



A Diabetic Prediction System Based on Mean Shift Clustering

Satyanarayana Murthy Teki^{1*}, Kuncham Venkata Sriharsha², Mohan Krishna Varma Nandimandalam³

¹ Computer Science and Engineering (CSE), Bapatla Engineering College, Bapatla 522102, Andhra Pradesh, India

² Computer Science and Engineering (CSE), Centurion University, Visakhapatnam 531173, Andhra Pradesh, India

³ Ocean IT Convergence Technology Research Lab, Hoseo University, Asan 31006, Korea

Corresponding Author Email: murthyteki@gmail.com

<https://doi.org/10.18280/isi.260210>

ABSTRACT

Received: 10 November 2020

Accepted: 23 February 2021

Keywords:

diabetes, mean-shift clustering, naive bayes, clustering

An abnormal rise in glucose levels may lead to diabetes. Around 30 million people are diagnosed with this disease in our country. In this perspective Indian Council of Medical Research funded by Registry of People with diabetes in India have taken an initiative and come up with numerous solutions but unfortunately neither of them has taken shape. Initially, the behavior of chemical reaction between glucose with chemical agent is estimated and tracked in the region of interest via mean shift algorithm using spatial and range information. This color change is related to plasma glucose concentration (plas), diastolic blood pressure, (pres.) Triceps skin fold thickness (skin), 2_hour serum insulin (insu), Body mass index and age. These features obtained from these 768 instances are classified using Naïve Bayes Algorithm. The results are compared with our previous work, an integrated system of K means and Naïve Bayes approach in terms of sensitivity, specificity, precision, and F-measure. It is worth noticing that our integration of mean-shift clustering and classification gives promising results with an utmost accuracy rate of 99.42% even after removing nearby duplicates in predefined clusters.

1. INTRODUCTION

Diabetes Mellitus (popularly known as Diabetes), a chronic medical condition caused due to abnormal sugar levels in blood. In normal cases, this glucose levels is controlled by insulin. But unfortunately in patient, insufficient production of insulin results Diabetes mellitus. This causes blindness, kidney failure and nerve damage. It also accelerates coronary heart diseases and other blood vessels in the body Porter and Green [1], Li et al. [2] and Panzarasa [3]. Preliminary results from a large community study conducted by the Indian Medical Research Council (ICMR) have shown that a lower proportion of the population is affected by diabetes in the northern states of our country than in southern India. In Hyderabad in particular, when compared with other metropolitan cities in India according to the National Urban Survey, the figure is approximately 16.6 per cent higher. Kuzuya et al. [4] work shows that pancreas islet cells can regenerate and early pathological changes in diabetes can be reversible. The analysis has already been carried out. It can contribute to appreciable regeneration of weakened Island Cells, and the importance of rapid diabetes development becomes evident when young people with more durable tissues are properly handled. As mentioned in Matheus et al. [5], data mining nowadays considered as statistical interface in predicting the trends of many kinds of diseases and illness. Rather than depending up on knowledge and experience, doctors are trying to invest their time much on knowledge discovery in database for predicting trends that would give better scope for patient diagnosis with advancements and innovations in the management system of medical databases, large volumes of medical data have evolved. This gives scope

for a researcher to assist health practitioners in research and diabetes prevention using data mining techniques. These data mining techniques are used in the analysis of various diseases including diabetes, cancer, heart disease, and kidney-related diseases for medical purposes. In order to achieve the best classification accuracy, abundant algorithms and diverse approaches like Naive Bayes, Artificial Neural Networks (ANN), Support Vector Machine (SVM), and Decision trees have been applied. In comparison to other diabetes prediction methods, when applied on larger datasets, Naive Bayes was considerably successful. Nevertheless, researchers are trying to improve the performance of Naive Bayes classification problems. As numbers of clusters are determined with respect to the data that is going to be handled, in our proposed work, we have used mean shift clustering as a preprocessing step to Naïve Bayes classification. This system is proposed based on the results of the Waikato Environment for Knowledge Analysis (WEKA) tool. The rest of the paper is organized as follows, section-2 contains literature study, and section-3 describes the knowledge discovery of database and existing classification approaches. Section-4 illustrates the proposed system architecture and its algorithm. Section-5 explains the experimental result of integrated Naive Bayes and mean shift clustering. Early diagnosis, early treatment and the new diagnostic criteria of diabetes mellitus. *British journal of nutrition*, 84(S2), S177-S181.

