

Received February 28, 2020, accepted March 30, 2020, date of publication April 6, 2020, date of current version April 17, 2020.

Digital Object Identifier 10.1109/ACCESS.2020.2985717

# Predicting Ayurveda-Based Constituent Balancing in Human Body Using Machine Learning Methods

VISHU MADAAN<sup>1,2</sup> AND ANJALI GOYAL<sup>1,3</sup>

<sup>1</sup>Department of Computer Science Engineering, IKG-Punjab Technical University, Kapurthala 144603, India

<sup>2</sup>Department of Computer Science Engineering, Lovely Professional University, Phagwara 144411, India

<sup>3</sup>Department of Computer Applications, Guru Nanak Institute of Management and Technology, Ludhiana 141002, India

Corresponding author: Vishu Madaan (vishumadaan123@gmail.com)

**ABSTRACT** Human Body constitution (*prakriti*) defines what is in harmony with human nature and what will cause to move out of balance and experience illness. *Tridosha* defines the three basic energies or principles that determine the function of our body on the physical and emotional levels. The three energies are known as *VATT*, *PITT* and *KAPH*. Each individual has a unique balance of all three of these energies. Some people will be predominant in one, while others will be a mixture of two or more. *Ayurveda-dosha* studies have been used for a long time, but the quantitative reliability measurement of these diagnostic methods still lags behind. A careful and appropriate analysis leads to an effective treatment. To collect a meaningful data set, a questionnaire with 28 different characteristics is validated by Ayurveda experts. Authors calculate Cronbach alpha of *VATT-Dosha*, *PITT-Dosha* and *KAPH-Dosha* as 0.94, 0.98 and 0.98, respectively to check the reliability of the questionnaire. Authors analyzed questionnaires of 807 healthy persons aged 20-60 years and found 62.1% men and 37.9% women. The class imbalance problem is resolved with oversampling and the equally distributed data set of randomly selected 405 persons is used for the actual experiment. Using computer algorithms, we randomly divide the data set (8:2) into a training set of 324 persons and a test data set of 81 persons. Model is trained using traditional machine learning techniques for classification analysis as Artificial Neural Network (ANN), K-Nearest Neighbour (KNN), Support Vector Machine (SVM), Naive Bayes (NB) and Decision Tree (DT). System is also implemented using ensemble of several machine learning methods for constitution recognition. Evaluation measures of classification such as root mean square error (RMSE), precision, recall, F-score, and accuracy is calculated and analyzed. On analyzing the results authors find that the data is best trained and tested with CatBoost, which is tuned with hyper parameters and achieves 0.96 precision, 0.95 recall, 0.95 F-score and 0.95 accuracy rate. The experimental result shows that the proposed model based on ensemble learning methods clearly surpasses conventional methods. The results conclude that advances in boosting algorithms could give machine learning a leading future.

**INDEX TERMS** Ayurveda, human body constituents, hyper parameter tuning, KAPH, optimized training model, PITT, VATT.

## I. INTRODUCTION

Ayurveda is one of the oldest medical sciences, which originated more than five thousand years ago on the Indian sub-continent [1]. The word “*Ayurveda*” is a combination of two words - “*Ayur*: life” and “*Veda*: know” thus defined as “science of life” [1], [2]. Ayurveda helps to maintain a healthy life by keeping the body, mind and soul of a person in balance. It is an art of life extension and a natural way of healing

The associate editor coordinating the review of this manuscript and approving it for publication was F. K. Wang.

and treatment. In today’s high-tech civilization, every human being is trapped by disease at a very early stage of life, so it is necessary to know the right way to live healthy. Ayurveda is such a path, which imparts the knowledge of natural harmony and helps to eliminate disharmony. Ayurveda finds the connection between the use of the senses and diseases. Incorrect use of the senses leads to disharmony between man and nature or imbalance in human nature. The equality between nature and self is the foundation upon which the principle of Ayurveda is built. Balancing the lifestyle among people is one of the powerful weapons of Ayurveda [1], [2]. A person who

is a follower of Ayurveda can adapt the rules and regulations of Ayurveda, and has more chances to live a healthy life without illness.

**A. PANCHMAHABHUTAS AND TRIDOSHA**

It is a belief in Ayurveda that each and every object in this universe consists of the five elements (space, water, fire, earth, air) called ‘*panchmahabhutas*’ [2]. These elements are represented in our body as three biological energies which regulate and control all life processes. These three forces like *VATT*, *PITT*, *KAPH* are known as ‘*Doshas*’ or ‘*Tridoshas*’ as shown in table 1 [3].

**TABLE 1. Panchmahabhutas and tridoshas.**

Ekadoshaja- dosha predominant	single	Dwandwaja- dosha predominant	Duel	Tridoshaja/ sannipataja- Tridosha predominant
VATT: Earth + Air		VATT - PITT		VATT - PITT - KAPH
PITT: Fire + Water		PITT - KAPH		-
KAPH: Earth + Water		KAPH - VATT		-

Human nature (*prakriti*) has two variants: the physical (*Shaaririka Prakriti*) and the intellectual (*Manasa Prakriti*). *VATT*, *PITT* and *KAPH* are *Sharirika prakriti* [1], [3].

‘*VATT*’ stands for the word wind. Wind moves and blows the cloud along. Wind is having the power to move the whole atmosphere in its blow. *VATT Dosh*a motivates the man to advance in life and it increases the will to live. Humans also need energy for physio- and biochemical processes. The power ‘*PITT*’ enhances focus and creates a glow in the human body. It regulates all metabolic processes in the body as well as body temperature and hormonal balance. ‘*KAPH*’ represents the element water, the fluid content in the “human body (as in tissues and organs) that lubricates the joints of the human body. These ‘*Tridosha*’ control all mental and physical processes in the living beings. ‘*Tridosha*’ characteristics are presented in table 2 [3]–[6]. 28 distinct features are selected for classifying an individual in suitable category of tridosha. Ayurveda says that each individual is different from others and unique in their relationship to each other [7]. Depending upon the relative predominance of the three physiological

**TABLE 2. Characteristics of tridoshas [5].**

Hair texture	Voice	Immunity	Temper level	Intelligence level
Skin Type	Sexual active- ness	Cold Tolerance	Patience level (Self-control)	Concentration Level
Physique type	Eating habits	Energy level	Talking behaviour	Grasping Power (level of under- standing)
Visibility of Tendons and veins	Walking speed	Sweating level	Friendship behavior	Memory recall and retention
Skin Tone	Appetite (De- sire to eat)	Sleeping habit	Jealous level	-
Eye Size	Stool type	Thirst Level	Activities performance level	-

factors (*dosha*) the psychosomatic constitution of an individual may be divided into seven categories namely *VATT*, *PITT*, *KAPH*, *VATT-PITT*, *VATT-KAPH*, *PITT-KAPH* and *samdosha* (balanced). An individual may have dominance of more than one *dosha* causing individual fall in any of the specific category of *dosha* type. Ayurveda science classifies man according to his physical structure and physiological characteristics. The study of (*prakriti*) helps the doctors to lead their patients into a healthy life. *Prakriti* is hereditary, which means that it is dissolved at the hour of its formation and is based on numerous parental, prenatal and postnatal elements [8], [9]. The personal analysis of *prakriti* helps to get to know the body and its requirements. *Prakriti* helps in maintenance of health, personal, family and professional life [10], [11].

