



PREDICTION OF DIABETES USING MACHINE LEARNING

Sanyam Jain¹, Uma Tomar², Jatin Chauhan³, Avas Gupta⁴, Nitin Chaudhary⁵

²Assistant Professor of Information Technology Engineering, Greater Noida Institute of Technology

^{1,3,4,5}Student, Department of Information Technology Engineering, Greater Noida Institute of Technology

Abstract: Diabetes is a common complaint caused by a set of metabolic affections where the sugar stages over drawn-out period is veritably high. It touches different organs of the mortal body which thus harm a huge number of the body's system, in precise the blood strains and jitters. Beforehand vaticination in similar complaint can be exact and save mortal life. To achieve the thing, this exploration work substantially discovers multitudinous factors associated to this complaint using machine literacy ways. Machine literacy styles give operative outgrowth to prize knowledge by erecting prognosticating models from individual medical datasets together from the diabetic cases. Scooping knowledge from similar data can

be precious to prognosticate diabetic cases. In this exploration, six popular used machine literacy ways, videlicet Random Forest (RF), Logistic Retrogression (LR), Naive Bayes (NB), C4.5 Decision Tree (DT), K-Nearest Neighbor (KNN), and Support Vector Machine (SVM) are compared in order to get outstanding machine literacy ways to read diabetes. Our new outgrowth shows that Support Vector Machine (SVM) achieved advanced delicacy compared to other machine literacy ways.

I. INTRODUCTION

Diabetic is a complaint that affects the hormone insulin, follow-on in abnormal metabolism of carbohydrates and advance way of sugar in the blood. This great blood sugar affects several organs of the mortal body which in turn complicates numerous source of the body, in precise the blood strains and jitters. The details of diabetic isn't nonetheless completely exposed, numerous experimenters supposed that both heritable rudiments and environmental goods are complex therein. As exposed by the International Diabetes Federation, extent of people having Diabetes stretched 382 million out of 2013 that takes 6.6 of the world's total adult population. According to the world healthcare medical data it has been probable that diabetic cases will be increased up to 490 billion within the time 2030. Likewise, diabetic is imaginably independent unproductive factor to micro-vascular snares. Diabetic cases are perhaps more unable against a hoisted threat of micro-vascular damage, in

this way long term difficulty of cardio-vascular complaint is the commanding reason of death. This micro-vascular detriment and hasty cardio vascular complaint eventually quick to retinopathy, nephropathy and neuropathy. Beforehand vaticination of similar complaint can be controlled over the conditions and save mortal life. To negotiate this thing, this exploration work substantially discovers the early vaticination of diabetes by taking into account colorful threat factors related to this complaint. For the restraint of the study we gathered individual dataset having 16 attributes diabetic of 2000 cases. These attributes are age, diet,

hyperactive- pressure, problem in vision, inheritable etc. In after part, we debate about these attributes with their conforming values. Grounded on these attributes, we figure vaticination model by means of colorful machine literacy ways to prognosticate diabetes mellitus. Machine literacy ways give well- organized result to excerpt knowledge by making prognosticating models from individual medical datasets composed from the diabetic cases. Haul out knowledge from similar data can be salutary to prognosticate diabetic cases. Innumerable machine literacy ways have the knack to prognosticate diabetes mellitus. Though it's veritably delicate to elect the stylish fashion to prognosticate grounded on similar attributes. Therefore for the determination of the study, we deal six popular machine learning

algorithms, videlicet Support Vector Machine (SVM), Naive Bayes (NB), K-Nearest Neighbor (KNN), Logistic Retrogression (LR), Random Forest (RF) and C4.5 decision tree (DT), on adult population data to prognosticate diabetes

