# Prediction Of Diabetes Using Machine Learning Classification Algorithms

**Naveen Kishore  G, V.Rajesh,  A.Vamsi Akki Reddy, K.Sumedh,  T.Rajesh Sai Reddy**

**Abstract :** Diabetes is a sickness that takes place whilst glucose content in your blood is too excessive. Insulin is a  hormone made through the pancreas, facilitates to separate glucose from meals get into your body- cells for power. On this we used category set of rules techniques of the device mastering to are expecting the diabetes. Five machine getting to know algorithms namely SVM, selection Tree, NaiveBayes , Logistic Regression and KNN are used to hit upon diabetes. This may be capable of predict the chance degrees of diabetes and gives the first-class getting to know set of rules with better accuracy comparatively different algorithms.

———————————————  ◆  ———————————————

## 1.  INTRODUCTION

Diabetes isn't a hereditary disorder however heterogeneous group of disorder which could ultimately result in an boom of glucose within the blood and lack of glucose inside the urine. Diabetes is typically resulting from genetics, way of life and surroundings. Eating an dangerous weight loss plan, being overweight play role in developing the diabetes. High blood sugar tiers can also result in kidney diseases, coronary heart illnesses. The excess of sugar in the blood can harm the tiny blood vessels in your frame. Signs of diabetes are blurry imaginative and prescient , extreme hunger, unusual weight reduction, common urination and thirsty. In this paper,  parameters used within the facts set to locate the diabetes are Glucose, Blood pressure, pores and skin thickness, Insulin, Age. Huge volumes of statistics units are generated by health care industries. Those facts sets is a collection of patient information about the diabetes from the hospitals. Big records analytics is the processing which it examines the information units and exhibits the hidden information. Pima Indians Diabetes Database (PIDD), this dataset is taken from the national Institute of Diabetes and Digestive diseases. The objective of the dataset is to predict whether or not the patient has diabetes or not, primarily based on diagnostic measurements in the dataset. Several constraints were taken from the massive database.
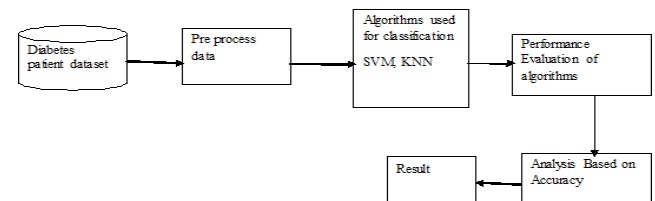
## 2.  LITERATURE SURVEY

SVM offers the great output whilst in comparison to different strategies for lots disorder prediction within the gift gadget studying algorithms [19][1].SVM is used for type and regression issues[17][29]. statistics mining is useful for getting significant information[14][5]. A massive data is generated from scientific institutes each 12 months[6] [2]. Many people have proposed one-of-a-kind structures for the prediction of diabetics. Orbi et al is one among them who proposed a machine for the prediction of diabetics [16][3].

_____

* *Naveen Kishore  G,  V.Rajesh, A.Vamsi Akki Reddy, K.Sumedh,*
* *Assoc.prof,  Professor,  T.Rajesh Sai Reddy  KLEF  KLEF naveengattim@live.com*

Many data sets are to be had for different diseases[9][21]. Mining of those datasets offers beneficial facts[18]. the principle goal of this device is to predict diabetes primarily based on the candidate struggling at unique age, with higher accuracy the use of decision Tree.[7][11][8]. usually choice timber are supervised studying strategies used for each category and regression.[20][26]. KNN is also used for the classification and regression.[31]. Genetic programming (GP) is utilized by Pradhan for testing and education of the database for the prediction of the diabetes[13][10]. however the output effects received within the Genetic programming approach has less accuracy whilst as compared to other techniques and additionally designed a prediction model for the diabetes disease with two models specifically ANN(synthetic Neural Networks) and the second one is FBS(Fasting Blood Sugar)[15][24]. The algorithms at the danger of diabetes mellitus become proposed by Nongyao[27]. He proposed 4 famend machine gaining knowledge of classifications techniques specifically decision Tree, artificial Neural Networks, Logistic Regression and Naive Bayes[22][30]. Bagging and Boosting strategies are used for growing the robustness of the designed version.[28][23]

