



Obesity Level Classification Based on Decision Tree and Naïve Bayes Classifiers

Salisu Garba^{1*}, Marzuk Abdullahi², Umar Abdullahi Umar³ and Nura Tijjani Wurnor⁴

^{1,2,3,4}Department of Computer Science, Sule Lamido University, Kafin Hausa

salisu.garba@slu.edu.ng¹, marzukabdullahi@gmail.com², umar.abdullahi@slu.edu.ng³ and nuruddeen.twurno@slu.edu.ng⁴

Abstract

This paper proposed an approach for obesity levels classification. The main contribution of this work is the use of boosting and bagging techniques in the decision tree (DT) and naïve Bayes (NB) classification model to improve the accuracy of obesity levels classification. This is achieved by introducing a boosting and bagging technique to further improve the recognition rate of obesity levels in the DT model, eliminating the correlated features, and eliminating the zero observations problem in the NB model. To validate the accuracy of the proposed approach, empirical evaluation was carried out using WEKA to determine the accuracy, precision, and recall. The results show that the DT classification model performs better in terms of accuracy and average precision. The proposed approach can help in software development for the classification of individuals with obesity.

Keywords: *Boosting and bagging technique, classification models, enhanced decision tree, naïve Bayes classifiers, obesity level.*

1. Introduction

Obesity is an anatomical condition characterized by an extreme growth of body fat. The obesity rate is increasing gradually; from prior research, obesity is a serious health disease around the globe. Some of the major complications caused by obesity include stroke, heart disease, diabetes, hypertension, dyslipidemia, and metabolic syndrome (Kim, Lim, & Kim, 2021). Many attributes such as height, level of sleep, calories intake, etc. are used to determine the level of obesity. Obesity can be classified into eight based on BMI Classification, that is insufficient weight, normal weight, overweight level i, overweight level ii, obesity type i, obesity type ii, and obesity type iii (Palechor & Manotas, 2019). The prediction of the level of obesity is very critical as early intervention programs and strategies that target the risks associated with each level of obesity can be developed. Moreover, the possible influencers of each level of obesity can also be identified (Kim et al., 2021).

According to Cheng et al., (2021), classification or prediction of obesity level is the process of assigning an individual or people to a predefined category according to biological (age, gender, height, weight) and social (meals, smoke, alcohol, e.t.c) attributes. For instance, an individual X can belong to any of the eight predefined categories. Before automated classification takes place in this modern technology, the traditional method has been used to analyze the individual parameters and manually classified by domain experts based on the biological and social attributes. Several classification techniques have been applied in text classification such as Decision Tree (DT) classifiers, Naïve Bayes (NB) probabilistic classifiers, Neural Network, K-Nearest Neighbours (KNN), Support Vector Machine (SVM), and Rocchio classifiers (Beunza et al., 2019; Dugan, Mukhopadhyay, Carroll, & Downs, 2015; Nimmala, 2019).

The most widely used methods for predicting obesity levels are DT (J48), and Bayesian networks (NB). This is due to their good performance and high accuracy in classifying obesity levels compared to SVM and KNN. DT is a high-performance classification algorithm that analysed data and uses a representation in a tree structure to gain a better insight into the data (Yadav & Thareja, 2019). NB is a classification technique based on the Bayes Theorem



with an assumption of independence among predictors. Naive Bayes classifier assumes that the presence of a particular feature in a class is unrelated to the presence of any other feature (Daud, Mohd Noor, Aljunid, Noordin, & Fahmi Teng, 2019; Priyanka & Ravikumar, 2017).

Although the usage of machine learning and other data mining techniques for classification has increased in the healthcare domain (obesity, covid-19, cancer e.t.c), it still has many challenges and limitations, especially in obesity prediction or classification (Kim et al., 2021). Several challenges are faced in the obesity level. Manual classification by domain experts is tedious and time-consuming. Clinical data are inherently noisy, though, and machine learning techniques provide more robust methods for handling missing and incorrectly recorded data. Besides, the high dimension of data is becoming a hindrance in obesity classification as it consists many of features of different sizes. However, the use of a few biological attributes may not always provide an accurate prediction of obesity (Palechor & Manotas, 2019). By using a broader range of biological and social attributes, the accuracy of the obesity level prediction model can be improved.

