



Performance comparison between K-Nearest Neighbor and Naive Bayesian Classifiers by using heart disease dataset

Yin Yin Htay¹, Ya Min²

¹Computer University, Magway, Myanmar

²Computer University, Lashio, Myanmar

ABSTRACT

Today, the diagnosis of diseases is a vital and intricate job in medicine. Medical diagnosis is regarded as an important yet complicated task that needs to be executed accurately and efficiently. Regrettably all doctors do not possess expertise in every sub specialty and moreover there is a shortage of resource persons at certain places. In this situation, an automatic medical diagnosis system is beneficial by bringing all of them together. For this diagnosis system, many classifiers are essential and needed for disease classification. So, this system is proposed as the performance comparison system about classifiers to know which classifier is more effective than other. To compare the performance, this system classifies the heart disease dataset by using K-Nearest Neighbor (KNN) and naive Bayesian (NB) classifiers.

Keywords— Comparison, KNN, NB

1. INTRODUCTION

Healthcare industry generates the large amount of data about patient, disease diagnosis etc. However, there is a lack of effective analysis tools to discover hidden relationships in data. Data mining provides a set of techniques to discover hidden patterns from data. Data mining is a process that uses a variety of data analysis tools to discover patterns and relationships in data that may be used to make valid predictions. A knowledge discovery process includes data cleaning, data integration, data selection, data mining, pattern evaluation and knowledge presentation. Data mining system can be classified according to the kinds of databases mined, the kinds of knowledge mined, the techniques used or the applications adapted. Data mining is the process of classification, association rule mining, clustering, etc. K-Nearest Neighbor (KNN) and Naive Bayesian (NB) classifiers are the most popular algorithms in the mining classification.

Major challenge facing Healthcare industry is quality of service. Quality of service implies diagnosing disease correctly and provides effective treatments to patients. Poor diagnosis can lead to disastrous consequences which are unacceptable.

So, this system is proposed to predict whether the patient is having heart disease or not by using KNN and NB classifier that are data mining techniques. After classifying according to these three classifiers, this system compares the accuracy and processing time of each classifier. According to the

comparison results, the healthcare industry can easily know which classifier is most effective for diagnosis system. In the remote areas like rural regions or country sides, the proposed system is also a user friendly, scalable and reliable system that can be implemented to imitate like human diagnosis expertise for treatment of heart disease.

2. RELATED WORK

In 2016, M. Panwar, A. Acharyya and R. A. Shafik [1] presented a new methodology based on novel preprocessing techniques, and K-nearest neighbor classifier. The effectiveness of the proposed methodology is validated with the help of various quantitative metrics and a comparative analysis, with previously reported studies using the same UCI dataset focusing on pima-diabetes disease diagnosis.

In 2016, D. VijayaKumar and V. J. R. Krishniah [2] used decision tree classification model for diagnosis of three brain diseases namely ischemic stroke, hemorrhage and hematoma, and tumor. This system helped the physicians to identify the type of human brain hemorrhage and hematoma and the type of brain tumors for further treatment.

In 2017, N. R. Gorrepati and N. R. Uppala [3] compared the performances of different classifiers on diagnosis of the Erythematous-Squamous disease. The classifiers examined here are support vector machine, discriminant classifier, K-nearest neighbor and decision tree. They have performed their analysis with two well-known multiclass implementation techniques. They demonstrated that the most reliable performance has been achieved using support vector machine classifier.

3. CLASSIFICATION

For decision making procedure, data mining is a very favorable and constructive method. Classification is a very simple and mostly used data mining technique. Knowledge of training data is mandatory for understanding of classification. There are two phases of classification procedure:

- Development of a model for training
- Evaluating the model using testing data.

For classification, there are various classifier. Bayesian classifier uses frequentist technique. The essence of frequentist technique is to apply probability to data. Bayesian calculations go straight for the probability of the hypothesis. K-nearest neighbor is a non-parametric method which depends on the use

of distance measurement. All available cases can be stored in it and whenever a new case entered, it can be classified based on the distance function [4].

4. KNN CLASSIFIER

When an unknown feature is introduced, the k-nearest neighbour (KNN) classifier finds k most similar training features that are closet to the unknown feature. The procedure of K-NN classifier is as follow [5]:

- Determine k.
- Calculate the similarity or distance between the testing data and all the training data.

$$(1) \quad d_{euc} = \sqrt{\sum_{j=1}^N |P_j - Q_j|^2}$$

where d_{euc} is the distance between the test and training data, P_j is the feature j of test data P , and Q_j is the feature j of training data Q .

