

Heart Failure Detection System using Random Forest Classification

Naveen P¹, Anjan Babu G²

Dept of Computer Science, SV University, Tirupati

Abstract— Coronary disease is maybe the most essential human ailments on earth and impacts human life harshly. Heart related disorders or Cardiovascular Diseases are the foremost defense innumerable passing in the world over the span of the several numerous years. In coronary disease, the heart can't push the fundamental proportion of blood to various bits of the body. Exact and on time investigation of coronary ailment is critical for cardiovascular breakdown evasion and treatment. Along these lines, there is a need of strong, exact and feasible system to dissect such infections on time for suitable treatment. The assessment is done on the Heart disappointment clinical records dataset taken from the University of California at Irvine Machine Learning Data Repository. The dataset contains a huge volume of feature estimations which are diminished using hereditary calculation-based component assurance methodology. The dataset contains a colossal rundown of capacities which is diminished using a further developed segment decision methodology named as covering procedure. The proposed covering procedure depends on an arbitrary woods and hereditary calculation estimation to pick the fundamental highlights from the given dataset. The picked subset of highlights then goes through a preprocessing step to introduce a consistency in the apportionment of data. Our proposed model has achieved 84 % precision.

I. INTRODUCTION

We see recently unique clinical affiliations are conveying colossal proportions of data which are difficult to manage. Facilities have assembled tremendous measures of information about patients and their clinical records. Data burrowing is searching for associations and models that could give important data to suitable dynamic. Clinical data mining is one of the primary concerns of dispute to get significant clinical data from clinical informational collections.

This is the mother support some associated clinical issues like coronary failure, liver disappointment, kidney dissatisfactions, nerves harms and vision incident. One of the critical real clinical issues is the area of diabetes at its starting stage.

Heart is the most key organ in human body accepting that organ gets influenced; it additionally impacts the other key bits of the body. Along these lines it is vital for individuals to go for a coronary infection examination [1].

The principle organ of the human body is heart. The limit of the heart is to siphon the blood and circles entire body [3]. The coronary disease (HD) has been considered as one of the complex and life deadliest human infections on earth. In this infection, for the most part the heart can't push the fundamental proportion of blood to various bits of the body to fulfill the conventional functionalities of the body, and along these lines, at last the cardiovascular breakdown occurs. As demonstrated by the World Health Organization (WHO), a normal 17 million people fail miserably consistently from cardiovascular disease, particularly coronary disappointments and strokes [10].

The results of coronary ailment fuse shortness of breath, deficiency of real body, swollen feet, and exhaustion with related signs, for example, raised jugular venous squeezing variable and periphery edema achieved by valuable heart or noncardiac abnormalities [11]. The assessment strategies in starting stages used to recognize coronary ailment were tangled, and its ensuing multifaceted design is one of the huge reasons that impact the standard of life. The coronary ailment assurance and treatment are amazingly multifaceted, especially in the farming countries, in view of the phenomenal availability of suggestive mechanical get together and absence of specialists and others resources which impact authentic assumption and treatment of heart patients. The exact and genuine investigation of the coronary disease peril in patients is imperative for diminishing their connected risks of genuine heart issues and further developing security of heart [11].

II. FEATURE SELECTION

Feature assurance has been extensively investigated and used by the AI and data mining neighborhood. In this interesting circumstance, a part, moreover called quality or variable, addresses a property of a cycle or structure than has been assessed or worked from the principal data factors. The goal of feature decision is to pick the most diminutive component subset given a particular hypothesis bungle, or then again tracking down the best component subset with k highlights, that yields the base theory botch [2][4]. Additional objections of highlight assurance are according to the accompanying: (I) further foster the hypothesis execution concerning the model amassed using the whole plan of features, (ii) give a more energetic theory and a faster response with unnoticeable data, and (iii) achieve a predominant and less perplexing understanding of the collaboration that makes the data. Feature assurance strategies are ordinarily described in three guideline social events: covering, embedded, and channel methods. Covers use the acknowledgment learning estimation as a segment of the limit surveying feature subsets [8]. The show is ordinarily estimated as far as the gathering rate got on a testing set, i.e., the classifier is used as a black box for assessing feature subsets [5]. Yet these methodologies might achieve a nice theory, the computational cost of setting up the classifier a combinatorial number of times becomes restrictive for high-dimensional datasets.