## 3.  METHODOLOGY
### 3.1



### 3.2    Algorithms used
* SVM
* Decision Tree
* KNN
* Logistic Regression
* Random Forest

### 3.2.1    SVM

SVM is the standard supervised learning algorithm. It does the complex data transformations and separates the data based on the outputs and it can be used for both classifications and regression challenges. In SVM there are different hyper planes which divides the data. In this

wehave to select the hyper plane which divides the class better. To find the better hyper plane you have to calculate the distance between the planes and the data which is called Margin. If the distance between the classes is low then the chance of miss conception is high and vice versa. So we need to select the class which has the high margin. The performance of SVM algorithm version for prediction of diabetes the use of confusion matrix is as follows

**Table 1.** *Confusion matrix for SVM*

| Tested Positive | 70 |
|---|---|
| Tested Negative | 141 |

### 3.2.2    Decision Tree

Decision tree set of rules is a supervised studying algorithm. It is used to remedy the type problems. In this algorithm whole facts is represented inside the shape of tree in which every leaf is corresponds to the class label and attribute are corresponds to inner node of the tree. The fundamental venture is to discover the foundation node in each stage. The overall performance of choice Tree set of rules model for prediction of diabetes is as follows

**Table 2.** *Confusion matrix for Decision Tree*

| Tested Positive | 79 |
|---|---|
| Tested Negative | 113 |

### 3.2.3    KNN

KNN is the one of the best set of rules used in gadget getting to know. In this set of rules the complete dataset is sorted. The data is divided into classes, if other data is want to classify then it finds the neighbours of that element based on the majority number of votes for the label . Initialize the data it calculates the distance between the classes and finding neighbours and voting for labels. The overall performance of KNN set of rules version for prediction of diabetes is as follows.

**Table 3.**  *Confusion matrix for KNN*

| Tested Positive | 44 |
|---|---|
| Tested Negative | 148 |

### 3.2.4    Logistic Regression:

Logistic regression is a machine getting to know classifier. This set of rules is used to split the observations for discrete classes. The outputs given by using the logistic regression is based totally on the opportunity feature. It uses the fee function that's known as as "sigma" characteristic. Sigma function is more complex than the normal linear function. Logistic regression limit the cost function value between 0 to 1.

**Table 4.**  *Confusion matrix for Logistic Regression*

| Tested Positive | 72 |
|---|---|
| Tested Negative | 120 |

### 3.2.5    Random Forest:

Random forest is a supervised system getting to know set of rules. It's also used to remedy classification and regression additionally. In this algorithm it consists of the trees. The number of tree structures present in the data is directly proportional to the accuracy of the result. Each internal node within the tree corresponds to an attribute and every leaf node represents class label.

**Table 5**.  *Confusion matrix for Random forest*

| Tested Positive | 68 |
|---|---|
| Tested Negative | 124 |

## 4 DATA SET

The statistics set is taken from Pima Indians Dataset Database (PIDD) that's to be had at UCI device studying. The statistics set has many impartial variables along with Glucose, Blood pressure, skin Thickness, BMI etc. Records set is trained to get the accurate end result and similarly it is tested.

**Table 4**. *Dataset illustration*

| Database | No. of Traits | No. of occurences |
|---|---|---|
| PIDD | 8 | 769 |

Accuracy is measured by the formula given below
$$ACCURACY = TP + TN \ / \ TP + FP + TN + FN$$
    Where,
    TP: True Positive
    TN: True Negative
    FP: False Positive
    FN: False Negative

## 5 RESULTS

In the given  Table we can see the algorithm which gives the better accuracy  on the data set.

**Table 5.** *Result Description*

| Algorithms | Accuracy | Misclassification |
|---|---|---|
| SVM | 73.43 | 26.5 |
| Decision Tree | 72.91 | 27.08 |
| Random Forest | 74.4 | 25.5 |
| KNN | 71.3 | 28.64 |
| Logistic Regression | 72.39 | 27.60 |

## 6 CONCLUSION

At last by using all these five machine learning algorithms we had measured different parameters within the dataset

aand we had came through better accuracy rate with random forest with nearly 75%. This work can be extended by adding any other algorithm which can give better accuracy than random forest.

