

Fungal Infection and Allergy related Disease prediction with the help of machine learning XGB Classifier and Decision Tree Algorithms

Ms. Divya Pachauri

Assistant Professor, Faculty of Computer Application, Manav Rachna International Institute of Research and Studies, Sector – 43, Delhi–Surajkund Road, Faridabad, Haryana-121004, India

Dr. Arvind Dangi

Associate Professor, Faculty of Computer Application, Manav Rachna International Institute of Research and Studies, Sector – 43, Delhi–Surajkund Road, Faridabad, Haryana-121004, India

Abstract:

The wide transformation of PC-based innovation in the medical care industry brought about the amassing of electronic information. Because of the significant information measures, clinical specialists need help in investigating side effects precisely and recognizing illnesses at the beginning phase. Nonetheless, managed AI (ML) calculations have displayed critical expectations in unparalleled standard frameworks for illness finding and supporting clinical specialists in the early location of high-risk sicknesses. The point is to perceive patterns across different managed ML models in sickness locations by assessing execution measurements. The occurrence of parasitic contaminations and sensitivity illnesses is expanding at an alarming rate, introducing a colossal test to medical care experts. This increment is straightforwardly connected with the developing populace of immuno-compromised people, coming about because of changes in clinical practice, like the utilization of concentrated chemotherapy and immunosuppressive medications. Shallow and subcutaneous contagious diseases influence the skin, keratinous tissues, and mucous films. Albeit seldom perilous, they can debilitating affect an individual's satisfaction and may, in certain conditions, spread to others or become obtrusive. Most shallow and subcutaneous contagious contaminations are handily analyzed and promptly manageable for treatment. Foundational parasitic diseases might be brought about by either a crafty organic entity that contaminates an in-danger or might be related to a more intrusive organic entity that is endemic to a particular geological region. Fundamental diseases can be hazardous and are related to high horribleness and mortality. Since determination is troublesome and the causative specialist is frequently affirmed exclusively at post-mortem, the specific occurrence of fundamental diseases takes time to decide. In this paper, we have anticipated the Contagious Contamination and Sensitivity related Illness forecast with AI XGB Classifier and Choice Tree Calculations. It is a lot of support in the well-being industry since this kind of sickness requires some investment to show its side effects and carve out an opportunity to be dealt with appropriately.

Keywords: *XGB classifier, Decision tree, supervised learning, unsupervised learning, Machine learning.*

1. Introduction:

Nowadays, People suffer from a different variety of diseases as a result of their living habits and the state of the environment. As a result, it becomes a difficult task to predict sickness at an early stage. As a doctor's perspective can say, "Health is Wealth." With the improvement of medical equipments, early detection of critical diseases is possible. In the case of a critical illness, the standard method of diagnosis may not be adequate. However, supervised machine learning (ML) algorithms have some significant potential in remarkable standard systems for the diagnosis of disease and helping medical experts in the early detection of high-risk diseases.

The mostly used Supervised ML algorithms were Decision Trees (DT), XGB Classifier, and Random Forest. With the rapid development of data and technology in the world, the healthcare domain is one of the most significant study fields in the contemporary era. The immeasurable amount of patient data is very difficult to manage. The Knowledge driven by big data analysis gives healthcare specialists clear insights about the disease that was not available before. In healthcare, machine learning is helpful in each stage of the process, from patient experience to medical research and outcomes. Machine Learning models approaches that assist in disease prediction and diagnosis. Here we explore how different supervised learning approaches are used to forecast diseases based on symptoms with the help of a decision tree and XGB classifier.

1.1. Decision tree classifier:

Choice trees are extraordinary instruments to assist anybody with choosing the best game plan. They create a profoundly important game plan where one can put choices and study the potential results of those choices. They likewise work with clients to make a fair thought of the upsides and downsides connected with every conceivable activity. A choice tree is utilized to address graphically the choices, the occasions, and the results connected with choices and occasions. Occasions are not entirely settled for every result.

A decision tree is a tree whose inside hubs can be taken as tests (on input information examples) and whose leaf hubs can be taken as classifications (of these examples). These tests are sifted down through the tree to get the right result for the info design. Choice Tree calculations can be applied and utilized in various fields. It may be utilized as a substitution for factual systems to find information, remove text, track down missing information in a class, and further develop web search tools. It likewise tracks down different applications in clinical fields. Numerous Choice tree calculations have been figured out. They have different exactness and cost-adequacy. We also need to know which calculation is ideal to utilize.

