C_q\}) = \beta = \beta_0 \times \Psi(x_i, x_t) \quad (4)$$

$$m_i(A) = 0 \quad \forall A \in 2^C \setminus \{C, \{C_q\}\} \quad (5)$$

$$m_i(C) = 1 - \beta \quad (6)$$

Now, we will discuss about two of our intuitions. First, a piece of evidence belonging to C_q will assign a larger belief to C_q when the evidence is more reliable which we call confidence. An evidence having higher posterior probability should get more confidence than the one which is in lower posterior probability region. The second intuition is that a neighbor will assign more belief to a specific class when the neighbor and the query instance are more proximate. The function defined in Eq. (7), $\Psi(\cdot)$ satisfies the two aforementioned intuitions where p_i is the confidence of x_i represented by the probability of class label (y_i) given x_i and $prx(x_i, x_t)$ represents the proximity between x_i and x_t .

$$\Psi(x_i, x_t) = prx(x_i, x_t) \times p_i \quad (7)$$

The procedure how PE k NN algorithm classifies a query instance is presented in Algorithm 1. The confidence assignment, proximity estimation and decision making steps are presented in detail in Sects. 3.1, 3.2 and 3.3 respectively.

3.1 Estimation of Confidence

The confidence (p_i) of an instance x_i ($x_i \in \mathbb{R}^l$) belonging to y_i is assigned in the following manner derived in Eq. (8).

$$p_i = P(y_i | x_i) = \frac{P(y_i) \times P(x_i | y_i)}{\sum_{j=1}^M P(C_j) \times P(x_i | C_j)} \quad (8)$$

where $y_i \in \{C_1, C_2, \dots, C_M\}$, $P(C_j)$ represents the prior of C_j in training space and $P(x_i | C_j)$ represents the likelihood in Bayes' theorem. Here, two approaches of estimating class-wise Probability Density Function (PDF) is presented. First one is using Single Gaussian Model (SGM) and another one is using Gaussian Mixture Model (GMM). When PE k NN uses confidence derived from SGM, we call it sPE k NN, and mPE k NN when it uses confidence derived from GMM.

Algorithm 1: PE k NN Algorithm

Input : Training data (X), Training data labels (Y), Neighborhood size (k),
Query instance (x_t)

Output: Predicted class label (y_t)

- 1 $conf, d_{max} \leftarrow \text{fitModel}(X, Y)$
- 2 $s \leftarrow$ indices of k nearest neighbors of x_t
- 3 Initialize a list bba of mass values
- 4 **for** $i = 1$ to k **do**
- 5 $index \leftarrow s[i]$
- 6 $confidence \leftarrow conf[index]$
- 7 $d \leftarrow \text{distance}(x_t, X[index])$
- 8 $proximity \leftarrow$ calculate proximity using Eq. (11) from d, d_{max}
- 9 $bba[i] \leftarrow$ assign mass value using Eq. (4), (5), (6) and (7) from $confidence,$
 $proximity$
- 10 **end**
- 11 $m \leftarrow$ combine mass values from bba using Eq. (2)
- 12 $P_{bet} \leftarrow$ calculate betting probabilities for all classes using Eq. (3) from m
- 13 $y_t \leftarrow$ calculate decision using Eq. (12) from P_{bet}
- 14 **function** $\text{fitModel}(X, Y)$:
- 15 Initialize Array, $conf$
- 16 $d_{max} \leftarrow 0$
- 17 **for** $i = 1$ to $|X|$ **do**
- 18 $conf[i] \leftarrow$ Calculate confidence using Eq. (8) from $X[i], Y[i]$
- 19 **for** $j = i+1$ to $|X|$ **do**
- 20 $d \leftarrow \text{distance}(X[i], X[j])$
- 21 $d_{max} \leftarrow \max(d, d_{max})$
- 22 **end**
- 23 **end**
- 24 **return** $conf, d_{max}$

Single Gaussian model assumes that all the features are independent and the continuous values associated with each class follow a normal distribution. Under these assumptions, the likelihood function can be represented as Eq. (9).

$$P(x) = \prod_{j=1}^l P(x_j) = \prod_{j=1}^l f(x_j; \mu_j, \sigma_j^2) = \prod_{j=1}^l \frac{1}{\sqrt{2\pi}\sigma_j} \times \exp\left(-\frac{(x_j - \mu_j)^2}{2\sigma_j^2}\right) \quad (9)$$

where x_j denotes the j -th feature of x and $f(\cdot)$ represents the normally distributed PDF parameterized by mean (μ) and variance (σ^2).