2. RELATED WORK

Lowanichchai et al. [6] have suggested a knowledge-based DSS for diabetic analyzer using the decision tree. This

Random Tree technique achieved high accuracy in the classification, i.e., 99.60 percent as compared with the medical diagnosis. Guo et al. [7] suggested seeking information from health repositories is important for a good diagnosis. An updated classifier for the Naive Bayes model was used here. Finally, the Waikato setting for Information Analysis tool was used, and the resulting model is 72.3% accurate. According to Guariguata et al. [8], Bellazzi et al. [9], Al Jarullah [10], National Public health on diabetes released that women reported highest with 9.6 million having diabetes. By 2050, this would rise from 17 million to 29 million. Six specific diabetic disease forecasting methods have been suggested by Krishnaveni and Sudha [11]. Those techniques are Naive Bayes, (k-Nearest Neighbor (KNN), SVM with Linear Kernel, SVM with Radial Basis Function (RBF) Kernel SVM with K-Means; et SVM with Fuzzy C-Means. Outcome is 76.3% with discrimination study, 71.1% with KNN, 76.1% with Naive Bay, 74.1% with Linear Kernel SVM, 74.1% with RBF Kernel SVM, 96.71% with K-Medium S VM and 94% with FuzzyC-Medium SVM. Several authors, however, such as Al Jarullah [10], Patil et al. [12], Raj and Rajan [13], Balpande and Wajgi [14] used different approaches to achieve the highest prediction efficiency. Our proposed integrated system (combination of K means clustering and Naïve Bayes classification is expected to achieve 95.416% accuracy rate even after removing nearby duplicates in predefined clusters.

3. DATASET DESCRIPTION

The Pima Indian Diabetes (PID) data set has been taken from the National Diabetes, Digestive and Cricket Disease Institute to determine whether an individual is affected or not by diabetes mellitus, which includes 768 samples along with 8 input attributes and 1 output attribute. (Class Label). Patients consist of women who are at least 21 years old from the Indian community of Pima. Dataset is saved in Attribute-Relation File Format (ARFF) format because it is a native method for storing data. Each format includes a list of rows and commas separate the attribute values for each row.

4. METHODOLOGY

This work utilizes four state-of-the-art. The following section gives a brief outline.

4.1 Knowledge discovery process

The Knowledge Discovery Process (KDD) process involves three stages. Data base information discovery (KDD). Data Preprocessing is the first step. The entire data will be formatted in this step in the database. Data collection, data cleaning, generation of attributes, normalization are the tasks involved in this stage. Step two is data mining where the use of various searching algorithms extracts useful information and/or patterns from the database. The final stage is the post-processing of data: formatted results are presented. Data mining tasks include differential analysis, generation of association rules, and identification of outlines, description, clustering and classification. We will discuss in detail about clustering and classification techniques. Classification techniques organize data into groups/classes. Naive Bayes, a well-known classification algorithm is derived Bayes theorem.

Let S be the training set of tuples, and D is the data sample, which do not know the class label. Suppose there are M classes say, $m_1, m_2, m_3, \dots, m_n$ and the classification is to derive maximum Posterior Probability i.e $Max[(P(m_i/D))]$.

$$P(m_i/D) = \frac{P\left(\frac{D}{m_i}\right)P(m_i)}{P(D)}$$

where, $P(m_i/D)$ is the Posterior probability of target class. $P(m_i)$ is called prior probability of class. $P\left(\frac{D}{m_i}\right)$ is the likelihood which is the probability of predictor of a given class. and $P(D)$ is the Prior probability of predictor of a given class.

4.2 Data preprocessing

Algorithm 1: Data Preprocessing

1. **Input:**
PID Dataset of instances (p).
 2. **Output:**
Preprocessed Dataset of instances (p) that are filtered.
 3. **Begin:**
 4. *Step 1: Load the dataset (p).*
 5. *Step 2: Find instance hardness values for each instance using five heuristics.*
 6. *Step 3: Misclassified instance(s)= instances with high instance hardness value.*
 7. *Step 4: remove the Misclassified instance(s).*
 8. *Step 4: Filter the inconsistent and noisy data using Naïve bayes.*
 9. *Step 5: return the preprocessed data.*
 10. **End**
-

