

# Disease Diagnosis using Machine Learning

## 1. Definition

### 1.1 Project Overview

Machine Learning (ML), a branch of Artificial Intelligence (AI), learns from the data using various algorithms and is a self-improving process in terms of performance as make adjustments during the learning process [1]. ML has been successfully applied to practically in every domain like robotics, education, travel to health care [2]. In healthcare domain, the ML approaches are mainly used with the purpose of disease diagnosis [3].

The machine learning approaches came into health sector domain in 1970s and an international AI journal *Artificial Intelligence in Medicine* was established in 1980 [4]. In the next two decades disease diagnosis domain adopted the classical ML approaches like Support Vector Machine, Naïve Bayes and some artificial neural networks [5]. The introduction of AlexNet in 2012 initiated the current wave of deep learning in this field as neural networks demonstrated superior performance [6]. Also, in this past decade, the investment in AI in healthcare applications has increased significantly.

The analysis of the clinical data can lead to the timely diagnosis of the disease which will help to start cure for the patient in time as well [3]. Traditional approach of diagnosing disease is generally costly and time consuming. And, the potential of time and cost proficient machine learning biased disease diagnosis approaches are proven by the researchers [7]. ML techniques have not only been able to diagnose the common diseases but are also equally capable of diagnosing the rare diseases [2], [8].

In general, a dataset table used to build a ML model for diagnosing a disease have columns for different attributes and a column variable for the class variable. Here, class variable indicates whether the instance in the table indicated is positively diagnosed with the disease under consideration. Usually, class values of 1 means positively diagnosed and 0 means negatively diagnosed. Supervised and unsupervised ML [9] approaches have been in practice for analyzing the health care data. In general, the disease diagnosis problems are based on supervised learning. We will present the detailed analysis of the dataset and ML algorithms used will be presented in Section 2.

### 1.2 Problem Statement

Although ML offers systematic and sophisticated algorithms of multi-dimensional clinical data, the accuracy of the ML in diagnosing the diseases is still a concern [10]. And the improvement in the performance of ML to diagnose disease is a hot topic in this domain. As different ML approach perform differently for different healthcare dataset, we are also in need to

find the way to apply many state of art algorithms to same dataset in reasonable time, so that the search of best ML method can be pursued to diagnose a particular disease.

The use of libraries like AutoGluon can help to find the best performing ML approach out of many approaches in diagnosing the disease for a given dataset in reasonable time and with optimal lines of code. This will decrease the probability of inaccurate diagnosis, which is a significantly important consideration while dealing with the health of the people. We will test the performance of ML approaches like Naïve Bayes, Support Vector Machine (SVM), K Nearest Neighbors (KNN), perceptron and robust deep neural networks in AutoGluon like LightGBM, XGBoost, MXNet etc.

### 1.3 Metrics

Disease diagnosis is a classification task. And, Classification ML Algorithms are evaluated using Classification Accuracy Measures like Accuracy, Precision, Recall and F1-score. Let us consider a value of 1 (having diabetes) be positive and value of 0 in class variable be negative in the considered dataset. Let True Positive (TP) be the correctly classified number of positive classes from a ML model. Similarly False Positive (FP) be the number of incorrectly classified as positive classes, True Negative (TN) be the correctly classified number of negative classes and False Negative (FN) be the number of classes incorrectly classified as Negative classes. Various classification accuracy measures are computed based on TP, FP, TN and, FN [11]. The Accuracy can be computed as

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

the Precision as

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

the Recall as

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

and the F1-Score as

$$\text{F1 - Score} = \frac{2 * \text{Precision} * \text{Recall}}{\text{Precision} + \text{Recall}} \quad (4)$$

Classification accuracy measures in Eq. (1), (2), (3) and (4) have been used to evaluate the performance of applied classifier algorithms.

In general, only one (mostly accuracy) evaluation metric is used to evaluate the performance of the ML algorithms. However, in our study we are using four evaluation metrics primarily

because the two reasons. The first reason is that in the used diabetes dataset Outcome class variables is highly imbalanced toward the value 0, and the accuracy measure from the imbalanced dataset can be misleading [12]. The next reason is that we are trying to avoid the case of accuracy paradox by considering four evaluation metrics [13], [14].

