

An Optimized Heart Disease Prediction Using Machine Learning

Shrey Jain¹, Shobhit Kumar²

¹Shrey Jain, B.Tech IVth Year, Galgotias University Greater Noida, Uttar Pradesh, India

²Shobhit Kumar, Assistant Professor, Galgotias University Grater Noida, Uttar Pradesh, India

Abstract - Heart is the accompanying critical organ diverging from cerebrum which has more noteworthy need in Human body. It siphons the blood and supplies to all organs of the entire body. Desire for occasions of heart ailments in clinical field is significant work. Gigantic Measure of patient related information is kept up on month to month premise. A touch of the information mining and AI strategies are utilized to imagine the coronary illness, for example, Artificial Neural Network (ANN), Decision tree, Fuzzy Logic, K-Nearest Neighbor (KNN), Naive Bayes and Support Vector Machine (SVM). This paper gives an information on the present count and it gives a general framework of the present work.

Cardiovascular ailment is one of the most deadly conditions in the current world.

Measurable information show the mortality of cardiovascular illness by uncovering the level of passing's overall brought about by heart ailments.

In this way, there is a verifiable need to foresee the condition at the soonest.

Key Words: Data mining, Heart disease, Machine learning, Medical Centre, Heart Ailments

1. INTRODUCTION

Coronary sickness is one of the prevalent affliction that can provoke lessen the future of individuals nowadays. Consistently 17.5 million people are kicking the basin in view of coronary sickness [1]. Life is reliant on part working of heart, since heart is fundamental piece of our body. Coronary ailment is a contamination that impacts on the limit of heart [2]. A check of a person's danger for coronary ailment is huge for certain pieces of prosperity headway and clinical prescription. A risk desire model may be traversed multivariate backslide examination of a longitudinal report [3]. In light of Cutting edge advancements are rapidly creating, social protection networks store tremendous proportion of data in their database that is marvelous and testing to assessment. Standard characteristics used for coronary sickness are Age, Sex, Fasting Blood Pressure, Chest Pain type, Resting ECG(test that measures the electrical development of the heart), Number of critical vessels tinted by fluoroscopy, Threst Blood Pressure (hypertension), Serum Cholestrol (choose the danger for making coronary ailment), Thalach (most outrageous heartbeat achieved), ST trouble (finding on an electrocardiogram, follow in the ST section is peculiarly low

underneath the example), painloc (chest torture zone (substernal=1, otherwise=0).

Table -1: Various sorts of coronary illness

| | |
|-------------------------------------|--|
| Arrhythmia | The heart beat is silly whether it may inconsistent, unnecessarily moderate or unreasonably snappy. |
| Cardiac arrest | A startling loss of heart capacity, cognizance and breathing happen out of nowhere. |
| Congestive cardiovascular breakdown | The heart doesn't siphon blood similarly as it should, it is the condition of relentless. |
| Congenital heart disease | The heart's variation from the norm which creates before birth. |
| Coronary artery disease | The heart's significant veins can harm or any infection happens in the veins. |
| Hypertension | It has a condition that the intensity of the blood against the flexibly course dividers is unnecessarily high. |
| Fringe course ailment | The constrained veins which lessen stream of blood in the limbs. |
| Stroke | Interference of blood flexibly happen harm to the mind. |