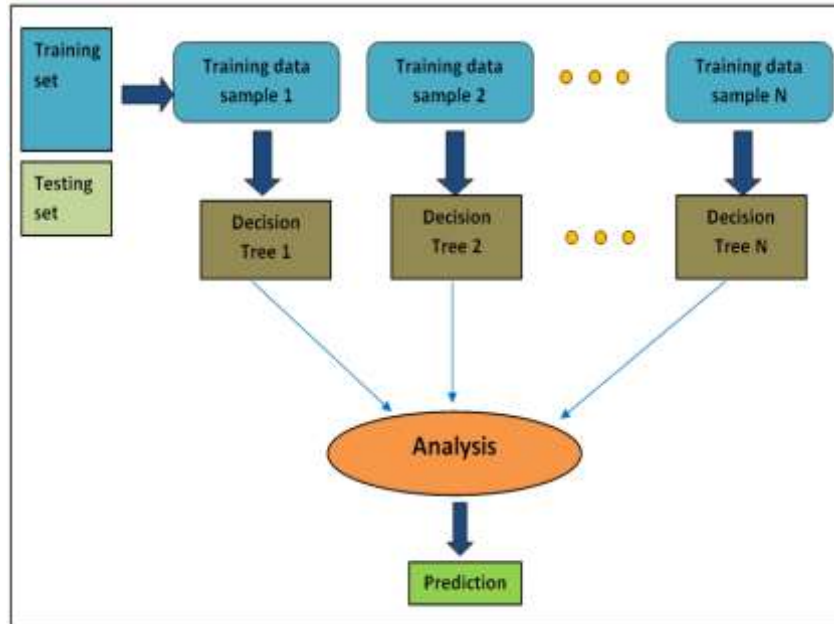


Fig1: Decision Tree Classifier Prediction Process Model

1.2. XGB Classifier:

Experts have made a couple of expert structures throughout the years to predict coronary disease early and assist cardiologists with further developing the tracking down the framework. This paper presents an illustrative structure that utilizes a smoothed XGBoost (Silly Slant Supporting) classifier to predict coronary disease.

2. Literature Review:

Choice tree characterization calculations comprise a few sorts that are utilized to create DT. This is by the control of both the consistent and occasional characteristics of the missing qualities. Hubs and branches are remembered for the DT. Every hub requires issues that depend on at least one property, for example, contrasting trait esteem and consistent or utilizing different capabilities to look at more than one property. With the end goal of the choice tree, the learning information assortment is, at times, alluded to as the result tree. To consolidate orders in AI and information mining utilizing the DT calculation. In the accompanying grouping, this calculation is applied iteratively, and the order requires a three-stage process: Develop Model (Learning), Assessment Model (Precision), and Model Use (Characterization). They arrive at the metric utilized to depict the test qualities for a hub in the tree, which is alluded to as the property determination scale (property). As a test capability for the ongoing hub, the best information property is determined. A few investigations proposed ways to deal with conquer the deficiencies of the DT issues so ideal trees can be determined, in view of a survey that was performed before, without definite subtleties and tests. DT techniques have demonstrated the way that such issues, as depicted above, can be kept away from. Moreover, it will give the predefined dataset a fitting arrangement. It was seen that many explorers were led with various informational indexes, and the DT approach was utilized to determine its shortcomings and to get better execution. A few enhancement procedures have been utilized in the review [4] to fortify the choice tree on the UCI ML datasets put away; in view of the evaluation discoveries, it was shown that the DT approach got the most elevated precision, which is 99.93% contrasting with different methods like KNN, LR, SVM, and NB which are less performing than the DT approach. In the division task, the review [7] utilized the DT way to deal with recognizing and removing the blood corridor for appropriate Optic Plate (OD) division, which brought about more noteworthy outcomes equivalent to 99.61%. In addition, in view of the review [6], it has been demonstrated the way that the DST strategy can likewise build the DT, where the two of them involved PT and MLT for DT in the UCI datasets; it has been shown that DST is more able to upgrade DT than different procedures. At last, by utilizing UCI AI Library datasets and CICIDS2017 datasets comprising of the most recent assaults among any remaining datasets, DT ended up being the most elevated, and their exactness was the best performing. While the review [3] used the DT(XGBoost) and RF on the datasets of the Smokers of the Chinese Community for Infectious prevention and Counteraction, it was seen that, once more, the DT approach accomplished the most elevated exactness, which is 84.11%. Besides, in light of studies [1], [6], [8] involving DT and KNN in the CICIDS2017, RNA-seq Jungle fever, and Wisconsin Bosom Disease datasets, it was tracked down that the DT Tijo and Abdulazeez/Diary of Applied Science and Innovation Patterns Vol. 02, No. 01, pp. 20 - 28 (2021) 26 methodology had the most noteworthy precision in each of the three examinations.