## Analysis

### 1.4 Data Exploration

For this project, we have chosen a healthcare dataset related to the diabetes. The dataset is Pima Indian Diabetes Dataset which is frequently used to evaluate the performance of developed ML techniques [15], [16]. We downloaded the dataset from [15]. This data set has 8 attributes and one class variable named Outcome. Outcome variable has possible value of 0 or 1, 1 being interpreted as tested positive for diabetes. The dataset has 768 instances, out of which 268 being those tested positive for diabetes.

Two the attributes in the dataset are continuous numerical variable the rest are discrete numerical integers.

*Table 1: Attribute Data Type*

Column	Non-Null Count	Dtype
-----	-----	-----
Pregnancies	768 non-null	int64
Glucose	768 non-null	int64
BloodPressure	768 non-null	int64
SkinThickness	768 non-null	int64
Insulin	768 non-null	int64
BMI	768 non-null	float64
DiabetesPedigreeFunction	768 non-null	float64
Age	768 non-null	int64
Outcome	768 non-null	int64

The detailed statistical description of each attribute is as shown below in Table 2.

*Table 2: Attributes statistical description*

	Pregnancies	Glucose	BloodPressure	SkinThickness	Insulin	BMI	DiabetesPedigreeFunction	Age
<b>count</b>	768.000000	768.000000	768.000000	768.000000	768.000000	768.000000	768.000000	768.000000
<b>mean</b>	3.845052	120.894531	69.105469	20.536458	79.799479	31.992578	0.471876	33.240885
<b>std</b>	3.369578	31.972618	19.355807	15.952218	115.244002	7.884160	0.331329	11.760232
<b>min</b>	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.078000	21.000000
<b>25%</b>	1.000000	99.000000	62.000000	0.000000	0.000000	27.300000	0.243750	24.000000
<b>50%</b>	3.000000	117.000000	72.000000	23.000000	30.500000	32.000000	0.372500	29.000000
<b>75%</b>	6.000000	140.250000	80.000000	32.000000	127.250000	36.600000	0.626250	41.000000
<b>max</b>	17.000000	199.000000	122.000000	99.000000	846.000000	67.100000	2.420000	81.000000

## 1.5 Exploratory Visualization

We performed exploratory visualization of the attributes with the histogram. The results are as shown in Figures 1. The idea behind the exploratory visualization was to check whether some variables are constant over the range. Such variables can be avoided while building the models. However, our exploratory visualization showed that every attribute can be important for the disease diagnosis with Machine Learning.