5. PROPOSED WORK

The input for the proposed system is a PID data set and the output is one class that represents Healthy or Diabetic. There are 2 steps to the proposed diabetes prediction system: preprocessing and classification based on the clusters. The proposed method uses Pima Indian Diabetes data set (PIDD), with classification results of 76.3% using a standard Naive Bay and 99.42% using an integrated Naive Bayes mean -shift classification system. In order to identify and remove the misclassified instances from database, we have used Smith and Martinez [15], PRISM method. Initially, instance hardness will be assigned to each instance to determine which instances are intrinsically difficult to correctly identify. The instance hardness for each instance p is defined as:

$$instance\ hardness(p) = \frac{\sum_i^N incorrect(LA_i, p)}{N} \quad (1)$$

where, p is a data instance, N is the number of learning algorithms. And $incorrect(LA, x)$ is a function that returns '1', if an instance p is misclassified by learning algorithms LA and 0 otherwise. Instead of using all the nine learning algorithms to check whether an instance is correctly classified or not, we use five heuristics proposed by Smith and Martinez [15] to predict the instance hardness for each instance p . The heuristics are as follows.

Heuristic 1: *k*-Disagreeing Neighbors (*kDN*):

$$kDN(p) = \frac{|\{y: y \in kNN(p) \wedge t(y) \neq t(p)\}|}{k} \quad (2)$$

where, $kDN(p)$ is the set of k nearest neighbors of p and $t(p)$ is the target associated with p .

Heuristic 2: *DisjunctSize* (*DS*):

$$\frac{|\mathit{disjunct}(p)| - 1}{\max_{y \in D} |\mathit{disjunct}(y)| - 1} \quad (3)$$

where the function $\mathit{disjunct}(p)$ returns the disjunct that covers instance p and D is the dataset that comprises instance p .

Heuristic 3 *Disjunct Class Percentage* (*DCP*):

$$DCP(p) = \frac{|\{z: z \in \mathit{disjunct}(p) \wedge t(z) = t(p)\}|}{|\mathit{disjunct}(p)|} \quad (4)$$

Heuristic 4 *Class Likelihood* (*CL*):

$$CL(p, t(p)) = \prod_i^{p_i} P(p_i | t(p)) \quad (5)$$

where, p_i is the value of instance p on its i^{th} attribute.

Heuristic 5 *Class Likelihood Difference* (*CLD*):

$$(p, t(p)) = CL(p, t(p)) - \underset{y \in Y - t(p)}{\operatorname{argmax}} CL(p, y) \quad (6)$$

Using these heuristics, instance hardness for each instance p , is identified and are filtered based on their high instance hardness values. Further using Naïve Bayes, noisy, empty and other inconsistent data is removed to boost the performance of decisions. Figure 1 illustrates the proposed system architecture.

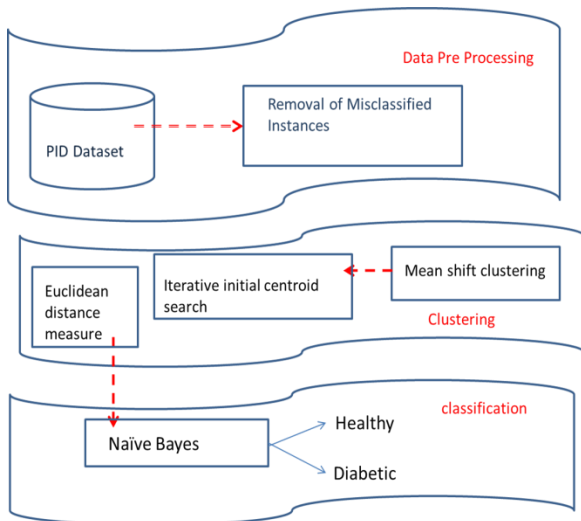


Figure 1. Proposed architecture

5.1 Pre-processing

There is no need to determine fixed number of clusters, automatically the mean shift cluster discovers the clusters and their corresponding center points by sliding window mechanism. And then based on the Euclidean distance

measure, Naïve Bayes classifier predicts the diabetic from healthy patient. However, the selection of initial cluster centers affects the results. Among those methods, Iterative Initial Centroid Search is used in this manuscript. Algorithm 2 illustrates the steps in mean shift clustering.