**B. PRAKRITI EXAMINATION**

Knowing the *dosha* type beforehand helps in planning the lifestyle and diet according to the body’s needs. This knowledge provides probable occurrence of qualitative and quantitative imbalances in the body. There are some methods in Ayurveda through which a thorough examination of the patient is done to gather a maximum of information about the patient before a treatment is prescribed. Recently, Shilpa et al worked on the development of a questionnaire for issuing *prakriti* and presented it as a satisfactory validity tool for *prakriti* prediction [5]. Various research studies based on questionnaires dealt with inter-rater variability, but a quantitative approach was missing [11]. In our study, we work on a questionnaire to collect significant data. Data collected using pilot study undergoes checks of internal consistency and then we develop a trained modal based on machine learning approach.

**C. ENSEMBLE LEARNING**

Ensemble learning is a way to solve complex computational intelligence problems. Ensemble methods use multiple machine learning algorithms to obtain predictive performance. Unlike ML models ensemble methods don’t pertain to any single learning model. Ensemble learning is used for assigning a confidence to the decision of a model, selecting optimal features, incremental learning for best results, and error-correcting [11], [12]. Ensemble based system is developed by combining the diverse models(classifiers) to improve the classification results. The procedural approach of ensemble learning is same as consulting number of physicians before agreeing for any medical treatment. This makes the confidence robust for a specific outcome. Commonly used ensemble learning algorithms are bagging (bootstrap aggregating), boosting, stacked generalization (stacking). Bagging and Boosting is a method of training the weak models. These methods produce class predictions by aggregating the predictions of several other sub-classifiers [33]. Other traditional machine Learning (ML) models are used for training the dataset as Artificial Neural Networks (ANN), Naive Bayes (NB) classifiers, K-Nearest Neighbour (KNN), Support Vector Machines (SVM), Decision Trees (DT), Logistic

Regression, Fuzzy Logic, Genetic Algorithm [13]–[16]. The researchers worked on various data mining techniques to find optimized solutions for medical problems. In this study authors use Boosting methods as XGboost and Catboost algorithms. To improve the performance of these algorithms an optimized selection of parameters is required. For the same we use hyper parameter selection methods. Other traditional methods of machine learning are also used for recognizing human body constitutions and their results are compared.

The rest of this paper is structured as follows. The section II highlights the related work. The section III explains the methodology. The section IV describes the implementation of machine learning models. The section V evaluates the results and analysis. Finally, the section VI concludes the complete research work.

## II. RELATED WORK

Many researchers researched in Ayurveda for better human health. Verma *et al.* [7] examine various aspects of Ayurveda and reviewed the literature on the history and application of Ayurveda. The authors discussed the current trend of Ayurveda and its important role in health care. There is a discussion of the aspects that need to be considered for the promotion and development of Ayurveda medicine.

Todkari and Lavekar *et al.* [10] discussed that Ayurveda consists of determining the constitution and imbalances of an individual through various modular approaches. The authors worked on three different questionnaires and tested their reliability to know the human *prakriti* level. They analyze and display the reliability of the tests in numerical values. They conclude that a standardized questionnaire-based tool for the investigation of *prakriti* is needed.

Dunlap *et al.* [11] assess *Tridoshas* i.e. *VATT*, *PITT*, and *KAPH* from psychological perspective in human beings on a personality scale. As per authors these *dosha's* are composed of the *Pancha-Mahabhutas*. One or multiple *dosha* can dominate at any point of time. There is no case when one or the other *Pancha-Mahabhutas* and the *Tridoshas* are totally absent. They explain in their research work that developed scale shows the psychometric properties. The scale assesses the psychological manifestation of the *Tridoshas*, which is the basic achievement. The standardization procedure developed by them involves the development of the Mysore Psychological *Tridosha* Scale for *prakriti* determination.

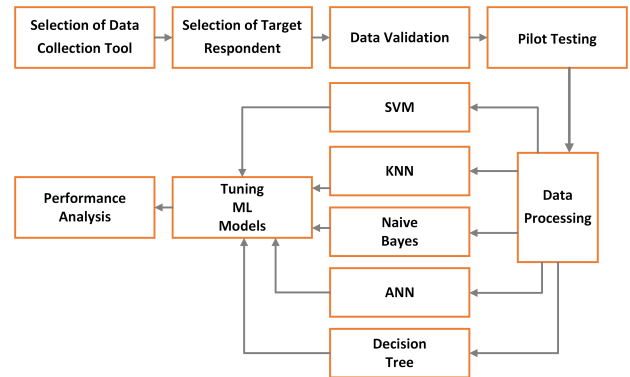
Ibrahim *et al.* [14] show in their study that the medical statements are full of ambiguities. Author have gone through a database of patient queries on medical websites are full of ambiguity and therefore the answers returned may not contain the desired information so. To avoid this, they propose end-to-end deep learning based medical diagnostic system (DL -MDS) to help diagnose diseases for users.

Woldaregay *et al.* [17] The authors discussed that the ability of machine learning to solve complex tasks with dynamic environment and knowledge has contributed to its success in their research. The work done so far in the field of medical and use of machine learning in health monitoring, health

healing has motivated for human *prakriti* determination with such decision-making expert tools.

## III. METHODOLOGY

In this section, we explain the complete research work process starting from data collection to performance analysis as shown in figure 1.



**FIGURE 1.** Work flow to achieve objective: After collecting the data from healthy individuals, data is validated and pilot testing is performed. Further dataset is collected and its internal consistency is checked, issue of class imbalance is resolved. Data undergoes training and testing with ML models as SVM, KNN, ANN, Naïve Bayes, Decision tree. Models are further tuned to improve the performances using hyper parameters and finally results of all ML models are compared.

### A. SELECTION OF DATA COLLECTION TOOL

The accurate and systematic data acquisition is a tedious task for training and accuracy of the model. The methods of data collection depends on the type of research. It may be document review, observation, interview, measurement or a combination of different methods. After identifying the scheme and characteristics of all the different categories, we prepare a questionnaire to collect the information on these characteristics from different people. Data collection by means of a questionnaire prepared by authors, is an economical way of obtaining information. It is easily accessible and provides wide coverage with less effort.

### B. SELECTION OF TARGET RESPONDENT

In order to exclude any influence of the participant's illness on the results, we ensure that healthy volunteers of both sexes are selected between the ages of 20 and 60. For this exercise we select students and employees of educational institutions to participate in this survey. This record contains 837 examples of 28 attributes listed in the table 2. We find 8 entries incomplete and remove these from the record.