XGBoost is a practical and versatile computerized reasoning classifier Chen and Guestrin upheld in 2016 [11]. Point supporting choice tree is the chief model of XGBoost, which joins different choice trees in a helping way. If all else fails, each new tree is made to decrease the extra of the past model by the inclination having an effect. The capabilities between the genuine and anticipated values give out the overabundance. The model has been organized until the number of choice trees exhibits the edge. Early regions and appropriate medicines manage to decrease passing rates accomplished by constant issues. These days, farsighted models are utilized in the early end and evaluation of smoking-related illnesses and afflictions [8-12]. In one review [13], the producers zeroed in on the association between normal factors and working on Crohn's illness among Japanese. Their outcomes recommend that isolated smoking history is associated with working on Crohn's problem. Another review [14] assessed the association between Parkinson's burden and public living, improvement, pesticide use, and cigarette smoking. The heaviness of the check and meta-evaluation showed a causal relationship between the bet of Parkinson's sickness and cigarette smoking, which has been constantly found in related creations. Inquisitively, everyday living, well-water use, improvement, and the utilization of pesticides, herbicides, bug showers, and fungicides were less strong in Parkinson's tainting.

Besides, one review [15] zeroed in on part risk pathways in the smoking-prompted cell breakdown in the lungs utilizing patient information from the Quality Verbalization Omnibus enlightening assortment. They advanced the capacities utilizing the inconsistency score and the recursive part expulsion (RFE) methodology. Then, at that point, the help vector with machining (SVM) based guess model was utilized. Their overview thought that smoking is the fundamental driver of cell breakdown in the lungs; stress and self-security in continuing with living things can be seen as turbulent elements. In another review [16], producers urged an altered classifier to build the precision of the constrained impacting strategy for an early finding of smoking-incited respiratory changes. They used a couple of man-made brainpower methodologies, for example, resolved direct classifiers, k closest neighbor (KNN), frontal cortex affiliations (NN), and SVM. Therefore, KNN and SVM classifiers accomplished a further expansion in precise accuracy.

3. Research Methodology:

The means in the proposed work process are shown, which include the pre-handling of preparing, testing information with determined models, assessment of results, and prediction of

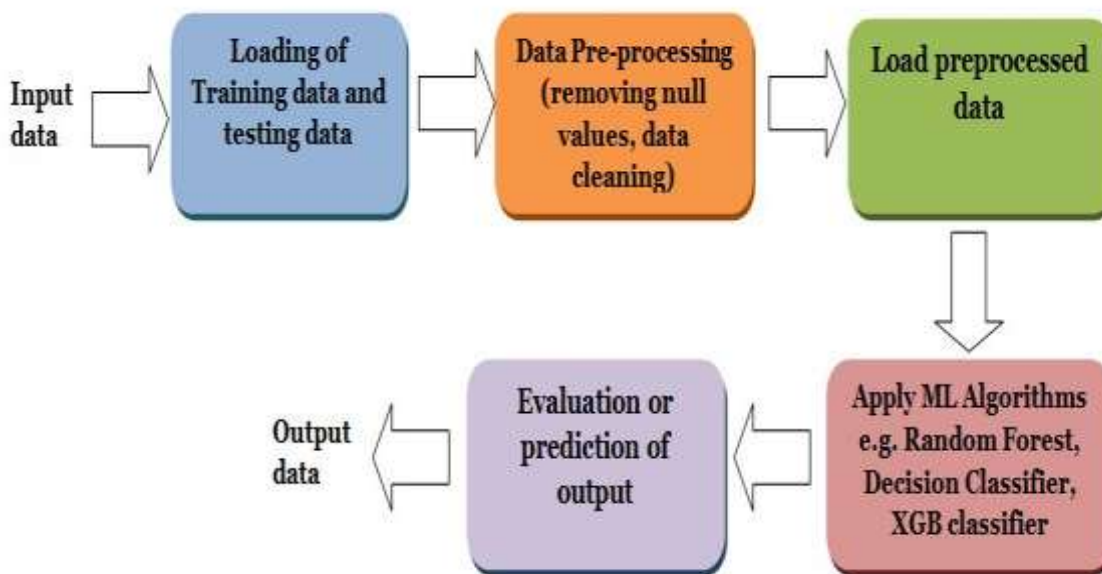


Fig2: Implementation of Disease prediction process

The first step of this model starts with collecting Kaggle data of 41 types of diagnosis diseases (i.e., fungal infection, Allergy, Chronic cholestasis) with different types of symptoms in the form of a CSV file. This set of input data is further divided into training and testing data in the next steps.

