



Development and Performance Analysis of Machine Learning Methods for Predicting Metabolic Syndrome Among Postmenopausal Women of India

Joyeta Ghosh^{1,3*}, Sudrita Roy Chaudhury², Khusboo Singh², Samarpita Koner²

¹Department of Dietetics and Applied Nutrition, Amity University Kolkata, Kolkata, India

²Department of Food and Nutrition, Swami Vivekananda University, West Bengal, India

³School of Computer Science and Mathematics, Liverpool John Moores University, Liverpool L3 3AF, UK

*Corresponding Author's Email: joyetaghosh01@gmail.com

Abstract

Aim: The objective of this study is to develop and evaluate the performance of machine learning methods for predicting metabolic syndrome among postmenopausal women in India. **Methods:** This work uses supervised machine-learning to construct a system that achieves notable accuracy. By assessing several factors, including as accuracy, sensitivity, specificity, precision, recall, F-Measure, Receiver Operating Characteristic (ROC), Precision–Recall Curve (PRC), and Area Under the Curve (AUC), several classification methods are used to identify the best-performing classifier. **Result:** The prevalence of MetS among postmenopausal women was found to be 40.17%, with 19.21% of respondents exhibiting a hyperglycaemic state and 57.86% having low HDL-C levels. In the Indian setting, among the machine learning algorithms tested, the Decision Tree and Random Forest classifiers emerged as the best-performing models, achieving an accuracy of 90.22%. These models utilized the six most essential features as identified by the International Diabetes Federation (IDF). Key predictive factors included waist circumference (WC), serum triglyceride levels (TG), and fasting blood sugar (FBS). **Conclusion:** This study will play a crucial role in predicting MetS and improving the quality of life for neglected post-menopausal women. Various software like web and mobile applications can adopt this paradigm. Swift diagnosis will lower the costs of diagnosis and further complications.

Keywords: Machine Learning, Metabolic Syndrome, India, Postmenopausal Women

Introduction

The metabolic syndrome (MetS) is a complex cluster of interconnected conditions encompassing central obesity, hypertension, dyslipidemia, and hyperglycaemia. These physiological disruptions significantly elevate the risk of developing Type 2 diabetes mellitus and cardiovascular diseases (CVD) in the future (Sattar *et al.*, 2003; Ghosh, 2023). Research demonstrates that women over 55 exhibit a heightened susceptibility to MetS and CVD, with risk factors substantially increasing during the postmenopausal phase (Mesh *et al.*, 2006; Lejsková *et al.*, 2011; Ghosh *et al.*, 2022a). The hormonal landscape transformation, characterized by declining estrogen levels and altered estrogen-testosterone ratios, has been directly linked to MetS development during menopausal transition (Mesh *et al.*, 2008; Janssen *et al.*, 2008). Beyond hormonal fluctuations, aging itself contributes to the clustering of cardio-metabolic risk factors, creating ongoing scientific discourse about whether

Received on :24th July 2024; Revised version received on :29th November 2024; Accepted: 21st December 2024

increased MetS incidence stems from menopausal changes or natural aging processes (Casiglia *et al.*, 1996). As menopause represents a universal female physiological experience (Jiang *et al.*, 2017), understanding its metabolic implications becomes crucial. Integrating contemporary research perspectives, Das *et al.* (2024) highlight the synergistic interactions of diet and genetic predispositions, while Shakil *et al.* (2024) underscore the importance of comprehensive nutritional assessment in understanding MetS vulnerability among elderly women, emphasizing the need for targeted, personalized risk prediction models. Artificial intelligence (AI) has become ubiquitous in many spheres of human civilization, including AI-powered services and professions that have been demonstrated to have significant negative effects on people, including in the healthcare industry (Ghosh *et al.*, 2024; Ghosh, 2024a & b). Given the growing health requirements of the population, especially the elderly, new approaches that maximize and improve the healthcare resources already in place are necessary to address the current staffing shortages in the global health and social system. The use of artificial intelligence (AI) in healthcare is presently the subject of extensive research, with a focus on prevention, diagnosis, novel drug discovery, and post-treatment support (Ghosh & Sanyal, 2024). The whole patient experience could undergo substantial changes because of these findings. It is crucial to remember, though, that AI is mostly used to treat conditions related to the heart, nervous system (including Parkinson's disease and stroke), and cancer (Jiang *et al.*, 2017; Ghosh, 2024a). Even if the technological advancements benefit people of all ages, it is critical to recognize that certain age groups may have unique health concerns. The use of a combined strategy that combines lifestyle changes and medication therapies effectively controls MetS, thereby reducing the risk of CVDs. The continued digital healthcare revolution has made the application and relevance of this specific technology increasingly apparent. The clinical care and research environment have seen major changes as a result of technological advancements including the widespread use of smart phones, wearable, embedded sensors, and the introduction of massive databases like electronic health records. AI technologies offer the capacity to produce insightful insights and dynamically evaluate complex data, hence enhancing treatment procedures and enhancing results. With the development and implementation of AI technology, metabolic healthcare and research could undergo a dramatic transformation using personalized predictions in clinical decision-making. More specifically, AI can help with the proactive and unbiased assessment of health problems, which can help with the diagnosis and delivery of therapy that is tailored to each patient's specific needs. This covers the administration of care as well as the supply of ongoing monitoring. Consequently, a system or model will make it easier to diagnose MetS in postmenopausal women. In addition to improving prevention, early detection will lower the cost of diagnosis. Given that India is frequently referred to as the "Diabetes Capital of the World," it is necessary for India. Considering this, we carried out the research to identify a machine learning-based model that postmenopausal women with MetS might utilize for their preliminary diagnostic. Research in India has not yet been conducted.

Methodology

Data and Study Design

This cross-sectional study merged and utilized two datasets to predict metabolic syndrome (MetS) risk among rural elderly women. The first dataset included 222 postmenopausal women aged 45–70 years, randomly selected from 30 villages in the Singur block, the rural field practice area of the All-India Institute of Hygiene and Public Health (AIIPH), West Bengal, India, from March 27, 2014, to August 1, 2016. The second dataset comprised 236 elderly women aged 60–70 years, randomly selected from Amdanga Block, North 24 Parganas District, West Bengal, India, collected between April 2014 and August 2018. Ethical clearance was obtained from the Ethics Committee of AIIPH, Kolkata, and informed written consent was secured (Ghosh *et al.*, 2020; Srimani *et al.*, 2017).

The dataset from Singur, comprising data on 222 postmenopausal women, was freely available online as part of the study published by Srimani *et al.*, 2017 and was utilized with appropriate acknowledgment. In contrast, the dataset from Amdanga, involving 236 elderly rural women, was collected as primary data by the author through field surveys and laboratory assessments, following

ethical guidelines. A preliminary report based on the Amdanga data has already been published, highlighting the findings and methodology (Ghosh *et al.*, 2020). The integration of these two datasets provided a robust framework for evaluating metabolic syndrome risk factors, ensuring a comprehensive analysis with both secondary and primary data sources. The dataset belongs to 7 attributes, which are mentioned in Table 1 in detail.