### C. DATA VALIDATION

A questionnaire used to sort information, is an exploratory tool consisting of a sequence of questions and various invitations to collect data from respondents. We consult with the Ayurveda experts *Dr. S.K. Bansal, Agra (Uttar Pradesh, India)* and *Dr. Ajay Grover, Fazilka (Punjab, India)* for a

detailed discussion and to assist in the preparation of the questionnaire. The questionnaire consists of questions in a closed format, where each question is answered with several options. The validity of the questionnaire is checked by Ayurveda experts.

**D. PILOT TESTING**

To ensure the consistency of data collection, we conduct a pilot study. First, we test the significance of the questionnaire with 50 of the randomly selected participants from the entire population under consideration. On obtaining the satisfactory accuracy from small scale implementation on this dataset, we proceed for further data collection with same procedures. On trusting the results of pilot study we further train and test complete model.

**E. DATA PRE-PROCESSING**

We collect data using unbiased surveying techniques and ask only healthy volunteers to participate. Using a simple random sample, we select the participants for the data collection. Simple Random Sampling is based on the selection of a number of people from a large population group [18]. We select each individual with the same chance or probability. While interacting with the participants, we explain and advise them about the role of Ayurveda in modern life. The volunteers, who are found motivated and given consents, are informed about the purpose of the study. We found that the participants are curious to know their Dosha type. For the data collection, we interact with the finite population, i.e. 807 persons aged 20 to 60 years. Of these, 306 are female and 501 male. The category wise details of the data is given in the table 3.

**TABLE 3. Age wise and region wise category for data collection.**

Gender		Age Group				Individuals From	
Male	Female	Adole-science	Young Adult	Adult	Elderly	North India	South India
501	306	40	388	305	74	498	309

**1) SAMPLE SIZE**

The sample size of the study depends on the type of population under consideration. For the finite population, we use Slovin’s formula as given in equation 1 to know the sample size for our research.

$$n = \frac{N}{1 + (N * e^2)} \tag{1}$$

n is the sample size for experiments, N is the total population, while e represents an error limit of 5%. We consider that the world population N in December 2019 is 7.75 billion people with an error rate of 0.05 at a confidence level of 95% and calculate the sample size at 400 [19]. It is important to pre-treat the information before it is shown for preparation. The data collected in raw format requires a number of technical improvements before further training can be provided. The main data problems which are resolved as follows:

**2) DATA MANAGEMENT (WRANGLING)**

We use Python libraries as pandas, numpy for data preprocessing. The data is manually converted or mapped from the “raw” form into a substantial format that allows for more convenient use of the data.

**3) IMBALANCED DATA**

We collect the data without knowing the results. A proportionate collection of data of different classes is considered as best for learning and good accuracy. Modest class imbalance can cause serious problems in training. [20] A number of different techniques are there to resolve this type of imbalanced class problem which focus either at the data level (sampling methods) or at the classifier level (modifying it internally). On analyzing the recorded data of 837 participants in detail, it was found that the collected data is non-uniformly distributed. Thus we select proportionate data that covers equal distribution of all classes.

**4) INCOMPLETE DATA**

The data collected for a particular task may have some incomplete entries, sometimes due to unknown features or technical errors. When formatting the data, we find some missing data entries. We solve this problem by approximating missing values and discovering a relationship between the known and unknown data [21].

**5) INTERNAL CONSISTENCY**

To check the internal consistency of the collected data and the reliability of the questionnaire, we consider a consistency measure, i.e. Cronbach Alpha is calculated. Cronbach alpha shows how closely a set of elements is connected as a group [22]. The value of Cronbach alpha for VATT, PITT and KAPH is given in the table 4.

**TABLE 4. Questionnaire reliability.**

Cronbach Alpha		
VATT	PITT	KAPH
0.94	0.98	0.98

**6) DATA SPLIT**

Significant dataset is divided into two groups as training and testing according to previous studies [23]–[25]. Dataset of 405 individuals is spilt into 8:2 ratios. 80% of the original dataset (324 samples) is used as training dataset and 20% (81 samples) is used as testing dataset. Same dataset is used for developing and evaluating all models. For developing all models, packages of python programming are used.

**IV. MACHINE LEARNING MODELS**

**A. SVM LEARNING**

SVM is a discriminatory classifier formally defined by a separating hyperplane. This hyperplane is a line that divides a plane into two parts, which in each class lie on either

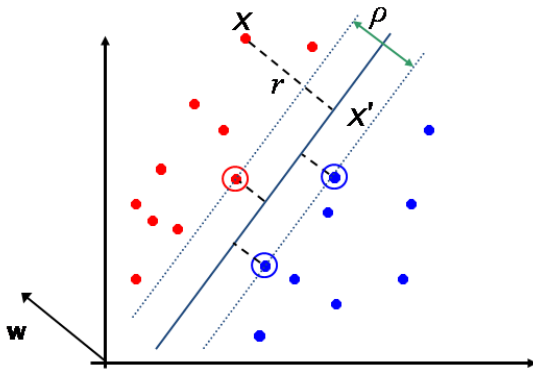


FIGURE 2. Hyper plane separating data points: Distance form data points to separator is  $r$ ,  $\rho$  is margin(width) between separator of classes.

side in two-dimensional space. It is implemented on pre-defined, labeled training data with a monitored learning method. The algorithm outputs an optimal hyperplane that categorizes the novel data points [26]. The algorithm of the support vector machine is supposed to find a hyperplane in an N-dimensional space that uniquely classifies the data points. Figure 2 shows the two different types of data points are separated by hyperplane. Equation 2 shows separator. Distance form data points to separator is  $r$  as shown in 3,  $\rho$  is margin(width) between separator of classes as shown in equation 4. The objective is to maximizes the distance between the hyperplane and the “difficult points” close to decision boundary.

$$W^T X_i + b \geq 1 \quad \text{if } Y_i = 1 \tag{2}$$

$$r = y \frac{W^T X + b}{\|W\|} \tag{3}$$

$$\rho = \frac{2}{\|W\|} \tag{4}$$

where  $X_i$  is a data point and  $Y_i$  is a particular class to which data point may belong. The reason behind high selection of SVM is because of its efficient implementations and performances proved to be excellent for high dimensional problems and small data sets.

**B. KNN METHOD**

KNN algorithm is a supervised machine learning algorithm used for solving classification and regression problems. This is based on the fact that data elements assigned the same class if these are closer in a feature space. This technique is based on Euclidean distance method by which we can calculate the distance between two points in a coordinate system as shown in equation 5.

$$ED = \sqrt{\sum_{i=1}^k (x_i - y_i)^2} \tag{5}$$

where ED is Euclidean distance. In a two dimensional, the distance is just between two points  $(x_i, y_i)$  in  $xy$  plane space where  $(i = 1, 2, \dots, k)$  be data points [27]. In similar way,

Manhattan formula and Minkowski formula is also used for calculating the distance. We train and tested  $k$  nearest neighbor algorithm on Gaussian distribution for pattern recognition.