II. RELATED WORK

Colorful experimenters have been shown variations in the area of diabetic by using machine literacy ways to prize knowledge from being medical data. For illustration, ALjumah et al. established a prophetic analysis model using support vector machine algorithm. In, Kavakiotis et al. used 10 fold cross confirmation as evaluation system in three different algorithms, including Logistic retrogression, Naive Bayes, and SVM, where SVM on condition that better performance and delicacy of 84 than other algorithms. In Zheng et al. applied Random Forest, KNN, Naive Bayes, SVM,



decision tree and logistic regression to prognosticate Diabetes at early stage, where drawing criteria can be bettered. Swarupa et al. applied J48, ANN, KNN, ZeroR and NB on colorful diabetes dataset. Pradeep et al. Applied Random Forest, KNN, SVM and J48 where J48 shows better performance than others. The bracket algorithms didn't assess using cross confirmation system. To prognosticate and control Diabetes Huang et al conversed three data mining styles, including IB1, Naïve Bayes and C4.5 in the time of 2000 to 2004. By smearing point selection fashion, the performance of IB1 and Naive Bayes handed better result. In Xue-Hui Meng et al. used three different data mining ways ANN, Logistic regression, and J48 to prognosticate the diabetic conditions using real world data sets.

Eventually it was concluded as J48 performs better delicacy than others. In this work, we examine real individual medical data grounded on multitudinous threat factors using popular machine learning bracket ways to assess their performance for prognosticating diabetes mellitus.

III. METHODOLOGY

In order to negotiate our thing study methodology includes of many stages, which are addendum of diabetes dataset with the applicable attributes of the cases, preprocessing the numeric value attributes, to smear several machine literacy ways and conforming prophetic analysis employing similar data. In the posterior, we fleetingly confer these stages.

Dataset and Attributes

In this work, datasets are collected among the Pima Indian feminine population close to Phoenix, Arizona. This explicit dataset has been wide employed in machine learning experiments and is presently obtainable through the UCI repository of normal datasets. This population has been studied unendingly by the National Institute of polygenic disorder, organic process and excretory organ. UCI repository contains 768 instances of observations and total nine attributes with no missing values rumored. information sets contains eight explicit variables that were thought-about high risk factors for the incidence of polygenic disorder, like range of times pregnant, plasma aldohexose concentration at a pair of hour in associate oral aldohexose tolerance take a look at (OGTT), pulsation force per unit area, a pair of hour bodily fluid hypoglycaemic agent, body mass index, polygenic disorder pedigree. All the patients during this datasets area unit feminine a minimum of twenty one years previous living close to Phoenix, Arizona. All attributes area unit numeric values except category is nominal sort. Attributes name and kinds area unit shown in table 1.

No	Nameofattributes	Type
1	Number of times pregnant	Numeric
2	Plasma glucose concentration a 2hours in an oral glucose to larence test	Numeric
3	Diastolic blood pressure	Numeric
4	Tricepsk in fold thickness	Numeric
5	2 hours eruminsulin	Numeric
6	Body mass index	Numeric
7	Diabetes pedigree function	Numeric

	Pregnancies	Glucose	BloodPressure	SkinThickness	Insulin	BMI	DiabetesPedigreefunction	Age	Outcome
0	2	138	82	35	0	35.6	0.127	47	1
1	0	84	82	31	125	30.2	0.233	23	0
2	0	145	0	0	0	44.2	0.630	31	1
3	0	135	68	42	250	42.3	0.365	24	1
4	1	139	82	41	450	40.7	0.535	21	0
5	0	173	78	32	255	46.5	1.199	58	0
6	4	99	72	17	0	25.6	0.254	28	0

Figure1 Sample dataset

It is sensible to see the correlations between the attributes. From the output graph below, the red round the diagonal suggests that attributes area unit related with one another. The yellow and inexperienced patches counsel some moderate correlation and therefore the blue boxes show negative correlations as shown below fig. 2.