On the other hand, GMM can also be used to estimate PDF from multivariate data. The class-wise PDF using m -component mixture model is given in Eq. (10).

$$P(x) = \sum_{i=1}^m \alpha_i P(x | Z_i) \quad (10)$$

The procedure of finding complete set of parameters ($Z_1, \dots, Z_m, \alpha_1, \dots, \alpha_m$) specifying the mixture model is briefly described in [14].

3.2 Estimation of Proximity

To capture the proximity between two instances, some distance measurement can be used. The proximity between two instances (x_i and x_j) from training samples will be maximum when x_i and x_j are identical. On the other hand, it will be lowest when they are the farthest two instances in the feature space. To measure this proximity, a normalization is applied as Eq. (11) so that $prox(x_i, x_j) \in [0, 1]$. Here, d_{max} is the distance between two farthest training instances.

$$prox(x_i, x_j) = 1 - \frac{d(x_i, x_j)}{d_{max}} \quad (11)$$

3.3 Decision Making

According to Eq. (7), $\Psi(\cdot)$ will return a larger value when a neighbor is more confident and more closer to the query instance. Now, for each of the k nearest neighbors, the BBAs are defined using Eq. (4), (5) and (6). In order to classify x_t , these BBAs are combined using DST. The betting probability (P_{bet}) for each singleton class from this combined decision will be then calculated using Eq. (3). Finally, the decision from this P_{bet} is taken using Eq. (12).

$$\hat{y} = \arg \max_{c \in \{C_1, \dots, C_M\}} P_{bet}(c) \quad (12)$$

where c is a singleton class so that the cardinality of c is 1.

Properties of β : Value of β is bounded between 0 to 1.

Proof. From Eq. (4), (7) and (8), it can be derived that,

$$\beta = \beta_0 \times P(y_i | x_i) \times prox(x_i, x_t) \quad (13)$$

Here, β_0 is a user given constant satisfying $0 < \beta_0 < 1$. The second term, $P(y_i | x_i)$, represents the posterior probability. The last term, $prox(x_i, x_t)$ is at most equal to 1 and at least equal to zero. As can be seen from Eq. (13), β is a product of three terms and all these terms are bounded between 0 to 1. It is sufficient to claim that, the value of β must be bounded between 0 to 1.

3.4 An Illustrative Example

Figure 1 shows the instances of a two-class imbalance problem where (+)s and (•)s represent the minority (Class-A) and majority class instances (Class-B) instances respectively. The class boundaries are represented as dotted lines and three query instances (t_1, t_2, t_3) are marked with (★)s. Here, first query instance t_1 is situated in a majority class region bounded by minority instances. Both k NN and PE k NN can successfully classify t_1 . Traditional algorithms such as C4.5 and NB face difficulties in this situation.

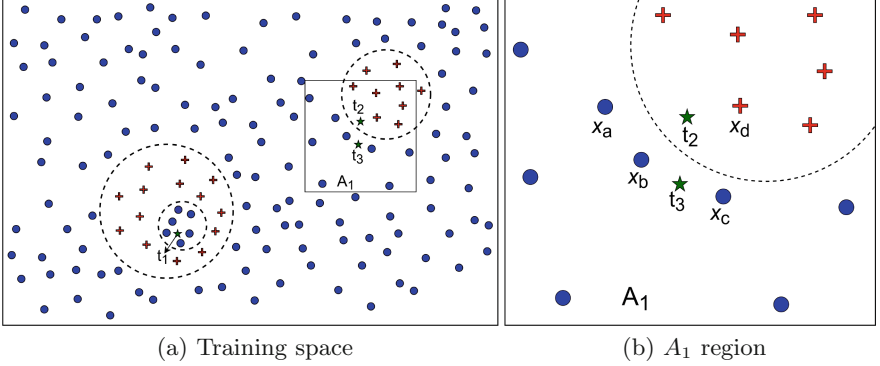


Fig. 1. A synthetic imbalanced dataset

The two other query instances t_2 and t_3 associated with a region namely A_1 (see Fig. 1b). Here, for both t_2 and t_3 , the four neighbors are x_a , x_b , x_c and x_d . Traditional k NN with $k = 4$, will classify both t_2 and t_3 as Class-B. PE k NN, on the other hand, considers the confidence of each neighbor. Here, x_d will provide a higher confidence compared to majority class instances (x_a , x_b and x_c). Assume, the confidence of x_a , x_b , x_c , and x_d are 0.30, 0.40, 0.30, and 0.75 respectively. And the proximity with respect to t_2 are 0.90, 0.95, 0.85 and 0.95 respectively. Then BBAs assigned by PE k NN for these neighbors are $m_a(\{B\}) = 0.2565$, $m_a(\{A, B\}) = 0.7435$, $m_b(\{B\}) = 0.3610$, $m_b(\{A, B\}) = 0.6390$, $m_c(\{B\}) = 0.2423$, $m_c(\{A, B\}) = 0.7577$ and $m_d(\{A\}) = 0.6769$, $m_d(\{A, B\}) = 0.3231$. Here, β_0 is set to 0.95. Now, combing these BBAs using DST, we get $P_{bet}(A) = 0.5325$ and $P_{bet}(B) = 0.4675$ which indicates that t_2 will be correctly classified as Class-A.

On the other hand, for the query instance t_3 , the proximity of x_a , x_b , x_c and x_d are 0.85, 0.95, 0.95 and 0.85 respectively. We, therefore, get $P_{bet}(A) = 0.4661$ and $P_{bet}(B) = 0.5339$ indicating that t_3 will be classified as Class-B. Therefore, t_3 is correctly classified as a majority class instance even though the neighbors of t_2 and t_3 are same.

Instead of DST, let us reconsider simpler techniques to combine evidences such as summing and taking the maximum of the proximity weighted confidences. If we simply sum class-wise proximity weighted confidences, both t_2 and t_3 get a higher value for Class-B as three of the four neighbors belong to that class. To avoid this bias, a query can be simply classified in the class for which it gets maximum proximity weighted confidence among the neighbors. But this method does not consider the local neighborhood priors. For which, it will classify both t_2 and t_3 as minority class which is not desired. PE k NN on the other hand using the DST framework successfully classifies both query instances.

4 Experiments and Results

Dataset description, implementation details and the performance metrics followed by the results obtained from the experiments with discussion are given in the following subsections.

4.1 Dataset Description

The characteristics of the 30 benchmark datasets are shown in Table 1 which are collected from UCI machine learning repository [3] and KEEL Imbalanced Datasets [2]. Imbalance Ratio (IR) between the samples of majority class and minority class of the datasets used in these experiment are at least 1.5 and values of all the features are numeric. A dataset is highly imbalanced when the value of IR is very high.

Table 1. Descriptions of Imbalanced Datasets. Idx, #Inst, #Cl and #Ftr represent index of a dataset, number of instances, classes and features respectively.

Idx	Name	#Ftr	#Cl	#Inst	IR	Idx	Name	#Ftr	#Cl	#Inst	IR
01	Appendicitis	7	2	106	4.05	16	Shuttle-c0-vs-c4	9	2	1829	13.87
02	Ecoli1	7	2	336	3.36	17	Vehicle0	18	2	846	3.25
03	Ecoli2	7	2	336	5.46	18	Vehicle1	18	2	846	2.90
04	Ecoli3	7	2	336	8.60	19	Vehicle2	18	2	846	2.88
05	Ecoli4	7	2	336	15.80	20	Vehicle3	18	2	846	2.99
06	Glass-0-1-2-3_vs_4-5-6	9	2	214	3.20	24	Yeast-1_vs_7	7	2	459	14.30
07	Glass1	9	2	214	1.82	21	Vowel0	13	2	988	9.98
08	Glass4	9	2	214	15.46	22	Wisconsin	9	2	683	1.86
09	Glass6	9	2	214	6.38	23	Yeast-0-5-6-7-9_vs_4	8	2	528	9.35
10	Haberman	3	2	306	2.78	25	Yeast-1-2-8-9_vs_7	8	2	947	30.57
11	Ionosphere	34	2	351	1.79	26	Yeast-2_vs_8	8	2	482	23.10
12	New-thyroid1	5	2	215	5.14	27	Yeast1	8	2	1484	2.46
13	Page-blocks0	10	2	5472	8.79	28	Yeast3	8	2	1484	8.10
14	Pima	8	2	768	1.87	29	Yeast5	8	2	1484	32.73
15	Segment0	19	2	2308	6.02	30	Yeast6	8	2	1484	41.40

4.2 Implementation Details and Performance Metrics

PE k NN is benchmarked against other algorithms including traditional learning algorithms (k NN, C4.5, NB), oversampling strategy (SMOTE), recent algorithms in the k NN family (EKNN, WKNN, CCWKNN, k ENN, GMDKNN) and few tree based recent algorithms for imbalanced classification (CCPDT, HDDT, iHDwDT). For PE k NN, we use $\beta_0 = 0.95$ in this experiment. For k ENN, the confidence level is set 0.1 and we set $p = 1$ for GMDKNN.

We have conducted 10-fold stratified cross validation to evaluate the performance of the proposed method. The Receiver Operating Characteristic (ROC) curve [17] is widely used to evaluate imbalanced classification. We use Area Under the ROC Curve (AUC) for evaluating the classifier performance.

For comparison, all the classifiers are ranked on each dataset in terms of AUC, with ranking of 1 is the best. We also perform Friedman tests on the ranks. After rejecting the null hypothesis using Friedman test that all the classifiers are equivalent, a post-hoc test called Nemenyi test [15] is used to determine the performance of which classifier is significantly better than the others.

4.3 Result and Discussion

Table 2 represents the comparison of 14 classifiers over 30 imbalanced datasets. The average ranks of these classifiers indicate that k NN is better performing algorithm compared to other traditional classifiers on imbalanced datasets. Though k NN performs better than C4.5, modifications of tree based algorithms for imbalanced datasets perform better than k NN. Moreover, k NN on SMOTE sampled datasets performs slightly better than k NN without sampling.

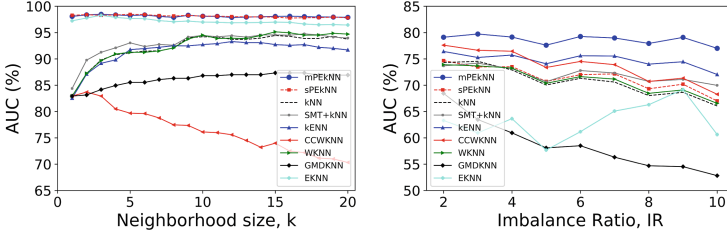
Now, if we compare k NN with its different variants, it can be observed that k ENN and WKNN improve the overall performance of traditional k NN although another variant CCWKNN fails to improve the performance in most cases over the experimented datasets. Moreover, it is investigated that, the recent generalized mean based k NN approach GMDKNN performs worse than k NN on imbalanced datasets. In contrast, we can observe from Table 2 that, EKNN performs better than all other classifiers except the proposed $sPEk$ NN and $mPEk$ NN. It indicates that, handling uncertainty can improve the performance of k NN on imbalanced datasets. Finally, average ranks show that $mPEk$ NN is the best performing classifier compared to others in the imbalanced datasets.

In addition, Table 2 summarizes the counts of Win-Tie-Loss (W-T-L) of $sPEk$ NN and $mPEk$ NN against other classifiers which indicates that $mPEk$ NN performs better than other classifiers in most cases. From Win-Tie-Loss, it is observed that $mPEk$ NN wins in at most 29 datasets with no loss against C4.5 and GMDKNN classifiers. In the least case, $mPEk$ NN performs better in 19 datasets and worse in 7 datasets compared to EKNN.

The results of Friedman test (Fr. Test) with two base classifiers ($sPEk$ NN and $mPEk$ NN) are shown in the last two lines of the Table 2. From Friedman test with 14 classifiers and 30 datasets, we can conclude that, all the fourteen classifiers are not equivalent. After rejecting that all fourteen classifiers perform equivalent, Nemenyi test is performed to determine which classifier performs significantly better than the others. A tick(✓) sign under a classifier indicates that Nemenyi test suggests the performance of that classifier is significantly different from the base classifier in pairwise comparison at 95% confidence level. Nemenyi test states that, $sPEk$ NN performs significantly better than all compared classifiers except EKNN, CCPDT and HDDT. More importantly, the test suggests that $mPEk$ NN is the best performing classifier among twelve classifiers.

4.4 Effects of Neighborhood Size and Imbalance Ratio

Here, we show the effects of neighborhood size and Imbalance Ratio (IR) on the performance of the proposed method compared to other k NN variants. Due to page limitations, only one dataset (Ionosphere) is used to present the comparison in terms of AUC with different the values of k ranging from 1 to 20. It is clear from Fig. 2a that sPE k NN and mPE k NN consistently perform better than the other algorithms and are less sensitive to the value of k .



(a) Effects of neighborhood size, k (b) Effects of Imbalance Ratio, IR

Fig. 2. Performance comparison among the algorithms belonging in k NN family

To visualize the effect of IR, we use a synthetic dataset of two-class problem in a two-dimensional space where instances of each class are taken from two Gaussian distributions. The characteristics of the dataset is given below where class-A is the minority class and Class-B is the majority class.

$$\eta_1^A = 0.6, \eta_2^A = 0.4, \mu_1^A = [3 \ 3]^T, \mu_2^A = [-2 \ -2]^T, \Sigma_1^A = 3I \text{ and } \Sigma_2^A = I$$

$$\eta_1^B = 0.9, \eta_2^B = 0.1, \mu_1^B = [0 \ 0]^T, \mu_2^B = [4 \ 3]^T, \Sigma_1^B = 8I \text{ and } \Sigma_2^B = I$$

Here η represents the mixture proportion and I is the identity matrix. Different datasets of 1500 samples are generated varying the class imbalance ratio ranging from 2 to 10. It is observable from Fig. 2b that, although the imbalance ratio increases, the performance of mPE k NN remains more steady compared to other k NN variants indicating less sensitivity of mPE k NN in these synthetic datasets.

5 Conclusion

This paper proposes an extended k NN algorithm to increase the performance of existing k NN by making it vigorous to imbalance class problem. In PE k NN, for a query instance, we calculate a confidence for each neighbor instance from the posterior probability of that instance which is then discounted by the proximity of that instance from the query instance. We show that this proximity weighted confidence increases the likelihood of classifying a minority class instance. To calculate the confidence we used two methods one using single Gaussian model

(sPE k NN) and other using Gaussian mixture model (mPE k NN). Results over 30 datasets provide the evidence that the proposed approach is better than twelve relevant methods in imbalanced datasets. However, one limitation of the proposed method is that we assume all the feature values as numeric. As future research direction, we have plan to extend the work for categorical features.

Acknowledgments. This research is supported by the fellowship from ICT Division, Ministry of Posts, Telecommunications and Information Technology, Bangladesh. Grant No - 56.00.0000.028.33.093.19-427; Dated 20.11.2019.