**C. NB CLASSIFIER**

Naive-Bayes is a classifier based on supervised learning or a statistical approach. Bayes’ theorem-based approach assumes that the presence of a feature in a class is completely independent of the presence of any other feature. The Naive Bayes consider problem instances as feature vectors, which classified through the method to specific classes [29]. These characteristics are not interdependent, that is, the value assigned to one characteristic does not influence the value of other characteristics. Equation 6 shows the relationship between the posterior probability  $P(T|A)$ ,  $P(T)$  Class Prior Probability),  $P(A)$  Predictor Prior Probability) and the probability  $P(A|T)$ .

$$P\left(\frac{T}{A}\right) = \frac{P\left(\frac{A}{T}\right) * P(T)}{P(A)} \tag{6}$$

where T- Goal, A- Attributes.

**D. ANN LEARNING**

The artificial neural network is purely similar to the human neural system [29]. It consists of a number of artificial neurons that are trained to output the classification. These neurons are connected to weights that are adjusted to modulate the effect of the input signals. Significant input data are fed into the network with their output. Figure 3 shows what a coherent network looks like and connectivity of neurons on different layers.

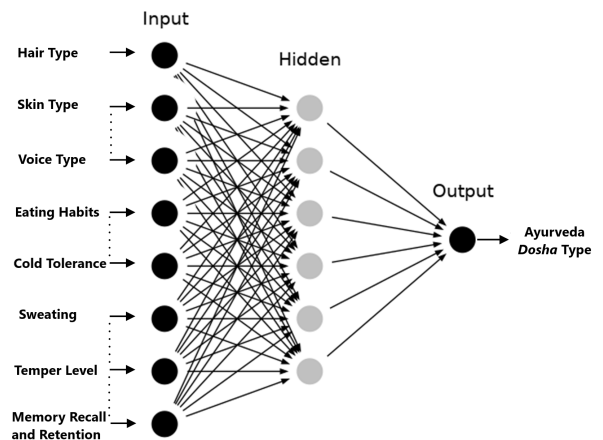


FIGURE 3. Basic Structure of ANN with 5 hidden layers is shown where 28 features are given at input layer and Ayurveda dsoha type is received at output layer.

For a defined training data set we select a suitable training model. Once the network is trained, we test ANN for its generalization performance. ANN gives feedback on whether the network has successfully classified the data or not. Many researchers use it ANN in medically related

applications [30], [31]. Two types of activation functions are used in artificial neural networks:

- (1) Rectified Linear Unit (ReLU), and
- (2) Softmax

1) ReLU

After the inputs from each layer have been passed, it is important to apply the activation function. ReLU is relatively simple, very useful and dynamic activation function. It offers the robustness to every small change like noise in the input. ReLU curve is half rectified. It means that for all negative input values, it turns the value into zero immediately as shown in equation 7 and 8. [32].

$$f(x) = \max(0, x) \tag{7}$$

derivative of ReLU is:

$$f'(x) = \begin{cases} x = 1, & \text{if } x > 0 \\ 0, & \text{otherwise} \end{cases} \tag{8}$$

2) SOFTMAX

We use the Softmax activation function when the data is to be divided into several classes. This function calculates the probability distribution of  $k$  output classes as shown in the equation 9.

$$p_j = \frac{e^{x_j}}{\sum_1^k e^{x_k}} \text{ for } j = 1, 2, \dots, k \tag{9}$$

where,  $x$  is input and output is  $p$ . Value of  $p$  lies between 0 and 1 [32]. We implement ANN in Python with Keras and TensorFlow libraries. The ANN training parameters are listed in table 5.

TABLE 5. ANN input parameters.

ANN Parameters	Value
Input Layer Neurons	28
Activation Function	ReLU, Softmax
Seed Value	7
Batch Size	128
Training Testing Ratio	8:2

E. DECISION TREE LEARNING

Decision tree is a non-parametric classifier and a predictive model based on the divide and conquer strategy. It is a classic example of soft computing and solves the purpose of the classifier. [33] A decision tree consists of the top node as root, lower child nodes, branches and internal nodes. The root connects the different classes of a tree. Leaf nodes represent the classes, the branches the results and the inner leaves the processes. The rules of classification rules form the paths from the root to the leaves. After a decision tree has been built from the data with the attributes  $A_1, A_2, \dots, A_n$  and the classes  $C_1, C_2, \dots, C_z$ , this decision tree can be used to classify a new data element with the attributes  $A_a, \dots, A_n$  into one of the classes  $C_1, C_2, \dots, C_z$ . For each new

data element to be classified, each node including the root is considered a question for the data sample. Based on the response suggested by any available branch, the next node is selected. When accessing the next node, another question about the data sample is answered until it reaches the leaf node. A leaf node is connected to one of the classes  $C_1, C_2, \dots, C_k$  [34], [35]. In this way, we assign a certain class to the data sample. A development of the tree-based algorithm is shown in figure 4.

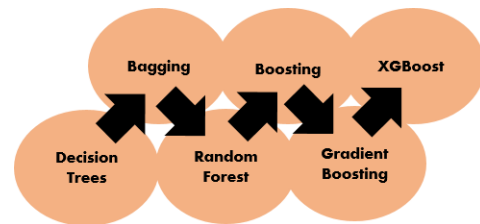


FIGURE 4. Evolution of decision tree algorithm.

1) BAGGING

Bagging or bootstrap aggregating involves combining inputs. A Bagging classifier fits base classifiers each on irregular subsets of the first dataset and afterward total their individual expectations might be by casting a ballot or by averaging to get a last prediction. Random Forest: It is an improvement over bagging algorithm. In random forest, only a subset of features is selected at random.

2) BOOSTING

It is a way to improve the predictability result using feedback. Earlier models are improved by minimizing the errors and by increasing the performance. Gradient Boosting: It is a special case of boosting approach where gradient descent algorithm is used to minimize the errors.

3) XGBoost

It is an Extreme gradient boosting algorithm. It is based on the software and hardware combination of optimization techniques to yield superior results using less computing resources in less time period. The focus of XGboost is on execution speed and performance of the training models [36].

4) CatBoost

It is also an Extreme gradient boosting algorithm. It outperforms the XGBoost algorithm in terms of execution time and memory limitations.

For datasets with a large number of features, XGBoost cannot run due to memory limitations, and Catboost converges to a good solution in the shortest time [37].

F. HYPER PARAMETER TUNING (HPT)

HPT is the process of searching and defining the optimized model architecture for accuracy. A hyperparameter is a parameter whose value is used to control the learning

process [38]–[40]. ‘Grid search’ and ‘Random Search’ are two methods of hyper parameter tuning. Grid search is most commonly used, applied on data. We use Inbuilt searching method ‘GridSearchCV’ from scikit–learn library of python to find optimal hyper-parameters and hence improve the accuracy/prediction results. We use max depth and min child weight parameters for tuning.