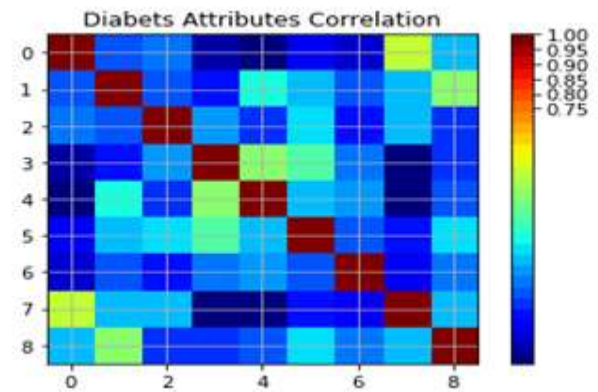


Figure2 Attributes Correlation

Next, we tend to visualize the information victimisation density plots to induce a way of the information distribution. From the outputs below, you'll see the information shows a general Gaussian distribution below fig 3.

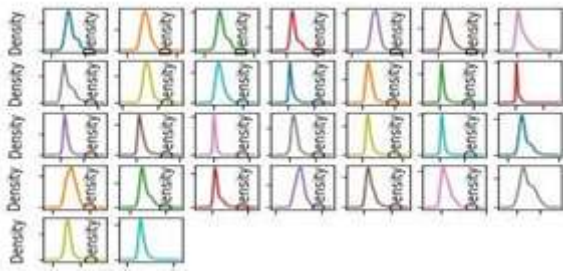


Figure 3 Density plots

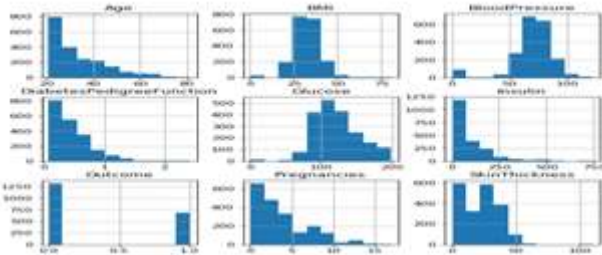


Figure 4 Histogram

According to Fig. 5 below, when process the matter we tend to collect the relevant information from the Diagnostic information Storage. we tend to then preprocess the information for the aim of building the prediction model. at the moment we tend to apply numerous machine learning techniques mentioned on top of on the coaching dataset. Finally,take a look at dataset is employed to live the performance of the techniques so as to settle on the most effective classifier for predicting diabetes.

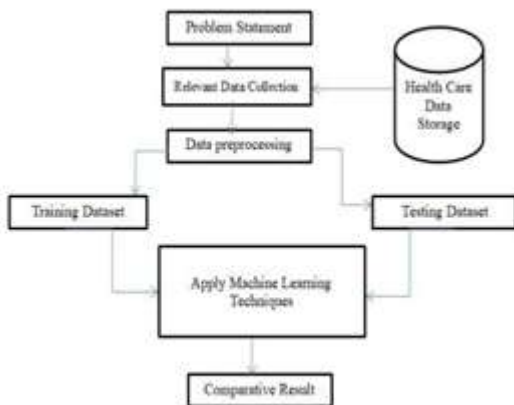


Figure 5 Showsan overview of the overall process of ourwork.

Apply Machine Learning Techniques

After the information has been prepared for demonstrating, we tend to use six common machine learning techniques to predict diabetes. Later we provide an summary of those

techniques. Support Vector Machines: is wide classification technique projected by J. Platt et. al. A Support Vector Machine is thought by characterizing the information by separating a hyper plane. SVM detaches entities in given categories. It may also acknowledge and classify instances that aren't supported by information. SVM isn't caring within the distribution of effort information of every category. The one delay of this formula is to execute multivariate analysis to supply a linear perform and another extension is learning to rank components to yield classification for individual components.

Naive Bayes: could be a common probabilistic classification technique projected by John et. al. . Naive Bayesn additionally known as Bayesian theorem could be a straightforward, effective and usually used machine learning classifier.

The formula calculates probabilistic results by investigation the regularity and combines the worth given in information set.nBy victimisation Bayesian theorem, it adopts that every one attributes area unit freelance and supported variable values of categories. In universe solicitation, the conditional independence hypothesis seldom holds true and offers well and a lot of sophisticate classifier results.