# 7 REFERENCES

[1] Yu, W., Liu, T., Valdez, R., Gwinn, M., Khoury, M.J., 2010. Application of support vector machine modeling for prediction of common diseases: The case of diabetes and pre-diabetes. BMC Medical Informatics and Decision Making 10. doi:10.1186/1472-6947-10-16, arXiv:arXiv:1011.1669v3

[2] .N.AdityaSundar, P.PushpaLatha, M.Rama Chandra "Performance Analysis of Classification Data Mining Techniques over Heart Disease Database ", International Journal of Engineering Science and Advanced Technology, Vol 2, Issue 3, p470-478,May-June 2012

[3] H. C. Koh and G. Tan, "Data Mining Application in Healthcare", Journal of Healthcare Information Management, vol. 19, no. 2,(2005).

[4] Andreas G. K. Janecek ,WilfriedN.Gansterer and Michael A.Demel,"On the Relationship Between Feature Selection and Classification Accuracy",JMLR: Workshop and Conference Proceedings 4: 90-105

[5] Arora, R., Suman, 2012. Comparative Analysis of Classification Algorithms on Different Datasets using WEKA. International Journal of Computer Applications 54, 21–25. doi:10.5120/8626-2492

[6] I. H. Witten, E. Frank, and M. A. Hall, Data Mining: Practical Machine Learning Tools and Techniques, 3rd ed. Morgan Kaufmann, 2011

[7] Dhomse Kanchan B., M.K.M., 2016. Study of Machine Learning Algorithms for Special Disease Prediction using Principal of Component Analysis, in: 2016 International Conference on Global Trends in Signal Processing, Information Computing and Communication, IEEE. pp. 5–10.

[8] Esposito, F., Malerba, D., Semeraro, G., Kay, J., 1997. A comparative analysis of methods for pruning decision trees. IEEE Transactions on Pattern Analysis and Machine Intelligence 19, 476–491. doi:10.1109/34.589207

[9] Aljumah, A.A., Ahamad, M.G., Siddiqui, M.K., 2013. Application of data mining: Diabetes health care in young and old patients. Journal of King Saud University - Computer and Information Sciences 25, 127–136. doi:10.1016/j.jksuci.2012.10.003.

[10] Kayaer K., Tulay, 2003. Medical diagnosis on Pima Indian diabetes using general regression neural networks, in: Proceedings of the international conference on artificial neural networks and neural information processing (ICANN/ICONIP), pp. 181-184

[11] Han, J., Rodriguez, J.C., Beheshti, M., 2008. Discovering decision tree based diabetes prediction model, in: International Conference on Advanced Software Engineering and Its Applications, Springer. pp. 99–109

[12] Garner, S.R., 1995. Weka: The Waikato Environment for Knowledge Analysis, in: Proceedings of the New Zealand computer science research students conference, Citeseer. pp. 57–64

[13] Bamnote, M.P., G.R., 2014. Design of Classifier for Detection of Diabetes Mellitus Using Genetic Programming. Advances in Intelligent Systems and Computing 1, 763–770. doi:10.1007/978-3-319-11933-5.

[14] Iyer, A., S, J., Sumbaly, R., 2015. Diagnosis of Diabetes Using Classification Mining Techniques. International Journal of Data Mining & Knowledge Management Process 5, 1–14. doi:10.5121/ijdkp.2015.5101, arXiv:1502.03774.

[15] Pradhan, P.M.A., Bamnote, G.R., Tribhuvan, V., Jadhav, K., Chabukswar, V., Dhobale, V., 2012. A Genetic Programming Approach for Detection of Diabetes. International Journal Of Computational Engineering Research 2, 91–94

[16] Perveen, S., Shahbaz, M., Guergachi, A., Keshavjee, K., 2016. Performance Analysis of Data Mining Classification Techniques to Predict Diabetes. Procedia Computer Science 82, 115–121. doi:10.1016/j.procs.2016.04.016

[17] Kumari, V.A., Chitra, R., 2013. Classification Of Diabetes Disease Using Support Vector Machine. International Journal of Engineering Research and Applications (IJERA) www.ijera.com 3, 1797–1801.