**Max depth:** Size of each new decision tree. Small size trees impose less complexity.

**Min child weight:** Minimum number of observations that must be in the children after a split. Small weight means more conservative. Gamma is tuned to make it less prone to overfitting [41], [42].

**V. RESULTS AND ANALYSIS**

In this section, we describe the implementation of ML models on a pre-processed data. It is important to analyze the performance of different learning techniques consistently as in machine learning no single algorithm works best for every problem. There are many factors that makes an impact, such as the size and structure of your dataset. This is the reason behind the implementation of multiple algorithms for our data, while using a hold-out “test set” of data to evaluate performance and select the winner [43]. Using traditional ML model, we are not finding any satisfactory results. For further improvement of the results, we apply hyper parameter tuning.

**A. SELECTION OF BEST LEARNING MODEL**

It’s not possible to design optimal model architecture only by looking at the data set. Each model is based on different performance characteristics. We can estimate the accuracy of each model on unseen data values by using cross-validation. These estimates provide a path to choose best models among all referred models. It is a good idea to visualize the data using different techniques from different perspectives. Similarly, various ways are used for model selection. We use different visualization methods in Python to know the average accuracy and other properties of the distribution of model accuracy.

Confusion matrix is a clean and unambiguous way to present the prediction results of a classifier [42]. To visualize the performance of classification models, we prepare a confusion matrix (also called as Error Matrix). Highlighted Matrix values shows correct classification while the other shows incorrect classification. We implement SVM on the data, in Python using Scikit-Learn the library. It contains the SVM library, which also contains built-in classes for various SVM algorithms [43]. To perform a classification task, the support vector classifier class is used, which is written as SVC in the Scikit-Learn’s svm library. The fit method of the SVC class is called to train the algorithm on the training data that is passed as parameters to the fit method. To select the best hyper parameters and training the dataset using SVM, Gridsearch and Cross Validation is used. Table 6 shows the classification results of model trained using SVM. Out of 81 sample tested, 69 samples are accurately classified,

**TABLE 6. Confusion matrix of data set trained using SVM.**

	VATT	PITT	KAPH	VATT-PITT	PITT-KAPH	VATT-KAPH	VATT-PITT-KAPH
VATT	<b>12</b>	0	0	0	0	0	1
PITT	0	<b>7</b>	0	1	0	0	3
KAPH	0	0	<b>9</b>	0	0	0	0
VATT-PITT	0	0	0	<b>11</b>	0	0	0
PITT-KAPH	0	1	3	0	<b>10</b>	0	0
VATT-KAPH	0	0	0	0	0	<b>13</b>	0
VATT-PITT-KAPH	1	0	0	2	0	0	<b>7</b>

while 21 samples are classified in incorrect classes. Accuracy rate of SVM classifier is 0.85.

**TABLE 7. Confusion matrix of data set trained using KNN.**

	VATT	PITT	KAPH	VATT-PITT	PITT-KAPH	VATT-KAPH	VATT-PITT-KAPH
VATT	<b>12</b>	0	0	0	0	0	1
PITT	0	<b>8</b>	0	1	0	0	2
KAPH	0	0	<b>10</b>	0	0	0	0
VATT-PITT	0	0	0	<b>10</b>	0	0	0
PITT-KAPH	0	1	3	0	<b>10</b>	0	0
VATT-KAPH	0	0	0	0	0	<b>13</b>	0
VATT-PITT-KAPH	1	0	0	1	0	0	<b>8</b>

We further apply KNN method on dataset, out of 81 samples, 71 samples are correctly classified while 10 samples are wrongly classified and achieve 0.87 accuracy rate as shown in table 7.

**TABLE 8. Confusion matrix of data set trained using Naive Bayes classification.**

	VATT	PITT	KAPH	VATT-PITT	PITT-KAPH	VATT-KAPH	VATT-PITT-KAPH
VATT	<b>9</b>	0	0	0	0	0	4
PITT	1	<b>7</b>	0	1	0	0	2
KAPH	0	0	<b>4</b>	0	1	1	3
VATT-PITT	3	0	0	<b>8</b>	0	0	0
PITT-KAPH	0	2	0	1	<b>8</b>	0	3
VATT-KAPH	0	0	0	0	0	<b>12</b>	1
VATT-PITT-KAPH	4	0	0	2	0	0	<b>4</b>

Similarly, we train system on data set using Naive Bayes method and achieve 0.642 accuracy rate as shown in table 8. Out of 81 testing samples, we obtain the correct classification

of 52 samples. We also implement Decision Tree method on dataset and find 72 samples are correctly classified out of total 81 testing samples. We achieve 0.88 accuracy rate, as shown in table 9.

**TABLE 9. Confusion matrix of data set trained using decision tree classification.**

	VATT	PITT	KAPH	VATT-PITT	PITT-KAPH	VATT-KAPH	VATT-PITT-KAPH
VATT	13	0	0	0	0	0	0
PITT	0	10	1	0	0	0	0
KAPH	0	1	5	0	0	1	2
VATT-PITT	0	0	0	11	0	0	0
PITT-KAPH	0	2	0	0	12	0	0
VATT-KAPH	0	0	1	0	0	12	0
VATT-PITT-KAPH	0	0	0	1	0	0	9

We further apply XGBoost to advance the decision tree method and improve the performance. We achieve 0.90 accuracy rate as model classifies 73 samples out of 81 as shown in table 10.

**TABLE 10. Confusion matrix of data set trained using XGBoost method.**

	VATT	PITT	KAPH	VATT-PITT	PITT-KAPH	VATT-KAPH	VATT-PITT-KAPH
VATT	13	0	0	0	0	0	0
PITT	0	9	1	1	0	0	0
KAPH	0	0	8	0	0	0	1
VATT-PITT	0	0	0	11	0	0	0
PITT-KAPH	0	0	2	0	12	0	0
VATT-KAPH	0	0	0	0	0	13	0
VATT-PITT-KAPH	0	0	0	2	1	0	7

System is further trained with optimized parameters and we get improved results as given in table 11. Model classifies 74 samples in correct classes out of total 81 testing samples of testing with 0.91 accuracy rate. We also implement Catboost algorithm to know the performance of training using other boosting algorithms and achieve 0.92 accuracy rate. Model classifies 75 sample in correct classes out of total 81 testing samples as shown in table 12. CatBoost algorithm is further tuned to know the best parameters of learning which are capable to provide best accuracy. On training it is found that accuracy rate is increased and it is 0.95 as the model is able to classify 77 samples out of 81 testing samples as shown in table 13. For further analysis, we consider other performance measures such as RMSE, Precision, Recall, F-score and Accuracy for comparing different classification models [24].