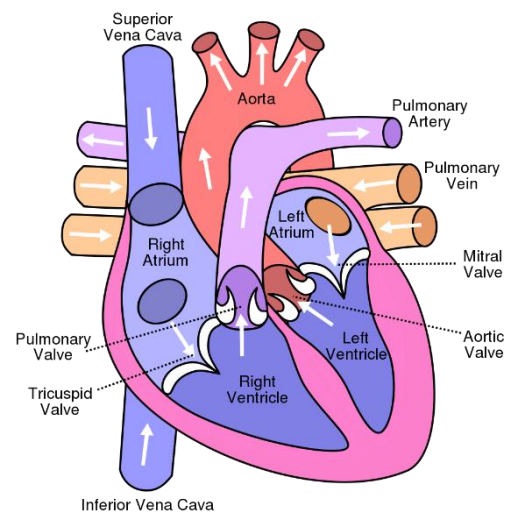


Fig - 1: Human Heart

2. Literature Review and inspiration

There are various works has been done identified with affliction gauge structures utilizing various information mining strategies and AI calculations in clinical center interests. K. Polaraju et al, proposed Prediction of Heart Disease utilizing Multiple Regression Model and it shows that Multiple Linear Regression is fitting for imagining coronary sickness credibility. The work is performed using planning educational assortment contains 3000 events with 13 particular qualities which has referenced previously. The informational rundown is allocated into two territories that is 70% of the information are utilized for preparing and 30% utilized for testing. Thinking about the outcomes, obviously the social occasion precision of Regression figuring is better stood apart from different calculations j48, SMO, and Bayes Net and Multilayer sharpness utilizing WEKA programming. Considering execution from various factor SMO and Bayes Net accomplish ideal execution than KStar, Multilayer recognition and J48 frameworks utilizing k-overlay cross support. The precision presentations achieved by those figurings are so far not attractive. Therefore, the exactness' acquaintance is improved more with give better choice to finding illness. S. Seema et al, spins around techniques that can foresee perpetual turmoil by mining the information containing in true blue thriving records utilizing Naïve Bayes, Decision tree, Support Vector Machine(SVM) and Artificial Neural Network(ANN). A relative report is performed on classifiers to gauge the better execution on an accurate rate. From this test, SVM gives most critical exactness rate, anyway for diabetes Naïve Bayes gives the most raised precision. Ashok Kumar Dwivedi et al, proposed various estimations like Naive Bayes, Classification Tree, KNN, Logistic Regression, SVM and ANN. The Logistic Regression gives better precision stood apart from different checks. MeghaShahi et al, proposed Heart Disease Prediction System utilizing Data Mining Techniques. The paper suggested SVM is persuading and furnishes more exactness as separated and other information mining figurings. Chala Beyene et al, suggested Prediction and Analysis the event of Heart Disease Using Data Mining Techniques. The proposed system is in like way crucial in human organizations relationship with experts that have no more information and twisted. It utilizes arranged clinical characteristics, for example, glucose and heartbeat, age, sex are a touch of the credits are joined to perceive if the individual has coronary illness or not.

3. Proposed System

3.1 Algorithms Used

Gullible Bayes Naive Bayes being a fundamental yet a suitable gathering system which relies upon the Naive Bayes Theorem. It expect freedom amongpointers, i.e., the characteristics or features should be not related to one another or should not, regardless, be related to each other. Whether or not there is dependence, still all of these features

or attributes openly add to the probability and that is the explanation it is called Naïve.

3.2 Dimensionality Reduction

Dimensionality Reduction includes choosing a numerical portrayal to such an extent that one can relate most of, yet not all, the change inside the given information, in this way including just most noteworthy data. The information considered for an assignment or an issue, May comprises of a great deal of traits or measurements, however not these properties may similarly impact the yield. Countless characteristics, or highlights, may influence the computational multifaceted nature and may even prompt over fitting which prompts poor outcomes.

3.3 Feature Extraction

In this, another arrangement of highlights is gotten from the first list of capabilities. Highlight extraction includes a change of the highlights. This change is regularly not reversible as few, or possibly many, helpful data is lost all the while. In and (PCA) is utilized for highlight extraction. HCA is a famously utilized direct change calculation. In the component space, it finds the headings that boost change and discovers bearings that are commonly symmetrical.

3.4 Bolster Vector Machine

Bolster Vector Machine is an inconceivably standard directed AI system (having a pre-portrayed objective variable) which can be used as a classifier similarly as a pointer. For gathering, it finds a hyper-plane in the component space that isolates between the classes. A BVM model addresses the readiness data centers as centers in the segment space, planned so that centers having a spot with free classes are detached by an edge as wide as could be normal the situation being what it is.

