

Using machine learning in chat box for heart disease prediction

Lingyun Yang
University of Nottingham
Psxly4@nottingham.ac.uk

ABSTRACT

Heart disease is one of the deadly diseases with high incidence. But the prediction or diagnosis of heart disease requires complex procedures in hospitals. Therefore, this article attempts to collect some data using simple conversations and use machine learning models to predict the likelihood of a user's heart disease, so as to try to help possible high-risk groups to prevent it as soon as possible. The core idea of this project is using a chat box to have conversation with the user and collect some data from the conversation, then use built-in machine learning models to make prediction. This paper contains two main contents, the design of the chat box and the design of the machine learning model.

1. INTRODUCTION

The CDC (Centers for Disease Control and Prevention, the national public health agency of the United States) reports heart disease as one of the leading causes of death for Americans of most races (African Americans, American Indians and Alaskan Natives, and white people) (PYTLAK, 2022). However, the data from NHS England shows for patients waiting to start treatment at the end of November 2021, the median waiting time was 11.5 weeks (Baker, 2022). Such long wait times can miss the best prevention and treatment period for early-stage patients. Therefore, this research tries to use a simple method to allow users to know whether they are a high-risk group of heart disease and help users to prevent it as soon as possible.

2. RELATED WORK

2.1. K_NN Algorithm

Why K-NN? The k-nearest neighbors algorithm (k-NN) is a non-parametric supervised learning method first developed by Evelyn Fix and Joseph Hodges in 1951 (Fix and Hodges, 1989). The KNN algorithm is a form of lazy learning in which the computation for the creation of predictions is postponed until classification. Although this approach increases computing costs compared to other methods, it is still the best option for applications where predictions are required seldom but accuracy is crucial. (Sharma and Sachdeva, 2016). In this project, the prediction is related to the user's health, so the accuracy rate is important, that's why K-NN is chosen in this research.

2.2. SVC

Why SVM(SVC)? Support Vector Machine (SVM) was originally mentioned in 1992, when Boser, Guyon, and Vapnik described it in COLT-92.(Cortes and Vapnik, 1995). In supervised learning, support vector machines (SVMs) are a collection of algorithms for categorising and predicting data. Support Vector Machines (SVMs) are classification and regression prediction tools that utilise machine learning techniques to increase predictive accuracy while reducing over-fitting(Bottou and Lin, 2007). This is exactly the goal that this research wants to achieve.

2.3. Chat Box

The chat box that used in this research is adapt from lab of module COMP4105 in University of Nottingham. A simple chatbot agent that has a store of knowledge, uses that knowledge to decide what questions to ask, and adds to that store of knowledge by interpreting the user's responses(Johnson, 2022).

3. EXPERIMENT

3.1. Tokenization using NLTK

According to the data from CDC. Nearly half of all individuals (47%) in the United States have at least one of three key risk factors for heart disease: high blood pressure, high cholesterol, and smoking. Other significant markers include diabetes, obesity (high BMI), insufficient physical exercise, and excessive alcohol consumption (PYTLAK, 2022). There are 14 factors selected as the basis for the design question. These factors are also 14 features in the machine learning models. Table -1 shows some typical factor and there corresponded questions and answers.

Factors	typical questions	Answer rescticeted
BMI	could you please tell me your BMI?	float <= 50
MentalHealth	how many days during the past 30 days was your mental health not good? (0-30 days)	float <=30
DiffWalking	Do you have serious difficulty walking or climbing stairs?	Binary [yes,no]
Sex	are you male or female?	Binary [male,female]
AgeCategory	could you please tell me your age category, like 20-29 , 40-49	categorical
Diabetic	have you ever had diabetic?	Binary [yes,no]
SleepTime	how many hour do you sleep in a 24-hour period?	float <=24

Table – 1 Question design and answer scope

After the user enters an answer use natural language toolkit NLTK to tokenize the answer and tagging each word. After this match the target word from the tagged word. If there is an target word in the answer, then convert it into the correct format to be used as input to the model to make predictions. If user's answer out of the restricted answer, ask their user input answer again. Please refer to Table -2, the Pseudo code of chate box:

```
# initial question list and create unknown question list
Create QAdic = questions and answer dictionary
unknown list = list QAdic
```

```

# if there is unknown questions in list, start ask questions
if unknown list is not empty:
    ask question loop:
        question = random pick from unknown list
        #get the input from user, tokenize and tag answer
        answer = get input from user
        tokenize answer and tag them
        # Match the target key word from tagged answer
        targetAns = matched [ tagged word list]
        if the tagged word list includes "bye":
            break the loop and say bye to user
        else:
            if(matched):
                remove question from unknown list
                transfer matched key words to right format
                continue loop
            else(user input out of answer scope):
                ask user input the answer again
                back to the step answer = get input from user
                ...

    If(unknown list is not empty):
        continue the loop
    else:
        end the loop and say bye to user