**TABLE 11. Confusion matrix of data set trained on XGBoost method with hyper parameter tuning.**

	VATT	PITT	KAPH	VATT-PITT	PITT-KAPH	VATT-KAPH	VATT-PITT-KAPH
VATT	13	0	0	0	0	0	0
PITT	0	9	1	1	0	0	0
KAPH	0	0	6	0	3	0	0
VATT-PITT	0	0	0	11	0	0	0
PITT-KAPH	0	0	0	0	14	0	0
VATT-KAPH	0	0	0	0	0	13	0
VATT-PITT-KAPH	0	1	0	1	0	0	8

**TABLE 12. Confusion matrix of data set trained using CatBoost method.**

	VATT	PITT	KAPH	VATT-PITT	PITT-KAPH	VATT-KAPH	VATT-PITT-KAPH
VATT	13	0	0	0	0	0	0
PITT	0	9	0	1	0	0	1
KAPH	0	0	8	0	0	0	1
VATT-PITT	0	0	0	11	0	0	0
PITT-KAPH	0	0	1	0	13	0	0
VATT-KAPH	0	0	0	0	0	13	0
VATT-PITT-KAPH	0	0	0	1	1	0	8

**TABLE 13. Confusion matrix of data set trained with hyper parameter tuning on CatBoost method.**

	VATT	PITT	KAPH	VATT-PITT	PITT-KAPH	VATT-KAPH	VATT-PITT-KAPH
VATT	13	0	0	0	0	0	0
PITT	0	9	2	0	0	0	0
KAPH	0	0	9	0	0	0	0
VATT-PITT	0	0	0	11	0	0	0
PITT-KAPH	0	0	0	0	14	0	0
VATT-KAPH	0	0	0	0	0	13	0
VATT-PITT-KAPH	0	1	0	1	0	0	8

1) ROOT MEAN SQUARE ERROR(RMSE)

It is a measure of the differences between the output predicted by a model and the output observed.

$$RMSE = \sqrt{(E - O)^2} \tag{10}$$

where E is expected value and O is observed value. RMSE can be calculated as shown in equation 10.

2) PRECISION (POSITIVE PREDICTED VALUE)

Precision is used to limit the number of false positives. It checks that how often a classifier predicts the positive results. It is calculated as number of correct positive



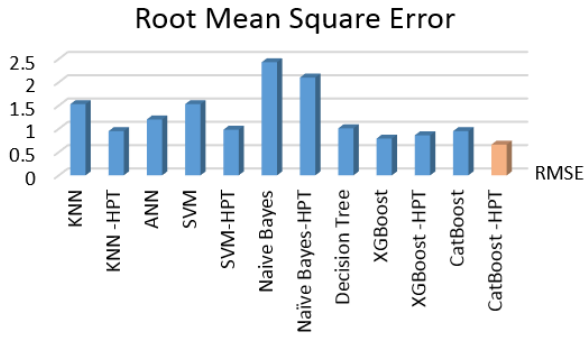


FIGURE 5. RMSE of implemented machine learning models.

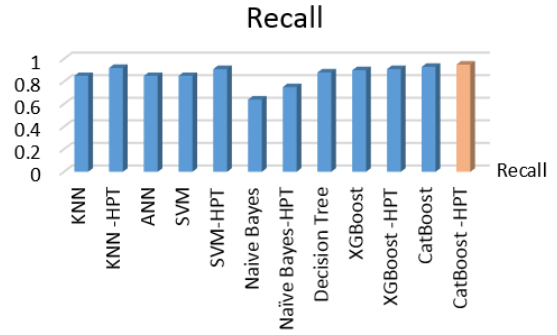


FIGURE 7. Recall values of implemented machine learning models.

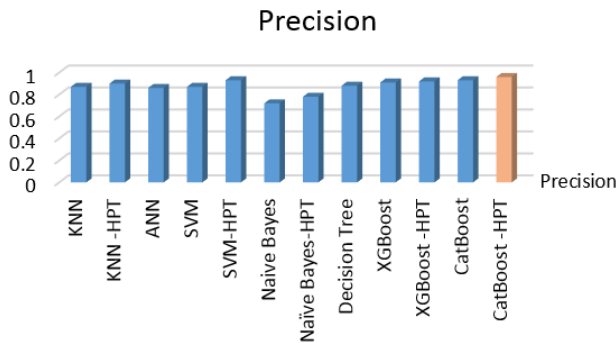


FIGURE 6. Precision of implemented machine learning models.

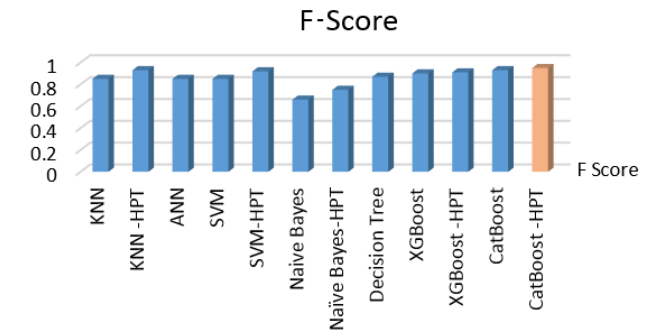


FIGURE 8. F-Score of implemented machine learning models.

predictions divided by total number of positive predictions can be calculated as shown in equation 11. It is an ability of a classification model to return only relevant instances.

$$Precision = \frac{TP}{TP + FP} \tag{11}$$

where TP is True Positive and FP is False Positive

### 3) RECALL (TRUE POSITIVE RATE)

It measures the sensitivity and is calculated as number of correct positive predictions divided by total number of positives as shown in equation 12. Best Sensitivity is 1 and worst is 0. It is an ability of a classification model to identify all relevant instances.

$$Recall = \frac{TP}{TP + FN} \tag{12}$$

where TP is True positive and FN is False Negative.

### 4) F-SCORE

F-Score is a measure of a test accuracy and is defined as the weighted harmonic mean of the precision and recall of the test. Single metric that combines recall and precision using the harmonic mean as shown in equation 13.

$$F - Score = 2 * \frac{(Precision * Recall)}{(Precision + Recall)} \tag{13}$$

### 5) ACCURACY

accuracy is a measure used for machine learning models to determine winning model for identifying relationships and patterns between variables in a dataset based supervised data. Accuracy can be calculated as shown in equation 14.

$$Accuracy = \frac{TP + TN}{TP + FP + TN + FN} \tag{14}$$

Correct predictions of model is True positives and True Negatives. All other predictions are *True Positives (TP)*, *True Negatives (TN)*, *False Positives (FP)* and *False Negatives (FN)*. Accuracy rate of all implemented models is shown in figure 9.

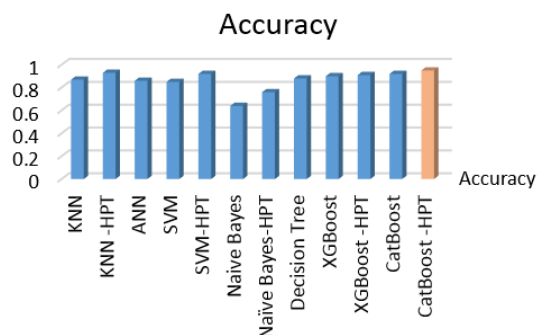


FIGURE 9. Accuracy of implemented machine learning models.