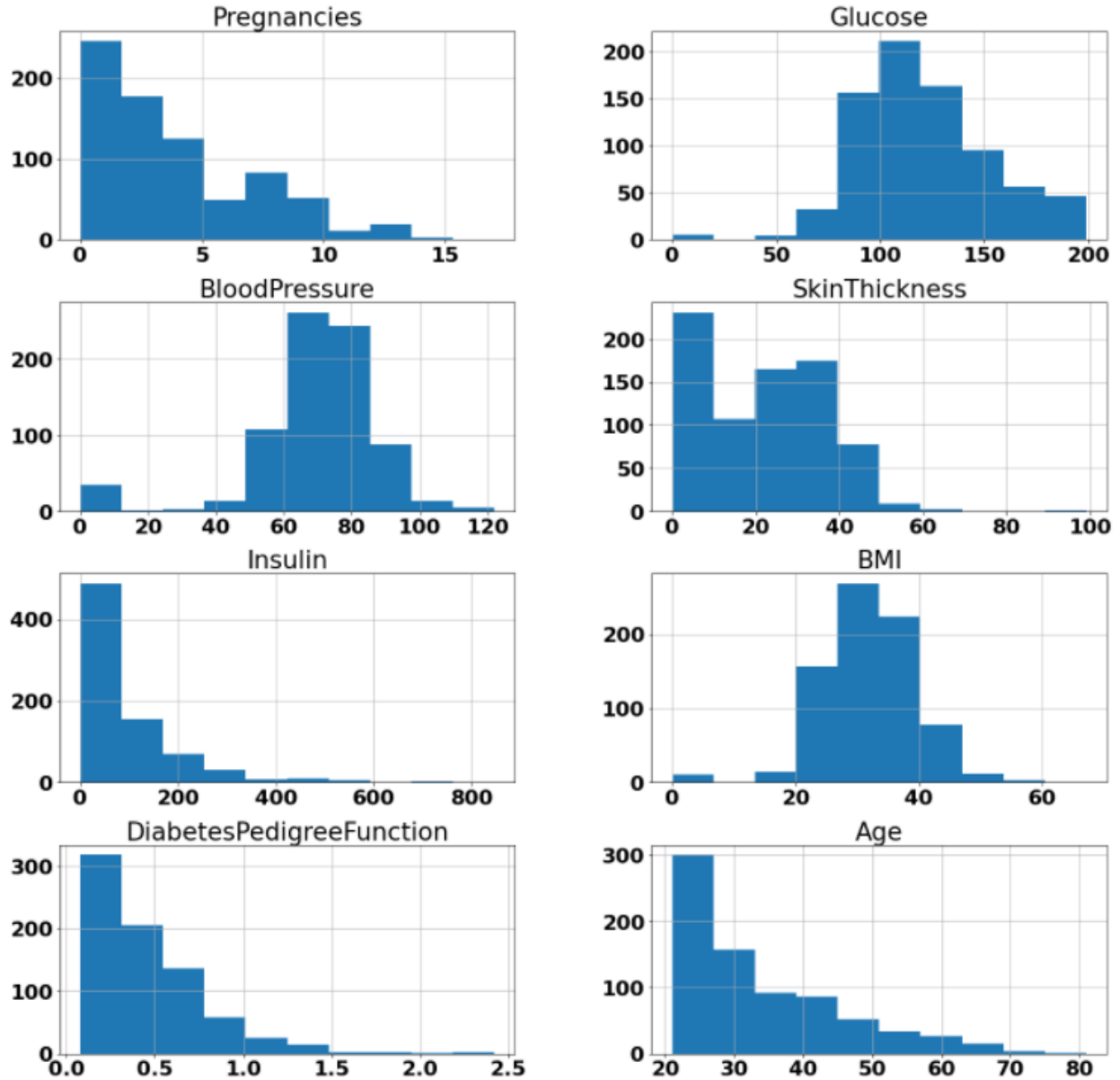


Figure 1: Histogram of Attributes

## 1.6 Algorithms and Techniques

Here, we will be applying classification algorithms from the scikit-learn library [17] and AutoGluon library [18] and checking the capacity of the algorithms to diagnose the diabetes disease. Scikit-learn is the most successful and robust library for machine learning in Python. This library is primarily written in Python and is based on the modules like NumPy [19], SciPy [20] and Matplotlib [21]. And the open source AutoML library AutoGluon-Tabular can train highly accurate different machine learning models with a single line of code [18]. The ML algorithms from scikit-learn library and Auto-Gluon library are implemented with AWS SageMaker [22]. The Amazon SageMaker is capable of building, training, and deploying state of art Machine Learning models with full managed infrastructure tools and workflows [23]. Some of the

classification ML used are Naïve Bayes, Support Vector Machine (SVM), K Nearest Neighbors (KNN), perceptron and robust deep neural networks in AutoGluon like LightGBM, XGBoost, MXNet etc. The list of ML algorithms evaluated for diabetes diagnosis are as below in Table 3 [18], [24]. The detail of the algorithm shadows the main goal of the project which is implementation of ML for disease diagnosis. Please visit the reference [18], [24], if the details of the Algorithms are of interest.

*Table 3 List of ML Algorithms Used*

S. N	ML Algorithm
1	Random Forest Classifier (Scikit-learn)
2	Decision Tree Classifier (Scikit-learn)
3	Naïve Bayes Classifier (Scikit-learn)
4	Perceptron (Scikit-learn)
5	Multilayer Perceptron (Scikit-learn)
6	Voting Classifier (Scikit-learn)
7	WeightedEnsemble_L2 (AutoGluon)
8	LightGBM_BAG_L1 (AutoGluon)
9	LightGBM_LARGE_BAG_L1 (AutoGluon)
10	NeuralNetFastAI_BAG_L1 (AutoGluon)
11	CATBoost_BAG_L1 (AutoGluon)
12	ExtraTreesGini_BAG_L1 (AutoGluon)
13	LightGBMXT_BAG_L1 (AutoGluon)
14	XGBoost_BAG_L1 (AutoGluon)
15	RandomForestEntr_BAG_L1 (AutoGluon)
16	RandomForestGini_BAG_L1 (AutoGluon)
17	ExtraTreesEntr_BAG_L1 (AutoGluon)
18	NeuralNetMXNet_BAG_L1
19	KNeighborsUnif_BAG_L1 (AutoGluon)
20	KNeighborsDist_BAG_L1 (AutoGluon)

## **1.7 Benchmark Model**

The author in [15] mentions that the baseline accuracy for diagnosing diabetes is about 65 percent when Pima Indian Diabetes Dataset is used for the classification problem and the best diagnosis accuracy obtained so far is about 77 percent. The decision trees algorithm, neural network, logistic regression has been able to give accuracy of about 77 percentage. We will compare the results of ML approaches applied in this study with these benchmark results.