2.1 Genetic Algorithm (GA)

Hereditary Algorithms for Feature Selection Among the various orders of feature assurance estimations, the groundbreaking computations particularly GAs are standard and by and large used. Hereditary estimations are search computations subject to the guidelines of ordinary assurance and innate characteristics, introduced by J Holland in the 1970's and moved by the normal advancement of living animals. Hereditary estimations hypothetical the issue space as a general population of individuals, and endeavor to explore the fittest individual by conveying ages iteratively. GA fosters a general population of starting individuals to a general population of extraordinary individuals, where each individual tends to an answer of the issue to be settled. The idea of every standard is assessed by a health function as the quantitative depiction of every standard's change to a particular environment. The strategy starts from a basic people of erratically created individuals. During each age, three fundamental inherited heads are sequentially applied to each individual with explicit probabilities, for instance assurance, half and half and change [6].

In GA, the request space contains strings, all of which tending to a contender answer for the issue and are named as chromosomes. The objective work worth of each chromosome is called its health regard. People are a lot of chromosomes close by their connected health. Ages are masses delivered in an accentuation of the GA [7]. Genetic estimation to glance through a space of up-and-comer answers for perceive the best one is according to the accompanying:

2.2 Genetic Algorithm Steps

1. [Start] Generate self-assertive people of n chromosomes (sensible responses for the issue).
2. [Fitness] Evaluate the wellbeing $f(x)$ of each chromosome x in the general population.
3. [New population] Create another general population by proceeding after steps until the new people is done
 - a) [Selection] Select two parent chromosomes from a general population as shown by their health (the better wellbeing, the more prominent chance to be picked).
 - b) [Crossover] With a half breed probability get over the watchmen to outline another successors (youths). If no half and half was performed, any kind of family down the line is an exact of watchmen.
 - c) [Mutation] With a change probability change new successors at each locus (position in chromosome).
 - d) [Accepting] Place new successors in another general population.
4. [Replace] Use new created people for a further run of computation.
5. [Test] If the end condition is satisfied, stop, and return the best game plan in current people.

6. [Loop] Go to organize 2.

III. RANDOM FOREST

Optional woodland region is a get-together learning framework subject to portrayal and break confidence trees. Each tree is ready on a bootstrap test, and optimal portions at each split are seen from a self-confident subset thing being what they are. Despite doubt, fearless trees can be used to assess variable importance measures to rank parts by reasonable importance. The irregular backwoods region is used to get the piece masterminding characteristics, and these traits are applied to pick which highlights are discarded in each accentuation of the appraisal [4][5]. The design combines the movement of an immense number of choice trees and inside strange trees; haphazardness is used in the going with ways: first thing, each choice tree is made using another bootstrap test. Likewise, during the improvement of each decision tree, each center split unites the conflicting certification of a subset of k portions, of which the best split is settled. It is especially valuable for enormous datasets several information highlights since it lessens the uproar, various nature and running time of the evaluation.

IV. EXPERIMENTAL RESULTS

The assessments have been coordinated by using Python programming language. It is an open-source programming language give stunning utilization of different data examination and Visualization methodologies. It is an earth-shattering library that gives numerous AI gathering estimations, capable mechanical assemblies for data mining and data assessment. The Python Scikit-learn is a pack for data request, backslide, bundling and portrayal. We have considered the Heart failure clinical records information from UCI Machine Learning Repository datasets [12]. The dataset subtleties are displayed in table-1 and the Statistical outline of the dataset as displayed in the figure-1 and figure-2. The standard dataset is parceled into two sets (70% and 30%), one for getting ready and another set for testing.