- Sort the distance and determine k nearest neighbors based on the K^{th} minimum distance.
- Gather the categories based on majority vote.
- Determine categories based on majority vote.

5. NAIVE BAYESIAN CLASSIFIER

Naive Bayesian (NB) classifier shows the relationship between the independent variables and the target variable [6]. The processing steps of NB classifier is as follows:

1. Each data sample is represented by n-dimensional feature vector, $X = (x_1, x_2, \dots, x_n)$ depicting n-measurements made on the sample from n- attributes, A_1, A_2, \dots, A_n .
2. Suppose that there are m classes, C_1, C_2, \dots, C_m . Given an unknown data sample X, Naïve Bayesian classifier assigns an unknown sample X to the class C_i if and only if

$$P(C_i | X) > P(C_j | X) \text{ for } 1 \leq j \leq m, j \neq i \quad (2)$$

The class C_i for which $P(C_i | X)$ that is maximized, called maximum posteriori hypothesis. By Bayes theorem,

$$P(C_i | X) = P(X | C_i) P(C_i) / P(X) \quad (3)$$

3. As $P(X)$ is constant for all classes, only $P(X | C_i) P(C_i)$ need to be maximized.
4. Given data sets with many attributes, the Naive assumption of class conditional independence is made. Thus,

$$P(X | C_i) = \prod_{k=1}^n P(x_k | C_i) \quad (4)$$

The probability $P(x_1 | C_i), P(x_2 | C_i), \dots, P(x_n | C_i)$ can be estimated from the data samples.

5. In order to classify an unknown sample X, $P(X | C_i) P(C_i)$ is evaluated for each class C_i . Sample X is then assigned to the class C_i if and only if

$$P(X | C_i) P(C_i) > P(X | C_j) P(C_j) \quad (5)$$

In other words, it is assigned to the class C_i for which $P(X | C_i) P(C_i)$ is the maximum [5].

6. PROPOSED SYSTEM DESIGN

In this system, the user must first put the patient information (the patient suffered disease symptom). For classification, this system extracts the training heart disease data into the system. Then, this system classifies the heart disease stage by using k-nearest neighbor (KNN) and naïve bayesian (NB) classifiers.

Proposed system design is shown in Figure 1. After classifying, this system produces and displays the heart disease stage. Then, this system measures the processing time of each classier and calculates the accuracy of each classifier by using Holdout method. Finally, this system displays the performance comparison result to the user.

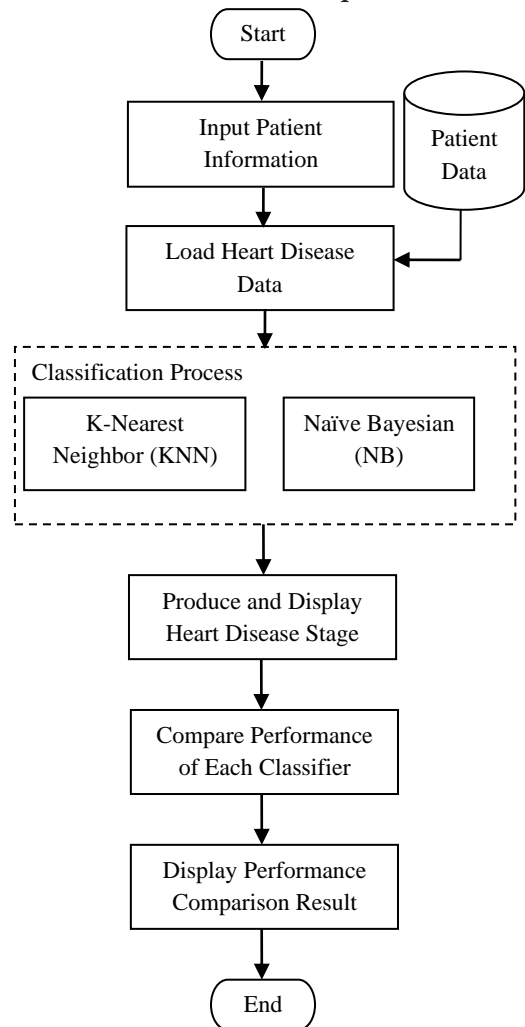


Fig. 1: Proposed System Design

7. EXPLANATION OF THE SYSTEM

For classification, this system uses the heart disease dataset. In this dataset, there are 13 (symptoms) attributes. As a sample, this system uses 10 records that are obtained from 10 patients who suffer heart disease. Heart disease dataset includes five class levels that are Normal (N), Level I (I), Level II (II), Level III (III) and Level IV (IV). Sample heart disease dataset is shown in Table I.