## **2. Methodology**

### **2.1 Data Preprocessing**

The exploratory analysis and visualization of the data did not suggest any preprocessing of the data for learning the ML models, as no anomaly was detected. Therefore, the process of evaluating a ML for diagnosing the disease was performed with no data preprocessing.

### **2.2 Implementation**

The implementation and evaluation of ML algorithms was performed in the notebook instance in the Amazon SageMaker. The six ML techniques were applied by importing the modules and the ML models directly as it was already installed in the Cuda Python 3 Kernel. However, the AutoGluon library is not pre-installed in the kernel. There, it had to be downloaded before importing the ML algorithms from it. The detailed implementation process is presented in the notebook project.ipynb which is kept in the author's GitHub repository. The results can be reproduced using the project.ipynb notebook.

### **2.3 Refinement**

We trained the 14 AutoGluon ML algorithms, first using the evaluation metric accuracy. As, the dataset we have an imbalanced dataset in terms of Outcome class, therefore we used the evaluation metric F1-score, which is a more favored evaluation metric while training with imbalanced data. We also tuned hyperparameters to check if better results are possible but the prediction accuracy with the tune hyperparameters came out to be lower than the untuned ones. Therefore, future research with extensive tuning of different hyper parameters is recommended to check the existence of better models with different set of hyper parameters.

## **3. Results**

### **3.1 Model Evaluation and Validation**

The evaluation of the different ML techniques in diagnosing diabetes from the given dataset is as shown in Table 4. The result shows that Naive Bayes Classifier ML algorithm performs better among the ones compared based on combined analysis of all the evaluation metrics.

Table 4: Evaluation of ML Algorithms

S. N	ML Algorithm	Accuracy	F1-score	Precision	Recall
1	Random Forest Classifier (Scikit-learn)	0.74	0.81	0.78	0.84
2	Decision Tree Classifier (Scikit-learn)	0.65	0.73	0.73	0.73
3	<b>Naïve Bayes Classifier (Scikit-learn)</b>	<b>0.77</b>	<b>0.83</b>	<b>0.80</b>	<b>0.86</b>
4	Perceptron (Scikit-learn)	0.49	0.47	0.71	0.35
5	Multilayer Perceptron (Scikit-learn)	0.68	0.76	0.75	0.77
6	Voting Classifier (Scikit-learn)	0.72	0.78	0.79	0.77
7	AutoGluon Best Performer	0.74	0.82	0.76	0.88

We present the accuracy performance of different AutoGluon ML algorithms when trained with accuracy as accuracy in Figure 3.