This study aims to classify the obesity level of individuals by categories, from underweight to obesity III level based on DT and NB classifiers. A recent, standard, and the publicly available dataset is used to validate the models and identify factors that could contribute to the obesity problems. DT is found to have higher accuracy in the classification of obesity levels. This can help in software development for the classification of individuals with obesity through integration between Weka, Java, and NetBeans. The structure of the research is: related works are discussed in Section 2, the description of the proposed approach is presented in Section 3, the experimental evaluations, and results obtained are discussed in Section 4, and finally in Section 5, the conclusions and future of the study are provided.

2. Related Works

Several research works have been explored using data mining techniques in healthcare and disease classification (Hossain, Mahmud, Hossin, Rashed, & Noori, 2018; Priyanka & Ravikumar, 2017). A limited number of publications exist that address using machine learning techniques to classify obesity (De-la-hoz-correa et al., 2019; Priyaa, Sathyapriya, & Arockiam, 2020). Machine learning is being used to examine large amounts of data to identify patterns and relationships that would otherwise go undetected (Cheng et al., 2021; Dugan et al., 2015; Nimmala, 2019). Some studies have applied these methods to build predictive models to understand the obesity problem.

(Priyaa et al., 2020) introduced an enhanced DT to improve the accuracy of the obesity level prediction model by considering unique health parameters such as calories consumed, pregnancy status, and bodybuilder which have proved to drastically improve the accuracy of the prediction model. In (De-la-hoz-correa et al., 2019), DTs (J48), NB, and Simple Logistic are used for the creation of the obesity level estimation model on which, software was built using the Weka library. (Dugan et al., 2015) employs ID3 (Iterative Dichotomiser 3) to generate a DT for the prediction of early childhood obesity based on data mined from the CHICA (Child Health Improvement through Computer Automation) system.

As has been previously reported in the literature, a comparison is used to ascertain the viability of classification techniques. For example, recent research by Maswadi *et al.*, (2021) shows that NB and DT classification algorithms are used to classify human activities such as: sitting, standing, walking, sitting down, and standing up which improved accuracy and effectiveness. Yadav and Thareja, (2019) compare the performance of naive Bayes which use an assumption that the pair of features being classified are independent and DT classification uses a divide and conquer by calculating the accuracy of each technique using the R language.

A closer look at the literature on the classification of the level of obesity, however, reveals a number of gaps and shortcomings. The dataset used in Daud et al., (2019) was based on three different types of data that were manually collected namely grocery data (food weight, quantity, description), demographic data (age, gender), and anthropometric data containing body height, weight, and physical activity level. The work was limited to the compatibility of only one algorithm and also the calorie value was assigned based on an approximation formula. In modern healthcare systems, there is a need for high precision and accuracy. In this way, any error in classifying the various level of obesity could negatively impact the accuracy of intervention programs, and strategies that target

the risks associated with each level of obesity can be developed. Therefore, performance parameters such as precision and accuracy must be as close to 100% as possible for obesity level classification.

The works closest to ours is the use of machine learning techniques to predict childhood obesity (Priyaa et al., 2020) and software for obesity data classification using DT classifier (De-la-hoz-correa et al., 2019). These works do an excellent job of comparing the performance metrics of NB (NB) and DT (DT) classification techniques. However, childhood obesity classification or prediction is a half-finished solution. Moreover, these works could be improved by considering a different set of attributes for the model. (Dugan et al., 2015) (Priyaa et al., 2020) (De-la-hoz-correa et al., 2019) analyzed the CHICA dataset, which is limited to basic demographics (like sex) and biometrics (like height, weight, and BMI). By using a more standard and comprehensive dataset with a broader range of attributes, the model accuracy can be improved, and some light can be shed on what parameters are most influential in determining the level of obesity in both adults and children (Palechor & Manotas, 2019). Given the shortcomings of previous research, the present study aims to propose an obesity level class based on enhanced DT and NB classifiers.

3. The proposed Approach

In this paper, an obesity level classification based on enhanced DT and NB classifiers is proposed to improve the classification accuracy of obesity using the dataset for estimation of obesity levels based on eating habits and physical condition. The overall approach is shown in Figure 1. The proposed approach comprises three stages (Data Cleaning and Partitioning, Modelling and Testing, Performance Evaluation). The proposed enhanced DT and NB classifiers are described in Sections 3.1, and 3.2, respectively.

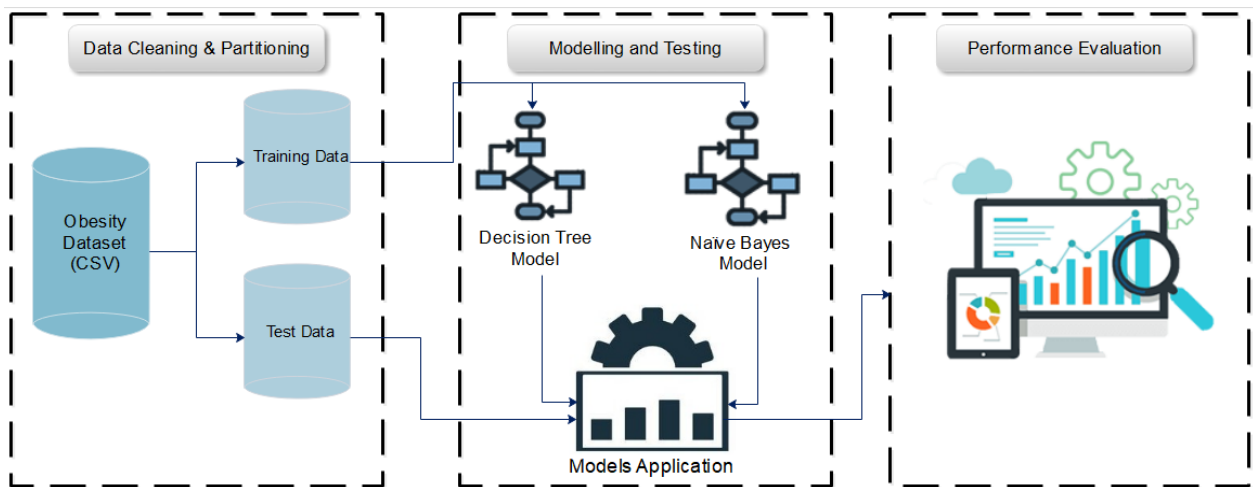


Figure 1. The three stages of the proposed obesity level classification approach

Data cleaning and Partition involve data cleaning, data transformation, and data reduction. In an effort to improve the accuracy of output, continuous data such as “age” are changed to discrete data, the number assigned to “Frequency of consumption of vegetables” is changed from “1, 2, 3” to LFC, FFC, and AFC where LFC= less, FFC= frequently and AFC= always. an alphabet. The pre-processing of the dataset described in Section 4.1 is carried out using Tableau Prep, Turbo Prep, and Rapid Miner.

The first step in the DT/NB Model development process is to read the cleaned dataset. Hence, cleaned data is loaded into the Rapid Miner for further analysis and building of the DT/NB Model. Instead of using the “Read CSV” operator, the data is loaded from a newly created repository. This data consists of 18 attributes and predictors are defined using 17 attributes out of which “NObesyedad” is defined independent variable as shown in Table 1. The second steps involve the use of a “split data” operator to partition the data into 2 (training and testing) to meet the requirements of machine learning modeling which requires a dataset to be broken down to support the accurate evaluation of the models later. Data is normally partitioned into 70% for training and 30% for testing.



3.1 Decision Tree Classification

The DT classifier is used to construct a model that predicts the value of a target variable, and the DT solves the problem by using the tree representation, where the leaf node corresponds to a class mark and attributes are expressed on the internal node of the tree. The ID3 (Iterative Dichotomiser 3) algorithm is widely used to generate a DT for classification or prediction. In this case, the ID3 is modified in execution efficiency and memory by introducing a boosting and bagging technique to further improve the recognition rate of the sample (obesity levels). Bagging requires that small changes to the training set should lead to different classifiers despite the instability while boosting requires that poor predictors produced in the training be accepted as long as their error on the given distribution is kept below 50%.

As such, the resulting DT model is smaller with higher accuracy, faster-running speed, and takes up less computer memory. These modifications also lower the complexity, easily while improving the adaptability of the model. Consequently, the main operator DT is used for the modeling process and the criterion 'gain ratio' is selected to reduce its bias on high-branch attributes. The gain ratio takes the number and size of branches into account when choosing an attribute. The last step is to apply the model to the dataset using the Apply Model operator and set the model performance check to the dataset using the performance operator. The empirical evaluation of the model is discussed in Section 4.

3.2 Naïve Bayes Classification

Naive Bayes classifier assumes that the absence or presence of a particular event from a group is not related or has nothing to do with the absence or presence of other events. The Naive Bayes classifiers do not offer an intrinsic method to evaluate feature importance. Naive Bayes methods work by determining the conditional and unconditional probabilities associated with the features and predicting the class with the highest probability. Thus, there are no coefficients computed or associated with the features used to train the model. Though Naive Bayes is a very simple and well-designed classification algorithm. However, it could be improved by eliminating issues that degrade the accuracy of the NB classification model.

Double counting leads to overestimating the importance of those features. So, correlated features are eliminated through feature selection based on probability, and the assumption of independent features is adhered to. Another issue that affects the accuracy of the NB classification model is the zero observations problem which normally occurs when the model comes across a categorical feature that is not present in the training set, the probability of 0 is assigned to that new category. This is very dangerous, as multiplying 0 with other features' probabilities will result in 0. Therefore, the zero observations problem is eliminated by using a smoothing parameter that assigns a very small probability estimate to such zero frequency occurrences to ensure that the probability value is never zero.

4. The proposed Approach

In this section, the experimental environment is described together with the datasets, and the evaluation metrics. The evaluation results for both DT and NB classifiers are presented and discussed. The experiments were carried out on a PC with Intel Core^(TM) CPU i5-5200U, @2.0 GHz processor (2 MB Cache, 800 MHz FSB), and 4GB RAM, running the Windows 10 Home Edition 64-bit OS. Tools used for cleaning the dataset (estimation of obesity levels based on eating habits and physical conditions) are Tableau Prep and Turbo Prep. The code for the basic versions of the DT and NB classifiers is adopted from RapidMiner, which is a data science software platform developed for data mining tasks. Rapid Miner contains tools for data pre-processing, classification, regression, clustering, and other predictive analytics.

4.1 Dataset Description

Given that Obesity is a problem that growing steadily and greatly influences the health of children, teenagers, and adults, an open and publicly accessible dataset entitled estimation of obesity levels based on eating habits and physical condition dataset was created (Palechor & Manotas, 2019). The performances of the enhanced DT and NB classifiers for obesity level prediction are tested on this dataset. This dataset includes data for the estimation of obesity levels in individuals from the countries of Mexico, Peru, and Colombia, with ages between 14 and 61, and diverse eating habits and physical conditions.

Dataset is stored in a comma-separated values file and the size of the dataset is 258KB, consisting of 17 attributes and 2111 records. There are no missing values in this dataset and all the data available are recorded in a consistent format. Table 1 describes the attributes of the datasets used in the experimental analysis. The attributes related to eating habits are FAVC, FCVC, NCP, CAEC, CH20, and CAL while the attributes related to the physical condition are SCC, FAF, TUE, MTRANS, and other variables obtained were gender, age, height, and weight. Class variable, NObesyedad was created using mass body index equation; $BMI = \text{Weight}/(\text{height}*\text{height})$ and compared with the information from WHO and Mexican Normativity. The categorization of obesity levels based on BMI calculated according to WHO and Mexican normativity is Underweight (< 18.5kg), Normal (18.5- 24.9kg), Overweight (25.0- 29.9kg), Obesity I (30.0-34.9kg), Obesity II (35.0- 39.9kg), and Obesity III (> 40.0kg).

Table 1. The attributes and definition of the estimation of obesity levels dataset

No.	Attributes	Definition	Data type
1	Gender	Gender of respondent	Categorical
2	Age	Age of respondent	Numerical (continuous)
3	Height	Height of respondent	Numerical (continuous)
4	Weight	Body weight of respondent	Numerical (continuous)
5	FHWOW	Family history with overweight	Categorical
6	FAVC	Frequent consumption of high caloric food	Categorical
7	FCVC	Frequency of consumption of vegetables	Categorical
8	NCP	Number of main meals	Numerical
9	CAEC	Consumption of food between meals	Categorical
10	SMOKE	Smoking behaviour of individuals	Categorical
11	CH2O	Consumption of water daily	Numerical (continuous)
12	SCC	Calories consumption monitoring	Categorical
13	FAF	Physical activity frequency	Numerical
14	TUE	Time using technology devices	Numerical
15	CALC	Consumption of alcohol	Categorical
16	MTRANS	Transportation used	Categorical
17	NObesyedad	Obesity level created based on mass body index	Categorical

4.2 Dataset Description

Performance of both DT and NB is measured by accuracy, precision, recall, and F- measure. Accuracy is the level of closeness between the prediction value and actual value of the model, the accuracy is calculated using Equation 1. Precision is the level of accuracy between the information requested by the user and the answer given by the model, the precision is calculated using Equation 2. A recall is the success rate of the model in rediscovering information, the recall is calculated using Equation 3. F- measure is the harmonic mean or weighted average of the precision and recall. The F-measure is calculated using Equation 4. There are four classifications that are used to measure the performance of NB and DT. These are True Positive (TP), False Positive (FP), False Negative (FN), and True Negative (TN).

$$\text{Accuracy} = \frac{TP+TN}{TP+TN+FP+FN} \quad (1)$$

$$\text{Precision} = \frac{TP}{(TP+FP)} \quad (2)$$

$$\text{Recall} = \frac{TP}{(TP+FN)} \quad (3)$$

$$\text{F-Measure} = 2 * \frac{\text{Precision} * \text{Recall}}{\text{Precision} + \text{Recall}} \quad (4)$$

4.2 Dataset Description

In this section, the results achieved by the analytical models are discussed based on the evaluation parameters such as accuracy, precision, recall, and F-measure. Moreover, a conclusion on which model is preferable for obesity level prediction (based on the comparison obtained) is provided with justification.

After the implementation of data pre-processing and analytical models (DT and NB), the results are obtained. Overall, the pre-processing is able to automatically remove noisy instances from training datasets for DT generation to avoid overfitting. It thus possesses more robustness and generalization capabilities. These machine learning algorithms are capable of identifying the most discriminative subset of attributes for classification. The evaluation results prove the efficiency of the proposed DT and NB algorithms for the classification of challenging real benchmark datasets (Estimation of obesity levels based on eating habits and physical conditions). Table 2 shows the summary of the results obtained from the implementation of the analytical models (DT and NB).

Table 2. The Precision and recall obtained from the DT and NB model

Class	Decision Tree		Naïve Bayes	
	Precision (%)	Recall (%)	Precision (%)	Recall (%)
Underweight	94.12	97.56	95.35	100.00
Normal	96.39	93.02	92.00	80.23
ObesityLevel1	99.01	95.24	99.00	94.29
ObesityLevel2	95.65	98.88	92.63	98.88
ObesityLevel3	97.98	100.00	98.97	98.97
Overweight	99.42	98.85	92.78	95.98

