

## **Application of Machine learning techniques in groundwater potential mapping along the west coast of India**

**Pankaj Prasad<sup>1,2\*</sup>, Victor Joseph Loveson<sup>1,3</sup>, Mahender Kotha<sup>2</sup>, Ramanand Yadav<sup>1,3</sup>**

<sup>1</sup>Geological Oceanography Division, CSIR- National Institute of Oceanography, Dona Paula 403004, Goa, India

<sup>2</sup>School of Earth, Ocean and Atmospheric Sciences, Goa University, Taleigao 403001, Goa, India

<sup>3</sup>Academy of Scientific and Innovative Research, Ghaziabad 201002, India

\*Corresponding Author:Pankaj Prasad

Email: ppankaj.earthscience@gmail.com;ppankaj@nio.org

Mobile No.: +919883329432

### **Abstract**

Groundwater potential mapping (GWPM) in the coastal zone is crucial for the planning and development of society and the environment. The current study is aimed to map the groundwater potential zones of Sindhudurg coastal stretch on the west coast of India, using three machine learning models: random forest (RF), boosted regression tree (BRT), and the ensemble of RF and support vector machine (SVM). In order to achieve the objective, fifteen groundwater influencing factors including elevation, slope, aspect, slope length (LS), profile curvature, plan curvature, topographical wetness index (TWI), distance from streams, distance from lineaments, lithology, geomorphology, soil, landuse, normalized difference vegetation index (NDVI), and rainfall were considered for inter-thematic correlations and overlaid with spring and well occurrences in a spatial database. A total of 165 spring and well locations were identified, which had been divided into two classes: training and validation, at the ratio of 70:30, respectively. The RF, BRT, and RF-SVM ensemble models have been applied to delineate the groundwater potential zones and categorized into five classes, namely very high, high, moderate, low, and very low. RF, BRT, and ensemble model results showed that 33.3%, 35.6%, and 36.8% of the research area had a very high groundwater potential zone. These models were validated with area under the receiver operating characteristics (AUROC) curve. The accuracy of RF (94%) and hybrid model (93.4%), were more efficient than BRT (89.8%) model. In order to further evaluate and validate, four different sites were subsequently chosen, and we obtained similar results, ensuring the validity of the applied models. Additionally, ground penetrating radar (GPR) technique was applied to predict the groundwater table and validated by measured wells. The mean difference between measured and GPR predicted groundwater table was 14 cm, which reflected the importance of GPR to guide the location of new wells in the study region. The outcomes of the study will help the decision-makers, government agencies, and private sectors for sustainable planning of groundwater in the area. Overall, the present study provides a comprehensive high-precision machine learning and GPR based groundwater potential mapping.

### **Keywords**

Groundwater potential; GIS; Machine learning; Ensemble model; GPR

## 1. Introduction

Groundwater is the most valuable resource on our planet. It reflects the socio-economic condition and development of an area (Naghibi et al., 2016). Since the last century, groundwater has been in high demand for domestic, agricultural, and industrial purposes worldwide (Mogaji et al., 2014; Chen et al., 2019). Extensive groundwater extraction has led to a continuous drop of water table (Naghibi et al., 2016; Das, 2019). Therefore, Groundwater management is necessary for sustainable use of water resources. Groundwater availability and its movement depend on topographical, hydrological, ecological, geological, and atmospheric factors (Oh et al., 2011; Golkarian et al., 2018). As groundwater is a hidden natural resource, the demarcation of groundwater potential zones is essential for planning, management, and sustainable development of an area.

Many studies have been done by researchers on GWPM using different methods. Earlier groundwater mapping was based on field surveys, which was more expensive and time-consuming (Ganapuram et al., 2009; Chen et al., 2018; Das, 2019). (Ganapuram et al., 2009; Chen et al., 2018; Das, 2019). In the present time, remote sensing (RS), geographic information system (GIS), statistical, and geophysical techniques have been applied to map the groundwater potentiality of a large area in time and cost-effective manner.

In earlier studies, GIS modeling is very successfully applied to identify the groundwater prospect region with a high prediction rate (Das, 2019; Chen et al., 2019). The combined use of RS and GIS techniques have been employed in different research for GWPM (Prasad et al., 2008; Oh et al., 2011; Magesh et al., 2012; Naghibi et al., 2016; Murasingh et al., 2018; Lee et al., 2018; Golkarian et al., 2018; Chen et al., 2018; Das, 2019). Besides, several statistical techniques have been employed along with RS and GIS techniques for GWPM such as frequency ratio (Guru et al., 2009; Oh et al., 2011; Pourtaghi and Porghasemi, 2014; Naghibi et al., 2015; Das, 2019), weights of evidence (Lee et al., 2012; Tahmassebi et al., 2016; Chen et al., 2018), and logistic regression (Nampak et al., 2014; Chen et al., 2018).

In this decade, approaches through machine learning models have continuously been increased for groundwater mapping, such as random forest (Naghibi and Pourghasemi, 2015; Rahmati et al., 2016; Naghibi et al., 2016; Golkarian et al., 2018; Naghibi et al., 2019; Arabameri et al., 2019), boosted regression tree (Naghibi et al., 2015; 2016; 2019; Kordestani et al., 2019), support vector machine (Naghibi et al., 2017b; Lee et al., 2018; Naghibi et al., 2018), C5.0 (Duan et al., 2016; Golkarian et al., 2018), classification and regression tree (Naghibi et al., 2015; 2016), and artificial neural network (Lee et al., 2018). Such models are also applied in numerous fields, namely landslide susceptibility mapping

(Pourghasemi et al., 2013; Youssef et al., 2016; Kim et al., 2018), gully erosion susceptibility mapping (Arabameri et al., 2019; Gayen et al., 2019), ecological study (Lek and Guegan, 1999; Recknagel, 2001; Elith et al., 2008; Crisci et al., 2012), flood susceptibility mapping (Tehrani et al., 2014; 2015; Khosravi et al., 2018; Shafizadeh-Moghadam et al., 2018), and land use land cover change detection (Friedl et al., 1999; Gislason et al., 2006; Rodriguez-Galiano and Chica-Rivas, 2014). The prediction rate of machine learning models is very high compared to statistical models as reported in different studies (Naghibi and Pourghasemi, 2015; Chen et al., 2018).

The aforementioned studies have used the advantages of a single model. However, the ensemble model is a combination of statistical and machine learning techniques (Naghibi et al., 2017a; Kordestani et al., 2019). Recently, Naghibi et al. (2017a), Kordestani et al. (2019), and Naghibi et al. (2019) have used ensemble models for GWPM with satisfactory accuracy. The present work was applied the ensemble model of RF and SVM for GWPM.

In the previous studies, the researchers used different models, but their applicability was restricted to a specific study region. The objective of the present research is to map the groundwater prospect zones along the Sindhudurg coastal stretch using machine learning techniques and evaluate the results in different regions for the validity of the models. Besides, GPR technology was introduced to identify the groundwater table. Groundwater potential maps of the study region can be helpful for better planning and management of groundwater resources.

## **2. Study area**

Sindhudurg coast stretches from 15°43'11.43"N to 16°33'45.63"N latitude and 73°18'36.53"E to 73°55'50.07"E longitude covering an area around 3177 sq km along the west coast of India. The study region is bounded in the north and south by rivers and in the east and west by Western Ghats and shoreline of the west coast of India (Fig. 1a). With these natural boundaries, the area under study represents a typical coastal environment as it is a transition zone between land and sea.

The geology of the study area shows formation from the Archean to the Recent age and the lithology of which mainly consists of granite, basalt, lateritic, and alluvial deposits. The Deccan basalts occupy more than 30% of the study area and the aquifers are mostly associated with fractures and joints. Whereas the aquifers associated with laterites are substantiated with porous nature. The elevation of the study area varies from 0 to 450 m above mean sea level from seashore to landward dissected hill ranges. From the geomorphological perspective, the study area is reclassified into six distinct regions viz., denudational origin-pediment pediplain complex, coastal origin-younger coastal plain, structural origin medium dissected plateau, denudational origin moderately dissected plateau, structural origin-

low dissected plateau, and others. Karli and Gad, the main river networks of the study area, originate from the Western Ghats and debouch into the Arabian Sea. Climatologically, the research area is subtropical with a minimum temperature of 15°C to a maximum temperature of 40°C having three seasons (rainy, winter, and summer) throughout a year. The annual average rainfall of the area is around 3000 mm, and the maximum rainfall occurs during the southwest monsoon season (June-September). Monsoon rainfall profoundly influences the groundwater in different parts of the study region. According to the groundwater surveys and development agency and central groundwater board reports, the groundwater level ranges between 0.20 and 21 m/bgl. During the field visit, it was observed that groundwater is the primary source for drinking and irrigation extracted through dug and pumping well. Hence, the groundwater management is a vital concern for the study region.

Apart from this primary study region, four secondary regions were selected, shown in Figure 1 (b1, b2, b3, and b4), for further assessment of the applied models. The geographical details of these areas are given in Table 1.

### **3. Materials and methods**

The methodology applied in the current research is presented in Figure 2, which involves the following stages: Firstly, different thematic layers and inventory map of spring and well were prepared and transformed into the spatial database. Secondly, the groundwater potential maps were produced using machine learning models. Finally, the accuracy of the models was examined by applying the AUROC curve.

#### **3.1. Spring and well inventory map**

Many researchers have used the locations of spring, well, and quant as inventory for groundwater potential mapping. In the present research, both spring and well points were considered for GWPM. The inventory map of the study region contains of 165 spring and well points, identified from numerous sources (Table 2), and field observation. Random partition algorithm was deployed to separate the spring and well points for training and validation purposes, where 116 (70%) points were preferred for training and the remaining 49 (30%) for validation of dataset. In the same way, the inventory maps of other regions were prepared.

#### **3.2. Groundwater conditioning factors**

It is essential to select the effective parameters for preparing the groundwater prospect map of an area. Based on the previous studies (Naghibi et al., 2016; Rahmati et al., 2016; Golkarian et al., 2018; Chen et al., 2018) and field examination, 15 groundwater conditioning factors (thematic maps) viz. elevation, slope, aspect, LS, profile curvature, plan curvature, TWI, distance from the streams, distance from the

lineaments, lithology, geomorphology, soil, land use, NDVI, and rainfall (Fig. 3a-3o) were taken into account for GWPM in the study area. These thematic maps were prepared using the ArcGIS 10 software from several data (Table 2). Each thematic layer was resampled into a uniform grid size of 30×30 m, and the grid of research area was prepared by 2201 columns and 3159 rows (3,530,426 pixels; 3177 km<sup>2</sup>). Similarly, thematic maps of the selected secondary regions were generated.

SRTM (30 m resolution) digital elevation model (DEM) was used to produce the following topographic factors. Elevation, one of the potential indicators of groundwater, plays a vital role for GWPM (Oh et al., 2011; Naghibi et al., 2016; Patra et al., 2018; Chen et al., 2019). The elevation map was created from the DEM. The slope is considered as the most relevant topographic variable for groundwater potentiality (Naghibi and Pourghasemi, 2015; Lee et al., 2018). The slope map was produced in the ArcGIS environment, and the slope values were grouped into six classes using the natural break method. Aspect is one of the important controlling factors for the GWPM. It defines the direction of the slope, which is exposed to sunlight, winds, lineament, and rainfall (Goudie, 2013; Chen et al., 2018). The aspect map was used to correlate the groundwater availability at the different directions of the slope. Slope length (LS) defines the length (L) and steepness (S) of the topography that influences the amount of groundwater storage. LS is calculated with the following equation (Moore and Burch, 1986).

$$LS = (fa \times cellsize/22.13)^{0.4} \times (\sin\theta/0.0896)^{1.3} \quad (1)$$

Where *fa* refers to flow accumulation and  $\theta$  represents the slope in degrees.

Curvature affects the surface and subsurface hydrology (Regmi et al., 2015). Profile curvature is parallel to the maximum slope in a particular direction. The negative value of profile curvature indicates the water flow decelerated in the surface, whereas the positive value indicates the water flow accelerated on the surface, and zero indicates the surface is linear. In another side, plan curvature defines the maximum slope in a perpendicular direction. It describes the convergence and divergence of water flow in the earth's surface. Negative values represent the concave slope of the surface, which causes the confluence of water flow. In contrast, positive values indicate the convex slope of the surface that determines the divergence of water flow in the region (ESRI, 2016). TWI expresses the effect of topography on the location, which is related to soil conditions of the area. TWI is calculated as follows (Moore et al., 1991).

$$TWI = \ln(fa/\tan\beta) \quad (2)$$

Here *fa* is the flow accumulation, and  $\beta$  is the slope angle at the point.

Distance from the stream is inversely related to the identification of groundwater prospects in an area. Lineaments are hydro-geologically very meaningful to control the groundwater movement and storage (Magesh et al., 2012). It is in a linear or curvilinear pattern on the earth's surface and identified from the satellite imagery, DEM, and field survey (Prasad et al., 2010; Magesh et al., 2012; Goudie, 2013; Rahmati et al., 2015). For this study, lineaments were extracted from the superimposed shaded relief maps at an interval of 45° azimuth angle. High lineament density indicates more groundwater potentiality in the area (Magesh et al., 2012). Euclidean distance method was adopted to examine the relationship between inventory and distance from streams and lineaments. In the study area, the age of the rock traced from the Archean to Recent, which controls the groundwater storage. Based on lithofacies and geological ages, the lithology of the study area was reclassified into six broad classes: Group 1: Archean schist and gneisses (Granite gneiss, quartzite, meta-gabbro, amphibole schist) (Ask), Group 2: Dharwar supergroup (Meta greywacke, metabasalt, granite) (Dsg), Group 3: Kalladgi supergroup (Sedimentary quartzite, shale) (Ksg), Group 4: Sahyadri supergroup (Unclassified flows, Aa flow, Mega cryst flow) (Ssg), Group 5: Laterite (Lat), Group 6: Alluvium (Alv). Geomorphology is another predisposing factor for predicting the potential of groundwater. In this study, the geomorphic unit was categorized into six major groups: 1. Denudational origin-pediment pediplain complex (DoPPc), 2. Coastal origin-younger coastal plain (CoYCP), 3. Structural origin moderately dissected plateau (SoMDp), 4. Denudational origin moderately dissected plateau (DoMDp), 5. Structural origin-low dissected plateau (SoLDp), 6. Others (Oth), using the topographical maps, and Landsat 8 OLI image. Soil is considered one of the most important indicators of the surface and sub-surface runoff, recharge, and infiltration processes (Mogaji et al., 2014; Rahmati et al., 2016). Soil map from the National Bureau of Soil Survey and Land use planning (NBSSLUP) was used to reclassify the soil into four categories: Ultic Typic Haplustalfs, Lithic Ustorthents, Ultic Haplustalfs, and Typic Ustropepts. Land use and land use changes influence the groundwater storage and aquifer yield (Ibrahim-bathis and Ahmed, 2016; Guru et al., 2017; Chen et al., 2018). Unsupervised classification technique was employed to produce the land use map from the Landsat 8 OLI image. The study area was classified into five major land use classes: forest, agricultural land, waste land, built-up, and water bodies. The accuracy of the land use classification was calculated as 86% using the Kappa index. NDVI provides the health of vegetation, and it is usually used to relate the vegetation density and groundwater potentiality (Pourghasemi et al., 2013; Chen et al., 2018). The range of the NDVI value lying from -1 to 1, with higher NDVI value indicates the healthy vegetation and vice versa. Based on the Landsat 8 OLI image, the NDVI map was created using the following equation:

$$NDVI = (NIR - R)/(NIR + R) \quad (3)$$

Here near-infrared (NIR) and red (R) bands represent the spectral reflectance measurement of these bands.

The distribution, duration, and intensity of rainfall are significant affecting factors for infiltration, runoff, and recharge conditions (Magesh et al., 2012). The rainfall map was prepared from the rain gauge data of Maharashtra government using the inverse distance weighted method and classified into five groups.

### 3.3 Methods

In this research, three machine learning models (RF, BRT, and ensemble of RF-SVM) were employed for GWPM. The relationship between different groundwater conditioning factors and inventory locations was calculated by the frequency ratio method. On the other side, the importance of these groundwater effective factors was measured by ‘variable importance’ function in R software. The raster values of 15 factors of each spring and well location were imported to R software; then, the models were applied using different packages in the R environment. To improve the classification accuracy and avoid the biasness, 10-fold cross-validation method (with five repetitions) was used in the models. The final output values of the models were transformed into a spatial dataset for GWPM using the ArcGIS software. At last, models were validated by the AUROC curve and also examined in different regions of the west coast of India. In addition, GPR technology was used to measure the groundwater table by RADAN and MATGPR software.

#### 3.3.1 Application of frequency ratio (FR) and variable importance function

FR is a bivariate statistical technique to explain the prospect of occurrence of a certain attribute (Bonham-Carter, 1994; Oh et al., 2011; Manap et al., 2014; Naghibi et al., 2016; Guru et al., 2017). It is defined by the relationships between dependent variables (spring and well location) and independent variables (groundwater conditioning factors) (Guru et al., 2017; Das, 2019). In this context, the FR model was used to show the quantitative relationship between the inventory and each sub-class of groundwater variables. FR is calculated as;

$$FR = (P_s/T_s)/(P_a/T_a) \quad (4)$$

Where  $P_s$  is the number of spring and well under each sub-class of the groundwater effective factors,  $T_s$  denotes the total spring and well of the study area;  $P_a$  is the number of pixels of each sub-class of the conditioning parameter,  $T_a$  is the total pixels of the study area. FR is the ratio of spring and well occurrences to the total area of each class of the affecting factors. So the FR value 1 indicates the average of the model. If the value is more than 1, it considers the high groundwater prospect and less

than 1 value indicates the low groundwater potential in each sub-class of the parameters (Lee and Pradhan, 2007; Pradhan and Lee, 2010; Oh et al., 2011; Manap et al., 2014; Nampak et al., 2014).

‘Variable importance’ function of RF model was used for measuring the effectiveness of each factor. It is a generic method for calculating the feature importance by trained methods (Khun et al., 2018). Each factor was evaluated individually using a filter approach. The importance values of the features range from 0 to 100. Higher the importance value, greater is the influence of the factor, and vice versa. The details of variable importance function are given in Khun et al. (2018).

### 3.3.2 Application of random forest (RF)

Random forest is an ensemble machine learning technique for both classification and regression tasks (Breiman, 2001; Youssef et al., 2016; Naghibi et al., 2017b; Kim et al., 2018). For classification, RF uses the resampling technique by randomly changing the predictive variables to increase the diversity in each tree (Youssef et al., 2016; Naghibi et al., 2017b). This method consists of multiple decision trees and merges them to explain the spatial relationship between controlling variables of groundwater and inventory of spring and well (Kim et al., 2018). The decision tree is generated by bootstrap samples and leaves few samples for validation to test the accuracy of the decision tree. The mean-squared error of each decision tree with their OOB samples ( $E_{oob}$ ) is used to calculate the learning error.  $E_{oob}$  is expressed as:

$$E_{OOB} = \frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2 \quad (5)$$

Here  $n$  denotes the total number of OOB samples;  $y_i$  is the observed output, and  $\hat{y}_i$  is the model output. The main advantages of this approach are; i) it can handle large datasets with high dimensionality as well as avoid the over-fitting of the datasets; ii) this does not need any assumption regarding the explanatory variables and response variables, and iii) it does not require any prior data to transformation and rescaling.

### 3.3.3 Application of boosted regression tree (BRT)

BRT model is advanced and different from classical regression methods. It uses statistical and machine learning techniques to enhance the performance of a single model by fitting many models and combining them for prediction (Schapire, 2003; Elith et al., 2008; Naghibi et al., 2016). It is a combination of two algorithms, namely boosting and regression tree (decision tree) (Elith et al., 2008; Youssef et al., 2016; Kim et al., 2018). In the BRT model, the initial decision tree (DT) reduces the loss function. At each iteration, the main target was to decrease the root mean square error and the residuals. Then, the next DT is fit for prediction residuals of the first tree. In this stagewise process,



the existing trees are unchanged as the model develops increasingly larger. At each step, the fitted value of each observation is reestimated to express the contribution of the recently added tree (Elith et al., 2008; Naghibi et al., 2016). DT was used for visualization and explicit decision-making. The advantages of the algorithm are that the predictive variable can be of any type (numeric, binary, and categorical, etc.), and the outcomes of the model are not affected by monotone transformations and different scales of measurement among predictors. It replaces the missing data in predictor variables using surrogates (Breiman, 2001; Elith et al., 2008; Youssef et al., 2016). Boosting is a technique to get higher accuracy from the predictive variables of regression trees. It is a sequential procedure to average many rough rules of thumb (Schapire, 2003; Elith et al., 2008).

By combining the algorithms, the BRT model establishes a binary tree with general classification and regression tree. The classified data split into two samples. Each sample defines the best point for the data partition, and it formulates the observed deviation and residuals at each partition (Kim et al., 2018). Finally, the model is capable of estimating the observed value.

In the BRT model, three parameters such as number of trees, shrinkage or learning rate, and interaction depth are required for tuning. Interaction depth defines the number of nodes in trees and the learning rate determines the importance of each tree in the built model. Based on these two parameters, the number of trees was decided for optimal prediction (Elith et al., 2008; Naghibi et al., 2016).

### **3.3.4 Application of ensemble of RF and support vector machine (SVM)**

SVM model is another popular supervised machine learning technique that is based on the concept of structural risk minimization and statistical learning theory (Tehrany et al., 2014; Mojaddadi et al., 2017; Naghibi et al., 2018). The principle of the method is to separate the hyperplane formation from the dataset. The hyperplane is defined as the centre of the maximum margin of separation (Marjanovic et al., 2011; Tehrany et al., 2015). On the basis of the hyperplane, point was classified as +1 or -1. For the case of linear separable data, a separating hyperplane can be computed with the help of following equation (Hong et al., 2017):

$$y_i (w \cdot x_i + b) \geq 1 - \zeta_i \tag{6}$$

Where  $w$  represents the coefficient vector which expresses the orientation of the hyperplane in the feature space,  $b$  is the offset of the hyperplane from the origin, and  $\zeta_i$  defines the positive slack variables. The following optimization problem can be solved by defining an optimal hyperplane (Samui, 2008).

$$\text{Minimize } \sum_{i=1}^n a_i - \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n a_i a_j y_i a y_j (x_i x_j) \quad (7)$$

$$\text{Subject to } \sum_{i=1}^n a_i y_i = 0, 0 \leq a_i \leq c \quad (8)$$

Here  $a_i$  indicates the lag range multiplier and  $C$  denotes the penalty. The precision of the successful classification in SVM model depends on the selection of kernel type (Yao et al., 2008). In the present research, radial basis function (RBF) was applied because of its higher capability in interpolation, as reported by many researchers (Tehrany et al., 2014, 2015; Gayen et al., 2019). RBF equation is calculated as follows

$$\text{RBF} : K(x_i, x_j) = \exp(-\gamma(x_i - x_j)^2) \quad (9)$$

Where  $K(x_i, x_j)$  is the kernel function,  $\gamma$  represents the RBF kernel function. The purpose of applying the algorithm was to reduce the error and model complexity (Naghibi et al., 2018).

The ensemble model was developed by combining two or more than two different predictive models. In recent studies (Naghibi et al., 2019; Kordestani et al., 2019; Chen et al., 2019), the ensemble model was applied to improve the results. In this work, RF and SVM models were ensembled in R software using the weighted average method. For the implementation of this machine learning models, R statistical software 3.6.1 version was used with the help of different packages.

### 3.4 Validation of groundwater potential maps

Validation is a fundamental step in modeling for the scientific significance of the research (Naghibi et al., 2016; Chen et al., 2019). In this research, area under the receiver operating characteristic (ROC) curve was opted for evaluation of the models. ROC is a graphical plot, which determines the performance of the models in a diagnostic test (Egan, 1975; Golkarian et al., 2018). The curve plots the true positive rate (sensitivity) on Y-axis and false positive rate (1- specificity) on X-axis (Youssef et al., 2016; Golkarian et al., 2018). Model prediction for occurrence and non-occurrence of springs and wells was evaluated using the area under the ROC curve. The area under the curve (AUC) represents the value between 0 and 1, and the higher value represents the better performance of the model (Youssef et al., 2016; Naghibi et al., 2016; Golkarian et al., 2018; Chen et al., 2018; 2019). To extract the generalized findings, it is necessary to apply the models in different regions. For this purpose, four different areas along the west coast of India were chosen (b1, b2, b3 and b4). The b1 and b4 regions show slightly different characteristics from the study area in terms of topography, lithology, and climatic conditions whereas b2 and b3 sites lying in close proximity to the study area appear to be

similar with respect to aforesaid criteria. Additionally, six GPR profiles were selected to examine the groundwater level with the corresponding water level of wells. GPR is a non-invasive geophysical technique based on propagation and reflection of the transmitted electromagnetic waves (Annan, 2003; Neal, 2004; Billy et al., 2014). A SIR-4000 model of GSSI (Geophysical Survey System Inc.) with 200 MHz antenna was used to identify the water table depth. GPR data was processed in Radan 7 and MatGPR 3.2 (compatible with Matlab software) software by applying the time zero removal, filtering, background removal, and migration.

## **4 Results**

### **4.1 Spatial relationship between groundwater conditioning factors and inventory of spring and well**

It is necessary to know the relationship between inventory with effective factors for GWPM in an area. The results of the above relation are shown in Table 3 using the FR model. The FR value ranges from 0-1.78 in different sub-classes of the groundwater conditioning variables. For elevation, FR value decreases with the higher altitude. In the case of slope, the class from 3°-6° is highly correlated with spring and well occurrences (FR=1.18). The northwest slope direction with FR value of 1.22 has more spring and well locations in compared to other directions of the slope. FR values of LS classes are decreasing with the increasing slope length. In the matter of profile and plan curvature, flat areas have the highest FR value of 1.43 and 1.41, respectively. In TWI class, the second class (8-12) has the maximum FR value (1.29). The FR values increase with decreasing distance from streams, indicating a high groundwater prospect. In relation to lineament, most of the springs and wells occur at 0 to 1.25 km distance from the lineaments. Regarding lithology, spring and well locations have mostly identified in the Kaladgi group of rock with FR value of 1.53. For geomorphology, FR of the denudational origin-pediment pediplain complex has a maximum value (1.55). In the soil groups, ultic haplustalfs soil has a strong correlation with spring and well occurrences with FR value 1.78. In relation to land use, agricultural land is the maximum FR value 1.50. In the case of NDVI, the highest FR value is 1.69 in the third class (0.25-0.34). In rainfall class, the highest rainfall (3900-4300 mm) shows the maximum FR value (1.22).

### **4.2 Groundwater potential models**

Groundwater potential maps were prepared using the three machine learning algorithms, namely RF, BRT, and the hybrid of RF-SVM models (Fig. 4). Based on Min-max normalization, the probability values of the models were normalized between 0 to 1 and subsequently classified into five zones: very low (0-0.20), low (0.20-0.45), moderate (0.45-0.70), high (0.70-0.85), and very high (0.85-1) with

same class range for comparison among the models. The higher value represents very good groundwater prospect of the area and vice versa.

The model-wise spatial variation of the groundwater potentiality is shown in Figure 4 and Table 4. According to the RF model, the high and very high prospective groundwater zones covered 16.73% and 33.31% of the study area, respectively. Low and very low groundwater prospect zones comprise 17.98% and 12.48% of the total area, respectively. In the BRT model, it was found that 19.13% and 35.60% of the entire research area was of high and very high groundwater potentiality. Besides, low and very low groundwater potentiality covered 12.93% and 26.37% of the area. Based on the ensemble model result, the study region was classified as very high (36.85%), high (19.18%), moderate (13.14%), low (13.99), and very low (16.84%) groundwater prospect zones (Table 4).

### **4.3 Validation of machine learning models and GPR profiles**

The validation of the model is crucial for the assessment of GWPM. In many studies, the ROC curve was used for quantitative validation of the models with high prediction rate (Golkarian et al., 2018; Chen et al., 2018). In the present study, the validity of RF, BRT, and ensemble models was confirmed using the ROC curve. The AUC values ranging from lower (0) to higher (1) represent the worst to the best model prediction for GWPM. Based on ROC result, the success rate of RF, hybrid of RF-SVM, and BRT models was measured as 94.0%, 93.4%, and 89.8%, respectively. The accuracy of the models from the selected sites was calculated. In the b1 area, the accuracy of RF, BRT, and RF-SVM were 82.2%, 72%, and 83.3%, respectively. The success rates of RF (94%, 96.5%); BRT (90.5%, 90.7%) and hybrid models (93.8%, 94.7%) were computed for the b2 and b3 regions, respectively. In case of b4 region, the precision of RF, BRT, and RF-SVM was evaluated as 82.8%, 80.7%, and 82.5%, respectively.

The water table from the GPR profiles (a, b, c, e, and f) was identified at the depth of 3.6, 3.4, 4.9, 4.8, 4.3, and 4.1 m respectively and the nearest water level of wells was observed at 3.6, 3.4, 4.9, 5.1, 4.1, and 4.4 m, respectively (Table 5). Only in the case of profile 'd', the water table could not be identified. From the five GPR profiles, the average and maximum difference between predicted and measured depths of groundwater were 14 and 30 cm, respectively. Overall, from six GPR profiles, five locations of water table were accurately determined, which means more than 80% accuracy achieved in detecting the groundwater table by GPR technology.

## **5 Discussion**

The presence of spring and well at the particular segments of the study area indicate the potentiality of high groundwater yield (Oh et al., 2011; Naghibi et al., 2017a). However, to assess the groundwater

prospect of the entire area, statistical and machine learning methods were used by many researchers with good results (Oh et al., 2011; Chen et al., 2019). The results of the present study are discussed as follows:

### **5.1 Important conditioning factors for GWPM**

The comparative importance of the fifteen influencing factors of groundwater was illustrated using 'variable importance' function of RF model. In this context, geomorphology had the highest importance followed by elevation, NDVI and, distance from the stream while soil was of lowest importance followed by slope length rainfall, and land use (Fig. 5). Geomorphology is the most effective factor since the geomorphic features of different landforms of the study area control the groundwater potentiality to a maximum extent. Various landforms on the earth's surface are associated with a different kind of groundwater storage (Deepika et al., 2013; Rajaveni et al., 2017). Structural hill, residual hill, and linear ridge represent the low groundwater potential, whereas pediplain and valley fill have the high groundwater potential due to high infiltration and groundwater recharge (Rajaveni et al., 2017; Berhanu and Hatiye, 2020). Almost fifty percent of the study area is associated with a pediment-pediplain complex, which indicates the good groundwater potentiality. Elevation is another critical factor for GWPM, which was an agreement with the results of Naghibi and Pourghasemi (2015), Naghibi et al. (2016), Rahmati et al. (2016), Naghibi et al. (2017b), and Naghibi et al. (2018). The lower elevation of the study area has the highest potentiality of groundwater due to hydraulic gradient and presence of low water table. The NDVI, which was the highest contributing factor for GWPM, found in the study of Naghibi et al. (2017a) and Chen et al. (2019), is too incorporated in the present study. It is also considered as significant for land use classification, where forest class contains high FR value (1.23). Distance from the stream was a dominant factor for GWPM in the research of Kordestani et al. (2019), which reflects in the current study. The parallel drainage pattern indicates the presence of a fault system (Deffontaines and Chorowicz, 1991) in the study area, which controls the groundwater movement and storage. The soil has the lowest controlling factor in GWPM because more than 30% of the area consists of hard rock, and which has negligible for primary porosity. Fractured rock, weathered basement, and depth of soil are favorable for the groundwater occurrence and movement (Prasad et al., 2008; Maiti et al., 2012; Das, 2017). In the study area, the groundwater occurrence is mainly controlled by secondary porosity (weathered and fracture rocks) rather than the primary porosity of soil.

In the other side, the FR model determines the importance of each sub-class of the groundwater affecting variables. The maximum values of the FR ratio were 1.78 (Ultic haplustalfts), 1.69 (0.25-0.34), 1.55 (Denudational origin-pediment pediplain complex), and 1.53 (Kalladgi), for soil, NDVI,

geomorphology, and geology factors, respectively. The importance of geomorphology and NDVI coincide with the result of variable importance index. In case of soil, the results of the variable importance index and FR were in contrast due to the consideration of a particular soil class that was conducive for high groundwater potentiality. The geological aspects (lithology, lineament) also influence the distribution and occurrence of groundwater (Berhanu and Hatiye, 2020). The Kaladgi supergroup mainly consists of sedimentary rocks having more porosity in comparison to hard rocks. In case of hard rock, groundwater occurrence is controlled by fault, boundaries between the different lithological units, weathered, and fracture zones. However, the properties of the study area and adopted methods have influenced the effective factors for GWPM.

## **5.2 Interpretation of the models and GPR profiles**

In different studies (Naghibi et al., 2017b; Golkarian et al., 2018) on GWPM, the RF model provides an excellent result. The results of the machine learning models from AUC exhibit that the RF is the best-fit model for the current study. The better performance of the RF model (AUC=94%) may be because of the model consists of the multiple decision trees with no overfitting of the data. Besides, the model provides the interaction ability between effective factors and non-linearity (Catani et al., 2013). In recent years, ensemble models are increasingly used in groundwater potential mapping with very high accuracy (Chen et al., 2019; Kordestani et al., 2019; Naghibi et al., 2019). It was observed that the RF-SVM model had performed well with 93.40% accuracy from AUC (Fig. 6). The advantages of RF and SVM models are the probable reason for high performance in the present work. The result of the BRT model has lesser accuracy compared to RF and ensemble models that may be due to the overfitting of the data. However, since the AUC values of the models are more than 0.7, the GWPM has been reliable for the study region (Naghibi et al., 2016; Golkarian et al., 2018). The RF and RF-SVM models were successfully applied in the other selected parts of India. In the b1 and b4 areas, the RF and RF-SVM models performed better than the BRT model (Fig. 6). On the other hand, the results from the b2 and b3 regions almost matched the result of the study area due to the homogenous hydro-geologic, topographic, and climatic properties. Moreover, the RF and hybrid models of RF-SVM are promising and sufficient to be advised as the method to prepare groundwater potential maps at the regional scale.

In many research works (Annan et al., 1991; Nakashima et al., 2001; Bano, 2006; Mahmoudzadeh et al., 2012; Manu and Preko, 2014), GPR was successfully used to identify the groundwater table for the better understanding of groundwater condition of an area. In this context, GPR technology was used to detect the groundwater table and advise the location for a new well in the research area. The strong radar reflection and amplitude variation from the groundwater table suggests the different

dielectric contrast of the earth materials (Shih et al., 1986; Doolittle et al., 2006; Manu and Preko, 2014). The results from GPR profiles revealed that radar reflection, amplitude variation, and high attenuation have prominent signatures of the water table. Groundwater table of the profile a, b, c, e, and f are precisely matched with the nearest water level data of the wells (Table 5). The water table from GPR profile 'd' was not identified due to the presence of different layers interpreted from the multiple radar reflections (Fig. 7).

### **5.3 Precision of the GWPM**

A better model defines that the area of high and very high classes of the models are precisely matched with minimal variation (Naghbi et al., 2017b). In this research, the outcomes of the five classes of groundwater potential area are consistent with the lowest percentage of variation for all the models. The validation results from the ROC suggest high preciseness of this model for GWPM. The cross-validation of the RF and RF-SVM models in all the selected sites has ensured the applicability of the models with good precision. For water table detection, the GPR has produced an excellent result which in turn validates the groundwater mapping in the area of interest.

### **6. Conclusion**

The increasing demand for groundwater made concern for groundwater potential mapping, especially on the west coast of India. In this research, machine learning algorithms have been applied to demarcate the groundwater potential zones in Sindhudurg coastal sector of the west coast of India. Based on literature review and field knowledge, fifteen groundwater-related thematic layers were superimposed with inventory location in the GIS environment and integrated with RF, BRT, and ensemble of RF-SVM models. According to the results from RF, BRT, and RF-SVM, the very high groundwater potential zone occupies 33.31%, 35.60%, and, 36.85%, respectively, of the research area. The prediction of the models was validated with AUROC curve. Based on AUROC curves, RF and RF-SVM models exhibited better performance than BRT model for GWPM in the research area. Likewise, in the other four selected regions, RF and RF-SVM models proved to be superior in comparison to BRT model. The most influencing factors of the groundwater prospect mapping were geomorphology, elevation, NDVI, distance from streams, and distance from lineament. The predicted depth of groundwater from GPR profiles and measured data during fieldwork on those wells are corroborating to each other, which helps to identify new potential well in the study region using GPR technology. The obtained results of the present study can be useful for government and private agencies in groundwater resource management, land use planning, and environmental protection in the study

region. Furthermore, the methodology of this research can be adopted to study other coasts and watersheds with more or less similar hydro-geologic, topographic, and climatic properties.

### **Acknowledgement**

We acknowledge the financial support from the Indian Rare Earth Limited (Project no. SSP-3232) and University Grant Commission (3160 NET-June 2015). Authors are thankful to the Director CSIR-NIO for support. The authors are also grateful to Prof. Jungho Im and anonymous reviewers for their critical comments and constructive suggestions to improve the manuscript. Field support from the survey team members KM Dubey, Lalit Arya, and Trilochan Kumar is thankfully acknowledged. The NIO contribution number is.....

**Conflict of Interest:** The authors declare no conflict of interest.

### **References**

- Annan, A.P., 2003. Ground penetrating radar principles, procedures, and applications. Sensors & Software Inc, Mississauga, ON, Canada.
- Annan, A.P., Cosway, S.W., Redman, J.D., 1991. Water table detection with ground-penetrating radar. Society of Exploration Geophysicists, 494-496. <https://doi.org/10.1190/1.1888793>
- Arabameri, A., Rezaei, K., Cerda, A., Lombardo, L., Rodrigo-Comino, J., 2019. GIS-based groundwater potential mapping in Shahroud plain, Iran. A comparison among statistical (bivariate and multivariate), data mining and MCDM approaches. Science of the Total Environment, 658, 160-177. <https://doi.org/10.1016/j.scitotenv.2018.12.115>
- Bano, M., 2006. Effects of the transition zone above a water table on the reflection of GPR waves. Geophysical Research Letters, 33. <https://doi.org/10.1029/2006GL026158>
- Berhanu, K.G. and Hatiye, S.D., 2020. Identification of Groundwater Potential Zones Using Proxy Data: Case study of Megech Watershed, Ethiopia. Journal of Hydrology: Regional Studies, 28, 100676. <https://doi.org/10.1016/j.ejrh.2020.100676>
- Billy, J., Robin, N., Hein, C.J., Certain, R., FitzGerald, D.M., 2014. Internal architecture of mixed sand-and-gravel beach ridges: Miquelon-Langlade Barrier, NW Atlantic. Marine Geology, 357, 53-71. <https://doi.org/10.1016/j.margeo.2014.07.011>
- Bonham-Carter, G.F., 1994. Geographic Information Systems for Geoscientists: Modelling with GIS, Elsevier, UK.
- Breiman, L., 2001. Random forests. Machine Learning, 45(1), 5-32. <https://doi.org/10.1023/A:1010933404324>
- Catani, F., Lagomarsino, D., Segoni, S., Tofani, V., 2013. Landslide susceptibility estimation by random forests technique: sensitivity and scaling issues. Natural Hazards and Earth System Sciences, 13(11), 2815-2831. <https://doi.org/10.5194/nhess-13-2815-2013>, 2013.



- Chen, W., Panahi, M., Khosravi, K., Pourghasemi, H.R., Rezaie, F., Parvinnezhad, D., 2019. Spatial prediction of groundwater potentiality using ANFIS ensembled with teaching-learning-based and biogeography-based optimization. *Journal of Hydrology*, 572, 435-448. <https://doi.org/10.1016/j.jhydrol.2019.03.013>
- Chen, W., Chen, W., Li, H., Hou, E., Wang, S., Wang, G., Peng, T., 2018. GIS-based groundwater potential analysis using novel ensemble weights-of-evidence with logistic regression and functional tree models. *Science of the Total Environment*, 634, 853–867. <https://doi.org/10.1016/j.scitotenv.2018.04.055>
- Crisci, C., Ghattas, B., Perera, G., 2012. A review of supervised machine learning algorithms and their applications to ecological data. *Ecological Modelling*, 240, 113-122. <https://doi.org/10.1016/j.ecolmodel.2012.03.001>
- Das, S., 2019. Comparison among influencing factor, frequency ratio, and analytical hierarchy process techniques for groundwater potential zonation in Vaitarna basin, Maharashtra, India. *Groundwater for Sustainable Development*, 8, 617-629. <https://doi.org/10.1016/j.gsd.2019.03.003>
- Das, S., 2017. Delineation of groundwater potential zone in hard rock terrain in Gangajalghati block, Bankura district, India using remote sensing and GIS techniques. *Modeling Earth Systems and Environment*, 3, 1589-1599. <https://doi.org/10.1007/s40808-017-0396-7>
- Deepika, B., Avinash, K., Jayappa, K.S., 2013. Integration of hydrological factors and demarcation of groundwater prospect zones: insights from remote sensing and GIS techniques. *Environmental Earth Sciences*, 70(3), 1319-1338. <https://doi.org/10.1007/s12665-013-2218-1>
- Deffontaines, B. and Chorowicz, J., 1991. Principles of drainage basin analysis from multisource data: application to the structural analysis of the Zaire Basin. *Tectonophysics*, 194(3), 237-263.
- Doolittle, J.A., Jenkinson, B., Hopkins, D., Ulmer, M., Tuttle, W., 2006. Hydropedological investigations with ground-penetrating radar (GPR): Estimating water-table depths and local ground-water flow pattern in areas of coarse-textured soils. *Geoderma*, 131(3-4), 317-329. <https://doi.org/10.1016/j.geoderma.2005.03.027>
- Duan, H., Deng, Z., Deng, F., Wang, D., 2016. Assessment of groundwater potential based on multicriteria decision making model and decision tree algorithms. *Mathematical Problems in Engineering*. <http://dx.doi.org/10.1155/2016/2064575>
- Egan, J.P., 1975. *Signal detection theory and ROC-analysis*. Academic Press, New York.
- Elith, J., Leathwick, J.R., Hastie, T., 2008. A working guide to boosted regression trees. *Journal of Animal Ecology*, 77(4), 802-813. <https://doi.org/10.1111/j.1365-2656.2008.01390.x>. ESRI 2016. ArcGIS for desktop. <http://desktop.arcgis.com>
- Friedl, M.A., Brodley, C.E., Strahler, A.H., 1999. Maximizing land cover classification accuracies produced by decision trees at continental to global scales. *IEEE Transactions on Geoscience and Remote Sensing*, 37(2), 969-977. DOI: 10.1109/36.752215
- Ganapuram, S., Kumar, G.V., Krishna, I.M., Kahya, E., Demirel, M.C., 2009. Mapping of groundwater potential zones in the Musi basin using remote sensing data and GIS. *Advances in Engineering Software*, 40(7), 506-518. <https://doi.org/10.1016/j.advengsoft.2008.10.001>

- Gayen, A., Pourghasemi, H.R., Saha, S., Keesstra, S., Bai, S., 2019. Gully erosion susceptibility assessment and management of hazard-prone areas in India using different machine learning algorithms. *Science of the Total Environment*, 668, 124-138. <https://doi.org/10.1016/j.scitotenv.2019.02.436>
- Gislason, P.O., Benediktsson, J.A., Sveinsson, J.R., 2006. Random forests for land cover classification. *Pattern Recognition Letters*, 27(4), 294-300. <https://doi.org/10.1016/j.patrec.2005.08.011>
- Golkarian, A., Naghibi, S.A., Kalantar, B., Pradhan, B., 2018. Groundwater potential mapping using C5.0, random forest, and multivariate adaptive regression spline models in GIS. *Environmental Monitoring and Assessment*, 190, 149. <https://doi.org/10.1007/s10661-6507-8>
- Goudie, A.S., 2013. *Encyclopedia of Geomorphology*. Routledge, London
- Guru, B., Seshan, K., Bera, S., 2017. Frequency ratio model for groundwater potential mapping and its sustainable management in cold desert, India. *Journal of King Saud University-Science*, 29(3), 333-347. <https://doi.org/10.1016/j.jksus.2016.08.003>
- Hong, H., Pradhan, B., Bui, D.T., Xu, C., Youssef, A.M., Chen, W., 2017. Comparison of four kernel functions used in support vector machines for landslide susceptibility mapping: a case study at Suichuan area (China). *Geomatics, Natural Hazards and Risk*, 8(2), 544-569. <https://doi.org/10.1080/19475705.2016.1250112>
- Ibrahim-Bathis., and K. Ahmed, S.A., 2016. Geospatial technology for delineating groundwater potential zones in Doddahalla watershed of Chitradurga district, India. *The Egyptian Journal of Remote Sensing and Space Science*, 19(2), 223- 234. <https://doi.org/10.1016/j.ejrs.2016.06.002>
- Khosravi, K., Pham, B.T., Chapi, K., Shirzadi, A., Shahabi, H., Revhaug, I., Prakash, I., Bui, D.T., 2018. A comparative assessment of decision trees algorithms for flash flood susceptibility modeling at Haraz watershed, northern Iran. *Science of the Total Environment*, 627, 744-755. <https://doi.org/10.1016/j.scitotenv.2018.01.266>
- Kim, J.C., Lee, S., Jung, H.S., Lee, S., 2018. Landslide susceptibility mapping using random forest and boosted tree models in Pyeong-Chang, Korea. *Geocarto International*, 33(9), 1000-1015. <https://doi.org/10.1080/10106049.2017.1323964>
- Kordestani, M.D., Naghibi, S.A., Hashemi, H., Ahmadi, K., Kalantar, B., Pradhan, B., 2019. Groundwater potential mapping using a novel data-mining ensemble model. *Hydrogeology Journal*, 27(1), 211-224. <https://doi.org/10.1007/s10040-018-1848-5>
- Kuhn, M., Wing, J., Weston, S., Williams, A., Keefer, C., Engelhardt, A., Cooper, T., Mayer, Z., Benesty, M., Lescarbeau, R., Ziem, A., Scrucca, L., Tang, Y., Candan, C., Hunt, T., 2018. Package “caret” Classification and Regression Training.
- Lee, S., Hong, S.M., Jung, H.S., 2018. GIS-based groundwater potential mapping using artificial neural network and support vector machine models: the case of Boryeong city in Korea. *Geocarto International*, 33(8), 847-861. <https://doi.org/10.1080/10106049.2017.1303091>

- Lee, S., Kim, Y.S., Oh, H.J., 2012. Application of a weights-of-evidence method and GIS to regional groundwater productivity potential mapping. *Journal of Environmental Management*, 96(1), 91-105. <https://doi.org/10.1016/j.jenvman.2011.09.016>
- Lee, S., and Pradhan, B., 2007. Landslide hazard mapping at Selangor, Malaysia using frequency ratio and logistic regression models. *Landslides*, 4(1), 33- 41. <https://doi.org/10.1007/s10346-006-0047-y>
- Lek, S., and Guegan, J.F., 1999. Artificial neural networks as a tool in ecological modelling, an introduction. *Ecological Modelling*, 120(2-3), 65-73. [https://doi.org/10.1016/S0304-3800\(99\)00092-7](https://doi.org/10.1016/S0304-3800(99)00092-7)
- Magesh, N.S., Chandrasekar, N., Soundranayagam, J.P., 2012. Delineation of groundwater potential zones in Theni district, Tamil Nadu, using remote sensing, GIS and MIF techniques. *Geoscience Frontiers*, 3(2), 189-196. <https://doi.org/10.1016/j.gsf.2011.10.007>
- Mahmoudzadeh, M.R., Francés, A.P., Lubczynski, M., Lambot, S., 2012. Using ground penetrating radar to investigate the water table depth in weathered granites-Sardon case study, Spain. *Journal of Applied Geophysics*, 79, 17-26. <https://doi.org/10.1016/j.jappgeo.2011.12.009>
- Maiti, S., Erram, V.C., Gupta, G. Tiwari, R.K., 2012. ANN based inversion of DC resistivity data for groundwater exploration in hard rock terrain of western Maharashtra (India). *Journal of Hydrology*, 464, 294-308. <https://doi.org/10.1016/j.jhydrol.2012.07.020>
- Manap, M.A., Nampak, H., Pradhan, B., Lee, S., Sulaiman, W.N.A., Ramli, M.F., 2014. Application of probabilistic-based frequency ratio model in groundwater potential mapping using remote sensing data and GIS. *Arabian Journal of Geosciences*, 7(2), 711-724. <https://doi.org/10.1007/s12517-012-0795-z>
- Manu, E., and Preko, K., 2014. Estimation of water table depths and local groundwater flow pattern using the Ground Penetrating Radar. *International Journal of Scientific and Research Publications*, 4(8).
- Marjanovic, M., Kovacevic, M., Bajat, B., Vozenilek, V., 2011. Landslide susceptibility assessment using SVM machine learning algorithm. *Engineering Geology*, 123(3), 225-234. <https://doi.org/10.1016/j.enggeo.2011.09.006>
- Mogaji, K.A., Lim, H.S., Abdullah, K., 2014. Regional prediction of groundwater potential mapping in a multifaceted geology terrain using GIS-based Dempster-Shafer model. *Arabian Journal of Geosciences*, 8(5), 3235-3258. <https://doi.org/10.1007/s12517-014-1391-1>
- Mojaddadi, H., Pradhan, B., Nampak, H., Ahmad, N., Ghazali, A.H.B., 2017. Ensemble machine-learning-based geospatial approach for flood risk assessment using multi-sensor remote-sensing data and GIS. *Geomatics, Natural Hazards and Risk*, 8(2), 1080-1102. <https://doi.org/10.1080/19475705.2017.1294113>
- Moore, I.D., Grayson, R.B., Ladson, A.R., 1991. Digital terrain modelling: a review of hydrological, geomorphological, and biological applications. *Hydrological Processes*, 5(1), 3-30. <https://doi.org/10.1002/hyp.3360050103>

- Moore, I.D., and Burch, G.J., 1986. Physical basis of the length-slope factor in the universal soil loss equation 1. *Soil Science Society of America Journal*, 50(5), 1294-1298. [doi:10.2136/sssaj1986.03615995005000050042x](https://doi.org/10.2136/sssaj1986.03615995005000050042x)
- Murasingh, S., Jha, R., Adamala, S., 2018. Geospatial technique for delineation of groundwater potential zones in mine and dense forest area using weighted index overlay technique. *Groundwater for Sustainable Development*, 7, 387-399. <https://doi.org/10.1016/j.gsd.2017.12.001>
- Naghibi, S.A., Dolatkordestani, M., Rezaei, A., Amouzegari, P., Heravi, M.T., Kalantar, B., Pradhan, B., 2019. Application of rotation forest with decision trees as base classifier and a novel ensemble model in spatial modeling of groundwater potential. *Environmental Monitoring and Assessment*, 191, 248. <https://doi.org/10.1007/s10661-019-7362-y>
- Naghibi, S.A., Pourghasemi, H.R., Abbaspour, K., 2018. A comparison between ten advanced and soft computing models for groundwater qanat potential assessment in Iran using R and GIS. *Theoretical and Applied Climatology*, 131, 967-984. <https://doi.org/10.1007/s00704-016-2022-4>
- Naghibi, S.A., Moghaddam, D.D., Kalantar, B., Pradhan, B., Kisi, O., 2017a. A comparative assessment of GIS-based data mining models and a novel ensemble model in groundwater well potential mapping. *Journal of Hydrology*, 548, 471-483. <https://doi.org/10.1016/j.jhydrol.2017.03.020>
- Naghibi, S.A., Ahmadi, K., Daneshi, A., 2017b. Application of support vector machine, random forest, and genetic algorithm optimized random forest models in groundwater potential mapping. *Water Resources Management*, 31(9), 2761-2775. <https://doi.org/10.1007/s11269-017-1660-3>
- Naghibi, S.A., