*Available Online through*       *Research Article*

www.ijptonline.com

# SOCIAL NETWORK ANALYSIS FOR DECISIVE THE ULTIMATE CLASSIFICATION FROM THE ENSEMBLE TO BOOST ACCURACY RATES

**B.Sundarraj*[1], Dr.K.P.Kaliyamurthie[2]**
Assistant Professor, Department of CSE, Bharath University, Chennai[1]
Professor, Department of Computer Science and Engineering, Bharath University, Chennai[2]
*Email: sundarrajboobalan@gmail.com*

**Abstract**

Classifier ensembles are used with success to boost accuracy rates of the underlying classification mechanisms. Through the utilization of collective classifications, it becomes doable to attain lower error rates in classification than by employing a single classifier instance. Ensembles area unit most frequently used with collections of call trees or neural networks because of their higher rates of error once used severally. during this paper, we are going to contemplate a novel implementation of a classifier ensemble that utilizes kNN classifiers. every categoryifier is ready-made to police investigation membership in a very specific class employing a best set choice method for variables. This provides the range required to with success implement associate ensemble. associate aggregating mechanism for decisive the ultimate classification from the ensemble is conferred and tested against many documented datasets.

**Keywords:** k Nearest Neighbor, Classifier Ensembles, Forward set choice

## 1. Introduction

**K-Nearest Neighbor Algorithmic Rule**

The k-Nearest Neighbors, or kNN algorithmic rule is documented to the information mining community, and is one in all the highest algorithms within the field [1]. The algorithmic rule achieves classification between m completely different categories. every instance to be classified is associate item that contains a set of r completely different attributes in set A= wherever a j corresponds to the jth attribute. Therefore, associate instance could be a vector p = of attribute values. for a few planned price of k, the closest k neighbors area unit determined through the utilization of a distance metric that is calculated mistreatment the distinction in distances between every of the attributes of the instance in question and its

neighbors. euclidian distance is far and away the foremost fashionable metric for scheming proximity. associate instance's membership among a given category may be computed either as a likelihood or by easy majority of the category with the foremost illustration within the nighest k neighbors. At the best level, this is often a tangle of binary classification, wherever information {is categoryified|is assessed|is classed} as being in a very bound class or not. because of completely different units of activity, there's conjointly a desire for standardisation across attribute variables so as to forestall one variable from dominating the classification mechanism [2]. one in all the issues with kNN is that while not some kind of coefficient theme for variables, every of the variables is treated as being equally necessary toward decisive similarity between instances. Combining completely different scales of activity across attributes once computing the space metric between instances will cause severe distortions within the calculations for decisive nearest neighbors. many completely different variable coefficient schemes and choice ways to beat this area unit mentioned by Wettschereck, Aha, Mohri [3]. Given the suggests that by that neighbors in kNN area unit calculated, unsuitable variables will have an outsized result on final classification. This becomes particularly problematic in cases wherever an outsized range of predictor variables area unit gift [4]. Closely associated with this downside is that the curse of spatiality whereby the typical distance between points becomes larger because the range of predictor variables will increase. one in all the advantages of correct variable choice is that it's the potential to assist mitigate the curse of spatiality.

It is usually control that kNN implementations area unit sensitive to the choice of variables, therefore alternative of the acceptable set of variables to be used in classification plays a crucial role [5]. one in all the ways is thru the utilization of forward set choice (FSS) with the kNN algorithmic rule [6]. FSS begins by distinctive the variable that results in the very best quantity of accuracy with regards to classifying associate instance. That attribute is then chosen for inclusion within the set of best variables. The remaining variables area unit then paired up with the set, and therefore the next variable for inclusion is once more calculated by decisive that one results in the best increase in classifier accuracy. This method of variable inclusion continues till no additional gains may be created in accuracy. Clearly, this is often a greedy methodology of decisive attributes for inclusion since the variable chosen at every step is that the one providing the most important gains in accuracy. Therefore, the set chosen at the conclusion of the algorithmic rule won't essentially be the most optimum since not all potential mixtures of variables were thought of. to boot, this algorithmic rule is kind of processor intensive.

Backward set choice (BSS) operates in a very similar manner, except {that all|that every one|that every one} variables area unit ab initio enclosed then a variable is discarded throughout each taste the attributes till no additional enhancements in accuracy area unit achieved. Work by Aha and Bankert [6] found that FSS of variables light-emitting diode to higher classification rates than BSS. They conjointly conjectured that BSS doesn't perform moreover with giant numbers of variables. kNN depends on forming a classification supported clusters of knowledge points. There area unit a spread of the way to think about kNN clusters for final classification. easy philosophical system is that the most typical, however there area unit alternative ways in which of coefficient the information [1]. Wettshereck, Aha, and Mohri [5] give a comprehensive summary of assorted choice and coefficient schemes employed in lazy learning algorithms, like kNN, wherever computation is deferred till classification. These modifications to the coefficient calculations of the algorithmic rule embody not solely international settings, however native changes to the weights of individual variables. The weights area unit adjustable betting on the composition of the underlying information. this enables for larger accuracy and adaptableness in bound parts of the information while not imposing international variable weightings.

## 2. Our Approach

Our approach begins with the assembly of associate ensemble of kNN classifiers. we have a tendency to selected to use kNN classifiers due to their ability to adapt to extremely nonlinear information, they're a reasonably mature technique, and there area unit variety of ways on the market for optimizing instances of kNN classifiers. every instance or object to be classified p could be a vector of values for r completely different attributes. This methodology works best for algorithms like kNN that produces activation as associate output to see category membership [22]. primarily every binary kNN classifier is that the analogue of a classification "stump", that could be a call tree that produces one categoryification of whether or not or not a given instance could be a member of a selected class. Classifiers that discriminate between all categories, like one model to see membership, have a slip-up rate determined by the quantity of misclassifications from the complete dataset. this is often as a result of the classifier is ready-made for and optimized over the gathering of m completely different categories. As a result, the parameters area unit adjusted in order that the error rate across all categoryifications is as low as doable while not deference to any explicit class. The set of variables that results in all-time low error rates once decisive membership in a very specific category area unit doubtless to be entirely completely different from the set of variables that area unit simplest in decisive membership in another category. the utilization of the FSS algorithmic rule permits every

individual binary classifier to tailor itself round the variables it deems most vital for decisive membership of associate instance. As a result, diversity amongst the kNN classifiers is achieved deterministically. the mandatory diversity is achieved by every individual categoryifier choosing the set of variables that area unit deemed most vital for distinctive specific class membership. this is often slightly {different|totally completely different|completely different} from the standard definition of diversity that stresses errors being created on different instances of knowledge. Since we have a tendency to use associate ensemble of individual kNN categoryifiers that area unit answerable for decisive membership in a very specific class, every individual categoryifier will have the parameters for variable weights adjusted to attain the very best classification rate for the particular class being analyzed. once employing a single classifier to differentiate between multiple categories, the variations within which variables area unit most necessary to agglomeration for identification of assorted categories becomes overshadowed.

In order to mix the individual votes of every member among the ensemble, we've 3 cases: one in all the individual classifiers identifies membership among the cluster, no membership is chosen, or there's a conflict concerning classification with 2 or additional classifiers presenting conflicting classifications. wherever classification is easy with one classification rising from the ensemble, we have a tendency to use that classification. within the latter 2 cases mentioned higher than, there should be how of achieving associate output. There area unit 2 doable approaches. the primary is to consider one overall kNN classifier that determines identification within the event of conflict. Therefore, if the ensemble is unsuccessful, the classification theme reverts back to one instance (master) classifier. The second approach is to use the classifier with the very best accuracy that chosen the instance for membership. A master categoryifier uses identical methodology however provides for classification between all doable categories within the dataset as opposition merely decisive membership in a very single class. This master classifier is employed to assign classification within the event that none of the members of the ensemble identifies associate item for sophistication membership.

## 3. Experimental Results

### 3.1. Datasets

The datasets that we have a tendency to utilised were from the UCI Machine Learning Repository with the exception of the IRIS dataset that is accessible within the R code package [23, 24, 25]. The statistics concerning every information set area unit conferred in Table one. we have a tendency to began with the IRIS information since it's one in all the foremost used

datasets in classification issues. what is more, it's an easy dataset with four predictors and provided an honest benchmark for initial results. we have a tendency to conjointly chosen the Low Resolution mass spectrometer (LRS) information since it contained an outsized range of variables and therefore the information needed no scaling before mistreatment the algorithmic rule. The dataset itself consists of header data for every entry, followed by intensity measurements at numerous spectral lengths. Finally, the ARRYTHMIA dataset was chosen because of the massive range of predictor variables that it offered. we have a tendency to were curious to check however well the FSS-kNN algorithmic rule performed at reducing the quantity of variables required to see category membership. there have been many instances within the ARRYTHMIA information set wherever missing information was problematic. These attributes were faraway from the dataset in order that classification may continue.

- ➢ Dataset:IrisLRSArrythmia
- ➢ Number of classes: three ten sixteen
- ➢ range of variables: four ninety three 263
- ➢ Number of knowledge points: one hundred fifty 532 442

## 3.2. Model Generation

Our methodology follows the we have a tendency to began by building the simplest categoryification model for every class within the dataset. The individual models were created mistreatment* FSS-kNN to see the simplest set of variables to use for decisive membership in every category. each set of variables was then tested mistreatment n-fold cross validation, wherever every part was foreseen mistreatment the remaining components within the kNN model over numerous k-values to see the foremost correct models for every category. This needed a modest quantity of processor time, however enabled America to use all of the on the market information for each coaching and testing that is one in all the advantages of n-fold cross validation. Following the generation of the individual classifiers, we have a tendency to engineered the master classifier.

After building our classifiers, we have a tendency to processed the information with the ensemble. the bulk of instances were chosen for membership by one in all the classifiers. within the event that quite one classifier categorised the instance as being a member of the category that it pictured, we have a tendency to reverted to the model accuracies of the individual classifiers, and appointed the item to the foremost correct classifier that known the item for sophistication membership.

Instances that weren't chosen for membership in a very category by any of the individual classifiers were processed by the master classifier.

We have a tendency to conducted n-fold cross-validation testing to see the accuracy of the ensemble.

The k-value and set of variables chosen for a private kNN classification model were the sole factors remaining identical between the classifications of instances.

**Construction Phase**:

- For every category within the information set algorithmic rule

- Build classifier c

- I that determines membership in school i mistreatment the forward set choice

- Compute the accuracy of this classifier

- Next category build a master classifier that considers membership amongst all categories

**Classification Phase:**

- For each item to be classified

- The item is evaluated by every classifier therefore request membership in individual categories

- If only 1 classifier known the item for membership

- Then assign the item to it category

- If quite one classifier known the item for membership

- Then assign category membership to the foremost correct classifier

- If no classifiers known the item for membership

- Then use the master classifier to assign a classification

- Exit item

## 4. Conclusions and Future Work

Our approach has incontestable that associate ensemble of categoryifiers trained to sight membership in a very given class are able to do high rates of classification. we've shown that we will reach larger categoryification rates by combining a series of classifiers optimized to sight class membership, than by mistreatment single instances of classifiers. Our model is best custom-made towards categoryification issues involving 3 or additional categories since a 2 class model may be promptly

handled by one classifier instance. We have not adjusted the importance of individual variables throughout the method of constructing individual classifiers for the ensemble. we've merely enclosed or excluded variables as being equally weighted while not scaling. whereas variable choice is useful in addressing a number of the issues made public, extra enhancements may be created to the kNN algorithmic rule by coefficient the variables that are chosen for inclusion into the model to account for variations in variable importance. Another weakness that has to be addressed is that the thought of incomplete datasets. Future work can specialize in developing extra classifiers to tell apart between instances that area unit chosen for sophistication membership by quite one classifier among the ensemble instead of reverting to the very best accuracy rate. a component of chance can be of hefty importance in medical specialty classifications. In larger datasets, there can be variety of cases wherever discerning membership amongst instances becomes tough. typically the determination happens between 2 categories that area unit terribly similar. In such cases wherever FSS- KNN leads to classifiers with comparatively low rates of classification, it would be necessary to look at the information to see whether or not the category in question is absolutely composed of many subclasses that would profit from their own individual binary classifiers among the ensemble. Finally, there remains the chance that we will use the predictor variables chosen as most vital for agglomeration by FSS to boost classification rates of alternative ways like neural nets and call trees.

## References

1. Xindong, W., Kumar, V., Quinlan, J. R., Ghosh, J., Yang, Q., Motoda, H., McLachlan, G. J., Ng, A., Liu, B., Yu, P. S., Zhou, Z. H., Steinbach, M., Hand, D. J. and Saul Steinberg, D. (2008) prime ten Algorithms in data processing. data and data Systems, 14 1: 1-37.

2. Xu, R., Wunsch, D., "Clustering", 2009, John Wiley &amp; Sons.

3. Wettschereck, D., Aha, D., &amp; Mohri, T. (1995). A Review and Empirical analysis of Feature coefficient ways for a category of Lazy Learning Algorithms. Tech. rept. AIC95-012. military service work, Navy Center for Applied analysis in AI, Washington, D.C.

4. Hand, D. , Mannila, H., Smyth, P., Principles of knowledge Mining, 2001.

5. Kotsiantis, S. "Supervised machine learning: a review of classification techniques", Informatica Journal (31), 2007, pp.249-268.

6.  Aha, D. W., &amp; Bankert, R. L. (1996). A Comparative analysis of consecutive Feature choice Algorithms. In D. Fisher &amp; H. H. Lenz (Eds.), AI and Statistics V. New York: Springer – Verlag.

7.  R. Caruana, A. Niculescu-Mizil, G. Crew, and A. Ksikes. "Ensemble choice from libraries of models." in Proceedings of the International Conference on Machine Learning (ICML), 2004.

8.  G. Brown, J. Wyatt, R. Harris, and X. Yao. Diversity creation methods: A survey and categorisation. Journal of data Fusion, 6(1):5–20, 2005.

9.  D. Opitz, R. Maclin, fashionable ensemble methods: associate empirical study, Journal of Articial Intelligence analysis eleven (1999), pp 169 - 198.

10. Tan, Pang-Ning, Steinbach, Michael, Kumar, Vipin, Introduction to data processing, Addison Wesley, 2006, pp 283-285.

11. N. El Gayar, associate Experimental Study of a Self-Supervised Classifier Ensemble, International Journal of data Technology, Vol. 1, No. 1, 2004.

12. A. Madabhushi, J. Shi, M.D. Feldman, M. Rosen, and J. Tomaszewski, "Comparing Ensembles of Learners: police investigation adenocarcinoma from High Resolution MRI," Proc. Second Int'l Workshop pc Vision Approaches to Medical Image Analysis (CVAMIA '06), pp. 25-36, 2006.

13. S. D. Bay. Combining Nearest Neighbor Classifiers Through Multiple Feature Subsets. Proc. 17th Intl. Conference on Machine Learning, pp. 37 – 45, Madison, WI, 1998.

14. C. Domeniconi and B. Yan. Nearest Neighbor Ensemble. In Proceedings of the seventeenth International Conference on Pattern Recognition, Cambridge, UK, pages 23–26, 2004.

15. L. S. Oliveira, M. Morita, R. Sabourin, and F. Bortolozzi. Multi-Objective Genetic Algoritms to form Ansemble of Classifiers. Proceedings of the Third International Conference on Evolutionary Multi-Criterion optimisation, Vol. 87, pp 592-606. 2005.

16. Fangfei Weng, Quingshan Jiang, Liang Shi, and Nannan Chinese., associate Intrusion Detection System supported the agglomeration Ensemble, IEEE International Workshop on 16-18 Apr 2007, pp. 121-124. 2007.

**Corresponding Author:**

**B.Sundarraj\*,**

*Email: sundarrajboobalan@gmail.com*