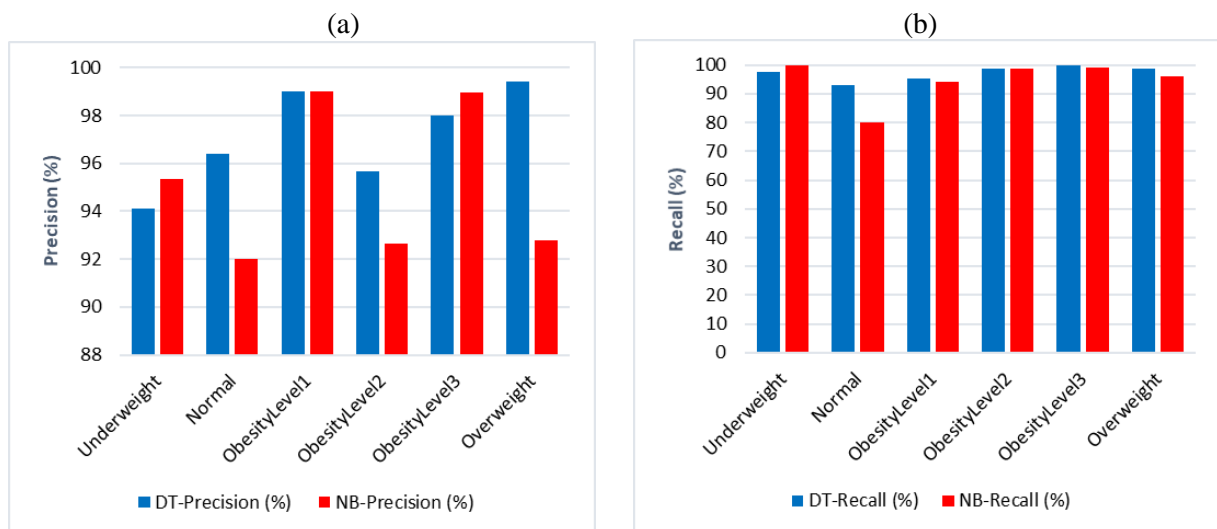


Figure 2. The precision and recall of the ND and DT models.

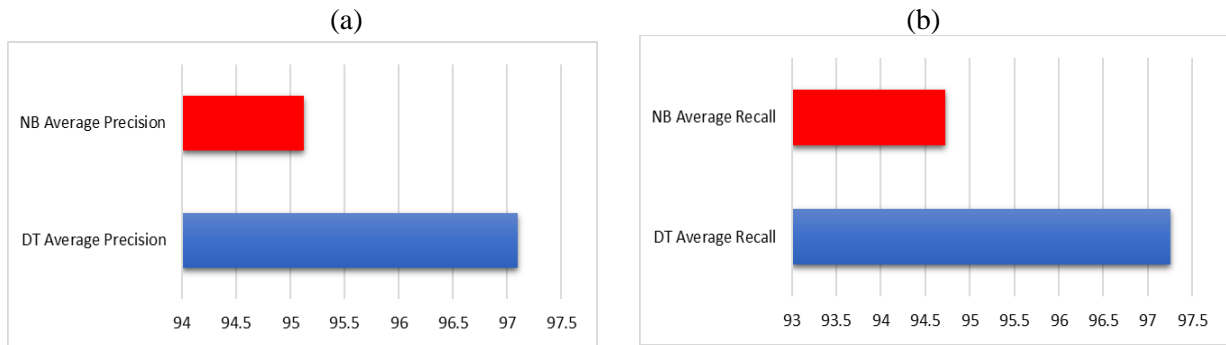


Figure 3. The average precision and recall of the ND and DT models.

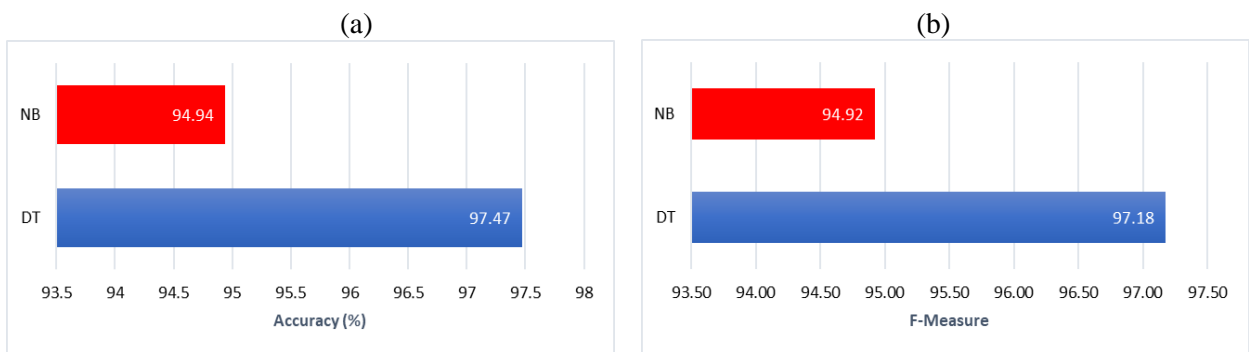


Figure 4. The accuracy and f-measure of the ND and DT models.