Table 1: Dataset Explanation (N=458)

Attribute Name	Description	Data Type
IDF MS	Presence of Metabolic Syndrome according to International Diabetes Federation (IDF),	Binary, Categorical
HDL	High Density Lipoprotein	Numeric, Continuous
TG	Triglyceride	Numeric, Continuous
FBS	Fasting Blood Sugar	Numeric, Continuous
WC	Waist Circumference	Numeric, Continuous
SBP	Systolic Blood Pressure	Numeric, Continuous
DBP	Diastolic Blood Pressure	Numeric, Continuous

Table 2: Distribution of post-menopausal women according to present and absent of metabolic syndrome (IDF MS) in relation to WC, FBS, HDL, TG, SBP & DBP (N=458)

Parameters	IDF MS		Total N (%)	Chi-square test (p)
	Yes N (%)	No N (%)		
WC (cm)				
<80	17 (0.09)	180 (99.91)	197 (100)	143.12 (0.00)
≥80	167(63.98)	94 (36.01)	261(100)	
FBS (mg/dl)				
<100	116 (31.35)	254 (68.64)	370 (100)	62.37 (0.00)
≥100	68 (77.27)	20(22.72)	88 (100)	
HDL (mg/dl)				
<50	69(26.04)	196 (73.96)	265(100)	52.29 (0.00)
≥50	115 (59.59)	78 (40.41)	193 (100)	
TG (mg/dl)				
<150	57(20.65)	219(79.35)	276 (100)	110.38 (0.00)
≥150	127 (69.78)	55 (30.22)	182 (100)	
DBP (mm of Hg)				
<85	47 (22.82)	159(77.18)	206 (100)	46.94 (0.00)
≥85	137(54.37)	115(45.63)	252(100)	
SBP (mm of Hg)				
<130	54 (22.31)	188(77.69)	242 (100)	68.10 (0.00)
≥130	130 (60.18)	86(39.81)	216 (100)	

Statistical Analysis

Every analysis was carried out on January 7, 2024. In the case of a normal distribution, continuous variables were represented by mean \pm standard deviation, while skewed distributions were represented by median and interquartile range (IQR). Percentages were used to express categorical variables. Operational definition and Selection of Predictors Systolic Blood Pressure (SBP), Diastolic Blood Pressure (DBP), Waist Circumference (WC), fasting blood glucose (FBG), serum triglyceride (TG) and high-density lipoprotein-cholesterol (HDL), were measured using standard procedure in both the project as mentioned (Ghosh *et al.*, 2022b; O'Keefe *et al.*, 2016; Diametra., 2019; Motamed *et al.*, 2015). Overnight fasting (10–12 h) blood specimens were collected early in the morning from the collection site for all biochemical estimations. The syndrome (MetS) occurs when a person has a combination of any three or more of the following metabolic factors: raised blood glucose or diabetes, high blood pressure, obesity, elevated triglycerides, and low levels of HDL cholesterol. Metabolic syndrome (MetS) was defined as per IDF, 2005 (for Asian Indians) criteria (Alberti *et al.*, 2005). Six factors were considered as potential predictors in our study as per IDF, 2005(for Asian Indians) criteria.

Using SPSS 22, the analysis above was conducted. Python 3.7 with the Scikit-learn software package finish the development and assessment of machine learning techniques. A bilateral p -value of less than 0.05 was deemed statistically significant.

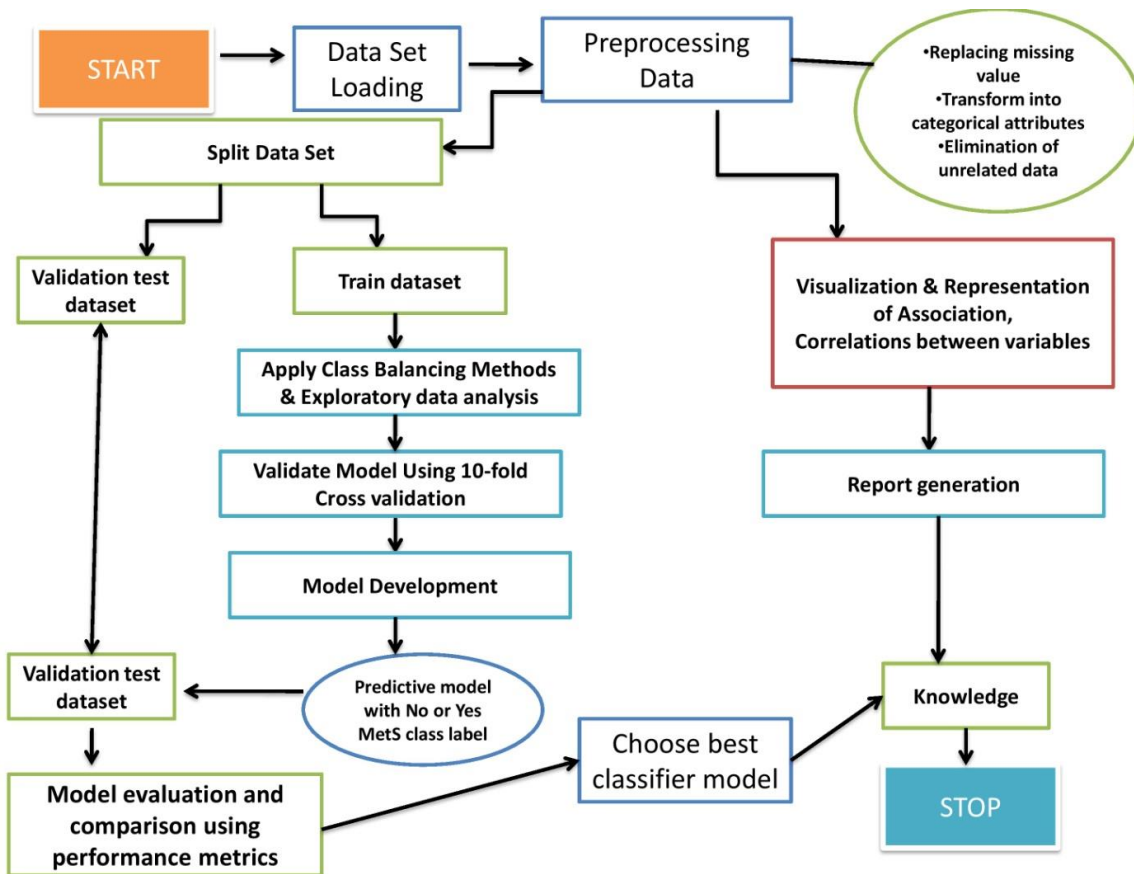


Figure 1: Framework of the MetS prediction model development

Model Construction and Evaluation

The research works technique is shown in figure 1. Its whole methodology is explained in the step-by-step order that follows. It offers a clear comprehension of the operational flowchart.

Prior to processing of data

Because of certain traits, a classifier may not always handle raw data efficiently (Ahmad et al., 2018). Sometimes raw data is noisy and incomplete. It can also be inconsistent at times (Ahmed et al., 2013). Raw data is highly susceptible to noise or inconsistency, which negatively impacts the analysis's progress and result. Thus, preparing data is crucial for data mining (Asaduzzaman et al., 2016). Preprocessing data is related to handling missing values, reducing dimension, selecting attributes, cleaning noisy and inconsistent data, etc. (Hasan et al., 2018). Therefore, preprocessing is crucial to prepare a dataset for machine learning and data mining so that better predictions can be made.