References

1. Akash, P.S., Kadir, M.E., Ali, A.A., Shoyaib, M.: Inter-node hellinger distance based decision tree. In: Proceedings of the Twenty-Eighth International Joint Conference on Artificial Intelligence, IJCAI (2019)
2. Alcalá-Fdez, J., et al.: KEEL data-mining software tool: data set repository, integration of algorithms and experimental analysis framework. *J. Multiple-Valued Logic Soft Comput.* **17**(2), 255–287 (2011)
3. Bache, K., Lichman, M.: UCI machine learning repository (2013)
4. Chawla, N.V., Bowyer, K.W., Hall, L.O., Kegelmeyer, W.P.: SMOTE: synthetic minority over-sampling technique. *J. Artif. Intell. Res.* **16**, 321–357 (2002)
5. Cieslak, D.A., Chawla, N.V.: Learning decision trees for unbalanced data. In: Daelemans, W., Goethals, B., Morik, K. (eds.) ECML PKDD 2008. LNCS (LNAI), vol. 5211, pp. 241–256. Springer, Heidelberg (2008). https://doi.org/10.1007/978-3-540-87479-9_34
6. Denoeux, T.: A k-nearest neighbor classification rule based on dempster-shafer theory. *IEEE Trans. Syst. Man. Cybern.* **25**(5), 804–813 (1995)
7. Dubey, H., Pudi, V.: Class Based Weighted K-Nearest Neighbor over Imbalance Dataset. In: Pei, J., Tseng, V.S., Cao, L., Motoda, H., Xu, G. (eds.) PAKDD 2013. LNCS (LNAI), vol. 7819, pp. 305–316. Springer, Heidelberg (2013). https://doi.org/10.1007/978-3-642-37456-2_26
8. Dudani, S.A.: The distance-weighted k-nearest-neighbor rule. *IEEE Trans. Syst. Man Cybern.* **SMC-6**(4), 325–327 (1976)
9. Fix, E., Hodges Jr., J.L.: Discriminatory analysis-nonparametric discrimination: consistency properties. Technical report, California University, Berkeley (1951)
10. Gou, J., Ma, H., Ou, W., Zeng, S., Rao, Y., Yang, H.: A generalized mean distance-based k-nearest neighbor classifier. *Expert Syst. Appl.* **115**, 356–372 (2019)
11. Li, Y., Zhang, X.: Improving k nearest neighbor with exemplar generalization for imbalanced classification. In: Huang, J.Z., Cao, L., Srivastava, J. (eds.) PAKDD 2011. LNCS (LNAI), vol. 6635, pp. 321–332. Springer, Heidelberg (2011). https://doi.org/10.1007/978-3-642-20847-8_27
12. Liu, W., Chawla, S.: Class confidence weighted k NN algorithms for imbalanced data sets. In: Huang, J.Z., Cao, L., Srivastava, J. (eds.) PAKDD 2011. LNCS (LNAI), vol. 6635, pp. 345–356. Springer, Heidelberg (2011). https://doi.org/10.1007/978-3-642-20847-8_29
13. Liu, W., Chawla, S., Cieslak, D.A., Chawla, N.V.: A robust decision tree algorithm for imbalanced data sets. In: Proceedings of the 2010 SIAM International Conference on Data Mining, pp. 766–777. SIAM (2010)
14. McLachlan, G., Peel, D.: Finite Mixture Models. Wiley, Hoboken (2004)

15. Nemenyi, P.: Distribution-free multiple comparisons. Ph.D. thesis, Princeton University (1963)
16. Shafer, G.: A Mathematical Theory of Evidence, vol. 42. Princeton University Press, Princeton (1976)
17. Swets, J.A.: Measuring the accuracy of diagnostic systems. *Science* **240**(4857), 1285–1293 (1988)
18. Trafalis, T.B., Alwazzi, S.A.: Support vector regression with noisy data: a second order cone programming approach. *Int. J. Gen Syst* **36**(2), 237–250 (2007)
19. Wu, X., et al.: Top 10 algorithms in data mining. *Knowl. Inf. Syst.* **14**(1), 1–37 (2008). <https://doi.org/10.1007/s10115-007-0114-2>
20. Yang, Q., Wu, X.: 10 challenging problems in data mining research. *Int. J. Inf. Technol. Decis. Making* **5**(04), 597–604 (2006)