K-Nearest Neighbor Algorithm: K-nearest neighbor is easy regression and classification formula that used non constant technique projected by Aha et. al. . The formula trains all usable attributes and classifies new question supported their likeness live. to work out the space from purpose of interest to points in coaching information set it uses tree like system. The attribute is classed by its neighbors.

Decision Tree: could be a tree that gives powerful classification techniques to predict diabetes. each distinct space and have of the domain is named a category. associate input feature of the category attribute is labeled with the interior node during a tree. The leaf node of the tree is labeled by attribute and every attribute related to a target price. There area unit some common call tree algorithms area unit accessible to classify diabetic information in machine learning techniques, together with ID3, C4.5 , J48, C5, CART and CHAID. C4.5 provides extended options of ID3 call tree formula projected by Ross Quinlan et. al. [14]. C4.5 call tree uses same coaching information as ID3, during which learned perform is introduced.the training technique are often accustomed diagnose medical information to predict the worth of the choice attribute.

Logistic regression: could be a probabilistic applied mathematics model for fact-finding a dataset during which there area unit one or a lot of freelance variables that govern a result. In logistical regression, the variable is binary meaning one as (TRUE, patient, etc.) or zero as (FALSE, healthy, etc.). logistical regression produces the coefficients of a procedure to predict a log it adjustment of the chance of incidence of the characteristic of inquisitiveness.

Random forest: is cooperative classification theme



supported call Tree. At the coaching stage, it produces a huge range of trees and creates a forest of call Trees. At the testing stage, every tree of the forest predicts a category label for every information. once every tree predicts a category label, then the ultimate call for every take a look at information depends on wide command balloting. that category label gets the bulk of votes this label accepts to be At the testing stage, each tree of the forest predicts a class label for each data. When each tree predicts a class label, then the final decision for each test data depends on widely held voting. Which class label gets the majority of votes this label accepts to be the correct label allotted to the test data. This process is continual for each of data in the dataset.

Experimental Results And Discussion

To conduct the experiment six common used machine learning techniques, particularly Random Forest (RF), provision Regression (LR), Naive mathematician (NB), C4.5 call Tree (DT), K- Nearest Neighbor (KNN), and Support Vector Machine (SVM) area unit were used. Machine learning techniques were enforced in pycharm three.6. associate experimental result shows that the performance of SVM is considerably superior to different machine learning techniques for the classification of diabetic information. The experimental results may assist health care to require early bar and create higher clinical selections to regulate polygenic disorder and so save human life. to require into consideration further attributes and analysis for additional analysis is our future work.

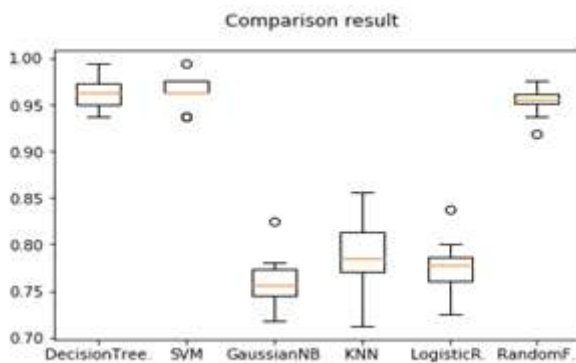


Figure6 Comparison result of Various Machine Learning Techniques

so as to assess the performance of varied machine learning techniques, we've got showed the prediction ends up in Fig. six on the premise of accuracy. The figure shows the results of studies like LR, SVM, NB, KNN, RF and C4.5 and SVM accomplishes higher results than different classifiers to predict DM. in keeping with Fig. 6, SVM achieves ninety six.4% on this dataset, that is larger than different learning techniques. This experimental result provides indication that SVM performs acceptable on medical datasets for the

determination of predicting polygenic disorder supported varied risk factors, deliberated within the earlier section.