Compute Missing Values

First, Weka (Data Mining Tools, Version 3.8.3) loads the dataset. It is then noted that there are some missing values in the dataset. To replace missing values, a filter named Replace Missing Values under the attribute of the unsupervised filter was applied. The filter uses the mean or mode to replace missing data (Gimpy et al., 2014). The term "mean" refers to numerical figures. Conversely, the mode is employed to swap out nominal values. This is the most widely used filter for replacing missing values (Peng et al., 2005).

Finding and eliminating outliers

An outlier is a data point that, as other data points demonstrate, does not represent typical behaviours (Rahman et al., 2019). Outliers have an impact on machine learning or data mining processes, which

help predict outcomes and identify effective solutions for pertinent problems (Kdnuggets.com,2018). Therefore, removing outliers and extreme values from the dataset is a crucial step in the machine learning process to get better prediction results. Weka has a filter called Inter Quartile Range that is used to identify extreme values and outliers. The similar filter can be used to identify outliers and extreme values. The dataset is divided into three quartiles to identify outliers. $Q1$ denotes the first quartile, $Q2$ denotes the second, and $Q3$ denotes the third. The formula used to calculate the Interquartile Range (IQR) is $IQR = Q3 - Q1$. Next, using the following equations (Satu et al.,2017) the values of B_{min} , lower boundary, and upper boundary were determined.

$$B_{min} = Q1 - 1.5 * IQR \quad (1)$$

$$B_{max} = Q3 + 1.5 * IQR \quad (2)$$

The figure displayed in this case, which is less than or more than P , is referred to as an outlier. A value that is either extremely little or extremely large inside a dataset is regarded as an extreme value. Following the removal of extreme values and outliers, data from either the positive or negative case will be adjusted more than that of the other groups. The class attributes become unbalanced as a result and poor accuracy results. Therefore, if the dataset is unbalanced, it is imperative that it be balanced. To balance the unbalanced dataset, a filter known as Synthetic Minority Oversampling Technique (SMOTE) is applied.

Implementation of Different Machine Learning Models

In present work, ten prediction models were used: Random Forest, a popular machine learning technique; Logistic regression; KNeighbors Classifier; MLP Classifier; Support Vector Machines; Perceptron; Linear SVC; Stochastic Gradient Decent; Gaussian NB; Decision Tree. The training set is used to train all ten models, and the test set is used to evaluate each model's prediction accuracy. We divided our data into training (70% of the total sample) and test (30%) datasets at random for each analysis. Furthermore, we utilized the verification approach (Krstajic et al., 2014), which calls for repeated nested cross-validation, to prevent performance estimations from being optimistically biased (overfitting). A two-stage procedure is hosted by repeated nested cross-validation. In step 1, the model was chosen hyperparameters were changed to optimize the accuracy of the model as determined by a validation data set. In contrast to regular model parameters (such as weights in a regression model), hyperparameters have an external configuration whose value cannot be inferred from the data. Hyperparameter tuning in the training dataset involves deploying a standard methodology with 10-fold cross-validation, for example, to determine the degree of model complexity. This process involves adjusting the model's learning from the data. Figure 1 illustrates this process. To determine the final forecasting performance of these models, the optimal parameter from stage 1 was used in stage 2. The accuracy, sensitivity, precision, and AUC of the discrimination index were chosen to assess the proposed model's prediction performance using the test datasets. We also determined each classifier's proportional relevance based on the predictors' contributions to prediction accuracy (Miche et al.,2020; Probst et al., 2018).

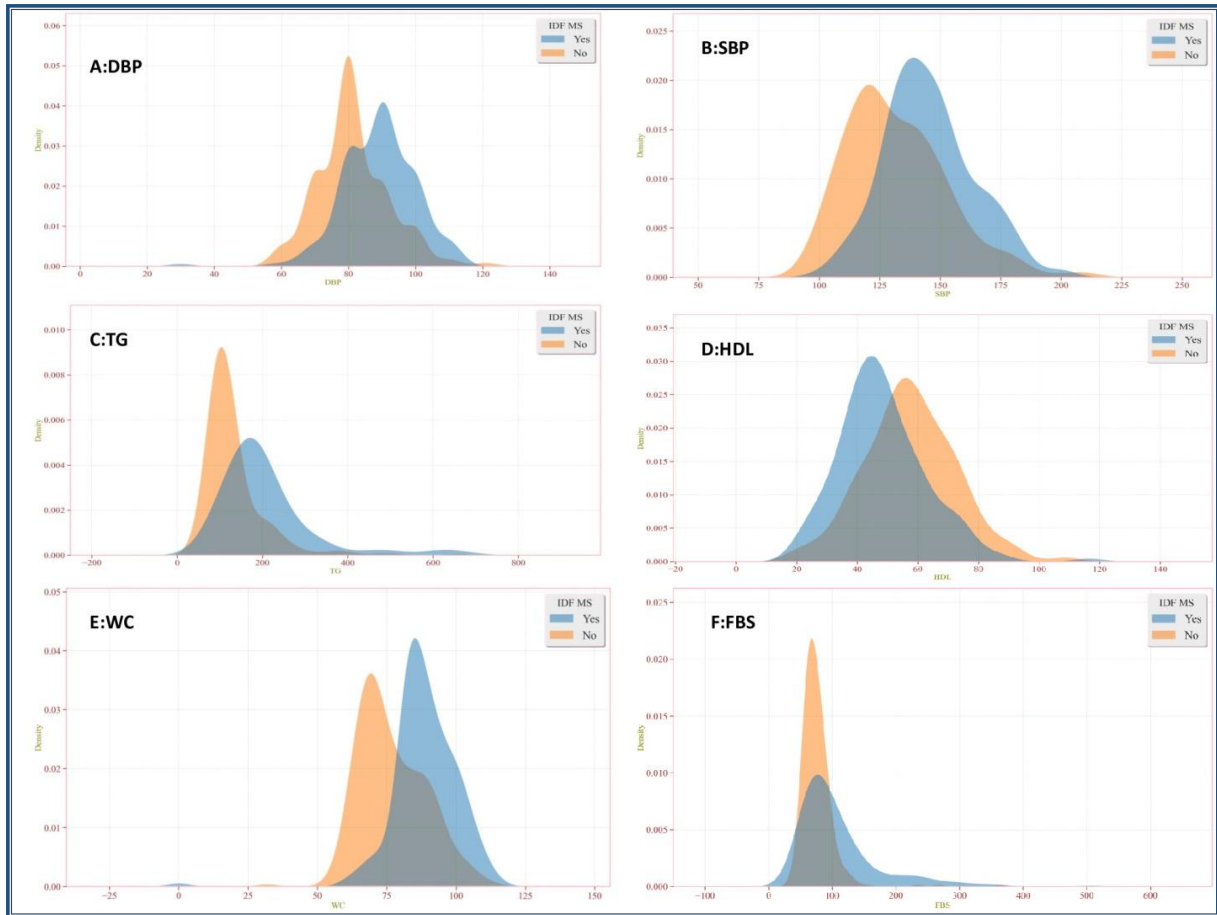


Figure 2: Kernel density estimate (KDE) figure for the participant density experiment based on MetS distribution (Yes, No) according to their (A)DBP, (B)SBP, (C)TG, (D)HDL, (E)WC, (F)FBS level

Results

The present study consists of 459 participant's data. The details descriptive analysis was shown in table 2. The prevalence of MetS was 40.17%, observed among post-menopausal women. Prevalence of hyperglycemic state remains 19.21%. Again, prevalence of low HDL-C was observed in 57.86% of respondents. Considering WC 56.98% population were at high risk, where considering TG level 39.73% were at high risk. Before the SMOTE filter was applied, figure 2 displays the kernel density estimate (KDE) figure for the participant density experiment based on the MetS distribution (Yes, No). Table 3 presents the results of several performance measures derived from the values of classification methods, including Random Forest, Logistic regression, Neighbors Classifier, MLP Classifier, Support Vector Machines, Perceptron, Linear SVC, Stochastic Gradient Decent, Gaussian NB, Decision Tree that were used to analyze the MetS dataset. In comparison to other classifiers, the appropriate classification method is assessed using these performance metrics. The best-performing classifier, according to Table 3, is Decision Tree, which has a 90.92% accuracy rate, 0.85 precision, and ROC as 0.895. Random forest is also showing significant result with 92.58% model accuracy score. The list of best fitting models, ranked by performance analysis score, is shown in figures 3 and 4. One of the key metrics used to assess a classification algorithm's performance is AUROC. Characteristics are mostly used to assess how well any classifier performs. Additionally, figure 5 displays the various Decision Tree classifier settings since it has the best performance rate out of all 10 models. The feature importance values for each used classification algorithm are shown in figure 6. With high values, it also indicates the most important risk factors. The most important risk factors are WC, HDL, TG, FBS SBP and DBP according to the figure.

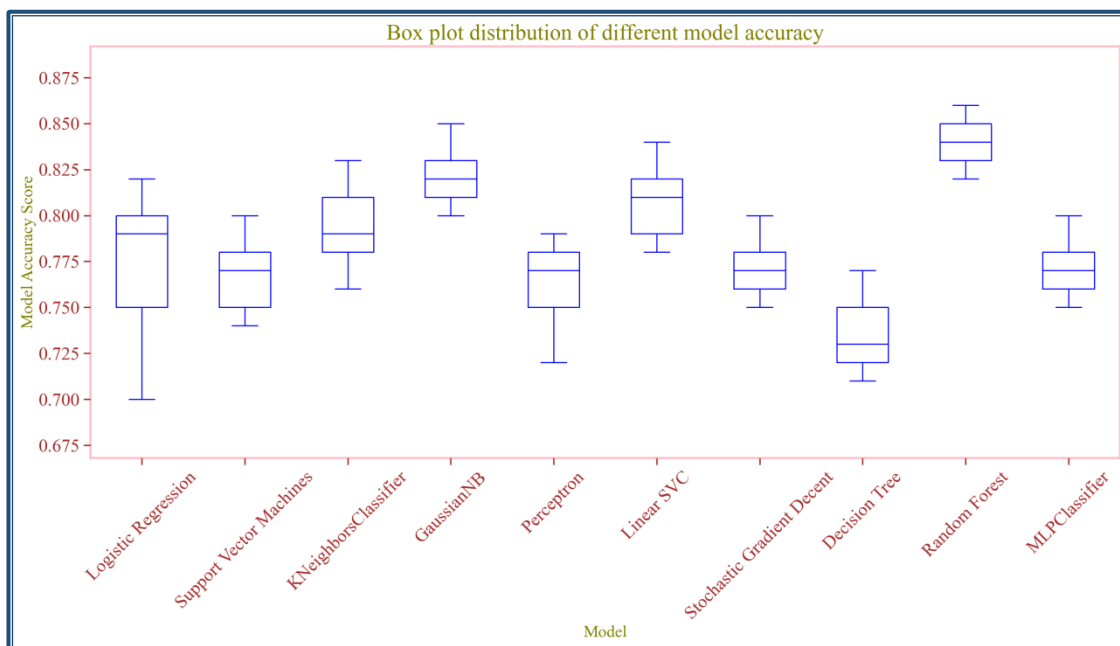


Figure 3: Box plot distribution of different model accuracy after application of SMOTE

The biggest risk variables are shown in table 4 & figure 6. The feature relevance score of each applied classification method identifies these risk variables. The most prevalent risk variables are gathered for additional machine learning examination. The table indicates that table 4 final columns include representations for WC, TG, FBS, HDL, and other parameters. According to the applied algorithms the most significant risk variables were WC, TG, and FBS.

Model	Training Accuracy	Model f1 Score	Model Accuracy Score
Decision Tree	100.00	86.15	90.22
Support Vector Machines	90.11	84.85	89.13
Random Forest	100.00	84.85	89.13
KNeighborsClassifier	92.03	81.25	86.96
Perceptron	85.16	81.16	85.87
MLPClassifier	87.91	80.60	85.87
Logistic Regression	84.89	79.41	84.78
GaussianNB	82.42	75.00	84.78
Linear SVC	87.36	80.00	84.78
Stochastic Gradient Decent	84.62	74.07	84.78

Figure 4: Overview of the performance estimates for each prediction model

Table 3: Application of different classification model and their performance in predicting MetS (After implementation of SMOTE and 10-fold Cross validation)

Classification Model Applied	Model Accuracy Score (Std. Dev) (%)	Kappa Statistics	Root Mean Squared Error	Relative Absolute Error	True Positive Rate	False Positive Rate	Precision	F-Measures	ROC Area
Random Forest	92.58(4.28)	0.811	0.294	0.000	0.91	0.093	0.852	0.878	0.911
Logistic Regression	84.08(6.32)	0.674	0.390	0.001	0.85	0.156	0.75	0.794	0.847
K Neighbors Classifier	84.89(4.98)	0.712	0.361	0.001	0.90	0.187	0.812	0.812	0.856
MLP Classifier	86.25(6.57)	0.695	0.375	0.001	0.84	0.156	0.771	0.805	0.855
Support Vector Machines	88.47(5.18)	0.763	0.329	0.001	0.90	0.125	0.823	0.848	0.887
Perceptron,	77.26(9.47)	0.699	0.375	0.001	0.85	0.125	0.756	0.811	0.862
Linear SVC	85.70(6.89)	0.678	0.390	0.001	0.83	0.125	0.736	0.799	0.854
Stochastic Gradient Decent,	81.55(9.42)	0.749	0.329	0.001	0.96	0.250	0.923	0.827	0.858
Gaussian NB	81.61(5.44)	0.643	0.390	0.001	0.95	0.343	0.875	0.750	0.803
Decision Tree	90.92(3.51)	0.785	0.312	0.001	0.91	0.125	0.848	0.861	0.895

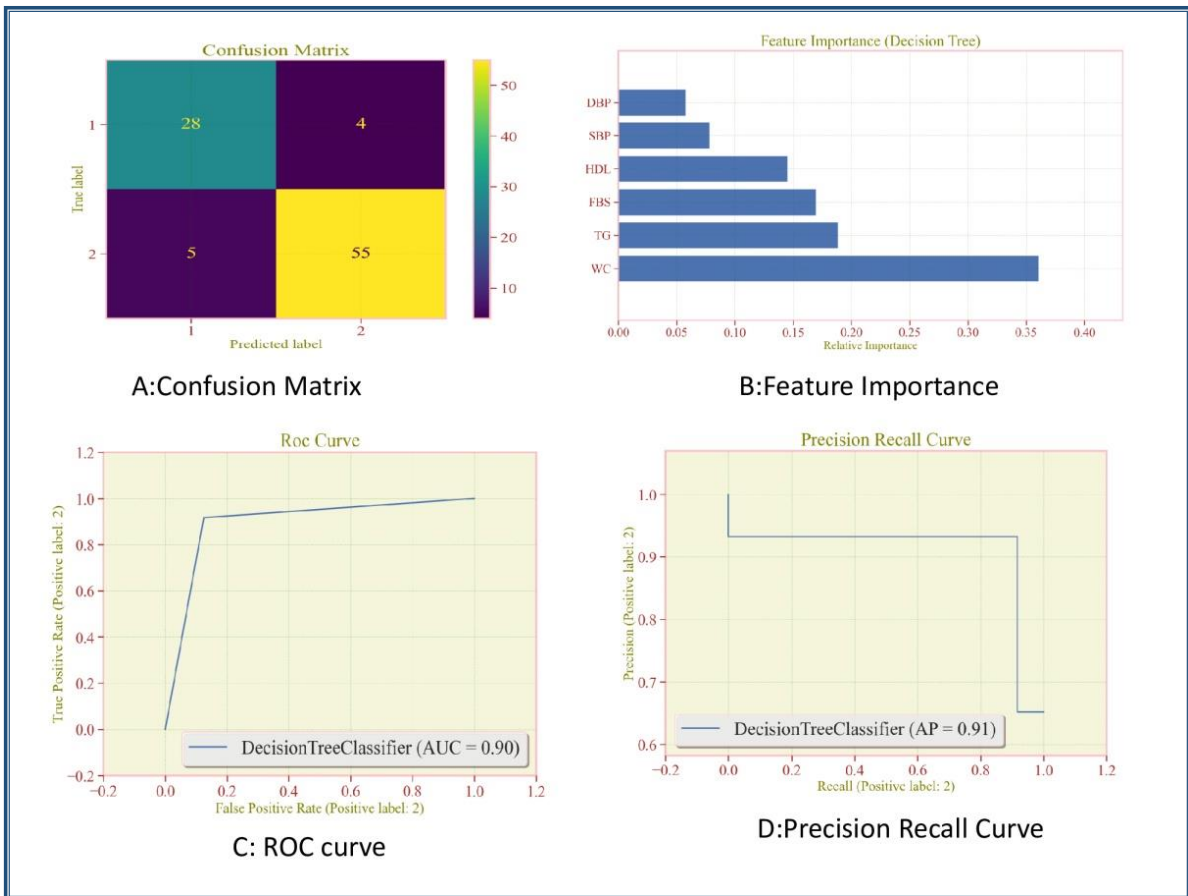


Figure 5: Different performance analysis of best fitted model Decision Tree Classifier according to accuracy comparison figure 4.

Table 4: Prediction of important risk factors for MetS among post-menopausal women according to different applied algorithm

Feature ranking	Random Forest	Decision Tree	MLP Classifiers	Linear SVC	Perceptron	Logistic Regression
1	WC	WC	TG	HDL	HDL	WC
2	TG	TG	FBS	DBP	DBP	TG
3	FBS	FBS	WC	SBP	SBP	FBS
4	HDL	HDL	HDL	WC	WC	SBP
5	SBP	SBP	SBP	FBS	FBS	DBP
6	DBP	DBP	DBP	TG	TG	HDL

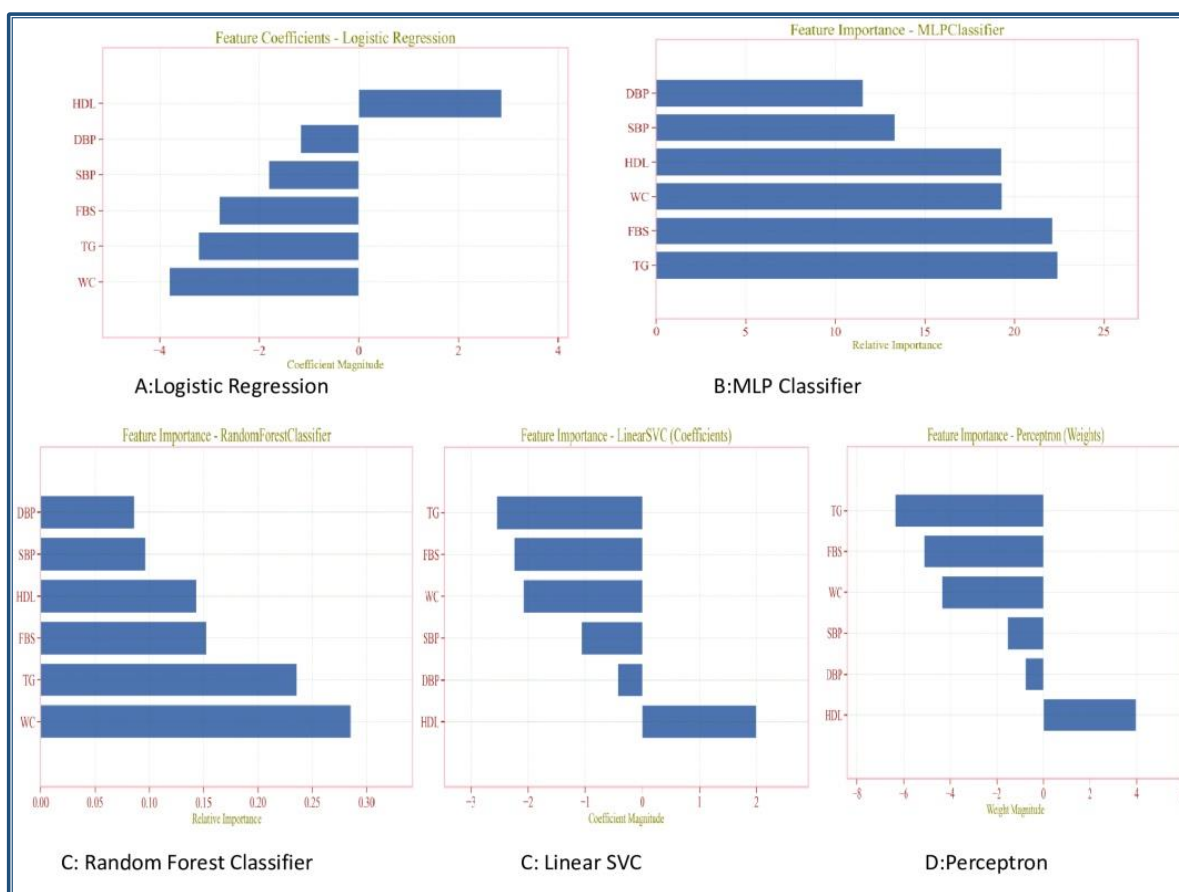


Figure 6: Feature importance score of the predictors according to different ML (Machine Learning) model applied

Discussion

Given the significance of Indian postmenopausal women to society, this is the first research of its kind. This section presents relevant works that use various datasets and use the ML models and approaches to forecast the incidence of metabolic syndrome as a reference point. Although their performance is not directly compared to the research results of the present submission, the outline of these studies aimed to (i) demonstrate the interest of the research community in this health condition, (ii) highlight the diversity in available datasets, particularly in the Indian context where IDF criteria is used to diagnose MetS, and (iii) identify the best-performing classifiers for metabolic syndrome risk prediction.