• Step 1 Loading/Reading of data:

This work is fully implemented in Python3. So NumPy and pandas library is used for working with numerical data and data frames included in Python3. Training and testing CSV file are converted into train data and test data set.

• Step 2 Data Preprocessing:

Train and test datasets contain some incomplete information, like null values and incorrect values. In the data preprocessing step, null values are converted to meaningful categorical data labels, and after that, LabelEncoder is used to convert these labels of categorical data into a numeric format to fit the classifier.

• **Step 3 Preprocessed data and ML Algorithm:**

Supervised classification algorithm including XGB Classifier & Decision Classification algorithm applied to data generated after preprocessing step. These three algorithms have their own individual correct diagnosis method with prediction and accuracy scores.

4. **Result and Implementation:**

In the next step, take the X and y datasets and break them into a training dataset according to the requirement of the algorithm and test (or validation) dataset need to train and test with the classifier `train_test_split()`, function from `scikit-learn`.

Out[4]:

..	scuring	skin_peeling	silver_like_dusting	small_dents_in_nails	inflammatory_nails	blister	red_sore_around_nose	yellow_crust_ooze	prognosis	Unnamed: 133
-	0	0	0	0	0	0	0	0	Fungal infection	NaN
-	0	0	0	0	0	0	0	0	Fungal infection	NaN
-	0	0	0	0	0	0	0	0	Fungal infection	NaN
-	0	0	0	0	0	0	0	0	Fungal infection	NaN
-	0	0	0	0	0	0	0	0	Fungal infection	NaN

Fig3: input training dataset with NULL values

Out[5]:

..	blackheads	scuring	skin_peeling	silver_like_dusting	small_dents_in_nails	inflammatory_nails	blister	red_sore_around_nose	yellow_crust_ooze	prognosis
-	0	0	0	0	0	0	0	0	0	Fungal infection
-	0	0	0	0	0	0	0	0	0	Allergy
-	0	0	0	0	0	0	0	0	0	GERD
-	0	0	0	0	0	0	0	0	0	Chronic cholestasis
-	0	0	0	0	0	0	0	0	0	Drug Reaction

Fig4: input testing dataset

Over the split, Null values loaded up for certain information values and the preparation information is put away in `X_train` and `y_train` and test information is put away in `X_test` and `y_test`. The `X_test` information isn't utilized during preparing model, rather being utilized subsequent to preparing stage to assess the model and survey its exactness utilizing some extraordinary exhibition assessment measurements.

5. **Evaluation and Prediction:**

In the evaluation stage, some statistical models are used to predict the output of the classification model.

5.1 **XGB Classifier:**

XGBClassifier is a scikit-learn API compatible class for classification approach.

xgboost=XGBClassifier()

XGBClassifier(base_score=0.5, booster='gbtree', colsample_bylevel=1, colsample_bynode=1, colsample_bytree=1, gamma=0, learning_rate=0.1, max_delta_step=0, max_depth=3, min_child_weight=1, missing=None, n_estimators=100, n_jobs=1, nthread=None, objective='multi:softprob', random_state=0, reg_alpha=0, reg_lambda=1, scale_pos_weight=1, seed=None, silent=None, subsample=1, verbosity=1) xgboost.fit(X_train,y_train)

Accuracy score is: 0.9761904761904762

Fig5: XGBClassifier Model outcome

	Actual	Predicted
0	15	15
1	4	4
2	16	16
3	9	9
4	14	14
5	33	33
6	1	1
7	12	12
8	17	17
9	6	6