The performance of DT and NB machine learning algorithms can be evaluated based on three criteria: Accuracy, precision, and recall. The accuracy of the model by DT is 97.47 % while the accuracy of the model for NB is 94.94% as shown in Figure 4. Hence, it is observed that the DT gains higher accuracy compared to NB. The level of precision and recall is obtained for different classes of obesity in the DT model as shown in Figure 2, this is because of the number of attributes available for classification in each class after pre-processing of the data.

The value of average precision and average recall is displayed in Figure 3. High precision means the model is better in prediction while high recall means the model is good in predicting specific classes in actual positive data. From the results obtained, both the models (DT and NB) have a comparatively small difference in error rate, this is because of the level of performance of both models (DT and NB) in terms of accuracy, precision, and recall. Therefore, both models are efficient in the estimation of obesity levels.

From the results achieved by the analytical models in Section 3 and the evaluation of performance parameters, it is clear that DT has a better accuracy percentage which is 97.47 % than NB classification whose accuracy rate is 94.94%. There are also other studies such as Maswadi et al., (2021) that get DT accuracy much better than the DT's accuracy. The structure of the tree also gives insight into the strongest predictors of obesity. Many of the strongest predictors seen in the DT modeling of the datasets (Estimation of obesity levels based on eating habits and physical conditions) have been independently validated in the literature as correlated with obesity, thereby supporting the validity of the model.

NB relies on the conditional independence assumption, in addition to the ability to handle large, continuous, and discrete datasets in comparison to DT, it can work on both binary and multi-class classification problems. However, we can conclude that the most preferable model (based on the comparison obtained) is DT. This is because the model can handle any type of data whether the calculation has qualitative or quantitative values. DT is also non-parametric. So, there is no compulsion for independent variables to follow any kind of probability distribution.



5. Conclusion

This work presented an obesity level classification based on NB and DT classifiers to improve the classification accuracy using the boosting and bagging technique. The classification (NB and DT) models are improved by introducing a boosting and bagging technique to further augment the recognition rate of the sample (obesity levels) in the DT classification model, this is achieved by eliminating the correlated features through feature selection based on probability, and eliminating the zero observations problem through a smoothing parameter to ensure that the probability value is never zero in the NB classification model. The proposed approach leveraged the improved DT and NB models to provide improved accuracy and average precision of more than 97%. It can be concluded that the structure of the tree gives insight into the strongest predictors of obesity as have been independently validated in the literature as correlated with obesity, thereby supporting the validity of the model. The future work aims at exploiting deep neural networks and diverse factors for the obesity levels classification model to further improve the precision and accuracy as close to 100% as possible for obesity level classification.