Specific Model Performance Studies

Gutiérrez-Esparza *et al.* (2020) ranked health factors, including clinical and anthropometric measurements, lifestyle information, and blood tests, from a dataset in Mexico City based on the National Cholesterol Education Program Third Adult Treatment Panel (ATP III) criteria. The random forest model achieved the highest performance in identifying abdominal obesity among individuals

with MetS, with sensitivity and specificity both at 0.93. Similarly, Karimi-Alavijeh *et al.* (2016) used two machine learning models, a decision tree and an SVM, to predict MetS incidence based on ATP III criteria. The SVM model demonstrated sensitivity, specificity, and accuracy of 0.774, 0.758, and 0.757, respectively, while the decision tree achieved 0.72 in sensitivity and 0.74 in specificity.

Diverse Datasets and Model Performances

Tavares *et al.* (2022) analyzed data from 17,182 adult participants in an annual screening program spanning 17 years, covering 37,999 visit pairs. For MetS prediction, the Light Gradient Boosting Machine (LGBM) model demonstrated superior performance with a sensitivity of 0.878, specificity of 0.702, and an AUC of 0.86. Similarly, Yu *et al.* (2020) evaluated the effectiveness of various decision tree-based machine learning algorithms in predicting MetS prevalence among self-paying patients using the FibroScan ultrasound instrument, where the Random Forest model achieved an AUC of 0.904. Furthermore, Choe *et al.* (2018) developed a MetS prediction model based on the genetic and clinical characteristics of non-obese Koreans. Among the models tested, the Naïve Bayes model performed best, with an AUC of 0.69, specificity of 0.80, and sensitivity of 0.42.

Non-Invasive and Ensemble Approaches

Datta *et al.* (2019) proposed an ML-based approach utilizing only non-invasive features for the early detection of MetS. Their ensemble classifier achieved an AUC of up to 0.90. Similarly, Cheong *et al.* (2015) investigated the predictive capabilities of BMI, waist circumference, and waist-to-hip ratio for identifying two or more non-adipose components of MetS, including high blood pressure, hypertriglyceridemia, low HDL-C, and high fasting plasma glucose. Receiver Operating Characteristic (ROC) curve analysis was used to evaluate the discriminative power of each anthropometric index, with AUC values indicating their effectiveness in distinguishing MetS from non-MetS cases.

Additional Modelling Techniques

Hosseini-Esfahani *et al.* (2021) utilized data-mining techniques to identify and prioritize key dietary and non-nutritional factors influencing the development of MetS. Their findings highlighted the Random Forest model's exceptional predictive ability, achieving a sensitivity of 0.97. In other studies, the eXtreme Gradient Boosting (XGBoost) model demonstrated strong performance, with AUC values of 0.88 (Lee *et al.*, 2017) and 0.93 (Li *et al.*, 2018). Logistic Regression was also employed in MetS prediction, yielding AUC values of 0.817 (Lee *et al.*, 2017) and 0.813 (Yang *et al.*, 2022). Additionally, Zou *et al.* (2018) developed and validated a new MetS risk score for forecasting MetS risk over three years, achieving an AUC of 0.68 with their machine learning model.

The development of metabolic syndrome is predisposed by several circumstances rather than being a distinct illness. It has a strong correlation with obesity, higher body weight, and sedentary behaviours. One of the study's limitations is that the dataset only includes the IDF parameters for MetS detection—HDL, WC, FBS, SBP, DBP, and TG. It is possible that adding demographic variables and additional anthropometric data to the model would have enhanced its MetS prediction accuracy. A smaller sample size may also be a contributing factor. Nevertheless, one important finding in this research is that it is one of the first studies of its kind to examine how an ML model may be used to predict MetS, as the IDF recommended five criteria are used to detect MetS in India. Because the illness progresses with time and becomes increasingly dangerous, early detection and prevention are crucial. Prioritizing early diagnosis and low-cost prevention is crucial since economic conditions play a significant role in the lives of this most disadvantaged segment of the Indian society. Clinicians will be able to anticipate MetS in Indian postmenopausal women with the use of this study. Clinicians may find it difficult to predict MetS in postmenopausal women because of milder symptoms. It is critical to diagnose MetS in postmenopausal women since it can lead to severe psychological and physical problems. In light of this, the study will be essential in predicting MetS and ensuring post-menopausal women's quality of life. Web apps, mobile applications, and other software can all use the suggested paradigm. With immediate diagnosis, the expense will be decreased.

Future work and limitations

This study explored predicting the risk of metabolic syndrome (MetS) using various machine learning models with and without SMOTE and 10-fold cross-validation. Models such as Random Forest, Logistic Regression, KNN, MLP, SVM, and Decision Tree were evaluated based on accuracy, precision, recall, F1 score, and AUC. Random Forest and Decision Tree performed best, achieving 90–92.58% accuracy using SMOTE. However, the exclusion of demographic factors, BMI, and waist-to-hip ratio limited prediction accuracy (Cheong *et al.*, 2015).

Future work will enhance explainability using Partial Dependency Plots (PDP) and Individual Conditional Expectation (ICE) to assess feature impacts (Molnar *et al.*, 2022). Dimensionality reduction techniques like t-SNE will also be applied to refine input features and improve classifier performance (Devassy *et al.*, 2020). An automated feature selection process and deep learning approaches will further expand the framework, enabling broader applicability and improved accuracy in MetS prediction.

Conclusion

Present research investigated metabolic syndrome (MetS) risk prediction using various machine learning models with SMOTE and 10-fold cross-validation. The experimental results demonstrated that Random Forest and Decision Tree models achieved the highest performance, with accuracy ranging from 90-92.58%. Future work will focus on enhancing feature selection, incorporating demographic factors, and exploring advanced machine learning techniques to improve predictive accuracy.

Conflict of Interest

Author declares no conflict of interests.

Acknowledgment

Author expresses her sincere gratitude and regards to all the researchers and professionals working tirelessly in the field of aging and nutrition, whose contributions have paved the way for advancements in understanding and addressing the complex nutritional needs of the elderly population.

References

- Ahmad, I., Basher, M., Iqbal, M. J., & Rahim, A. (2018). Performance comparison of support vector machine, random forest, and extreme learning machine for intrusion detection. *IEEE Access*, 6, 33789-33795. <https://doi.org/10.1109/ACCESS.2018.2841987>
- Ahmed, K., Emran, A. A., Jesmin, T., Mukti, R. F., Rahman, M. Z., & Ahmed, F. (2013). Early detection of lung cancer risk using data mining. *Asian Pacific Journal of Cancer Prevention*, 14(2), 595-598. <https://doi.org/10.7314/apjcp.2013.14.1.595>
- Alberti, K. G., Zimmet, P., Shaw, J., & IDF Epidemiology Task Force Consensus Group. (2005). The metabolic syndrome—a new worldwide definition. *The Lancet*, 366(9491), 1059-1062. [https://doi.org/10.1016/s0140-6736\(05\)67402-8](https://doi.org/10.1016/s0140-6736(05)67402-8)
- Asaduzzaman, S., Chakraborty, S., Hossain, M. G., Bashar, M. I., Bhuiyan, T., Chandan, S. S., Ahmed, K., & Paul, B. K. (2016). Hazardous consequences of polygamy, contraceptives and number of children on cervical cancer in a low income country: Bangladesh. *Cumhuriyet Üniversitesi Fen-Edebiyat Fakültesi Fen Bilimleri Dergisi*, 37(1), 74-84. <https://doi.org/https://dergipark.org.tr/tr/download/article-file/230736>
- Casiglia, E., d'Este, D., Ginocchio, G., Colangeli, G., Onesto, C., Tramontin, P., Ambrosio, G. B., & Pessina, A. C. (1996). Lack of influence of menopause on blood pressure and cardiovascular risk profile: A 16-year longitudinal study concerning a cohort of 568 women. *Journal of Hypertension*, 14(6), 729-736. <https://doi.org/10.1097/00004872-199606000-00008>
- Cheong, K. C., Ghazali, S. M., Hock, L. K., Subenthiran, S., Huey, T. C., Kuay, L. K., Mustapha, F. I., Yusoff, A. F., & Mustafa, A. N. (2015). The discriminative ability of waist circumference, body mass index and waist-to-hip ratio in identifying metabolic syndrome: Variations by age, sex and race. *Diabetes & Metabolic Syndrome: Clinical Research & Reviews*, 9(2), 74-78. <https://doi.org/10.1016/j.dsx.2015.02.006>