Fig6: XGB classifier Actual vs predicted output

Predicted	0	1	2	3	4	5	6	7	8	9	...	31	32	33	34	35
Actual	0	1	2	3	4	5	6	7	8	9	...	31	32	33	34	35
0	1	0	0	0	0	0	0	0	0	0	...	0	0	0	0	0
1	0	1	0	0	0	0	0	0	0	0	...	0	0	0	0	0
2	0	0	1	0	0	0	0	0	0	0	...	0	0	0	0	0
3	0	0	0	1	0	0	0	0	0	0	...	0	0	0	0	0
4	0	0	0	0	1	0	0	0	0	0	...	0	0	0	0	0
5	0	0	0	0	0	1	0	0	0	0	...	0	0	0	0	0
6	0	0	0	0	0	0	1	0	0	0	...	0	0	0	0	0
7	0	0	0	0	0	0	0	1	0	0	...	0	0	0	0	0
8	0	0	0	0	0	0	0	0	1	0	...	0	0	0	0	0
9	0	0	0	0	0	0	0	0	0	1	...	0	0	0	0	0
10	0	0	0	0	0	0	0	0	0	0	...	0	0	0	0	0
11	0	0	0	0	0	0	0	0	0	0	...	0	0	0	0	0
12	0	0	0	0	0	0	0	0	0	0	...	0	0	0	0	0
13	0	0	0	0	0	0	0	0	0	0	...	0	0	0	0	0
14	0	0	0	0	0	0	0	0	0	0	...	0	0	0	0	0
15	0	0	0	0	0	0	0	0	1	0	...	0	0	0	0	0
16	0	0	0	0	0	0	0	0	0	0	...	0	0	0	0	0
17	0	0	0	0	0	0	0	0	0	0	...	0	0	0	0	0
18	0	0	0	0	0	0	0	0	0	0	...	0	0	0	0	0
19	0	0	0	0	0	0	0	0	0	0	...	0	0	0	0	0
20	0	0	0	0	0	0	0	0	0	0	...	0	0	0	0	0
21	0	0	0	0	0	0	0	0	0	0	...	0	0	0	0	0
22	0	0	0	0	0	0	0	0	0	0	...	0	0	0	0	0
23	0	0	0	0	0	0	0	0	0	0	...	0	0	0	0	0
24	0	0	0	0	0	0	0	0	0	0	...	0	0	0	0	0
25	0	0	0	0	0	0	0	0	0	0	...	0	0	0	0	0
26	0	0	0	0	0	0	0	0	0	0	...	0	0	0	0	0
27	0	0	0	0	0	0	0	0	0	0	...	0	0	0	0	0
28	0	0	0	0	0	0	0	0	0	0	...	0	0	0	0	0
29	0	0	0	0	0	0	0	0	0	0	...	0	0	0	0	0
30	0	0	0	0	0	0	0	0	0	0	...	0	0	0	0	0
31	0	0	0	0	0	0	0	0	0	0	...	1	0	0	0	0
32	0	0	0	0	0	0	0	0	0	0	...	0	1	0	0	0
33	0	0	0	0	0	0	0	0	0	0	...	0	0	0	0	0
34	0	0	0	0	0	0	0	0	0	0	...	0	0	0	0	1
35	0	0	0	0	0	0	0	0	0	0	...	0	0	0	0	0
36	0	0	0	0	0	0	0	0	0	0	...	0	0	0	0	0
37	0	0	0	0	0	0	0	0	0	0	...	0	0	0	0	0
38	0	0	0	0	0	0	0	0	0	0	...	0	0	0	0	0
39	0	0	0	0	0	0	0	0	0	0	...	0	0	0	0	0
40	0	0	0	0	0	0	0	0	0	0	...	0	0	0	0	0

Fig7: Predicted output value of XGBClassifier

5.2 Decision Tree Classifier:

Decision Tree is a Directed AI based classifier that utilizes a bunch of rules to decide, likewise as the human take choice. In this calculation, in test_size contention 40% of the information to utilize for test, with the other 70% utilized for preparing.
 dtc= DecisionTreeClassifier(random_state=42)
 model_dtc = dtc.fit(xtrain, ytrain)