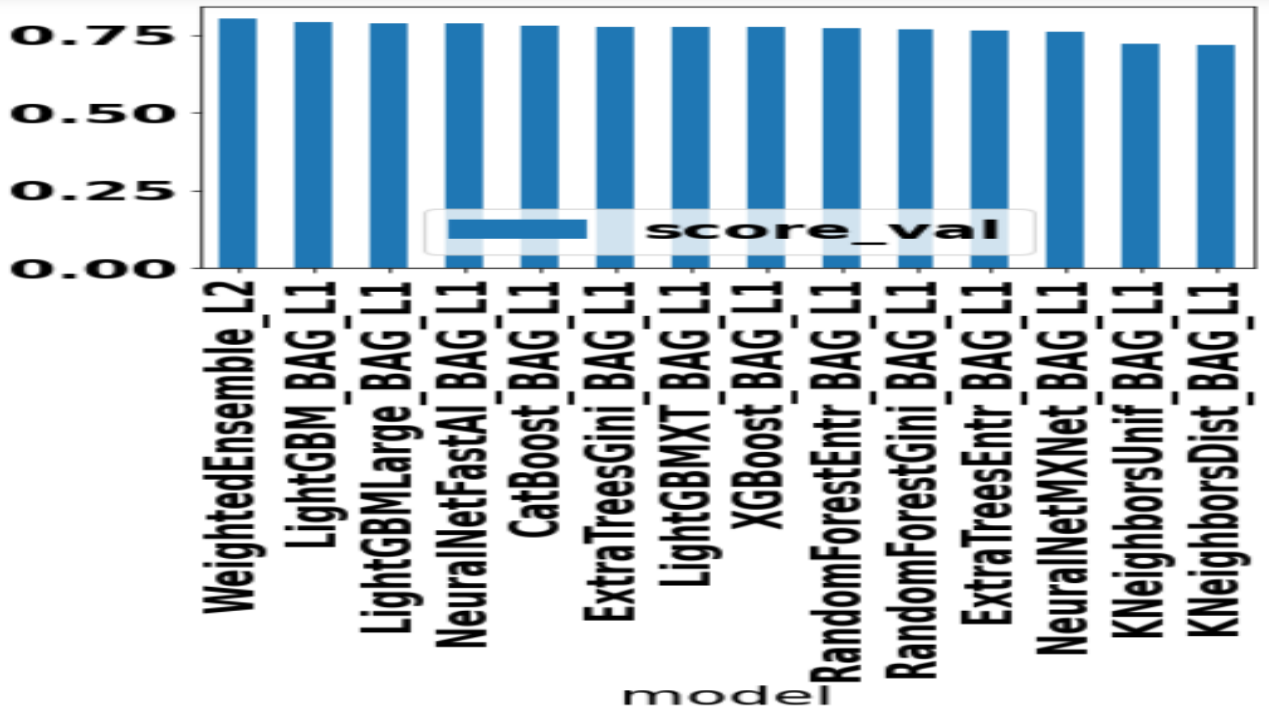


Figure 2: Evaluation of AutoGluon ML algorithms when trained with accuracy as evaluation metric



Similarly, performance in terms of F1-scores is shown when trained with F1-scores as evaluation metric as in Figure 4. It is seen that the Weighted Ensemble ML technique performs better for both the cases and KNN based ML has the least performance.

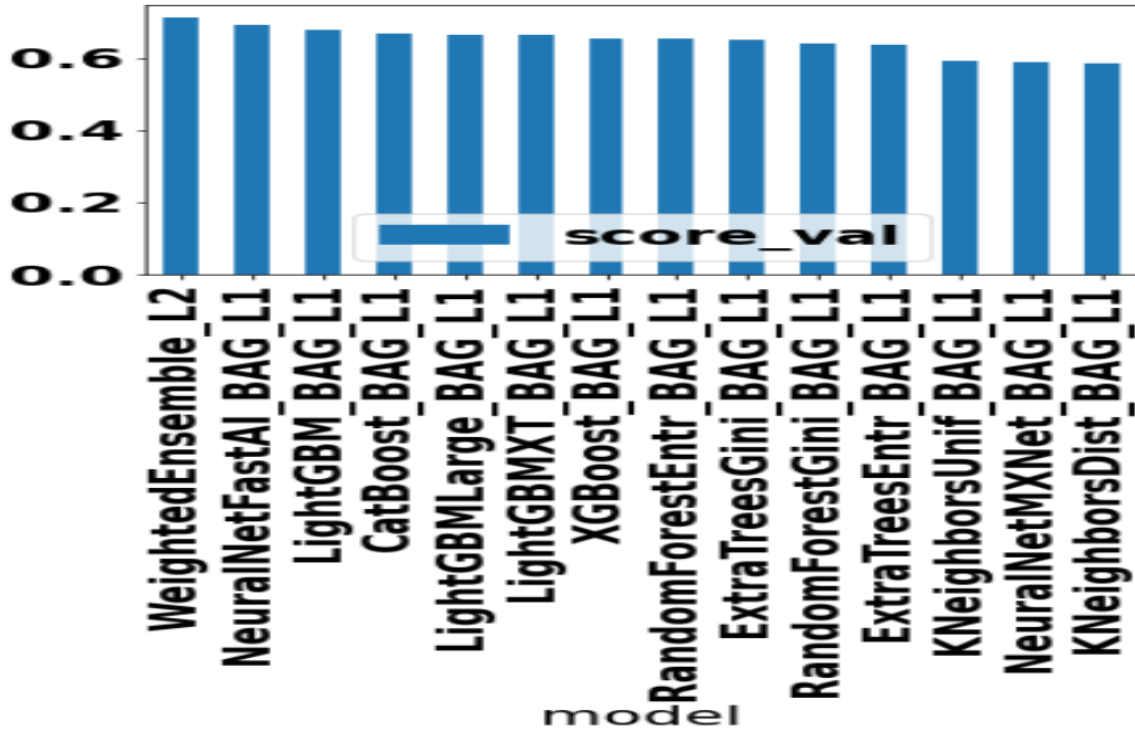


Figure 3: Evaluation of AutoGluon ML algorithms when trained with F1-score as evaluation metric