- Choe, E. K., Rhee, H., Lee, S., Shin, E., Oh, S. W., Lee, J. E., & Choi, S. H. (2018). Metabolic syndrome prediction using machine learning models with genetic and clinical information from a nonobese healthy population. *Genomics & Informatics*, 16(4), e31. <https://doi.org/10.5808/gi.2018.16.4.e31>
- Das, P., Banka, R., Ghosh, J., Singh, K., Choudhury, S. R., & Koner, S. (2024). "Synergism of Diet, Genetics, and Microbiome on Health. In S. Patnaik, A. Hamad, D. Paul, P. Dutta, & M. Shafiq (Eds.), *Nutrition Controversies and Advances in Autoimmune Disease*. IGI Global Scientific Publishing 131-189. <https://doi.org/10.4018/979-8-3693-5528-2.ch006>
- Datta, S., Schraplau, A., Da Cruz, H. F., Sachs, J. P., Mayer, F., & Böttinger, E. (2019). A machine learning approach for non-invasive diagnosis of metabolic syndrome. In Proceedings of the 2019 IEEE 19th International Conference on Bioinformatics and Bioengineering (BIBE) (pp. 933-940). IEEE. <http://dx.doi.org/10.1109/BIBE.2019.00175>
- Devassy, B. M., & George, S. (2020). Dimensionality reduction and visualisation of hyperspectral ink data using t-SNE. *Forensic Science International*, 311. <https://doi.org/10.1016/j.forsciint.2020.110194>
- Ghosh J., & Sanyal P.(2024) Development and Evaluation of Machine Learning Models for Predicting Constipation and Its Risk Factors Among College-Aged Females. *Nutrition and Food Science*, 12(3). Available at: <https://bit.ly/3MPo6eH>.
- Ghosh, J. (2023). A review on understanding the risk factors for coronary heart disease in Indian college students. *International Journal of Noncommunicable Diseases*, 8(3), 117. https://doi.org/10.4103/incd.incd_68_23
- Ghosh, J. (2024a). Recognizing and predicting the risk of malnutrition in the elderly using artificial intelligence: A systematic review. *International Journal of Advancement in Life Sciences Research*, 7(3). <https://doi.org/10.31632/ijalsr.2024.v07i03.001>
- Ghosh, J. (2024b). Leveraging data science for personalized nutrition. In Nutrition controversies and advances in autoimmune disease (pp. 572-605). IGI Global. <https://www.igi-global.com/chapter/leveraging-data-science-for-personalized-nutrition/353809>
- Ghosh, J., Chaudhuri, D., Saha, I., & Chaudhuri, A. N. (2020). Prevalence of metabolic syndrome, vitamin D level, and their association among elderly women in a rural community of West Bengal, India. *Medical Journal of Dr. DY Patil Vidyapeeth*, 13(4), 315-320. <https://journals.lww.com/mjdy/toc/2020/13040>
- Ghosh, J., Chaudhuri, D., Saha, I., & Nag Chaudhuri, A. (2022b). Association of conicity index with different cardiovascular disease risk factors among the rural elderly women of West Bengal, India. *Indian Journal of Community Medicine*, 47, 18–22. <https://pubmed.ncbi.nlm.nih.gov/articles/PMC8971865/>
- Ghosh, J., Singh, K., Choudhury, S. R., Koner, S., Basu, N., & Maity, S. (2022a). Impact of diet and nutrition on memory T cell development, maintenance, and function in the context of a healthy immune system. *Acta Scientific Nutritional Health*, 6(8), 142–154. <http://dx.doi.org/10.31080/ASNH.2022.06.1108>
- Gimpy, M. D. R. V. (2014). Missing value imputation in multi attribute data set. *International Journal of Computer Science and Information Technology*, 5(4), 1-7. <https://api.semanticscholar.org/CorpusID:14050728>
- Gutiérrez-Esparza, G. O., Infante Vázquez, O., Vallejo, M., & Hernández-Torruco, J. (2020). Prediction of metabolic syndrome in a Mexican population applying machine learning algorithms. *Symmetry*, 12(4). <https://doi.org/10.3390/sym12040581>
- Hasan, S. M. M., Mamun, M. A., Uddin, M. P., & Hossain, M. A. (2018). Comparative analysis of classification approaches for heart disease prediction. In 2018 International Conference on Computer, Communication, Chemical, Material and Electronic Engineering (IC4ME2) (pp. 1-4). IEEE. <https://api.semanticscholar.org/CorpusID:52299844>
- Hosseini-Esfahani, F., Alafchi, B., Cheraghi, Z., Doosti-Irani, A., Mirmiran, P., Khalili, D., & Azizi, F. (2021). Using machine learning techniques to predict factors contributing to the incidence of metabolic syndrome in Tehran: Cohort study. *JMIR Public Health and Surveillance*, 7(4). <https://doi.org/10.2196/27304>
- Janssen, I., Powell, L. H., Crawford, S., Lasley, B., & Sutton-Tyrrell, K. (2008). Menopause and the metabolic syndrome: The Study of Women's Health Across the Nation. *Archives of Internal Medicine*, 168(14), 1568-1575. <https://doi.org/10.1001/archinte.168.14.1568>
- Jiang, F., Jiang, Y., Zhi, H., Dong, Y., Li, H., Ma, S., Wang, Y., Dong, Q., Shen, H., & Wang, Y. (2017). Artificial intelligence in healthcare: Past, present and future. *Stroke and Vascular Neurology*, 2(4), 230-243. <https://doi.org/10.1136/svn-2017-000101>
- Karimi-Alavijeh, F., Jalili, S., & Sadeghi, M. (2016). Predicting metabolic syndrome using decision tree and support vector machine methods. *ARYA Atherosclerosis*, 12(3), 146-152. <https://pubmed.ncbi.nlm.nih.gov/27752272/>
- Krstajic, D., Buturovic, L. J., Leahy, D. E., & Thomas, S. (2014). Cross-validation pitfalls when selecting and assessing regression and classification models. *Journal of Cheminformatics*, 6(1), 10. <https://doi.org/10.1186/1758-2946-6-10>