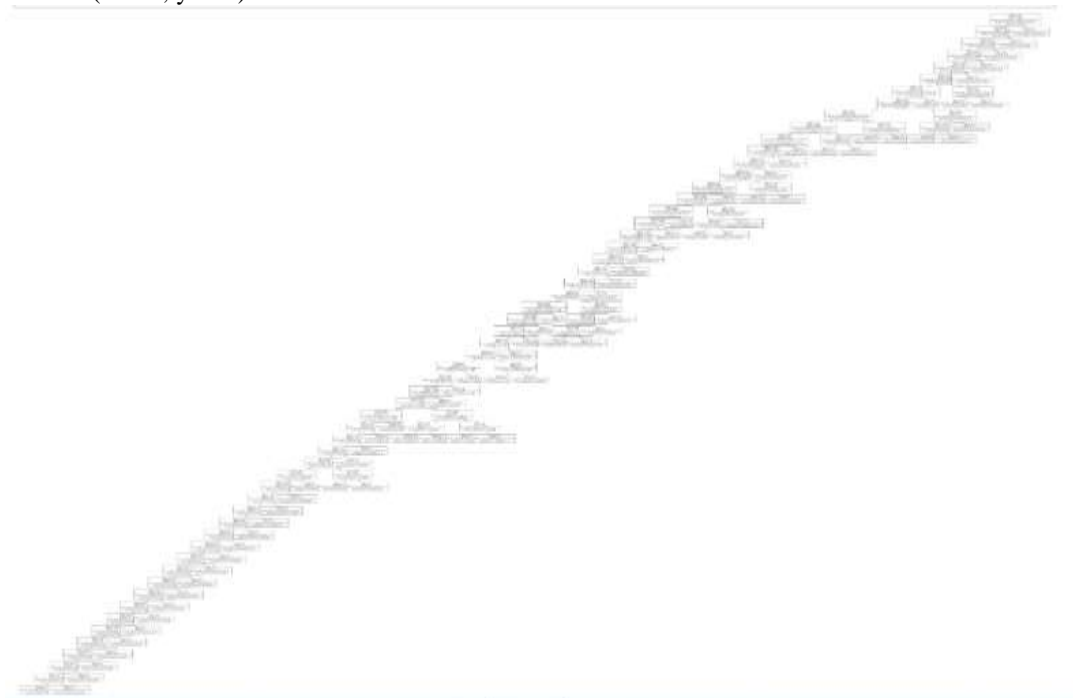


Fig8: Decision Tree-Based Classifier Graph

```
training accuracy is: 1.0
testing accuracy is: 0.9969512195121951
```

Fig9: Decision Tree classifier Model outcome

5.3 Comparison Table:

Classifier	Prediction Accuracy
XGBClassifier	0.9761904761904762
Decision Tree	0.9969512195121951

Table1: XGBClassifier Vs Decision Tree Classifier

A decision Tree-Based algorithm predicts the above output for 40% of testing data, and if dividing testing data in 30% with 70% of training data, it gives 100% accuracy with the same dataset.

6. Conclusion and future work:

As a result of our research, we have applied XGB Classifier and Decision tree-based classifier algorithms under supervised machine learning on the same training and testing datasets.

After comparison of the above two classifiers, as a result, we find out different prediction accuracy scores. The decision Tree-based algorithm predicts 0.99 accuracies for 40% of testing data, and the XGB classifier predicts 0.97. On the other hand, if we divided testing data by 30% and training data by 70%, then the accuracy score of the decision tree gives 100% accuracy with the same dataset.

As an outcome of this research, we have identified the Decision tree-based approach as the best-preferred algorithm for predicting diseases with multiple symptoms.

As we have applied this approach for small datasets, but if we vary the data size, the prediction accuracy rate may also be different from the current scenario. For future research, there is scope for applying the Random-Forest algorithm with an increased dataset of different types of diseases with their symptoms.

References:

- [1] I. Ramadhan, P. Sukarno, and M. A. Nugroho, "Comparative Analysis of K-Nearest Neighbor and Decision Tree in Detecting Distributed Denial of Service," in 2020 8th International Conference on Information and Communication Technology (ICoICT), Yogyakarta, Indonesia, Jun. 2020, pp. 1–4, doi: 10.1109/ICoICT49345.2020.9166380.
- [2] V. M. E. Batitis, M. J. G. Caballes, A. A. Ciudad, M. D. Diaz, R. D. Flores, and E. R. E. Tolentin, "Image Classification of Abnormal Red Blood Cells Using Decision Tree Algorithm," in 2020 Fourth International Conference on Computing Methodologies and Communication (ICCMC), Mar. 2020, pp. 498–504, doi: 10.1109/ICCMC48092.2020.ICCMC-00093.
- [3] Y. Zhang, J. Liu, Z. Zhang, and J. Huang, "Prediction of Daily Smoking Behavior Based on Decision Tree Machine Learning Algorithm," in 2019 IEEE 9th International Conference on Electronics Information and Emergency Communication (ICEIEC), Jul. 2019, pp. 330–333, doi: 10.1109/ICEIEC.2019.8784698. [76] S. Nandhini and J. M. K.S, "Performance Evaluation of Machine Learning Algorithms for Email Spam Detection," in 2020 International Conference on Emerging Trends in Information Technology and Engineering (ic-ETITE), Feb. 2020, pp. 1–4, doi: 10.1109/ic-ETITE47903.2020.312.
- [4] A. I. Taloba and S. S. I. Ismail, "An Intelligent Hybrid Technique of Decision Tree and Genetic Algorithm for E-Mail Spam Detection," in 2019 Ninth International Conference on Intelligent Computing and Information Systems (ICICIS), Dec. 2019, pp. 99–104, doi: 10.1109/ICICIS46948.2019.9014756.
- [5] M. O. Arowolo, M. Adebisi, A. Adebisi, and O. Okesola, "PCA Model For RNA-Seq Malaria Vector Data Classification Using KNN And Decision Tree Algorithm," in 2020 International Conference in Mathematics, Computer Engineering and Computer Science (ICMCECS), Mar. 2020, pp. 1–8, doi: 10.1109/ICMCECS47690.2020.240881.
- [6] S. Pathan, P. Kumar, R. Pai, and S. V. Bhandary, "Automated detection of optic disc contours in fundus images using decision tree classifier," *Biocybernetics and Biomedical Engineering*, vol. 40, no. 1, pp. 52–64, 2020.
- [7] A. A. Nagra et al., "Hybrid self-inertia weight adaptive particle swarm optimisation with local search using C4. 5 decision tree classifier for feature selection problems," *Connection Science*, vol. 32, no. 1, pp. 16–36, 2020.
- [8] A. Ahmim, L. Maglaras, M. A. Ferrag, M. Derdour, and H. Janicke, "A novel hierarchical intrusion detection system based on decision tree and rules-based models," in 2019 15th International Conference on Distributed Computing in Sensor Systems (DCOSS), 2019, pp. 228–233.
- [9] M. Li, H. Xu, and Y. Deng, "Evidential decision tree based on belief entropy," *Entropy*, vol. 21, no. 9, p. 897, 2019.
- [10] P. Sathiyarayanan, S. Pavithra, M. S. SARANYA, and M. Makeswari, "Identification of Breast Cancer Using The Decision Tree Algorithm," in 2019 IEEE International Conference on System, Computation, Automation and Networking (ICSCAN), 2019, pp. 1–6.
- [11] Chen, T.; Guestrin, C. Xgboost: A scalable tree boosting system. In Proceedings of the 22nd Acm Sigkdd International Conference on Knowledge Discovery and Data Mining, New York, CA, USA, 13–17 August /2016; pp. 785–794.
- [12] Al-Obaide, M.A.; Ibrahim, B.A.; Al-Humaish, S.; Abdel-Salam, A.S.G. Genomic and bioinformatics approaches for analysis of genes associated with cancer risks following exposure to tobacco smoking. *Front. Public Health* 2018, 6, 84.
- [13] Kondo, K.; Ohfuji, S.; Watanabe, K.; Yamagami, H.; Fukushima, W.; Ito, K. Japanese Case-Control Study Group for Crohn's disease. The association between environmental factors and the development of Crohn's disease with focusing on passive smoking: A multicenter case-control study in Japan. *PLoS ONE* 2019, 14, e0216429.
- [14] Breckenridge, C.B.; Berry, C.; Chang, E.T.; Sielken Jr, R.L.; Mandel, J.S. Association between Parkinson's disease and cigarette smoking, rural living, well-water consumption, farming and pesticide use: Systematic review and meta-analysis. *PLoS ONE* 2016, 11, e0151841.
- [15] Chen, R.; Lin, J. Identification of feature risk pathways of smoking-induced lung cancer based on SVM. *PLoS ONE* 2020, 15, e0233445.
- [16] Amaral, J.L.; Lopes, A.J.; Jansen, J.M.; Faria, A.C.; Melo, P.L. An improved method of early diagnosis of smoking-induced respiratory changes using machine learning algorithms. *Comput. Methods Programs Biomed.* 2013, 112, 441–454.
- [17] Narin, A.; Isler, Y.; Ozer, M. Early prediction of Paroxysmal Atrial Fibrillation using frequency domain measures of heart rate variability. In Proceedings of the 2016 Medical Technologies National Congress (TIPTEKNO), Antalya, Turkey, 27–29 October 2016. <https://doi.org/10.1109/tiptekno.2016.7863110>.
- [18] Shah, D.; Patel, S.; Bharti, S.K. Heart Disease Prediction using Machine Learning Techniques. *SN Comput. Sci.* 2020, 1, 345. <https://doi.org/10.1007/s42979-020-00365-y>.
- [19] Alotaibi, F.S. Implementation of Machine Learning Model to Predict Heart Failure Disease. *Int. J. Adv. Comput. Sci. Appl.* 2019, 10, 261–268. <https://doi.org/10.14569/ijacsa.2019.0100637>.
- [20] Hasan, N.; Bao, Y. Comparing different feature selection algorithms for cardiovascular disease prediction. *Health Technol.* 2020, 11, 49–62. <https://doi.org/10.1007/s12553-020-00499-2>.
- [21] Ouf, S.; ElSeddawy, A.I.B. A proposed paradigm for intelligent heart disease prediction system using data mining techniques. *J. Southwest Jiaotong Univ.* 2021, 56, 220–240. <https://doi.org/10.35741/issn.0258-2724.56.4.19>.
- [22] Khan, I.H.; Mondal, M.R.H. Data-Driven Diagnosis of Heart Disease. *Int. J. Comput. Appl.* 2020, 176, 46–54. <https://doi.org/10.5120/ijca2020920549>.
- [23] Kaggle Cardiovascular Disease Dataset. Available online: <https://www.kaggle.com/datasets/sulianova/cardiovascular-diseasedataset> (accessed on November 1, 2022).
- [24] Han, J.A.; Kamber, M. *Data Mining: Concepts and Techniques*, 3rd ed.; Morgan Kaufmann Publishers: San Francisco, CA, USA, 2011.