Algorithm 2: Mean Shift based Naive Bayes

-
- 1. Input:**
 2. Preprocessed data set(D)
 - 3. Output:**
A Reduced dataset that returns classification result.
 - Begin:**
 4. Step 1: Apply ten-fold cross-validations.
 5. Step 2: Segregate reduced data set into trained and testing data sets.
 6. Step 3: Train the classifier using trained data and apply it to the testing data.
 7. Step 4: Attribute selection for clustering.
 8. Step 4: Select the window in random fashion over data points centered at ‘C’ and radius ‘r’ as kernel.
 9. Step 5: **For** $i=1$ to $N(\text{iteration}(s))$.
 10. Calculate the center of gravity (mean vector $m_r(p)$): dataset instance.
 11. Shift the search window to the mean vector.
 12. Repeat step 10 until convergence. ($\nabla f(p) = 0$).
 13. Best centroid for each cluster is obtained.
 - 14. End**
-

Preprocessing represents that:

1. Input data points.
2. Data points after clustering.
3. Removal of Noise and Inconsistent Data.
4. Reduced Dataset.

5.2 Proposed architecture

The model consists of a pipeline of generator and discriminator networks. The generator is a U-Net network that has an encoder and decoder network with skip connections. The down sampling occurs in the encoder and the up sampling occurs in the decoder network. The features from same level of encoder is passed to the decoder output during up sampling.

5.3 Initial centroid selection method

Our intention is to find the final centroid and their corresponding groups as part of clustering. In this perspective, we use a non-parametric clustering technique, a mean shift algorithm. This technique neither requires prior knowledge of clusters nor the shape of clusters. Given p instances $p_i, i=1, 2 \dots n$ on a d dimensional space, R^d . For given set of datapoints (p instances) we chose a starting point and window randomly in the first iteration.

The multivariate kernel density estimate obtained with kernel $K(p)$ and window radius ‘r’ is,

$$f(p) = \frac{1}{nr^d} \sum_{i=1}^n K\left(\frac{p - p_i}{r}\right) \quad (7)$$

To evaluate the kernel $k(x)$ profile that is suitable for the radially symmetric kernels, it should satisfy the relation:

$$K(p) = c_{k,d} k(\|p\|^2) \quad (8)$$

where, $c_{k,d}$ is normalization constant.

The gradient of the density estimator is:

$$\begin{aligned} \nabla f(p) &= \frac{2c_{k,d}}{nr^{d+2}} \sum_{i=1}^n (p_i - p) g\left(\left\|\frac{p-p_i}{r}\right\|^2\right) \\ &= \frac{2c_{k,d}}{nr^{d+2}} \left[\sum_{i=1}^n \left(\left\|\frac{p-p_i}{r}\right\|^2\right) \right] \left[\frac{\sum_{i=1}^n p_i g\left(\left\|\frac{p-p_i}{r}\right\|^2\right)}{\sum_{i=1}^n g\left(\left\|\frac{p-p_i}{r}\right\|^2\right)} \right] \end{aligned} \quad (9)$$

where, $g(s) = -k'(s)$. The first term is proportional to the density estimate at p computed with kernel $G(p) = c_{g,d} g(\|p\|^2)$ and the second term.

$$m_r(p) = \frac{\sum_{i=1}^n p_i g\left(\left\|\frac{p-p_i}{r}\right\|^2\right)}{\sum_{i=1}^n g\left(\left\|\frac{p-p_i}{r}\right\|^2\right)} - p \quad (10)$$

The mean shift vector always points toward the direction of the maximum increase in the density. The mean shift procedure, obtained by successive computation of mean shift vector $m_r(p^t)$ and translation of the window $p^{t+1} = p^t + m_r(p^t)$ until it is converged. At this point the gradient of density function becomes zero. Initially we have used kNN (k nearest neighbor) algorithm to find window size (r). If $p_{i,k}$ is the k nearest neighbor of p_i then bandwidth is calculated as:

$$r_i = \|p_i - p_{i,k}\| \quad (11)$$

As mentioned in Zhao et al. [16] choice of window size always influences convergence rate and number of clusters. In subsequent iterations, in order to fit the major cluster size, we keep updating the window size and shape in each iteration without any provisional estimation. When using kNN to determine 'r', choice of k influences the value of 'r'. For achieving better results, it is always essential to increase the value of 'r' when dimension of data increases. But at the same time Mean shift might not work well with higher dimensions. In higher dimensions, number of local maxima is pretty goes to high and as a result it converges to local optima soon.

As we are dealing with low dimensional data, we are able to achieve optimal convergence rate. The results are tabulated in Table 1. It is inferred that From 268 entries, around 34.89% diabetic patients are effected and from 500 entities, 65.11% patients are not affected.