```

Table -2 pseudo code of Chat Box

3.2. Machine learning Models

In this research, I build two models based on two different machine learning methods: the K_NN algorithm and the SVM algorithm. I also compare these two models based on the same dataset. This research uses a machine learning tool called scikit-learn to build and optimize the model. Scikit-learn (formerly scikits. learn and known as sklearn) is a Python library for machine learning. Support vector machines, random forests, gradient boosting, k-means and DBSCAN, are included, as well as various classification, regression and clustering algorithms. The NumPy and SciPy Python libraries are integrated with it.

3.2.1 Machine learning Models

The dataset used for this research is part of Personal Key Indicators of Heart Disease from Kaggle.com [Kaggle]. The original datasets from Kaggle.com contains 319794 instances. In

order to reduce the computational load, I just randomly picked 1000 instances as a new datasets for this research. All data was split into training set and test set, the proportion of test size is 16%.

The Original dataset comes from the CDC and is an important component of the Behavioural Risk Factor Surveillance System (BRFSS), which conducts annual telephone surveys to collect data on the health of U.S. citizens. Established in 1984, the BRFSS now collects information from all 50 states, the District of Columbia, and three U.S. territories. The BRFSS performs approximately 400,000 adult interviews yearly, making it the largest continuously conducted health survey system. The most current dataset comprises information from 2020. (PYTLAK, 2022).

3.2.2 Model Experiment

1) K_NN model

KNN manipulates the training data and classifies the fresh test data using distance measures. It determines the k-nearest neighbours of the test data, and then classifies the data based on the majority of class labels(Band, 2020). K value indicates the count of the nearest neighbours, and there are no pre-defined statistical methods to find the most favourable value of K. So first, randomly pick a K value and start the computing. After this I derive a plot between accuracy and K values in a defined range from 1 to 14(because the number of features is 14). Then choose the K value as having a maximum accuracy.

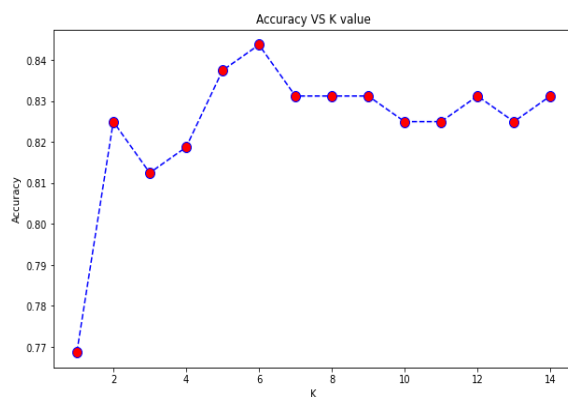


Figure – 1 Visualize plot between accuracy and K value

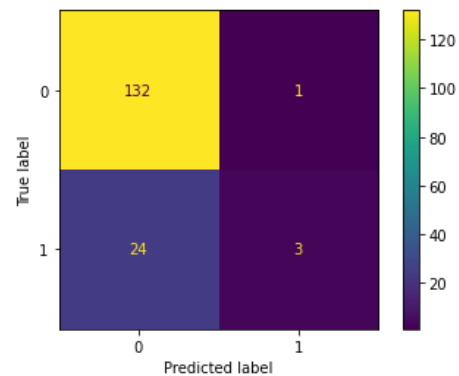


Figure – 2 Visualize plot confusion matrix , K = 6

The result is shown in Figure – 1. From the plot, the maximum accuracy got in this experiment is 0.84 at K = 6. Therefore, set the K = 6 and calculate the confusion metrics. The result of confusion metrics is shown below. And visualization of confusion metrics is shown in Figure -2.

```
confusion metrics:
[[132  1]
 [ 24  3]]
```

2) SVC Model

Another model built in this research is based on machine learning method SVC(support vector classifier).