- [25] Rivero, R.; Garcia, P. A Comparative Study of Discretization Techniques for Naive Bayes Classifiers. *IEEE Trans. Knowl. Data Eng.* 2009, 21, 674–688. *Algorithms* 2023, 16, 88 15 of 15
- [26] Khan, S.S.; Ning, H.; Wilkins, J.T.; Allen, N.; Carnethon, M.; Berry, J.D.; Sweis, R.N.; Lloyd-Jones, D.M. Association of body mass index with lifetime risk of cardiovascular disease and compression of morbidity. *JAMA Cardiol.* 2018, 3, 280–287. <https://doi.org/10.1001/jamacardio.2018.0022>.
- [27] Kengne, A.-P.; Czernichow, S.; Huxley, R.; Grobbee, D.; Woodward, M.; Neal, B.; Zoungas, S.; Cooper, M.; Glasziou, P.; Hamet, P.; et al. Blood Pressure Variables and Cardiovascular Risk. *Hypertension* 2009, 54, 399–404. <https://doi.org/10.1161/hypertensionaha.109.133041>.
- [28] Yu, D.; Zhao, Z.; Simmons, D. Interaction between Mean Arterial Pressure and HbA1c in Prediction of Cardiovascular Disease Hospitalisation: A Population-Based Case-Control Study. *J. Diabetes Res.* 2016, 2016, 8714745. <https://doi.org/10.1155/2016/8714745>.
- [29] Huang, Z. A Fast Clustering Algorithm to Cluster Very Large Categorical Data Sets in Data Mining. *DMKD* 1997, 3, 34–39.
- [30] Maas, A.H.; Appelman, Y.E. Gender differences in coronary heart disease. *Neth. Heart J.* 2010, 18, 598–602. doi:0.1007/s12471-010-0841-y.
- [31] Bhunia, P.K.; Debnath, A.; Mondal, P.; D E, M.; Ganguly, K.; Rakshit, P. Heart Disease Prediction using Machine Learning. *Int. J. Eng. Res. Technol.* 2021, 9.
- [32] Mohanty, M.D.; Mohanty, M.N. Verbal sentiment analysis and detection using recurrent neural network. In *Advanced Data Mining Tools and Methods for Social Computing*; Academic Press: 2022; pp. 85–106. <https://doi.org/10.1016/b978-0-32-385708-6.00012-6>.
- [33] Menzies, T.; Kocagüneli, E.; Minku, L.; Peters, F.; Turhan, B. Using Goals in Model-Based Reasoning. In *Sharing Data and Models in Software Engineering*; Morgan Kaufmann: San Francisco, CA, USA, 2015; pp. 321–353. <https://doi.org/10.1016/b978-0-12-417295-1.00024-2>.
- [34] Fayez, M.; Kurnaz, S. Novel method for diagnosis diseases using advanced high-performance machine learning system. *Appl. Nanosci.* 2021. <https://doi.org/10.1007/s13204-021-01990-6>.
- [35] Hassan, C.A.U.; Iqbal, J.; Irfan, R.; Hussain, S.; Algarni, A.D.; Bukhari, S.S.H.; Alturki, N.; Ullah, S.S. Effectively Predicting the Presence of Coronary Heart Disease Using Machine Learning Classifiers. *Sensors* 2022, 22, 7227. <https://doi.org/10.3390/s22197227>.
- [36] Subahi, A.F.; Khalaf, O.I.; Alotaibi, Y.; Natarajan, R.; Mahadev, N.; Ramesh, T. Modified Self-Adaptive Bayesian Algorithm for Smart Heart Disease Prediction in IoT System. *Sustainability* 2022, 14, 14208. <https://doi.org/10.3390/su142114208>.

Authors Profile



Ms. Divya Pachauri, the author has completed his B.Tech(CSE), and M.Tech(CSE) degrees in the year 2013 and 2016 respectively. Her research area includes Big Data analysis, Machine Learning, medical image processing.



Dr. Arvind Kumar, is currently working as an Associate Professor in the Faculty of Computer Applications at Manav Rachna International Institute of Research and studies. He has more 13 years of teaching experience. He is a Gold medalist in PGDIS Course from IGNOU University, New Delhi. He did his doctorate from Jamia Milia Islamia, M.Tech from MDU, MCA & BCA from GGSIPU, New Delhi. He has worked in the Computer Science and Cyber Security Branch at PCTI, TISS.