Summarized performance of all implemented models is shown in table 14. Results shows that CatBoost is

**TABLE 14. Performance of different machine learning models according to different measure.**

Classification Model	Performance Measure				
	RMSE	Precision	Recall	F-Score	Accuracy
KNN	1.53	0.87	0.85	0.85	0.87
KNN (with hyper parameter tuning)	0.95	0.90	0.92	0.93	0.93
ANN	1.2	0.86	0.85	0.85	0.86
SVM	1.53	0.87	0.85	0.85	0.85
SVM (with hyper parameter tuning)	0.98	0.93	0.91	0.92	0.92
Naive Bayes	2.43	0.72	0.64	0.66	0.64
NaiveBayes (with hyper parameter tuning)	2.1	0.78	0.75	0.75	0.76
Decision Tree	1.01	0.88	0.88	0.87	0.88
XGBoost	0.79	0.91	0.90	0.90	0.90
XGBoost (with hyper parameter tuning)	0.86	0.91	0.90	0.90	0.90
CatBoost	0.95	0.93	0.93	0.93	0.93
CatBoost (with hyper parameter tuning)	<b>0.66</b>	<b>0.96</b>	<b>0.95</b>	<b>0.95</b>	<b>0.95</b>

implemented with least output of 0.66 rmse, 0.96 precision, 0.95 recall, 0.95 f-score, and 0.95 accuracy rate.

## VI. CONCLUSION

The main goal of this work was to predict Ayurveda based human body constituents using different machine learning models and analyse their performance on various parameters. To achieve this, we first consulted the medical experts to design the questionnaire and collected the records from end users. After that, we again consulted the experts to validate and pre-process the recorded data then we implemented K-nearest neighbor, artificial neural networks, support vector machine, Naive Bayes, decision tree, XG-Boost and CatBoost methods with and without hyper parameter tuning. In this paper, we present performance results comparison among these machine learning methods on various performance parameters to predict human body constituents. We achieved 0.95 accuracy rate using CatBoost implemented with optimized parameters. The novelty of our work is the application based advanced algorithms of machine learning to recognize the human body constituents (Ayurveda Dosha). This proposed system may be proven as a supportive tool for Ayurvedic medical practitioners for recognizing human constituents. Every healthy or unhealthy individual can use this system without consulting an Ayurvedic practitioner to know about constituents imbalances in ones body.

## REFERENCES

- [1] B. Patwardhan, "Bridging Ayurveda with evidence-based scientific approaches in medicine," *EPMA J.*, vol. 5, no. 1, pp. 1–7, Dec. 2014.
- [2] *Ayurveda: A Brief Introduction and Guide*. Accessed: Jun. 23, 2019. [Online]. Available: <https://www.ayurveda.com/resources/general-information>
- [3] M. M. Mathpati, S. Albert, and J. D. H. Porter, "Ayurveda and medicalisation today: The loss of important knowledge and practice in health?" *J. Ayurveda Integrative Med.*, vol. 11, no. 1, pp. 89–94, Jan. 2020, doi: 10.1016/j.jaim.2018.06.004.
- [4] R. V. Rao, "Ayurveda and the science of aging," *J. Ayurveda Integrative Med.*, vol. 9, no. 3, pp. 225–232, Jul. 2018.
- [5] R. Chinthala, S. Kamble, A. S. Baghel, N. N. L. Bhagavathi, "Ancient archives of Deha-Prakriti (human body constitutional traits) in ayurvedic literature : A critical review," *Int. J. Res. Ayurveda Pharmacy*, vol. 10, no. 3, pp. 18–26, May/Jun. 2019.
- [6] S. Bhalerao and K. Patwardhan, "Prakriti-based research: Good reporting practices," *J. Ayurveda Integrative Med.*, vol. 7, no. 1, pp. 69–72, Mar. 2016.
- [7] V. Verma, S. Agrawal, and S. Gehlot, "Possible measures to assess functional states of Tridosha: A critical," *Int. J. Health Sci. Res.*, vol. 8, no. 1, pp. 219–231, Jan. 2018.
- [8] V. V. L. Prasuna, B. K. Sharma, and A. Narayana, "Comparative study of personality with Ayurvedic Prakriti," *Int. J. Ayurveda Pharmacy Res.*, vol. 2, no. 1, pp. 124–136, 2014.
- [9] B. A. Wani, S. K. Mandal, and P. Godatwar, "Prakriti analysis and its clinical significance," *Int. J. Ayurveda Pharmacy Res.*, vol. 5, no. 9, pp. 86–90, Sep. 2017.
- [10] D. P. Todkari and G. S. Lavekar, "Critical appraisal of Panchamahabhuta Siddhant," *Int. Ayurvedic Med. J.*, vol. 3, no. 5, pp. 1454–1461, 2015.
- [11] C. Dunlap, D. Hanes, C. Elder, C. Nygaard, and H. Zwickey, "Reliability of self-reported constitutional questionnaires in ayurveda diagnosis," *J. Ayurveda Integrative Med.*, vol. 8, no. 4, pp. 257–262, Oct. 2017.
- [12] S. Shilpa and C. Venkatesha Murthy, "Development and standardization of Mysore Tridosha scale," *AYU (Int. Quart. J. Res. Ayurveda)*, vol. 32, no. 3, p. 308, 2011.
- [13] X. Liu, Z. Liu, G. Wang, Z. Cai, and H. Zhang, "Ensemble transfer learning algorithm," *IEEE Access*, vol. 6, pp. 2389–2396, 2018.
- [14] Y. Ibrahim, S. Kamel, A. Rashad, L. Nasrat, and F. Jurado, "Performance enhancement of wind farms using tuned SSSC based on artificial neural network," *Int. J. Interact. Multimedia Artif. Intell.*, vol. 5, no. 7, p. 118, 2019.
- [15] Q. Xue and M. C. Chuah, "Explainable deep learning based medical diagnostic system," *Smart Health*, vol. 13, Aug. 2019, Art. no. 100068.
- [16] *Modern Machine Learning Algorithms: Strengths and Weaknesses*. Accessed: Jan. 15, 2020. [Online]. Available: <https://elitedatascience.com/machine-learning-algorithms>
- [17] A. Z. Woldaregay, E. Årsand, S. Walderhaug, D. Albers, L. Mamykina, T. Botsis, and G. Hartvigsen, "Data-driven modeling and prediction of blood glucose dynamics: Machine learning applications in type 1 diabetes," *Artif. Intell. Med.*, vol. 98, pp. 109–134, Jul. 2019.
- [18] P. W. West, "Simple random sampling of individual items in the absence of a sampling frame that lists the individuals," *New Zealand J. Forestry Sci.*, vol. 46, no. 1, p. 15, 2016.
- [19] *Slovin's Formula for Sampling*. Accessed: Aug. 30, 2019. [Online]. Available: <https://www.statisticshowto.datasciencecentral.com/how-to-use-slovins-formula/>
- [20] B. Krawczyk, "Learning from imbalanced data: Open challenges and future directions," *Prog. Artif. Intell.*, vol. 5, no. 4, pp. 221–232, Nov. 2016.
- [21] E. Karanja, J. Zaveri, and A. Ahmed, "How do MIS researchers handle missing data in survey-based research: A content analysis approach," *Int. J. Inf. Manage.*, vol. 33, no. 5, pp. 734–751, Oct. 2013.
- [22] S. Keith Taber, "The use of Cronbach's alpha when developing and reporting research instruments in science education," *Res. Sci. Educ.*, vol. 48, no. 1, pp. 1273–1296, 2018.
- [23] G. P. Herrera, M. Constantino, B. M. Tabak, H. Pistori, J.-J. Su, and A. Naranpanawa, "Data on forecasting energy prices using machine learning," *Data Brief*, vol. 25, Aug. 2019, Art. no. 104122.
- [24] A. M. Ibrahim and B. Bennett, "The assessment of machine learning model performance for predicting alluvial deposits distribution," *Procedia Comput. Sci.*, vol. 36, no. 1, pp. 637–642, Nov. 2014.
- [25] P. Sharada Mohanty, P. David Hughes, and M. Salathé, "Using deep learning for image-based plant disease detection," *Frontiers Plant Sci.*, vol. 7, p. 1419, Sep. 2016.
- [26] E. Sadrfaridpour, T. Razzaghi, and I. Safro, "Engineering fast multilevel support vector machines," *Mach. Learn.*, vol. 108, no. 11, pp. 1879–1917, Nov. 2019.
- [27] W. Cherif, "Optimization of K-NN algorithm by clustering and reliability coefficients: Application to breast-cancer diagnosis," *Procedia Comput. Sci.*, vol. 127, no. 1, pp. 293–299, Mar. 2018.
- [28] A. M. Mansour, "Texture classification using Naive Bayes classifier," *Int. J. Comput. Sci. Netw. Secur.*, vol. 18, no. 1, pp. 112–121, Jan. 2018.