Support Vectors Classifier attempts to identify the optimal hyperplane for classifying data by optimising the distance between sample points and the hyperplane (Bottou and Lin, 2007). There are 3 main parameters in the Support Vector Classifier: Kernel, gamma and C. The kernel settings determine the type of hyperplane utilised to split the data. Using 'linear' will utilise a linear hyperplane (a line in the case of 2D data). 'rbf' and 'poly' exploit a nonlinear hyper-plane (Cortes and Vapnik, 1995). Gamma is a parameter for nonlinear hyperplanes. The greater the gamma value, the closer it attempts to match the training data set. C is the mistake term's penalty parameter. It regulates the trade-off between a smooth decision boundary and the proper classification of training points. Set the parameter as default setting, and computing an result. After that, use GridSearchCV method to optimize those parameters. Parameter range of C is [0.1,1,10,1000]. After calculation, the best parameters are { 'C': 0.1 } with a score of 0.7875. So set the C equals 0.1 get the confusion metrics shown below, and the visualization of confusion matrix shown in Figure -3.

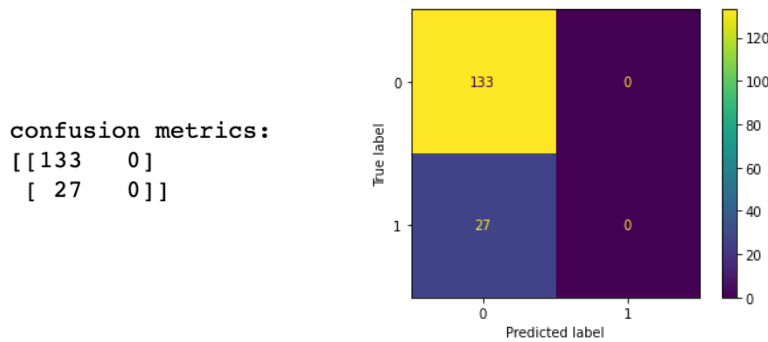


Figure - 3 Visualize plot confusion matrix

4. MODEL EVALUATION

In this research, two models were compared using Accuracy, Recall, Precision and AUC. The summary of Accuracy, Recall Precision and AUC is shown in Table - , and the visualized plot is shown in Figure -

Accuracy is the simplest basic performance metric and is simply the ratio of properly predicted observations to the total number of observations (Solutions, 2016). The calculation formula of Accuracy is:

$Accuracy = \frac{TP+TN}{TP+FP+FN+TN}$. In this research, the accuracy of SVM model is 0.7875, the accuracy of K_NN model is 0.8438.

Precision is the ratio of correctly predicted positive observations to the total predicted positive observations (Solutions, 2016). The calculation formula is: $Precision = \frac{TP}{TP+FP}$. This formula directly shows what precision talks about. It about how precise/accurate the model is out of those predicted positive, how many of them are actual positive. Precision is a good measure to determine, when the costs of False Positive is high. The precision of SVM model is 0.18, the precision of K_NN model is 0.75.

Recall actually calculates how many of the Actual Positives that model capture through labeling it as Positive (True Positive). It is the ratio of correctly predicted positive observations to the all observations in actual class. The Recall of SVM model is 0.074, the recall of K_NN model is 0.111.

	SVM	KNN
Accuracy	0.7875	0.8438
Precision	0.1818	0.7500
Recall	0.0741	0.1111
AUC	0.4684	0.6331

An ROC curve (receiver operating characteristic curve) is a graph showing the performance of a classification model at all classification thresholds. This curve plots two parameters: True Postive Rate and False Positive Rate. Calculation formula are: $TPR = \frac{TP}{TP+FN}$, $FPR = \frac{FP}{FP+TN}$. ROC curves plot TPR versus FPR at different classification thresholds. False positives and True Positives both increase when the classification threshold is lowered(Solutions, 2016). The following Figure -3 shows a typical ROC curve.

AUC stands for "Area under the ROC Curve." That is, AUC measures the entire two-dimensional area underneath the entire ROC. The AUC calculation results of these two models are 0.498 for SVC model and 0.633 for K_NN model From the point of view of AUC, the model built with K_NN algorithm have a better performance than the model built with SVC.