TABLE 1
HEART FAILURE CLINICAL RECORDS DATASET

S. No	Dataset	No. of Features	No. of Instances	Class Division
1	Heart failure clinical records	13	299	Dead_96 Live_203

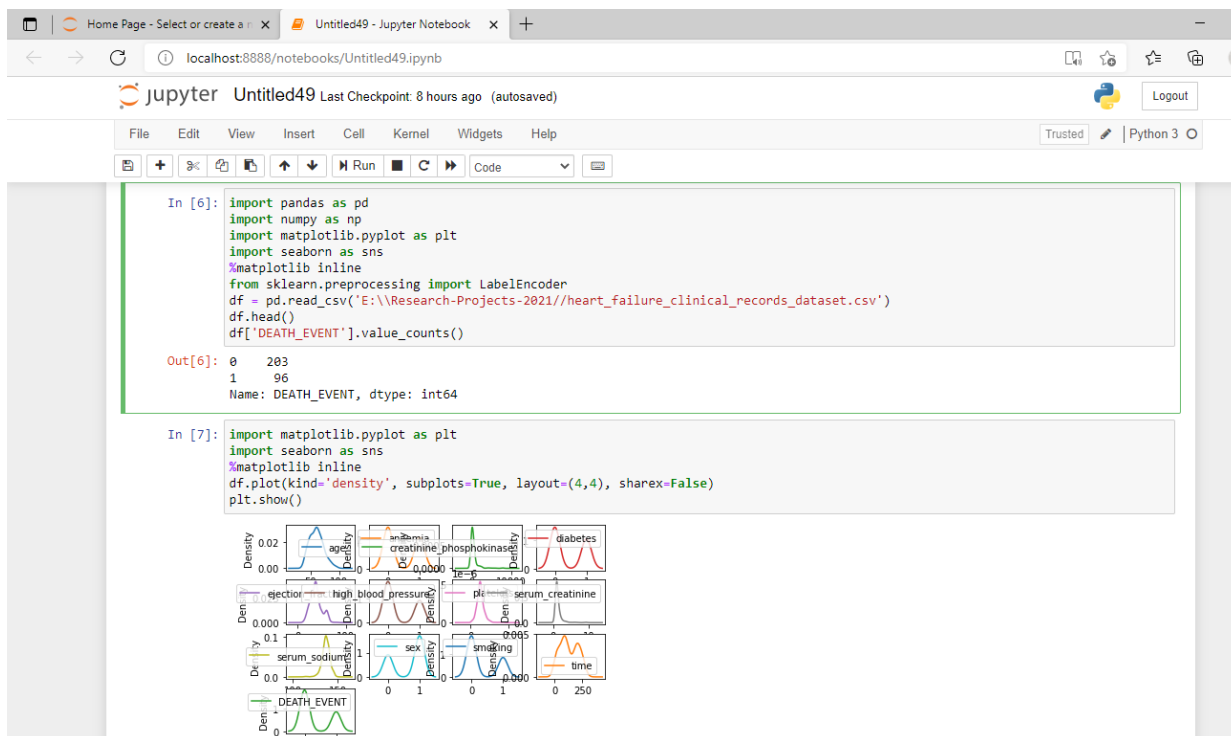


FIGURE 1: Detailed information of the dataset

```

X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.30)
print ('The train data has {0} rows and {1} columns'.format(X_train.shape[0],X_train.shape[1]))
print ('-----')
print ('The test data has {0} rows and {1} columns'.format(X_test.shape[0],X_test.shape[1]))

Total No.of Records
(299, 13)
The train data has 209 rows and 12 columns
-----
The test data has 90 rows and 12 columns

In [10]: df.describe()
Out[10]:

```

	age	anaemia	creatinine_phosphokinase	diabetes	ejection_fraction	high_blood_pressure	platelets	serum_creatinine	serum_sodium
count	299.000000	299.000000	299.000000	299.000000	299.000000	299.000000	299.000000	299.000000	299.000000
mean	60.833893	0.431438	581.839465	0.418060	38.083612	0.351171	263358.029264	1.39388	136.625418
std	11.894809	0.496107	970.287881	0.494067	11.834841	0.478136	97804.236869	1.03451	4.412477
min	40.000000	0.000000	23.000000	0.000000	14.000000	0.000000	25100.000000	0.500000	113.000000
25%	51.000000	0.000000	116.500000	0.000000	30.000000	0.000000	212500.000000	0.900000	134.000000
50%	60.000000	0.000000	250.000000	0.000000	38.000000	0.000000	262000.000000	1.100000	137.000000
75%	70.000000	1.000000	582.000000	1.000000	45.000000	1.000000	303500.000000	1.400000	140.000000
max	95.000000	1.000000	7861.000000	1.000000	80.000000	1.000000	850000.000000	9.400000	148.000000

```

In [11]: clf = DecisionTreeClassifier()
clf = clf.fit(X_train,y_train)
y_pred = clf.predict(X_test)
pd.crosstab(y_test, y_pred, rownames=['Actual Disease'], colnames=['Predicted Disease'])

```

FIGURE 2: Statistical summary of the dataset

V. RESULTS

In the first stage Random Forest algorithm was trained on the original set of features was used in the experiment. In the second stage we implement a GA algorithm for obtaining the adequate number of features to identify the features selected. The results that we got for forest without GA based feature selection and with GA based feature selection are shown below in the table-2 and same as shown in the figure-3 with their corresponding values.

**TABLE 2
PERFORMANCE OF CLASSIFIERS**

Algorithm	Precision	Recall	f1-score
Random Forest without GA	80	82	80
Random Forest with GA	84	85	84

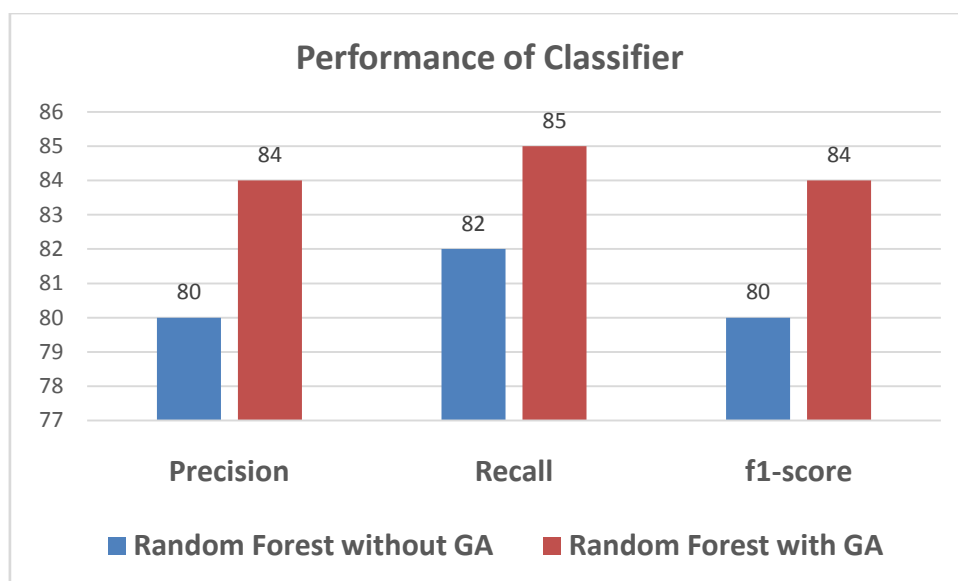


FIGURE 3: Performance of classifier

From the figure-3, we observe the performance of random forest without GA based on accuracy has got 80%, whereas the performance of random forest with GA feature selection based on accuracy has achieved 84%. However, there is an improvement in the accuracy with feature selection. The accuracy rate is increased 4% with feature selection.

In our experimental result the random forest with GA feature selection algorithm shows the highest accuracy compared without GA. With the improvement the accuracy, the proposed model demonstrated that it performs well after selecting relevant features. This result provided new insight using a classification learning algorithm and reduction technique to selection relevant and important feature in order to improve the accuracy of the system and to identify possible features which may contribute to this improvement. Most of the proposed research system could effectively utilize feature selection process to improve detection rate of their system and minimize considerably the false alarm rate.

