

Diabetes Disease Prediction Model Deployment on Heroku-based Cloud Computing Platforms using Homogeneous Ensemble Machine Learning Algorithms

Belayneh Endalamaw Dejene*
belayneh.endalamaw@uog.edu.et

¹University of Gondar, Gondar-Ethiopia

Abstract

Background: Diabetes is a chronic disease that continues to be a major and global problem since it affects the health of the entire population. According to a study published in 2017 by the International Diabetes Federation, there were 425 million diabetics in the world at the time, with that number anticipated to climb to 625 million by 2045. This research aims to create a diabetes disease prediction model that uses homogenous ensemble machine learning methods and then deploys it on the cloud utilizing Heroku-based cloud computing.

Methods: This study was conducted by following a design science approach. The data were extracted from the github, and the data was collected based on the UCI repository and PIDD features. We have preprocessed to get quality data that are suitable for the machine learning algorithm to develop a model that predicts diabetes disease. Random forest, cat boost, extreme gradient boosting, and extra tree classifier were conducted to develop the predictive model with a total of 1393 instances with 8 features with the class label, and a training and testing dataset split ratio of 80/20 using random splitting methods. After developing the predictive model we have also evaluated using objective-based evaluation metrics like accuracy, precision, recall, f1-score, ROC-AUC, and K-fold cross-validation. Finally, we have deployed the best-performed model for the potential users on the Heroku-based cloud computing platform by designing the interface using HTML language as a front end and integrating it with the predictive model using the Flask framework.

Results: The overall accuracy of Random forest, cat boost, extra tree classifier, and extreme gradient boosting, is 87.81%, 90.32%, 87.81%, and 86.38%, respectively.

Conclusion: Finally, the researcher decided to use cat boost algorithms for further use in the development of artifacts, model deployment, risk factor analysis, and generating rules because it has registered better performance with 90.32% accuracy. The most determinant risk factors of diabetes were identified

using feature importance. Some of them are Glucose, BMI, and Diabetes pedigree function.

Keywords: Ensemble, Diabetes, Flask, Heroku, Machine learning

I. Background

According to [1] Diabetes mellitus is a collection of metabolic disease marked by hyperglycemia caused by insulin production, insulin action, or both. Diabetes is a disease caused by elevated blood glucose levels in the body [2]. Diabetes is a chronic condition defined by an increase in glucose or blood sugar levels caused by the body's inability or insufficiency to generate insulin, or insulin's inability to function on the organism's cells [3] [4]. Diabetes is a hormonal illness in which the body's failure to create insulin causes improper sugar metabolism in the body, elevating blood glucose levels in the body of a specific individual [5]. Certain risk factors, like age, BMI, glucose levels, blood pressure, and so on, play a significant role in the disease's impact [5]. It's a common chronic disease that puts people's health at risk [6]. Diabetes is defined by blood glucose levels that are greater than usual, which is caused by either faulty insulin secretion or its biological effects, or both [6]. Diabetes is a huge global concern because it has an impact on everyone's health [6]. According to [4] Diabetes is a metabolic disease that causes high blood sugar levels and a plethora of other complications such as stroke, kidney failure, and heart and nerve problems. It is becoming a highly common disease that affects several organs in the human body [7]. Diabetes Mellitus is one of the most serious diseases, and it affects a large number of people [8]. Diabetes mellitus can be caused by a variety of factors such as age, obesity, lack of exercise, inherited diabetes, lifestyle, poor food, high blood pressure, and so on [8]. Diabetes is one of the leading causes of death in underdeveloped and developing countries [4]. Diabetics are unable to convert carbohydrates consumed into glucose sugar, which provides energy as a fuel for daily activities [6]. As a result, glucose remains in the bloodstream and does not reach every cell in the body [6]. According to the World Health Organization (WHO), approximately 422 million people suffer with diabetes, primarily in low- and middle-income nations, and this number might rise to 490 billion by 2030 [2]. According to [3] Diabetes is a leading cause of divergence, and the biggest number of diabetics are found in adults over 65 years of age, due to genetic, family history, health, and