Table 1. Naive Bayes classification results

CLASS	Entries	No of Persons (in %)
Diabetic (P)	268	34.89
Non-diabetic (N)	500	65.11

6. EXPERIMENTAL RESULTS

6.1 Experiment setup

The PID data set has been taken from the National Institute of Diabetes, Digestive and Kidney Diseases. The dataset comprises 768 samples with each of 8 input attributes and one output attribute (class Label). The output Label is to predict whether a person is affected with diabetes mellitus or not based on the symptoms. Patients include females aged at

least 21 years from Pima Indian Heritage. The dataset is stored in ARFF format. The format comprises list of rows, and the attribute values for each row are separated by commas. The PID dataset contains 9 attributes which represent the symptoms such as Number of times Pregnant (PGY), Glucose levels or concentration in plasma (PLSM), Blood Pressure (BP), Triceps skin Fold thickness (SKIN), Serum-Insulin levels (INS), Body mass index (BMI), Diabetes Pedigree Function (PDF), Age of the Women (AGE) and the Class label (CLASS), which is the only output attribute. In this data set, there are no missing values, but there were some zeros included as missing values. Among these missing values, 5 patients had a sugar level of 0, 11 more had body mass index record as 0, 28 patients had a diastolic blood pressure of 0, 140 others had serum insulin levels at 0, and 192 others had a skin fold thickness value as 0. Performance can be measured with tenfold cross-validation. The dataset is divided into 10 equal partitions during the cross validation process, each method runs 10 times. Every time a different partition is used as a dataset testing, the other 9 partitions are used as data set training. Precision, accuracy, sensitivity, specificity, recall, F-measure and error rate metrics are considered for performance assessment.

Error Rate: calculated as the ratio of all incorrectly predicted instances to the total number of instances.

$$\text{Error rate} = \frac{(FP+FN)}{TOTAL} \text{ OR } \frac{(FP+FN)}{P+N} \quad (12)$$

Accuracy: calculated as ratio of all correctly predicted instances to the total number of instances.

$$\text{Accuracy} = \frac{(TP+TN)}{TOTAL} \text{ OR } \frac{(TP+TN)}{P+N} \quad (13)$$

Sensitivity or Recall: classifies positive instances. It is also called a Recall or True Positive rate (TPR).

$$\text{Sensitivity} = \frac{(TP)}{P} \quad (14)$$

Specificity: classifies negative instances. It is also called the True Negative rate (TNR).

$$\text{Specificity} = \frac{(TN)}{N} \quad (15)$$

Precision: defined as the ratio of the number of correct positively predicted instances to the total positive instances. It is also called Positive Predictive values (PPV).

$$\text{Precision} = \frac{(TP)}{TP+FP} \quad (16)$$

F-measure: calculates the harmonic mean of Precision and recall.

6.2 Confusion matrix

It identifies the correctly and incorrectly classified instances. The 2x2 matrix representation of confusion matrix as shown in below Table 2.

True positive (TP): The set of positive instances that were correctly classified.

True Negative (TN): The set of negative instances that are correctly classified.

False Positive (FP): These are the set of negative instances that are misclassified as positive.

False Negative (FN): These are the set of positive instances that are misclassified as

Table 2. Comparison of proposed model with recent works

Method Used	Accuracy %
Naïve Bayes	76.3
SVM	65.1
Decision Tree	73.82
SVM+ K Means Clustering	96.71
SVM+Fuzzy C Means Clustering	94
Proposed Approach: Naïve Bayes+ Mean Shift Clustering	99.45

6.3 Results

Total time taken to build this classification model was very low. This integrated method increased Naive Bayes accuracy in the diabetes forecast as shown in results. The proposed method in combination of SVM with K-means and Fuzzy C-means, was also higher precision than the integrated method.

7. CONCLUSIONS

In this paper, a Naïve Bayes+ Mean Shift Clustering algorithm has been proposed for diabetes prediction system. The experimental results depict that the integrated means shift cluster-based naive Bayes classification enhanced the traditional naive Bayes accuracy in predicting diabetes patients. The best accuracy of 99.45% is achieved.