- [29] M. Borhani and N. Ghasemloo, "Soft computing modelling of urban evolution: Tehran metropolis," *Int. J. Interact. Multimedia Artif. Intell.*, vol. 6, no. 1, p. 7, 2019.
- [30] N. A. Almansour, H. F. Syed, N. R. Khayat, R. K. Altheeb, R. E. Juri, J. Alhiyafi, S. Alrashed, and S. O. Olatunji, "Neural network and support vector machine for the prediction of chronic kidney disease: A comparative study," *Comput. Biol. Med.*, vol. 109, pp. 101–111, Jun. 2019.
- [31] H. Omrani, A. Tayyebi, and B. Pijanowski, "Integrating the multi-label land-use concept and cellular automata with the artificial neural network-based land transformation model: An integrated ML-CA-LTM modeling framework," *GIScience Remote Sens.*, vol. 54, no. 3, pp. 283–304, May 2017.
- [32] X. Jiang, Y. Pang, X. Li, J. Pan, and Y. Xie, "Deep neural networks with elastic rectified linear units for object recognition," *Neurocomputing*, vol. 275, pp. 1132–1139, Jan. 2018.
- [33] U. R. Acharya, S. L. Oh, Y. Hagiwara, J. H. Tan, and H. Adeli, "Deep convolutional neural network for the automated detection and diagnosis of seizure using EEG signals," *Comput. Biol. Med.*, vol. 100, no. 1, pp. 270–278, Sep. 2018.
- [34] S. Deepak and P. M. Ameer, "Brain tumor classification using deep CNN features via transfer learning," *Comput. Biol. Med.*, vol. 111, Aug. 2019, Art. no. 103345.
- [35] E. M. Karabulut and T. Ibriki, "Analysis of cardiocogram data for fetal distress determination by decision tree based adaptive boosting approach," *J. Comput. Commun.*, vol. 02, no. 09, pp. 32–37, 2014.
- [36] V. A. Dev and M. R. Eden, "Formation lithology classification using scalable gradient boosted decision trees," *Comput. Chem. Eng.*, vol. 128, pp. 392–404, Sep. 2019.
- [37] S. V. Murty and R. K. Kumar, "Accurate liver disease prediction with extreme gradient boosting," *Int. J. Eng. Adv. Technol.*, vol. 8, no. 6, pp. 2249–8958, Aug. 2019.
- [38] S. Lim and S. Chi, "Xgboost application on bridge management systems for proactive damage estimation," *Adv. Eng. Informat.*, vol. 41, Aug. 2019, Art. no. 100922.
- [39] M. Luckner, B. Topolski, and M. Mazurek, "Application of XGBoost algorithm in fingerprinting localisation task," in *Proc. IFIP Int. Conf. Comput. Inf. Syst. Ind. Manage.* Cham, Switzerland: Springer, pp. 661–671, 2017.
- [40] P. Schratz, J. Muenchow, E. Iturritxa, J. Richter, and A. Brenning, "Hyperparameter tuning and performance assessment of statistical and machine-learning algorithms using spatial data," *Ecol. Model.*, vol. 406, pp. 109–120, Aug. 2019.
- [41] Y. Yoo, "Hyperparameter optimization of deep neural network using univariate dynamic encoding algorithm for searches," *Knowl.-Based Syst.*, vol. 178, pp. 74–83, Aug. 2019.
- [42] R. Andonie, "Hyperparameter optimization in learning systems," *J. Membrane Comput.*, vol. 1, no. 4, pp. 279–291, Dec. 2019.
- [43] R. Hidayati, K. Kanamori, L. Feng, and H. Ohwada, "Implementing majority voting rule to classify corporate value based on environmental efforts," in *Proc. ICDMBD*, Bali, Indonesia, 2016, pp. 59–66.



**VISHU MADAAN** received the B.Tech. and M.Tech. degrees in computer science engineering from Lovely Professional University, Phagwara, India. She is currently pursuing the Ph.D. degree in computer science with IKG-Punjab Technical University, Punjab. She is also working as an Assistant Professor with Lovely Professional University. She has published more than 20 research articles in peer-reviewed conferences and journals. Her area of research includes soft computing, expert systems, pattern recognition, and machine learning. She is also a Reviewer of many international conferences and journals of high repute.



**ANJALI GOYAL** received the bachelor's degree in electronics from Kurukshetra University, in 1993, the master's degree in computer applications from Panjab University, Chandigarh, in 1996, and the Ph.D. degree from Punjab Technical University (PTU), Jalandhar, India, in 2013. Her academic achievement includes university merit position in graduation. She has served in various colleges of Kurukshetra University and PTU. She has a teaching experience of 23 years. She is also guiding many Ph.D. students. She is currently working as an Assistant Professor with the Department of Computer Application, Guru Nanak Institute of Management and Technology, Ludhiana, affiliated to PTU. Her research interests include content base image retrieval, digital watermarking, and pattern recognition. She has a number of international journal and conference publications to her credit. She is a Reviewer of many reputed international journals.

...