environmental factors in developed countries. While the majority of diabetics in developing nations, including our country, are between the ages of 45 and 64, type 2 diabetes has become more widespread in recent years among those between the ages of 30 and 40 [3]. Diabetes should not be removed, if it is untreated, it leads to serious complications such as heart-disease, kidney disease, stroke, blood pressure, eye damage, hypertension, cardiovascular disease, and it can also affect other organs of human body [2][6][8][9]. Diabetes is becoming more common over the world as a result of environmental and genetic causes [5]. The numbers are rapidly increasing due to a variety of factors, including bad meals, lack of physical activity, and many others [5]. Individuals and the government are investing money in research and programs to find a cure for this deadly disease [6]. To predict diabetes disease, several types of research have been conducted. For example, [1][3][4][6][8][10][11][12][13] and [14] investigated the prediction model which predicts diabetes disease using supervised machine learning algorithms. Besides this [7] tries to develop a predictive model with a relevant artifacts, that shows the performance of the model when the data were uploaded. Most of these studies, however, used Pima Indian diabetes dataset (PIDD), limited amount of datasets, only focused on developing predictive models, in existing method, the prediction accuracy is not so high. The previous studies were not generate relevant rules which is used for the decision makers to decide or develop policies and strategies about the prevention and control of diabetes disease, did not identify the most determinant risk factors of diabetes disease, doesn't develop a prototype or artifacts to interact potential users into the system. The previous studies weren't deploy the predictive model over the cloud for the potential users. This study, hence, aims to develop a predictive model that predicts the diabetes disease using homogeneous ensemble machine learning algorithms and deploy the predictive model using Heroku-based cloud computing platforms by investigating the following research questions: (1) Which homogeneous ensemble machine learning algorithms is suitable for predicting diabetes diseases? (2) What are the most determinant risk factors that influence the occurrence of diabetes disease? (3) What are the important rules that may use to develop strategies and policies towards preventing and/or reducing diabetes disease? The rest of this document is organized as follows: Section II presents related works, Section III discusses materials and methods used, Section IV mentions experimental setup and result discussion, and Section V presents the conclusion.

II. Related works

Several researchers have attempted to construct an

accurate diabetes prediction model over the years[1][4][7][8][10]. However, this subject still faces significant open research issues due to a lack of appropriate data sets and prediction approaches, which pushes researchers to use machine learning (ML)-based methods [4]. R. Krishnamoorthi et al.[4] Aimed to develop a novel diabetes healthcare disease prediction framework using machine learning techniques, the researcher uses PIDD data set for the prediction of diabetes. The researcher uses logistic regression, support vector machine, random forest, and k-nearest neighbor machine learning algorithms and the logistic regression model performs the best with a ROC of 86%. The researcher concludes that the logistic regression model is a suitable model to predict diabetes, and they recommend future work to develop a diabetes prediction model and deploy it with different mechanisms. U. Butt et al.[9] Investigates for the classification, early-stage identification, and prediction of diabetes disease using a machine learning-based approach, the researcher uses PIDD to conduct this experiment. For diabetes classification, random forest (RF), multilayer perceptron (MLP), and logistic regression (LR) classifiers have been employed. For predictive analysis, the researcher have employed long short-term memory (LSTM), moving averages (MA), and linear regression (LR). In this study, the researcher observed that MLP outperforms other classifiers with 86.08% of accuracy and LSTM improves the significant prediction with 87.26% accuracy of diabetes. M. Soni [2] the researcher aimed to develop a diabetes prediction model using machine learning techniques. The researcher uses PIDD to conduct this study. The proposed approach uses various classification and ensemble learning methods in which SVM, Knn, Random Forest, Decision Tree, Logistic Regression, and Gradient Boosting classifiers from these classifiers RF classifier performs best with the accuracy of 77%. Furthermore, [1][3][6][8][10][11][12][13] and [14]aimed to construct a predictive model, but they did not identify risk factors, they didn't extract rules which are important to make evidence-based strategies, and interventions towards preventing and/or reducing diabetes, they didn't develop information technology artifacts to interact potential users with the system, and didn't deploy the predictive model for the potential users. This study, hence, motivated to fill these gaps by constructing a predictive model, identifying risk factors, extracting relevant rules, designing an innovative artifact and deploying the predictive model for potential users on Heroku-based cloud computing platforms by designing prototypes or artefacts using HTML as a front end and integrating with the predictive model using Flask framework.

III. Materials and methods

A. Dataset Collection

The data used in this research was extracted from github repository, we have no further data collection method conducted. The data was 2000 records with 9 features and class levels, but after preprocessing the final data that we have used to develop the predictive model was 1393 records with 8 features and class level. The dataset description is given below

Table 1: Data set descriptions

No.	Attributes	Values
1.	Pregnancy	{0-19}
2.	Glucose	{0-199}
3.	Blood pressure(mm Hg)	{0-138}
4.	Skin thickness	{0-110}
5.	insulin	{0-846}
6.	Body Mass Index(BMI)	{0-80.6}
7.	Diabetes pedigree function	{0.078-2.42}
8.	Age	{21-81}
9.	Diabetes status	{0 & 1}

B. Data understanding method

In the data understanding phase, we have to understand the data such as types of features, fields, kinds of values in each feature have and also identify which features are continuous and which features are discrete-valued [15] This phase is useful for, knowing about the basic statistics of the data, is important to make it easier to fill the missing values, smooth noisy values, and spot outliers during data preprocessing [15][16] To understand the dataset we used data visualization techniques like box plot, bar plot, and histogram to show the outliers and we use manual methods to see in detail about the dataset.

C. Data Pre-processing

Data preparation involves data selection, data cleaning, data integration, feature selection, handling imbalances, and data transformation to make available to extract value from those data [17]. The extracted datasets consist of a total of 2000 records with 8 features including the class label. By nature, the collected data has missing values for some features, inconsistent values, redundant records, features that require encoding, noise, and outliers. As all these features are not relevant for developing a predictive model that can predict diabetes disease, data preprocessing techniques such as data cleaning, data transformation, and data encoding methods were applied. Data cleaning is a way of removing noise, inconsistencies, redundancy, and missing values to carefully develop the model [16]. Without cleaning the collected data, we can't get an accurate result [16]., so the researcher cleans the data appropriately by using appropriate methods [16]. The missing values were handled using median imputation techniques for continuous data features such as BMI, Insulin, and skin thickness. And mean imputation methods for the categorical data features such as Glucose and Blood Pressure. Redundant data were removed manually, or we have just removed them for achieving

better performances. Features filled with inconsistent values are also filled with mean and mode imputation methods by making a missing value. We have not applied any features selection methods, because the data by nature were collected with limited and relevant features. After conducting all the required data preprocessing tasks, a total of 1393 instances with 8 features were considered for further analysis and prediction model development. Finally, the dataset was divided into training and testing datasets following an 80/20% ratio.

D. Model Building

To conduct this study, we used four homogeneous ensemble machine learning algorithms namely random forest, Cat Boost, extra tree classifier, and extreme gradient boosting. To build a predictive model that predicts diabetes disease, homogeneous ensemble machine learning algorithms such as extreme gradient boosting, random forest, extra tree classifier, and cat boost algorithms were selected for an experiment. This method was chosen because of its ability to provide the best results while avoiding bias, variation, and over-fitting [18]. Random-forest was chosen for numerous types of datasets because of its excellent prediction accuracy, efficiency, interpretability, and nonparametric nature [19]. Cat boost is a method that supports numerical, categorical, and text characteristics, as well as having a good categorical data handling technique [20]. Cat boost is an example of a homogeneous ensemble approach that can increase model performance while lowering over-fitting and tuning time [20]. The cat boost technique is simple to develop, has good performance, is quick to train and predict, and is robust [20]. Data scientists use extreme gradient boosting (Xgboost) because of its high out-of-core computing speed [21]. The researcher employed homogeneous ensemble machine learning algorithms because they are more powerful in terms of accuracy than individual algorithms [22]. For developing the predictive models, all 8 features of the dataset were used. To get better performance by removing the overfitting, and preventing the bias, a grid search was implemented to tune the hyperparameters of each algorithm, as the performance of the algorithm highly depends on the selection of Hyperparameter, which has always been a crucial step in the process of machine learning model development [23][24][25]. The performance of each predictive model was evaluated using accuracy, precision, recall, F1- score, K-fold cross-validations, and ROC-AUC. Fig 1 above represents the proposed model architecture that was implemented in this study to develop a predictive model.

E. Parameter tuning

To conduct these experiments, we have used grid search parameter tuning mechanisms for each algorithm, and we have used those tuned parameters for developing the final model.

Table 2: List of tuned parameters

algorithms	Cat boostalgorithms	XGB Classifier	Random Forest Classifier	Extra Tree Classifier
Parameters tuned by grid search	depth=10, iterations=300, l2_leaf_reg=1 learning_rate=0.15	objective = 'reg:linear', colsample_bytree = 0.3, learning_rate = 0.01,max_depth = 5,alpha = 10, n_estimators = 100	criterion='entropy',max_features='sqrt', min_samples_split=3,n_estimators=500, random_state=0,max_depth=20, max_leaf_nodes=400, n_jobs=-1	n_estimators=500, max_features=8, max_depth=200, max_leaf_nodes=400, criterion='entropy'

F. Model evaluation

After developing a predictive model, we have evaluated the prediction results using various objective evaluation metrics like classification accuracy, precision, recall, f1-score, k-fold cross-validations, and ROC-AUC.

G. Artifact development

Design science research aims to create useful information technology artifact that solves organizational problems by integrating users with the system and deploying the system to the potential set of users [26]. To design those information technology artifacts, we used a flask web development framework. Web Framework represents a collection of libraries and modules that enable web application developers to write applications without worrying about low-level details such as protocol, thread management, and so on [27]. To develop information technology artifacts, we used HTML for the front

end and python with Flask framework for the back. Flask is one of the machine learning model deployment web development frameworks written in python, that doesn't require a particular library, supports extensions that can add application features, lightweight, and is suitable for small size model deployment [27]. After developing the diabetes disease prediction model and evaluating the performance of the model, we have developed a user interface with information technology artifacts that interacts with the potential user and the system. The artifacts were designed by using HTML languages as a front end. To display the results of the predictive model on the designed artifacts, we have used Flask API as a back end to extract the predictive model result on the front end.

Figure 1: The designed artifacts

Diabetes disease prediction model

Pregnancies	<input type="text" value="Select Pregnancies....."/>	Glucose	<input type="text" value="select Glucose....."/>
BloodPressure	<input type="text" value="BloodPressure....."/>	SkinThickness	<input type="text" value="SkinThickness....."/>
Insulin	<input type="text" value="select Insulin....."/>	BMI	<input type="text" value="select BMI....."/>
Dia_PedigreeFun	<input type="text" value="Diabetes Pedigree....."/>	Age	<input type="text" value="select Age....."/>

H. Domain expert evaluation

Before deploying the best-performed model, we have evaluated the best-performed model and the designed artifacts from the University of Gondar specialized hospital staff as domain experts by preparing questionnaires. See appendix I. Besides the questioner, the domain expert or health specialists who work in the University of Gondar specialized hospital uses the model for consequent three months, and they evaluate the model using test cases. Using the test case approach the domain expert tests, and treat the patient based on this predictive model artifact, and 94% of the predictive model match with the domain experts.

I. Model deployment

Deployment is a way of designing a framework to be applied in public-policy settings as an external partner or internal analytic team trying to use machine learning methods in a production environment [28]. To deploy a final model for the potential users, we used a Heroku-based cloud computing platform. Heroku is so easy to use, easy to scale, secure, free to use, and that it's a top choice for many development projects [29]. Heroku is a container-based cloud Platform as a Service is used to deploy, manage, and scale modern apps [29]. After evaluating the developed innovative artifacts with the best-performed model using Flask framework using domain experts we have deployed the best-performed model over the cloud using Heroku-based cloud computing platforms. After developing the predictive model end evaluating the result, we have prepared an innovative artifacts using HTML as a front end and Flask API python module as a back end. After developing those artifacts we have

evaluated it using test case. Based on the domain experts idea we have deployed it for the potential set of users over the cloud using Heroku-based cloud computing plat forms. All potential users can access (<https://diabetes-model1.herokuapp.com/>) the predictive model to evaluate the person is whether diabetic or not.

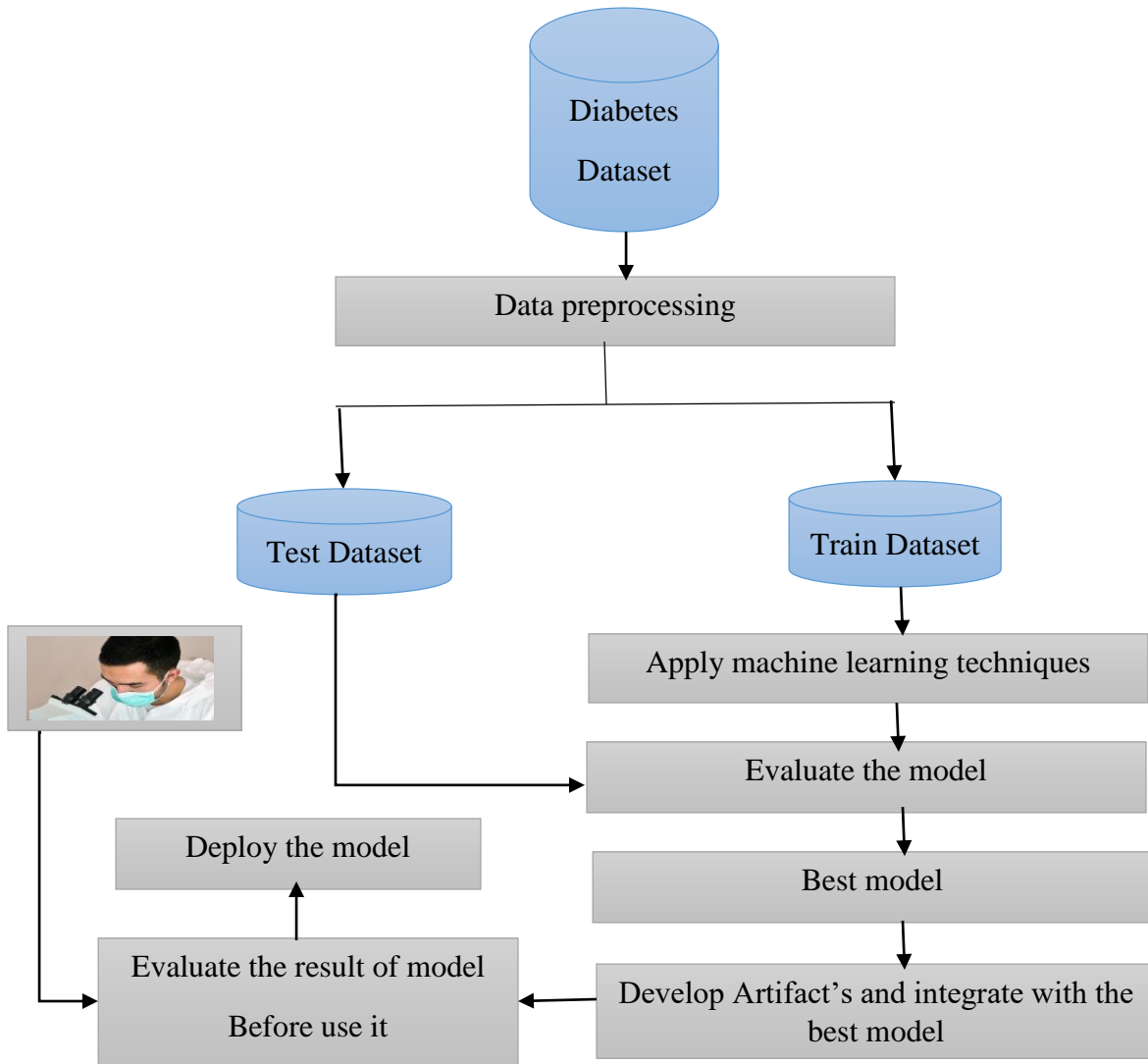


Figure 2: Proposed model architecture

IV. Experimental setups and result discussion

Here below we have discussed all the results based on the research questions raised above.

1) Which homogeneous ensemble machine learning algorithms are suitable for predicting diabetes diseases?

To answer this question, we have conducted four experiments using homogeneous ensemble machine learning algorithms

namely random forest, extreme gradient boost, extra tree classifier, and cat boost algorithms. The experiments showed that the model that was developed using the cat boost algorithm performs better in predicting diabetes disease with 90.3% of accuracy, 88.41% of precision, 89.91% of recall, 89.08 % of f1_score, 93.05% of cross-validation, and, 89.91% of ROC-AUC see table 3 below, the model was developed using all the tuning parameters see table 2 above.

Table 3: Model performance

	Evaluation metrics	Result
Extra Tree classifier	Accuracy	87.81%
	Precision	86.12%
	Recall	85.68%
	F1_score	85.9%
	Cross-validation	90.4%
	ROC_AUC	85.68%
Random forest	Accuracy	87.81%
	Precision	85.98%
	Recall	85.98%
	F1_score	85.98%
	Cross-validation	90.61%

	ROC_AUC	85.98%
Cat Boost	Accuracy	90.32%
	Precision	88.41%
	Recall	89.91%
	F1_score	89.08%
	Cross-validation	93.05%
	ROC_AUC	89.91%
Extreme gradient Boost	Accuracy	86.38%
	Precision	84.22%
	Recall	84.62%
	F1_score	84.42%
	Cross-validation	91.54%
	ROC-AUC	84.62%

2) What are the most determinant risk factors that influence the occurrence of diabetes disease?

To answer this question, feature importance analysis was performed using the model that was developed with the best performing algorithm which is cat boost. See in Table 4 here below the Glucose attribute has the highest feature importance
Table 4: Feature importance values

No.	Attributes	Feature importance values
1.	Glucose	17.30552168
2.	Diabetes pedigree function	15.23109942
3.	Body Mass Index(BMI)	15.0465794
4.	Blood pressure(mm Hg)	11.94622687
5.	Age	11.7521685
6.	Pregnancy	10.89282892
7.	Skin thickness	9.1230587
8.	Insulin	8.7025165

3) What are the important rules that can be generated from the predictive model?

To answer this question, we used all the features that we used to develop the predictive model and generate all the important rules by using the best-performed algorithms (cat boost algorithms) for diabetes disease.

I. Conclusion

Diabetes is a chronic disease that continues to be a major and global problem since it affects the health of the entire population. Diabetes is a hormonal disorder in which the inability of the body to produce insulin causes the metabolism of sugar in the body to be abnormal, thereby, raising the blood glucose levels in the body of a particular individual. Certain risk factors such as age, BMI, Glucose Levels, Blood Pressure, etc., play an important role in the contribution of the disease. This study aimed to develop a predictive model for diabetes disease by using homogeneous ensemble machine learning algorithms and deploying it over the cloud using Heroku-based cloud computing platforms and Flask API python module. This study was conducted using design science methodology. To construct a predictive model, we used homogeneous ensemble machine learning algorithms namely random forest, extreme gradient boosting, extra tree classifier, and cat boost algorithms. Of those algorithms, the cat boost

value of 17.30552168, i.e. it is the most determinant risk factor for diabetes. Insulin has a feature importance value of 8.7025165, it is the least feature importance value of the rest features, and i.e. it is the least determinant risk factor when we compare it with other features.

algorithm has registered the highest performance with 90.32% of accuracy, 88.41% of precision, 89.91% of recall, 89.08 % of f1_score, 93.05% of cross-validation, and, 89.91% of ROC-AUC. Due to the better performances of cat boost algorithms we have used them for identifying the determinant risk factors using feature importance analysis, generating the most important rules using the best fit model for developing policies and interventions towards maintaining diabetes. Finally, we recommend that future researchers conduct a predictive model for diabetes disease and develop it with mobile applications.