REFERENCES

- [1] Porter, T., Green, B. (2009). Identifying diabetic patients: a data mining approach. *AMCIS 2009 Proceedings*, 500.
- [2] Li, L., Tang, H., Wu, Z., Gong, J., Gruidl, M., Zou, J., Tockman, M., Clark, R.A. (2004). Data mining techniques for cancer detection using serum proteomic profiling. *Artificial Intelligence in Medicine*, 32(2): 71-83. <https://doi.org/10.1016/j.artmed.2004.03.006>
- [3] Panzarasa, S. (2010). Data mining techniques for analyzing stroke care processes. *Proceedings of the 13th World Congress on Medical Informatics*, pp. 939-943. <https://doi.org/10.3233/978-1-60750-588-4-939>
- [4] Kuzuya, T., Nakagawa, S., Satoh, J., Kanazawa, Y., Iwamoto, Y., Kobayashi, M., Kashiwagi, A., Araki, E., Ito, C., Inagaki, N., Iwamoto, Y., Kasuga, M., Hanafusa, T., Haneda, M., Ueki, K., Committee of the Japan Diabetes Society on the Diagnostic Criteria of Diabetes Mellitus. (2002). Report of the Committee on the classification and diagnostic criteria of diabetes mellitus. *Diabetes Research and Clinical Practice*, 55(1): 65-85. [https://doi.org/10.1016/s01688227\(01\)00365-5](https://doi.org/10.1016/s01688227(01)00365-5)
- [5] Matheus, C.J., Piatetsky-Shapiro, G., McNeill, D. (1996). 20 selecting and reporting what is interesting: The kefir application to healthcare data. <http://citeseerx.ist.psu.edu/viewdoc/summary?doi=10.1.1.200.7107>.
- [6] Lowanichchai, S., Jabjone, S., Puthasimma, T. (2006). Knowledge-based DSS for an analysis diabetes of elder using decision tree. Faculty of Science and Technology Nakhon Ratchasima Rajabhat University, Nakhonratchasima, 30000. <https://doi.org/10.18231/2454-9150.2018.0637>
- [7] Guo, Y., Bai, G., Hu, Y. (2012). Using Bayes network for prediction of type-2 diabetes. *2012 International Conference for Internet Technology and Secured Transactions*, London, UK.
- [8] Guariguata, L., Whiting, D., Weil, C., Unwin, N. (2011). The international diabetes federation diabetes atlas methodology for estimating global and national prevalence of diabetes in adults. *International Diabetes Federation*, 94(3): 322-332. <https://doi.org/10.1016/j.diabres.2011.10.040>
- [9] Bellazzi, R. (2008). Telemedicine and diabetes management: Current challenges and future research directions. *J. Diabetes Sci. Technol.*, 2(1): 98-104. <https://dx.doi.org/10.1177%2F193229680800200114>
- [10] Al Jarullah, A.A. (2011). Decision tree discovery for the diagnosis of type II diabetes. In *2011 International Conference on Innovations in Information Technology*, pp. 303-307. <https://doi.org/10.1109/INNOVATIONS.2011.5893838>
- [11] Krishnaveni, G., Sudha, T. (2017). A novel technique to predict diabetic disease using data mining classification techniques. *International Journal of Advanced Scientific Technologies, Engineering and Management Sciences (IJASTEMS)*, 3.
- [12] Patil, B.M., Joshi, R.C., Toshniwal, D. (2010). Hybrid prediction model for type-2 diabetic patients. *Expert Systems with Applications*, 37(12): 8102-8108. <https://doi.org/10.1016/j.eswa.2010.05.078>
- [13] Raj, S., Rajan, G.V. (2013). Correlation between elevated serum ferritin and HbA1c in type 2 diabetes mellitus. *Int J Res Med Sci*, 1(1): 12-15.
- [14] Balpande, V., Wajgi, R. (2017). Review on prediction of diabetes using data mining technique. *International Journal of Research and Scientific Innovation (IJRSI)*. <https://www.semanticscholar.org/paper/Review-on-Prediction-of-Diabetes-using-Data-Mining-Balpande-Wajgi/a4e56f035cc20b725fea47d4b768b0b77076b35a>.
- [15] Smith, M.R., Martinez, T. (2011). Improving classification accuracy by identifying and removing instances that should be misclassified. In *the 2011 International Joint Conference on Neural Networks*, pp. 2690-2697. <http://dx.doi.org/10.1109/IJCNN.2011.6033571>
- [16] Zhao, Q., Li, H., Wang, X., Pu, T., Wang, J. (2019). Analysis of users' electricity consumption behavior based on ensemble clustering. *Global Energy Interconnection*, 2(6): 479-488. <https://doi.org/10.1016/j.gloi.2020.01.001>