The experimental results screen shots are shown in the figure-4 and figure-5.

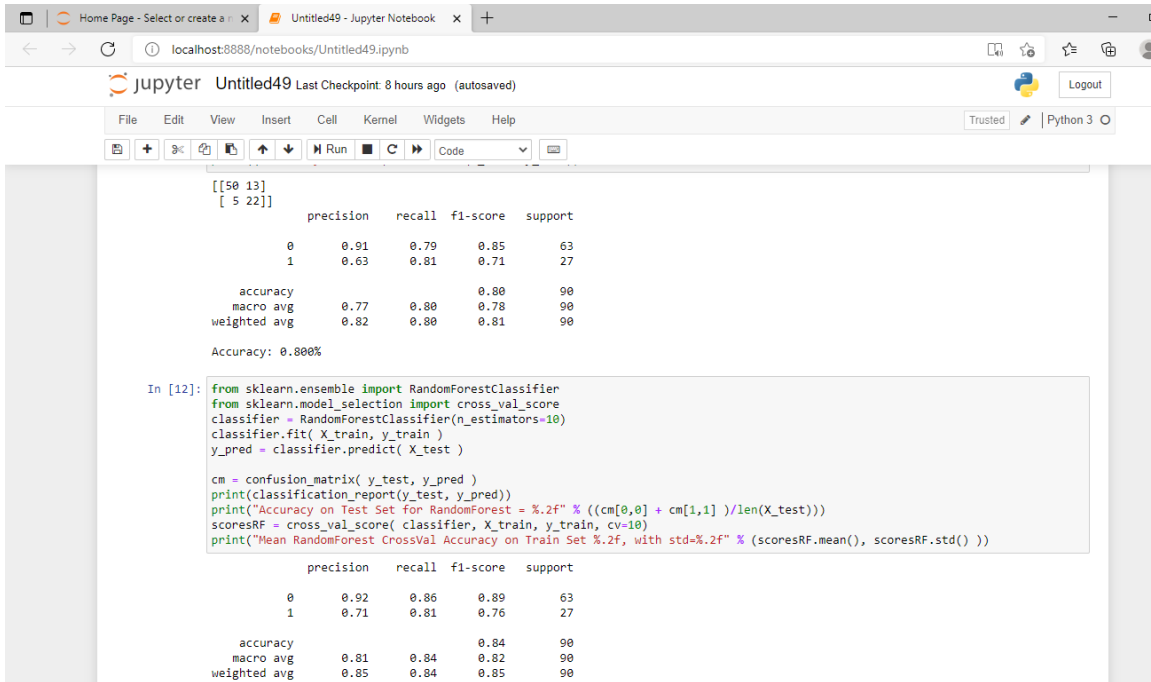


FIGURE 4: Experimental result Screen shot

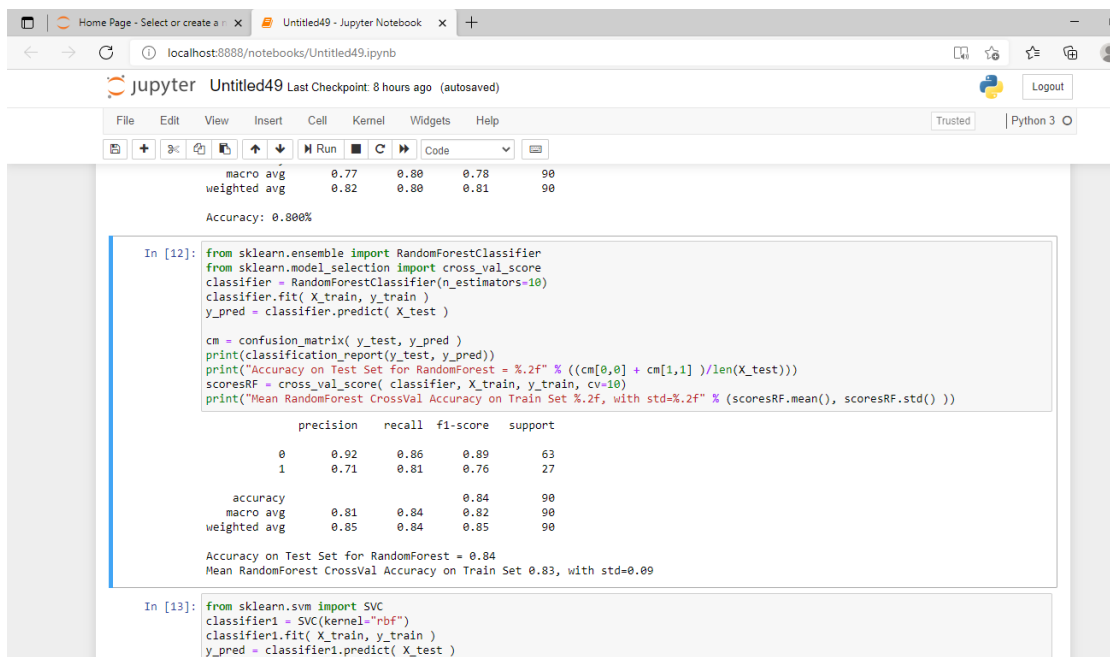


FIGURE 5: Experimental result Screen shot

VI. CONCLUSION

This paper has explored the ways to deal with tackle the significant characterization issue of the element choice. A show and recommendation of a component determination strategy which comprise of a GA highlight disposal utilizing an irregular timberland classifier to recognize significant highlights have been finished. The element choice, preprocessing, and grouping strategies have delivered a blend which gives promising outcomes to arrangement. The assessment the adequacy of the strategy utilizing distinctive order metric estimation has been made and it has been demonstrated that by diminishing the quantity of highlights, the exactness of the model was improved. To recognize class from enormous dataset, identification calculation, and highlight determination strategy have excessively more productive.

REFERENCES

- [1] Blake, C.L., Mertz, C.J.: "UCI Machine Learning Databases", <http://mllearn.ics.uci.edu/databases/heartdisease/>, 2004.
- [2] Guyon I, Elisseeff A (2003) An introduction to variable and feature selection. *J Mach Learn Res* 3:1157–1182