3.5 Choice Tree

Choice tree is an directed learning count. These strategies are generally utilized in arrangement issues. It performs easily with ceaseless and straight out characteristics. This count disconnects the people into in any event two similar sets reliant on the most critical indicators. Choice Tree count, first processes the entropy of each and every characteristic. By then the dataset is part with the help of the components or pointers with most extraordinary information increment or least entropy.

3.6 Problem Statement

- Not ready to productively perform deals arranging around cardiovascular field.
- Spending much time in distinguishing an area for cardiovascular field.

• No estimating framework accessible around cardiovascular.

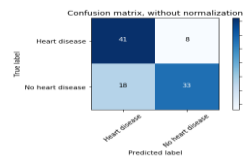
3.7 Solution Component

- Python Script
- Pandas: data analysis
- LibraryNumpy:scientific computing
- Library Sklearn: machine learning
- library Itertools : library contains functions for creating iterators to use for efficient looping
- Matplotlib: 2D plotting
- Library Seaborn: statistical data visualization

3.8 Work Flow

- The dataset used in this project contains 14 features from Cleveland Clinic Foundation for heart disease.
- Clean the data and find any missing values
- Evaluate two models linear Stochastic Gradient Descent and decision tree
- Train the data using both models and predict the probabilities of disease belonging to a particular class
- Applying cross validation on the training and test set for validating our Models

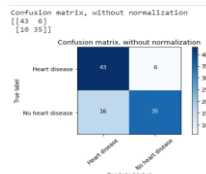
```
Decision tree Cross-Validation scores:
[ 0.81  0.84  0.9  0.77  0.73  0.73  0.67  0.77  0.73  0.76]
('Mean Decision tree Cross-Validation score = ', 0.77145346681497962)
Decision tree Cross-Validation scores:
[ 0.87  0.84  0.84  0.84  0.7  0.7  0.57  0.77  0.77  0.72]
('Mean Decision tree Cross-Validation score = ', 0.75112347852289316)
('The best parameters for model are ', 'gini')
('The Cross-Validation score = ', 0.77145346681497962)
```



Decision tree Test score:
0.74

Fig -2: Model Evaluation – Decision tree

```
('Mean Linear regression SGD Cross-Validation score = ', 0.736576760554518)
('Parameters for model', ('log', 'l1', 0.01, 1000))
Linear regression SGD Cross-Validation scores:
[ 0.80222301  0.87860774  0.50864516  0.4536129  0.3  0.56666667]
[ 0.46666667  0.86666667  0.73333333  0.82758621]
('Mean Linear regression SGD Cross-Validation score = ', 0.706737015131187)
('The best parameters for model are ', ('log', 'l1', 0.01, 1000))
('The Cross-Validation score = ', 0.82460140897293288)
```



Linear regression SGD Test score:
0.78

Fig -3: Model Evaluation – SGD

```
('The best parameters for model are ', ('log', 'l1', 0.01, 1000))
('The Cross-Validation score = ', 0.82460140897293288)
Linear regression SGD Test score:
0.8
```