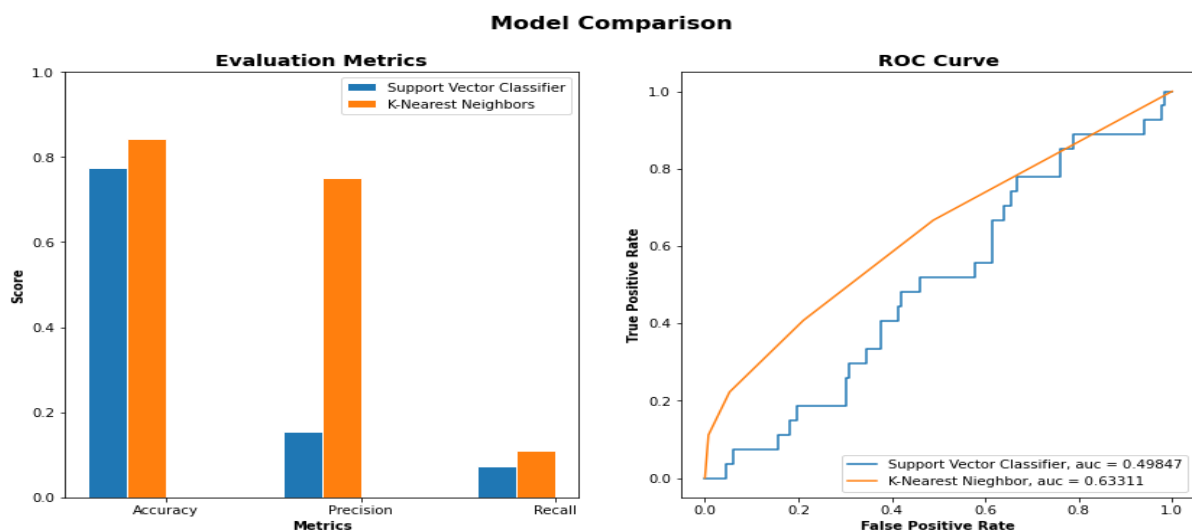


Figure – 3 Model comparison

5. FUTURE WORK AND CONCLUSION

In this research, a chat box and two machine learning model were applied in order to built a simple agent to predict the probability of someone have heart disease, so that users can be reminded to change their lifestyles as soon as possible for prevention, or get treatment. The

model built on the K_NN algorithm performs better in both Accuracy, Precision, Recall and AUC than the model built on the SVM algorithm.

In the future, there are a few things that can be improved and developed if more time is available. First, the chat box now just have 14 questions built in with almost no interactive between users and the agent. So a more interactive chatbox might be better at engaging users to chat, increasing the likelihood of collecting more data. In addition, during the construction of the machine learning model, more training data sets and more features can also increase the accuracy of prediction.

REFERENCES

- BAKER, C. 2022. NHS Key Statistics: England, February 2022.
- BAND, A. 2020. *How to find the optimal value of K in KNN?* [Online]. Available: <https://towardsdatascience.com/how-to-find-the-optimal-value-of-k-in-knn-35d936e554eb> [Accessed].
- BOTTOU, L. & LIN, C.-J. 2007. Support vector machine solvers. *Large scale kernel machines*, 3, 301-320.
- CORTES, C. & VAPNIK, V. 1995. Support-vector networks. *Machine learning*, 20, 273-297.
- FIX, E. & HODGES, J. L. 1989. Discriminatory analysis. Nonparametric discrimination: Consistency properties. *International Statistical Review/Revue Internationale de Statistique*, 57, 238-247.
- JOHNSON, C. 2022. COMP3004/COMP4105 Designing Intelligent Agents Coursework. In: NOTTINGHAM, U. O. (ed.).
- PYTLAK, K. 2022. *Personal Key Indicators of Heart Disease* [Online]. Kaggle.com: Kaggle.com. Available: <https://www.kaggle.com/datasets/kamilpytlak/personal-key-indicators-of-heart-disease> [Accessed 2022].
- SHARMA, J. & SACHDEVA, K. 2016. Offline signature verification using nn, KNN and SURF. *International Journal in IT & Engineering*, 4.
- SOLUTIONS, E. 2016. *Accuracy, Precision, Recall & F1 Score: Interpretation of Performance Measures* [Online]. Available: <https://blog.exsilio.com/all/accuracy-precision-recall-f1-score-interpretation-of-performance-measures/> [Accessed].