- Lee, S., Lee, H., Choi, J. R., & Koh, S. B. (2020). Development and validation of prediction model for risk reduction of metabolic syndrome by body weight control: A prospective population-based study. *Scientific Reports*, 10(1), 10006. <https://doi.org/10.1038/s41598-020-67238-5>
- Lee, S., Lee, S. K., Kim, J. Y., Cho, N., & Shin, C. (2017). Sasang constitutional types for the risk prediction of metabolic syndrome: A 14-year longitudinal prospective cohort study. *BMC Complementary and Alternative Medicine*, 17(1), 438. <https://doi.org/10.1186/s12906-017-1936-4>
- Lejsková, M., Alušik, S., Suchánek, M., Zecová, S., & Piha, J. (2011). Menopause: Clustering of metabolic syndrome components and population changes in insulin resistance. *Climacteric*, 13(1), 83-91. <https://doi.org/10.3109/13697131003692745>
- Li, G., Esangbedo, I. C., Xu, L., Fu, J., Li, L., Feng, D., Han, L., Xiao, X., Li, M., Mi, J., Li, M., Gao, S., & Hou, Y. (2018). Childhood retinol-binding protein 4 (RBP4) levels predicting the 10-year risk of insulin resistance and metabolic syndrome: The BCAMS study. *Cardiovascular Diabetology*, 17(1), 69. <https://pubmed.ncbi.nlm.nih.gov/29759068/>
- Mesch, V. R., Boero, L. E., Siseles, N. O., Royer, M., Prada, M., Sayegh, F., Schreier, L., Benencia, H. J., & Berg, G. A. (2006). Metabolic syndrome throughout the menopausal transition: Influence of age and menopausal status. *Climacteric*, 9(1), 40-48. <https://doi.org/10.1080/13697130500487331>
- Mesch, V. R., Siseles, N. O., Maidana, P. N., Boero, L. E., Sayegh, F., Prada, M., Royer, M., Schreier, L., Benencia, H. J., & Berg, G. A. (2008). Androgens in relationship to cardiovascular risk factors in the menopausal transition. *Climacteric*, 11(6), 509-517. <https://doi.org/10.1080/13697130802416640>
- Miche, M., Studerus, E., Meyer, A. H., Gloster, A. T., Beesdo-Baum, K., Wittchen, H. U., & Lieb, R. (2020). Prospective prediction of suicide attempts in community adolescents and young adults, using regression methods and machine learning. *Journal of Affective Disorders*, 265, 570-578. <https://doi.org/10.1016/j.jad.2019.11.093>
- Molnar, C., König, G., Herbringer, J., Freiesleben, T., Dandl, S., Scholbeck, C. A., Casalicchio, G., Grosse-Wentrup, M., & Bischl, B. (2022). General pitfalls of model-agnostic interpretation methods for machine learning models. In *Proceedings of the xxAI-Beyond Explainable AI: International Workshop, Held in Conjunction with ICML 2020* (pp. 39-68). Springer. [General Pitfalls of Model-Agnostic Interpretation Methods for Machine Learning Models | SpringerLink](https://doi.org/10.1007/978-3-030-84166-4_3)
- Motamed, N., Perumal, D., Zamani, F., Ashrafi, H., Haghjoo, M., Saeedian, F. S., Maadi, M., Akhavan-Niaki, H., Rabiee, B., & Asouri, M. (2015). Conicity index and waist-to-hip ratio are superior obesity indices in predicting 10-year cardiovascular risk among men and women. *Clinical Cardiology*, 38(9), 527-534. <https://doi.org/10.1002/clc.22437>
- O'Keefe, E. L., DiNicolantonio, J. J., Patil, H., Helzberg, J. H., & Lavie, C. J. (2016). Lifestyle choices fuel epidemics of diabetes and cardiovascular disease among Asian Indians. *Progress in Cardiovascular Diseases*, 58(5), 505-513. <https://doi.org/10.1016/j.pcad.2015.08.010>
- Peng, L., & Lei, L. (2005). A review of missing data treatment methods. *Intelligent Information Management Systems and Technologies*, 1(3), 412-419.
- Probst, P., & Boulesteix, A. L. (2018). To tune or not to tune the number of trees in random forest. *Journal of Machine Learning Research*, 18, 1-18. [17-269.pdf](https://arxiv.org/abs/1702.06999)
- Rahman, M. R., Islam, T., Zaman, T., Shahjaman, M., Karim, M. R., Huq, F., Quinn, J. M., Holsinger, R. D., Gov, E., & Moni, M. A. (2019). Identification of molecular signatures and pathways to identify novel therapeutic targets in Alzheimer's disease: Insights from a systems biomedicine perspective. *Genomics*, 112(2), 1290-1299. <https://doi.org/10.1016/j.ygeno.2019.07.018>
- Sattar, N., Gaw, A., Scherbakova, O., Ford, I., O'Reilly, D. S., Haffner, S. M., Isles, C., Macfarlane, P. W., Packard, C. J., Cobbe, S. M., & Shepherd, J. (2003). Metabolic syndrome with and without C-reactive protein as a predictor of coronary heart disease and diabetes in the West of Scotland Coronary Prevention Study. *Circulation*, 108(4), 414-419. <https://doi.org/10.1161/01.cir.0000080897.52664.94>
- Satu, M. S., Atik, S. T., & Moni, M. A. (n.d.). A Novel Hybrid Machine Learning Model To Predict Diabetes Mellitus. https://doi.org/10.1007/978-981-15-3607-6_36
- Shakil S., Ghosh J., Singh K., Chaudhury S.R. (2024) Comparative analysis of nutritional status among institutionalized and community-dwelling elderly women and its association with mental health status and cognitive function. *J Fam Med Prim Care*, 13(8):3078-3083. DOI: 10.4103/jfmpc.jfmpc_1932_23. <https://pmc.ncbi.nlm.nih.gov/articles/PMC11368306/>
- Srimani, S., Saha, I., & Chaudhuri, D. (2017). Prevalence and association of metabolic syndrome and vitamin D deficiency among postmenopausal women in a rural block of West Bengal, India. *PLoS One*, 12(11), e0188331. <https://doi.org/10.1371/journal.pone.0188331>
- Tavares, L. D., Manoel, A., Donato, T. H. R., Cesena, F., Minanni, C. A., Kashiwagi, N. M., da Silva, L. P., Amaro, E., Jr., & Szlejf, C. (2022). Prediction of metabolic syndrome: A machine learning approach to help

primary prevention. *Diabetes Research and Clinical Practice*, 191, 110047. <https://doi.org/10.1016/j.diabres.2022.110047>

Yang, H., Yu, B., OUYang, P., Li, X., Lai, X., Zhang, G., & Zhang, H. (2022). Machine learning-aided risk prediction for metabolic syndrome based on 3 years study. *Scientific Reports*, 12(1), 2248. <https://doi.org/10.1038/s41598-022-06235-2>

Yu, C. S., Lin, Y. J., Lin, C. H., Wang, S. T., Lin, S. Y., Lin, S. H., Wu, J. L., & Chang, S. S. (2020). Predicting metabolic syndrome with machine learning models using a decision tree algorithm: Retrospective cohort study. *JMIR Medical Informatics*, 8(3), e17110. <https://doi.org/10.2196/17110>

Zou, T. T., Zhou, Y. J., Zhou, X. D., Liu, W. Y., Van Poucke, S., Wu, W. J., Zheng, J. N., Gu, X. M., Zhang, D. C., Zheng, M. H., Pan, X. Y., & Gao, F. (2018). MetS risk score: A clear scoring model to predict a 3-year risk for metabolic syndrome. *Hormone and Metabolic Research*, 50(9), 683-689. <https://doi.org/10.1055/a-0677-2720>