Acknowledgment

We would like to acknowledge the University of Gondar specialized hospital staff for their support to evaluate everything from the data understanding to model evaluations.

Abbreviations

API: application interface, BMI: body mass index, IDF: International Diabetes Federation, HTML: hypertext markup language, PIDD: Pima Indian dataset, WHO: world health organization

Data Availability

The data that support the findings of this study are available on request from the corresponding author.

Ethics approval and consent to participate

All methods used in this study followed guidelines and regulations. Health care professionals who work in the University of Gondar specialized hospital approved this study.

Conflicts of Interest

The authors of this manuscript declare that they do not have any conflicts of interest.

Consent for publication

Not applicable.

Competing interests

The authors report that they have no conflicts.

Author Details

The authors of this study are affiliated with the University of Gondar, College of Informatics, Gondar, Ethiopia.

References

- [1] L. F. Aparicio, J. Noguez, L. Montesinos, and J. A. G. García, "Machine learning and deep learning predictive models for type 2 diabetes : a systematic review," *Diabetol. Metab. Syndr.*, 2021, doi: 10.1186/s13098-021-00767-9.
- [2] M. Soni, "Diabetes Prediction using Machine Learning Techniques," vol. 9, no. 09, pp. 921–925, 2020.
- [3] O. Llaha and A. Rista, "Prediction and Detection of Diabetes using Machine Learning."
- [4] R. Krishnamoorthi *et al.*, "A Novel Diabetes Healthcare Disease Prediction Framework Using Machine Learning Techniques," *J. Healthc. Eng.*, vol. 2022, 2022, doi: 10.1155/2022/1684017.
- [5] C. Lyngdoh, N. A. Choudhury, and S. Moulik, "Diabetes Disease Prediction Using Machine Learning Algorithms," no. May, 2021, doi: 10.1109/IECBES48179.2021.9398759.
- [6] Q. Zou, K. Qu, Y. Luo, D. Yin, Y. Ju, and H. Tang, "Predicting Diabetes Mellitus With Machine Learning Techniques," vol. 9, no. November, pp. 1–10, 2018, doi: 10.3389/fgene.2018.00515.
- [7] P. Model *et al.*, "Predict Diabetes Mellitus Using Machine Learning Algorithms Predict Diabetes Mellitus Using Machine Learning Algorithms," 2021, doi: 10.1088/1742-6596/2089/1/012002.
- [8] A. Mujumdar and V. Vaidehi, "ScienceDirect ScienceDirect ScienceDirect Diabetes Prediction using Machine Learning Aishwarya Mujumdar Diabetes Prediction using Machine Learning Aishwarya Mujumdar Aishwarya," *Procedia Comput. Sci.*, vol. 165, pp. 292–299, 2019, doi: 10.1016/j.procs.2020.01.047.
- [9] U. M. Butt *et al.*, "Machine Learning Based Diabetes Classification and Prediction for Healthcare Applications," vol. 2021, 2021.
- [10] H. Zhou, R. Myrzashova, and R. Zheng, "Diabetes prediction model based on an enhanced deep neural network," 2020.
- [11] J. Wang and Z. Xu, "Cluster file systems: A case study," *Futur. Gener. Comput. Syst.*, vol. 18, no. 3, pp. 373–387, 2002, doi: 10.1016/S0167-739X(01)00057-7.
- [12] M. D. Dithy and V. Krishnapriya, "Anemia selection in pregnant women by using random prediction (Rp) classification algorithm," *Int. J. Recent Technol. Eng.*, vol. 8, no. 2, pp. 2623–2630, 2019, doi: 10.35940/ijrte.B3016.078219.
- [13] P. Anand, R. Gupta, and A. Sharma, "Prediction of Anaemia among children using Machine Learning Algorithms," no. June, pp. 469–480, 2020.
- [14] S. S. Yadav and S. M. Jadhav, "Machine learning algorithms for disease prediction using Iot environment," *Int. J. Eng. Adv. Technol.*, vol. 8, no. 6, pp. 4303–4307, 2019, doi: 10.35940/ijeat.F8914.088619.
- [15] D. Skillicorn, *Understanding complex datasets: Data mining with matrix decompositions*. 2007.
- [16] K. Natarajan, J. Li, and A. Koronios, "Data mining techniques for data cleaning," *Eng. Asset Lifecycle Manag. - Proc. 4th World Congr. Eng. Asset Manag. WCEAM 2009*, no. September, pp. 796–804, 2009, doi: 10.1007/978-0-85729-320-6_91.
- [17] N. H. Son, "Data cleaning and Data preprocessing," 2011, [Online]. Available: <http://www.mimuw.edu.pl/~son/datamining/DM/4-preprocess.pdf>.
- [18] H. Inoue and R. Inoue, "A very large platform for floating offshore facilities," *Coast. Ocean Sp. Util. III. Proc. Symp. Genoa, 1993*, pp. 533–551, 1995.
- [19] J. Ali, R. Khan, N. Ahmad, and I. Maqsood, "Random Forests and Decision Trees," *Int. J. Comput. Sci. Issues*, vol. 9, no. 5, pp. 272–278, 2012.
- [20] "How CatBoost Algorithm Works In Machine Learning." <https://dataaspirant.com/catboost-algorithm/> (accessed Aug. 16, 2021).
- [21] A. Ibrahim, A. Osman, A. Najah, M. Fai, Y. Feng, and A. El-shafie, "Extreme gradient boosting (Xgboost) model to predict the groundwater levels in Selangor Malaysia," *Ain Shams Eng. J.*, vol. 12, no. 2, pp. 1545–1556, 2021, doi: 10.1016/j.asej.2020.11.011.
- [22] O. Kramer, "Ensemble Learning," no. May, pp. 25–32, 2013, doi: 10.1007/978-3-642-38652-7_3.
- [23] M. J. Healy, "Statistics from the inside. 15. Multiple regression (1).," *Arch. Dis. Child.*, vol. 73, no. 2, pp. 177–181, 1995, doi: 10.1136/adc.73.2.177.
- [24] R. G. Mantovani, A. L. D. Rossi, E. Alcobaça, J. C. Gertrudes, S. B. Junior, and A. C. P. de L. F. de Carvalho, "Rethinking Defaults Values: a Low Cost and Efficient Strategy to Define Hyperparameters," 2020, [Online]. Available: <http://arxiv.org/abs/2008.00025>.
- [25] M. M. RAMADHAN, I. S. SITANGGANG, F. R. NASUTION, and A. GHIFARI, "Parameter Tuning in Random Forest Based on Grid Search Method for Gender Classification Based on Voice Frequency," *DEStech Trans. Comput. Sci. Eng.*, no. cece, 2017, doi: 10.12783/dtcse/cece2017/14611.
- [26] S. Y. R. Esearch, B. A. R. Hevner, S. T. March, J. Park, and S. Ram, "D ESIGN S CIENCE IN I NFORMATION," vol. 28, no. 1, pp. 75–105, 2004.

- [27] “What is Flask Python - Python Tutorial.”
<https://pythonbasics.org/what-is-flask-python/>
(accessed Aug. 16, 2021).
- [28] K. Ackermann *et al.*, “Deploying machine learning models for public policy: A framework,” *Proc. ACM SIGKDD Int. Conf. Knowl. Discov. Data Min.*, no. August, pp. 15–22, 2018, doi: 10.1145/3219819.3219911.
- [29] “What is Heroku and What is it Used For? - MentorMate.”
<https://mentormate.com/blog/what-is-heroku-used-for-cloud-development/> (accessed Sep. 13, 2021).