Table 1: Sample Heart Disease Dataset

Age	Sex	Chest Pain Type	Trestbps	Chol	Fasting Blood Sugar	Rest ECG	Thalach	Exang	Oldpeak	Slope	Ca	Thal	Class
63	1	4	0	260	140	1	112	1	3	2	0	0	II
44	1	4	0	209	130	1	127	0	0	0	0	0	N
60	1	4	0	218	132	1	140	1	1.5	3	0	0	II
55	1	4	0	228	142	1	149	1	2.5	1	0	0	I
66	1	3	1	213	110	2	99	0	1.3	2	0	0	N
66	1	3	0	0	120	1	120	0	0.5	1	0	0	N
65	1	4	1	236	150	1	105	1	0	0	0	0	III
62	1	3	0	0	180	1	140	1	1.3	2	0	0	N
60	1	3	0	0	120	0	141	1	2	1	0	0	III
60	1	2	1	267	160	1	157	0	0.5	2	0	0	I

The patient input heart disease symptom into the system. The inputted information includes age (60), sex (1), chest pain type

(4), trestbps (0), chol (260), fasting blood sugar (140), restECG (1), thalach (140), exang (1), oldpeak (1.5), slope (3), ca (0) and thal (0). Then, this system calculates and classifies the heart disease stage that the patient suffers.

7.1. KNN Classification Process

By using training and testing data, this system calculates the distance between each data. After calculating, this system chooses the most similar training class for the testing data. KNN classifier results are shown in Table 2.

Table 2: KNN Classifier Results

ID	Test and Training Data	Distance Result
1	Test and Training Data 1	24.5
2	Test and Training Data 2	95.5
3	Test and Training Data 3	50
4	Test and Training Data 4	27
5	Test and Training Data 5	113.2
6	Test and Training Data 6	299
7	Test and Training Data 7	47.5
8	Test and Training Data 8	220.2
9	Test and Training Data 9	282.5
10	Test and Training Data 1	40

According to the KNN classifier, this system produces the heart disease stage “**Level II**” to the patient.

7.2 NB Classification Process

According to Naïve Bayesian classifier, this system calculates probability of each attribute. NB classifier result is shown in Table 3.

Table 3: NB Classifier Results

ID	$P(X C_i) P(C_i)$	Probability Result
1	$P(X Class = I) P(Class = I)$	0
2	$P(X Class = II) P(Class = II)$	0.00312
3	$P(X Class = III) P(Class = III)$	0
4	$P(X Class = N) P(Class = N)$	0

The heart disease level that the patient suffered is “**Level II**”.

8. EXPERIMENTAL RESULT

To compare the performance of the system, this system uses the heart disease dataset from the UCI website. This system is tested 500 records. By using hold out method, this system calculates the accuracy of each classifier. The experimental result of the system is shown in Table 4.

Table 4: Experimental Result of the System

Test Data	Accuracy	
	NB	KNN
150	90%	91%
220	91%	92%
270	82%	85%
300	84%	85%
330	80%	81%

After calculating the performance, KNN classifier is more precise than NB classifier. Sometimes, NB classifier can face problem about probability calculation of each attribute. Finally, KNN is precise and takes more processing time than NB and KNN.

9. CONCLUSION

This system is proposed an effective heart disease prediction system by using k-nearest neighbor and naïve bayesian classifiers. This system is also implemented to show which classifier is more than another. This system points out the KNN classifier that can quickly produce the heart disease result. Finally, this system is helpful for practice to confirm his/ her decision during heart disease prediction.

10. REFERENCES

[1] M. Panwar, A. Acharyya and R. A. Shafik, "K-Nearest Neighbor Based Methodology for Accurate Diagnosis of Diabetes Mellitus", IEEE, 2016.

[2] D. VijayaKumar and V. J. R. Krishniah, "An Automated Framework for Stroke and Hemorrhage Detection using Decision Tree Classifier", IEEE, 2016.

[3] N. R. Gorrepati and N. R.Uppala, "Comparative Performance Analysis of Different Classifiers on Diagnosis of Erythmato-Squamous Diseases", International Conference on Innovations in Information, Embedded and Communication Systems, IEEE, 2017.

[4] C. Jalota and R. Agrawal, "Analysis of Educational Data Mining using Classification", International Conference on Machine Learning, Big Data, Cloud and Parallel Computing, pp. 243-247, 2019.

[5] H. Jiawei and K. Micheline, "Data Mining Concepts and Techniques", Simon Fraser University, United States of America, 2001.

[6] M.S. Basarslan and I. D. Argun, "Classification of A Bank Data Set on Various Data Mining Platforms", IEEE, 2018.