Fig -4: Combination of Models

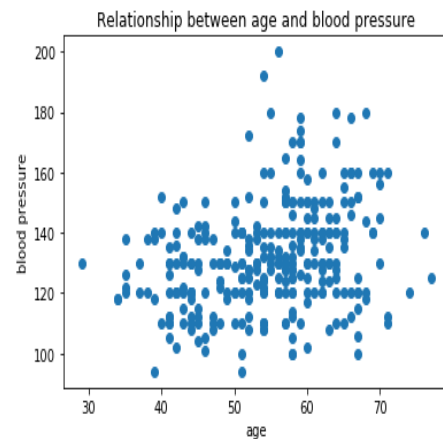


Chart -1: Relationship between age and blood pressure

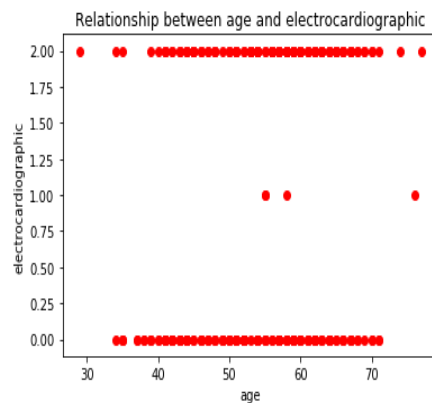


Chart -2: Relationship between age and electrocardiographic



Chart -3: Relationship between age and max heart rate

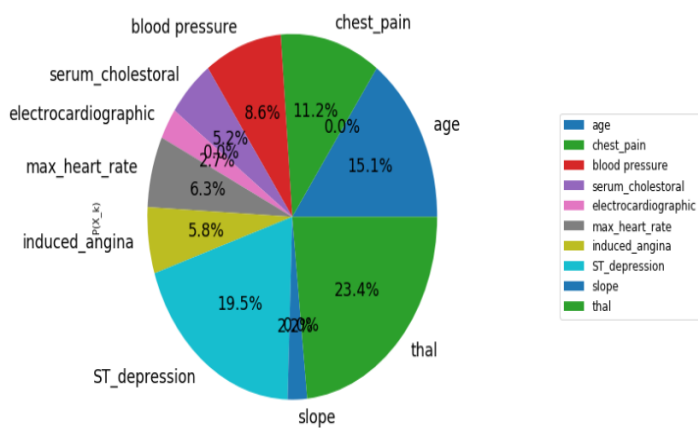


Fig -5: Most Deciding Parameters

4. Conclusion & Future Scope

By using different kinds of data mining and AI strategies to envision the occasion of coronary sickness have summarized. There are a couple of treatment procedures for calm, if they once resolved to have the particular kind of coronary disease. Data mining can be of very data structure such sensible dataset. Taking everything into account, as distinguished through the writing overview, accept just a negligible achievement is achieved really taking shape of judicious model for coronary disease patients and accordingly there is a necessity for combinational and progressively complex models to grow the accuracy of the envisioning the early phase of coronary ailment. Might want to make use of testing distinctive discretization strategies, various classifier casting a ballot procedure and diverse choice tree types.

REFERENCES

- [1] Ramadoss and Shah B et al."A. Reacting to the danger of incessant infections in India". Lancet. 2005; 366:1744–1749.doi: 10.1016/S0140-6736(05)67343-6.
- [2] Global Atlas on Cardiovascular Disease Prevention and Control. Geneva, Switzerland: World Health Organization, 2011
- [3] Dhomse Kanchan B and Mahale Kishor M. et al. "Investigation of Machine Learning Algorithms for Special Disease Prediction utilizing Principal of Component Analysis", 2016 International Conference on Global Trends in Signal Processing, Information Computing and Communication.
- [4] R.Kavitha and E.Kannan et al. "An Efficient Framework for Heart Disease Classification using Feature Extraction and Feature Selection Technique in Data Mining ", 2016
- [5] Shan Xu ,Tiangang Zhu, Zhen Zang, Daoxian Wang, Junfeng Hu and Xiaohui Duan et al. "Cardiovascular Risk Prediction Method Based on CFS Subset Evaluation and Random Forest Classification Framework", 2017 IEEE second International Conference on Big Data Analysis.
- [6] Manpreet Singh, Levi Monteiro Martins, Patrick Joanis and Vijay K. Mago et al. " Building a Cardiovascular Disease Predictive Model utilizing Structural Equation Model and Fuzzy Cognitive Map", 978-1-5090-0626-7/16/\$31.00 c 2016 IEEE.