- [3] HeonGyu Lee, Ki Yong Noh, KeunHoRyu, "Mining Biosignal Data: Coronary Artery Disease Diagnosis using Linear and Nonlinear Features of HRV," *LNAI 4819: Emerging Technologies in Knowledge Discovery and Data Mining*, pp. 56-66, May 2007.
- [4] Ian H. Witten and Eibe Frank. *Data Mining: Practical machine learning tools and techniques*. 2nd ed. San Francisco: Morgan Kaufmann, 2005.
- [5] J.Han and M.Kamber, "Data Mining concepts and Techniques", the Morgan Kaufmann series in Data Management Systems, 2nd ed. San Mateo, CA; Morgan Kaufmann, 2006.
- [6] Kohavi R, John GH (1997) Wrappers for feature subset selection. *ArtifIntell* 97(1-2):273–324
- [7] Noraini Mohd Razali, John Geraghty "A genetic algorithm performance with different selection strategies", *Proceedings of the World Congress on Engineering Vol II*, 2011.
- [8] N. Michael, "Artificial Intelligence - A Guide to Intelligent Systems", 2nd edition, Addison Wesley, 2005.
- [9] P.N.Tan, M.Steinbach and V.Kumar "Introduction to Data Mining", A: Addison-Wesley, 2005.
- [10] Sitar-Taut, V.A., et al., Using machine learning algorithms in cardiovascular disease risk evaluation. *Journal of Applied Computer Science & Mathematics*, 2009.
- [11] "The Atlas of Heart Disease and Stroke", [online]. http://www.who.int/cardiovascular_diseases/res_ources/atlas/en/
- [12] UCI Machine Learning Repository. <https://archive.ics.uci.edu/ml/>.