Table 2

AccuracyResultsofVariousMachineLearningTechniques

Classification Algorithms	Accuracy
Support Vector Machine (SVM)	96.4%
Decision Tree(DT)	95.8%
Random Forest(RF)	95.3%
Naïve Bayes (NB)	76.0%
Logistic Regression(LR)	77.6%
K-Nearest Neighbors(KNN)	78.8%

Overall, we've got designated the most effective machine learning technique to predict polygenic disorder to realize high performance, supported the analysis criteria discuss on top of. All the techniques mentioned over area unit calculable on associate unseen testing diabetic dataset. The technique that accomplishes the best performance in terms of accuracy is taken into account to be the most effective alternative. supported Fig. 6, it is ascertained that SVM achieved the higher accuracy of ninety six.4 you must predict polygenic disorder utilizing a given medical dataset.

IV. CONCLUSION

In this work, we've got investigated the first prediction of polygenic disorder} by taking into consideration many risk factors associated with this disease exploitation machine learning techniques To predict polygenic disorder with efficiency, we've got done our investigation exploitation six common machine learning algorithms, particularly Support Vector Machine (SVM), Naive mathematician (NB), K-Nearest Neighbor (KNN), provision Regression (LR.), Random Forest (RF.) and C4.5 call tree, on adult population information to predict DM. The technique that accomplishes the best performance in terms of accuracy is taken into account to be the most effective alternative. supported Fig.6, it is ascertained that SVM achieved the higher accuracy of ninety six.4 you must predict polygenic disorder utilizing a given medical dataset.

V. REFERENCES

- [1]. V., A. K. and R., C. 2013. Classification of polygenic disorder malady exploitation Support Vector Machine. International Journal of Engineering analysis and Applications. 3, (April. 2013), 1797-1801.
- [2]. Carlo, B G., Valeria, M. and Jesús, D. C. 2011. The impact of polygenic disorder on health care prices in European country. skilled review of pharmacoeconomics & outcomes analysis. 11, (Dec. 2011),709-19.



- [3]. Nahla B., St. Andrew et al. 2010. Intelligible support vector machines for designation of DM. data Technology in Biomedicine, IEEE Transactions. 14, (July. 2010), 1114-20.
- [4]. Abdullah A. Aljumah et al., Application of knowledge mining: polygenic disorder health care in young and recent patients, Journal of King Saud University - pc and knowledge Sciences, Volume 25, Issue 2, July 2013, Pages 127-136
- [6]. Rani, A. Swarupa, and S. Jyothei. "Performance analysis of classification algorithms underneath totally different datasets." In Computing for property international Development (INDIACom), 2016 third International Conference on, pp. 1584-1589. IEEE, 2016.
- [7]. Kandhasamy, J. Pradeep, and S. Balamurali. "Performance analysis of classifier models to predict DM." Procedia computing forty seven (2015): 45-51.
- [8]. Y. Huang, P. McCullagh, N. Black, R. Harper, Feature choice and classification model construction on kind a pair of diabetic patients' data, computing in medication forty one (3) (2015) 251–262.
- [9]. Meng, X. H., Huang, Y. X., Rao, D. P., Zhang, Q., & Liu, Q. (2013). Comparison of 3 data processing models for predicting polygenic disorder or prediabetes by risk factors. The Kaohsiung journal of medical sciences, 29(2), 93-99.
- [10]. Platt, John C. "12 quick coaching of support vector machines exploitation consecutive lowest improvement." Advances in kernel strategies (1999): 185-208.
- [11]. John, George H., and Pat Langley. "Estimating continuous distributions in Bayesian classifiers." Proceedings of the Eleventh conference on Uncertainty in computing. Morgan Kaufmann Publishers INC., 1995.
- [12]. Aha, David W., Dennis Kibler, and Marc K. Albert. "Instance-based learning algorithms." Machine learning six.1 (1991): 37-66.
- [13]. Ross Quinlan (1993). C4.5: Programs for Machine Learning. Morgan Kaufmann Publishers, San Mateo.