[18] Nai-Arun, N., Moungmai, R., 2015. Comparison of Classifiers for the Risk of Diabetes Prediction. Procedia Computer Science 69, 132–142. doi:10.1016/j.procs.2015.10.014

[19] Orabi, K.M., Kamal, Y.M., Rabah, T.M., 2016. Early Predictive System for Diabetes Mellitus Disease, in: Industrial Conference on Data Mining, Springer. Springer. pp. 420–427

[20] Priyam, A., Gupta, R., Rathee, A., Srivastava, S., 2013. Comparative Analysis of Decision Tree Classification Algorithms. International Journal of Current Engineering and Technology Vol.3, 334–337. doi:JUNE 2013, arXiv:ISSN 2277 – 4106

[21] Kavakiotis, I., Tsave, O., Salifoglou, A., Maglaveras, N., Vlahavas, I., Chouvarda, I., 2017. Machine Learning and Data Mining Methods in Diabetes Research. Computational and Structural Biotechnology Journal 15, 104–116. doi:10.1016/j.csbj.2016.12.005.

[22] Ray, S., 2017. 6 Easy Steps to Learn Naive Bayes Algorithm (with code in Python)

[23] Sisodia, D., Shrivastava, S.K., Jain, R.C., 2010. ISVM for face recognition. Proceedings - 2010 International Conference on Computational Intelligence and Communication Networks, CICN 2010 , 554–559doi:10.1109/CICN.2010.109.

[24] .Bhargavi V.R., Senapati R.K., Curvelet fusion enhacement based evaluation of diabetic retinopathy by the identification of exudates in optic color fundus images ,2016, Biomedical Engineering - Applications, Basis and Communications, Vol: 28, Issue: 6, ISSN 10162372

[25] Bhargavi V R., Senapati R.K., Swain G., Prasad P.M.K., Computerized diabetic patient's fundus

image screening for lesion regions detection and grading ,2016, Biomedical Research (India), Vol: 2016, Issue: Special Issue 1, pp: S443 - S449, ISSN 0970938X

[26] Kumar K.V.V., Kishore P.V.V., Indian classical dance mudra classification using HOG features

[28] impedance analysis in health care systems,2017 Journal of Medical Imaging and Health Informatics, Vol:7, issue:6, pp: 1126-1138, DOI: 10.1166/jmihi.2017.2211, ISSN: 21567018

[29] . .Rao G.A., Syamala K., Kishore P.V.V., Sastry A.S.C.S. .," Deep convolutional neural networks for sign language recognition ", 2018, International Journal of Engineering and Technology(UAE) ,Vol: 7 ,Issue: 1.5 Special Issue 5 ,pp: 62 to:: 70 ,DOI: ,ISSN: 2227524X

[30] . .Reddy S.S., Suman M., Prakash K.N. .," Micro aneurysms detection using artificial neural networks ", 2018, Lecture Notes in Electrical

[33] earch, ISSN:16740440, Vol No:46, 2019, pp:72 - 77.

and SVM Classifier,2017 International Journal of Electrical and Computer Engineering, Vol:7, issue:5, pp: 2537-2546, DOI: 10.11591/ijece. v7i5. pp2537-2546, ISSN: 20888708

[27] . Mirza S.S., Rahman M.Z.U., Efficient adaptive filtering techniques for thoracic electrical bio-Engineering ,Vol: 471 ,Issue: ,pp: 273 to:: 282 ,DOI: 10.1007/978-981-10-7329-8_28 ,ISSN: 18761100 9.78981E+12

[31] 209. Putluri S., Ur Rahman M.Z., Fathima S.Y. .," Cloud-based adaptive exon prediction for DNA analysis ", 2018, Lecture Notes in Electrical Engineering ,Vol: 434 ,Issue: ,pp: 409 to:: 417 ,DOI: 10.1007/978-981-10-4280-5_43 ,ISSN: 18761100 9.78981E+12

[32] Rajesh V., Saikumar K., Ahammad S.K.H., "A telemedicine technology for cardiovascular patients diagnosis feature using knn-mpm algorithm", Journal of International Pharmaceutical Res