### 3.2 Justification

Our study shows that most of the ML methods perform better than the benchmark of baseline accuracy of 65 percent set for this dataset while diagnosing the diabetes. However, about 77 percent of the accuracy seems to be the best case for the state of art ML algorithms. Considering the case of having the imbalanced data, we can emphasize the capability of Naïve Bayes method to perform better among the rest considering the combined analysis of all the evaluation metrics.

## 4. Conclusion

### 4.1 Free-Form Visualization

The boxplots of the attributes are as presented as the free form visualization. Here, the boxplots indicate that the attributes vary differently in terms of their distribution, This could be the reason for reaching a good F1-score of about 0.83, in spite of having a imbalanced dataset. The next important point is regarding the box plot of Body Mass Index (BMI). It shows the mean BMI of the collected data is more than 30, however the dataset does have significantly less number of

instances diagnosed with the diabetes, which is against the general assumption. Thus, the BMI cannot only account for high probability of having the diabetes.

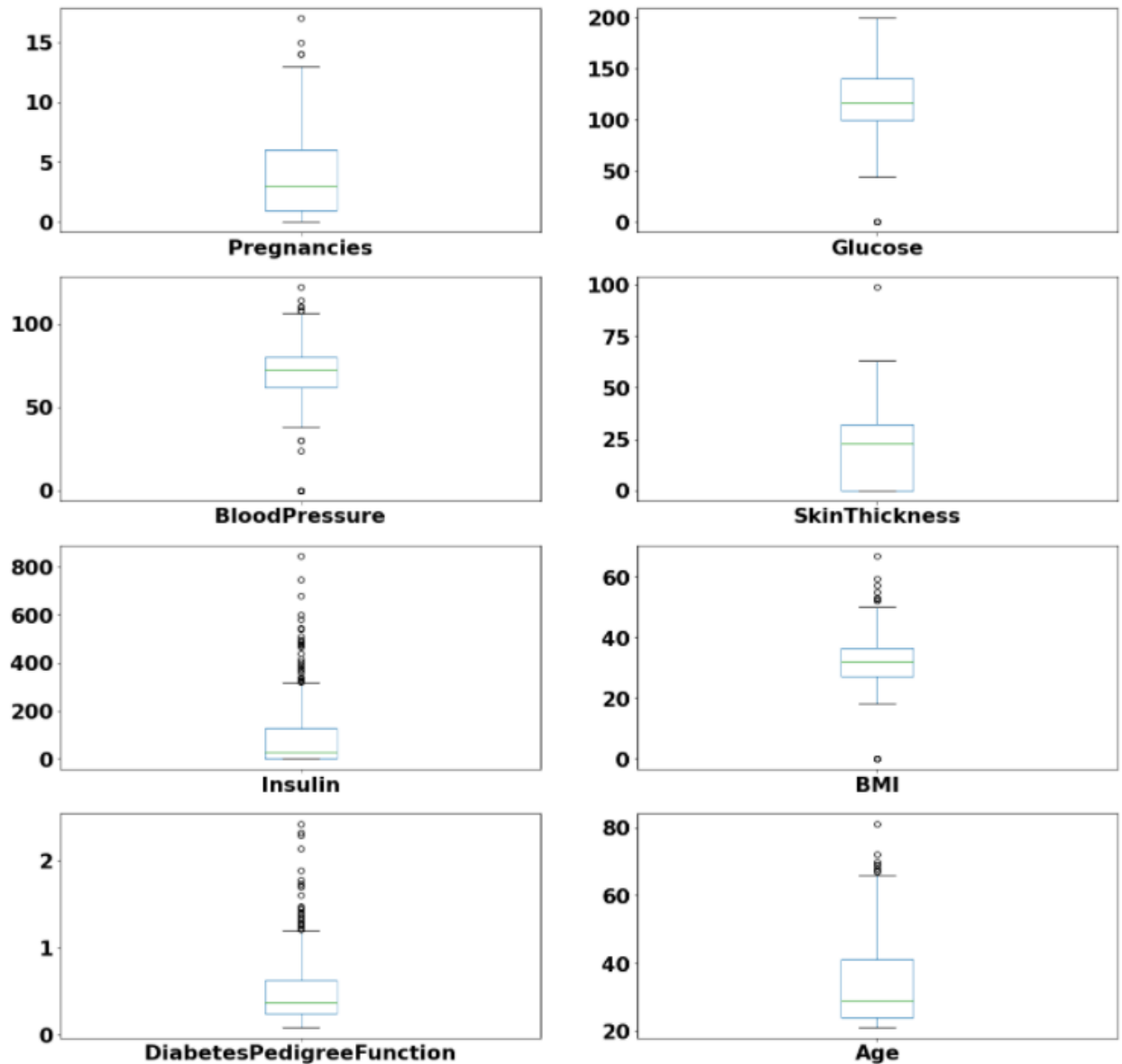


Figure 4: Box-plot of Attributes

## 4.2 Reflection

The project initiates by downloading the Pima Indian Diabetes Dataset [15]. Then, we created a notebook instance in the Amazon SageMaker. When the notebook instance was successfully created, a Jupyter Notebook was opened and downloaded dataset were uploaded to the notebook instance. We used the most cost-effective instance ml.t2.medium as our data set was not significantly large. After that we imported and downloaded (if needed) the required module and libraries. The next step was data exploration and exploratory visualization result of which suggested the direct use of current dataset without preprocessing for evaluation ML models to

diagnosing diabetes. Then, we applied different state of art classification ML algorithms to Diabetes Dataset using scikit-learn and AutoGluon libraries and computed the performance of each approach in diagnosing the diabetes. We made sure that same training and test set were used for each of the ML algorithm by defining the parameter *seed* with same number=42. Our results show that the that Naive Bayes Classifier ML algorithm performs better among the ones compared.

### 4.3 Improvement

The possibility of the improvement in the performance of ML models in future can be started by finding the correlation among each attribute and dropping the highly correlated attributes. Because the highly correlated attributes can confuse a model in the learning phase.

## 5. Reference

- [1] T. M. Mitchell, *Machine Learning*. New York: McGraw-Hill, 1997.
- [2] M. Fatima and M. Pasha, "Survey of Machine Learning Algorithms for Disease Diagnostic," *J. Intell. Learn. Syst. Appl.*, vol. 09, no. 01, pp. 1–16, 2017, doi: 10.4236/JILSA.2017.91001.
- [3] Neural Designer, "Machine learning use cases | Neural Designer." <https://www.neuraldesigner.com/solutions> (accessed Jan. 13, 2022).
- [4] Q. Danial, "Demystifying AI in Healthcare: Historical Perspectives and Current Considerations." <https://www.physicianleaders.org/news/demystifying-ai-in-healthcare-historical-perspectives-and-current-considerations> (accessed Jan. 13, 2022).
- [5] I. Kononenko, "Machine learning for medical diagnosis: history, state of the art and perspective," *Artif. Intell. Med.*, vol. 23, no. 1, pp. 89–109, Aug. 2001, doi: 10.1016/S0933-3657(01)00077-X.
- [6] A. Krizhevsky, I. Sutskever, and G. E. Hinton, "ImageNet Classification with Deep Convolutional Neural Networks," *Adv. Neural Inf. Process. Syst.*, vol. 25, 2012, Accessed: Jan. 13, 2022. [Online]. Available: <http://code.google.com/p/cuda-convnet/>.
- [7] P. Sajda, "Machine learning for detection and diagnosis of disease," *Annu. Rev. Biomed. Eng.*, vol. 8, pp. 537–565, 2006, doi: 10.1146/ANNUREV.BIOENG.8.061505.095802.
- [8] J. Schaefer, M. Lehne, J. Schepers, F. Prasser, and S. Thun, "The use of machine learning in rare diseases: a scoping review," doi: 10.1186/s13023-020-01424-6.
- [9] J. Schmidt, M. R. G. Marques, S. Botti, and M. A. L. Marques, "Recent advances and applications of machine learning in solid-state materials science," *npj Comput. Mater.* 2019 51, vol. 5, no. 1, pp. 1–36, Aug. 2019, doi: 10.1038/s41524-019-0221-0.
- [10] T. Freeman, "Deep learning for disease diagnosis confounded by image labels – Physics World." <https://physicsworld.com/a/deep-learning-for-disease-diagnosis-confounded-by-image-labels/> (accessed Jan. 13, 2022).
- [11] P. Galdi and R. Tagliaferri, "Data Mining: Accuracy and Error Measures for Classification and Prediction Neonatal MRI View project Computational methods for omics data View project Data Mining: Accuracy and Error Measures for Classification and Prediction," doi: 10.1016/B978-0-12-809633-8.20474-3.
- [12] J. M. Johnson and T. M. Khoshgoftaar, "Survey on deep learning with class imbalance,"

*J. Big Data*, vol. 6, no. 1, pp. 1–54, Dec. 2019, doi: 10.1186/S40537-019-0192-5/TABLES/18.

- [13] Wikipedia, “Accuracy paradox - Wikipedia.” [https://en.wikipedia.org/wiki/Accuracy\\_paradox](https://en.wikipedia.org/wiki/Accuracy_paradox) (accessed Jan. 14, 2022).
- [14] F. J. Valverde-Albacete and C. Peláez-Moreno, “100% Classification Accuracy Considered Harmful: The Normalized Information Transfer Factor Explains the Accuracy Paradox,” *PLoS One*, vol. 9, no. 1, Jan. 2014, doi: 10.1371/JOURNAL.PONE.0084217.
- [15] J. Brownlee, “10 Standard Datasets for Practicing Applied Machine Learning.” <https://machinelearningmastery.com/standard-machine-learning-datasets/> (accessed Jan. 12, 2022).
- [16] J. W. Smith, J. E. Everhart, W. C. Dickson, W. C. Knowler, and R. S. Johannes, “Using the ADAP Learning Algorithm to Forecast the Onset of Diabetes Mellitus,” *Proc. Annu. Symp. Comput. Appl. Med. Care*, p. 261, 1988, Accessed: Jan. 13, 2022. [Online]. Available: [/pmc/articles/PMC2245318/?report=abstract](https://pubmed.ncbi.nlm.nih.gov/2245318/).
- [17] F. Pedregosa FABIANPEDREGOSA *et al.*, “Scikit-learn: Machine Learning in Python Gaël Varoquaux Bertrand Thirion Vincent Dubourg Alexandre Passos PEDREGOSA, VAROQUAUX, GRAMFORT ET AL. Matthieu Perrot,” *J. Mach. Learn. Res.*, vol. 12, pp. 2825–2830, 2011, Accessed: Jan. 12, 2022. [Online]. Available: <http://scikit-learn.sourceforge.net>.
- [18] N. Erickson *et al.*, “AutoGluon-Tabular: Robust and Accurate AutoML for Structured Data,” Mar. 2020, Accessed: Jan. 12, 2022. [Online]. Available: <https://arxiv.org/abs/2003.06505v1>.
- [19] C. R. Harris *et al.*, “Array programming with NumPy,” *Nat.* 2020 5857825, vol. 585, no. 7825, pp. 357–362, Sep. 2020, doi: 10.1038/s41586-020-2649-2.
- [20] P. Virtanen *et al.*, “SciPy 1.0: fundamental algorithms for scientific computing in Python,” *Nat. Methods*, vol. 17, no. 3, pp. 261–272, Mar. 2020, doi: 10.1038/S41592-019-0686-2.
- [21] J. D. Hunter, “Matplotlib: A 2D graphics environment,” *Comput. Sci. & Eng.*, vol. 9, no. 3, pp. 90–95, 2007, doi: 10.1109/MCSE.2007.55.
- [22] SageMaker, “Amazon SageMaker – Machine Learning – Amazon Web Services.” <https://aws.amazon.com/sagemaker/> (accessed Jan. 13, 2022).
- [23] AWS, “Amazon SageMaker Amazon Sagemaker API Reference Amazon SageMaker: Amazon Sagemaker API Reference.”
- [24] Scikit-learn, “1. Supervised learning — scikit-learn 1.0.2 documentation.” [https://scikit-learn.org/stable/supervised\\_learning.html](https://scikit-learn.org/stable/supervised_learning.html) (accessed Jan. 14, 2022).