Reference

- Beunza, J. J., Puertas, E., García-Ovejero, E., Villalba, G., Condes, E., Koleva, G., ... Landecho, M. F. (2019). Comparison of machine learning algorithms for clinical event prediction (risk of coronary heart disease). *Journal of Biomedical Informatics*, 97(July), 103257. <https://doi.org/10.1016/j.jbi.2019.103257>
- Cheng, X., Lin, S., Liu, J., Liu, S., Zhang, J., Nie, P., ... Xue, H. (2021). Does Physical Activity Predict Obesity — A Machine Learning and Statistical Method-Based Analysis, 1–11.
- Daud, N., Mohd Noor, N. L., Aljunid, S. A., Noordin, N., & Fahmi Teng, N. I. M. (2019). Predictive Analytics: The Application of J48 Algorithm on Grocery Data to Predict Obesity. In *2018 IEEE Conference on Big Data and Analytics, ICBDA 2018* (pp. 1–6). <https://doi.org/10.1109/ICBDAA.2018.8629623>
- De-la-hoz-correa, E., Mendoza-palechor, F. E., De-la-hoz-manotas, A., Morales-ortega, R. C., Hernández, S., & Adriana, B. (2019). Obesity Level Estimation Software based on Decision Trees. *Journal of Computer Science*, 15(1), 67–77. <https://doi.org/10.3844/jcssp.2019.67.77>
- Dugan, T. M., Mukhopadhyay, S., Carroll, A., & Downs, S. (2015). Machine learning techniques for prediction of early childhood obesity. *Applied Clinical Informatics*, 6(3), 506–520. <https://doi.org/10.4338/ACI-2015-03-RA-0036>
- Garba, S. (2021a). Self-adaptive mobile web service discovery approach based on modified negative selection algorithm. *Neural Computing and Applications*, 6, 1–23. <https://doi.org/10.1007/s00521-021-06486-6>
- Garba, S. (2021b). Social Media and Covid-19 Vaccine Hesitancy: A Data Story Telling Approach. In *1st National Conference of the School of Preliminary Studies, Sule Lamido University Kafin Hausa* (p. 0).
- Garba, S., Mohamad, R., & Saadon, N. A. (2020). Context Ontology for Smart Healthcare Systems. In *The 5th International Conference of Reliable Information and Communication Technology (IRICT 2020)* (pp. 199–206).
- Hossain, R., Mahmud, S. M. H., Hossain, A., Rashed, S., & Noori, H. (2018). Predicting Risk Factor of Obesity among Middle-Aged People Using Data Mining Techniques. In *Procedia Computer Science* (Vol. 132, pp. 1068–1076). Elsevier B.V. <https://doi.org/10.1016/j.procs.2018.05.022>



- Kim, H., Lim, D. H., & Kim, Y. (2021). Classification and Prediction on the Effects of Nutritional Intake on Overweight/Obesity, Dyslipidemia, Hypertension and Type 2 Diabetes Mellitus Using Deep Learning Model: 4--7th Korea National Health and Nutrition Examination Survey. *International Journal of Environmental Research and Public Health*, 18(11), 5597.
- Maswadi, K., Ghani, N. A., Hamid, S., & Rasheed, M. B. (2021). Human activity classification using Decision Tree and Naïve Bayes classifiers. *Multimedia Tools and Applications*, 80(14), 21709–21726. <https://doi.org/10.1007/s11042-020-10447-x>
- Muhammad, J., & Garba, S. (2019). Web-based Clinic Management System (CMS). *International Journal of Science and Engineering Applications*, 8(05), 131–135.
- Nimmala, S. (2019). Machine Learning Approaches to Classify Diabetes Patients based on Age, Obesity level and Cholesterol level. *CVR Journal of Science & Technology*, 16(1), 97–101. <https://doi.org/10.32377/cvrjst1617>
- Palechor, F. M., & Manotas, A. de la H. (2019). *Dataset for estimation of obesity levels based on eating habits and physical condition in individuals from Colombia, Peru and Mexico. Data in Brief* (Vol. 25). Elsevier Ltd. <https://doi.org/10.1016/j.dib.2019.104344>
- Priyaa, P. K., Sathyapriya, S., & Arockiam, L. (2020). An Enhanced Decision Tree Ensemble Technique For Obesity Prediction. *International Journal of Scientific & Technology Research*, 9(02).
- Priyanka, N., & Ravikumar, P. (2017). Usage of data mining techniques in predicting the heart diseases - Naïve Bayes & decision tree. In *Proceedings of IEEE International Conference on Circuit, Power and Computing Technologies, ICCPCT 2017* (pp. 1–7). <https://doi.org/10.1109/ICCPCT.2017.8074215>
- Yadav, K., & Thareja, R. (2019). Comparing the Performance of Naive Bayes And Decision Tree Classification Using R. *International Journal of Intelligent Systems and Applications*, 11(12), 11–19. <https://doi.org/10.5815/ijisa.2019.12.02>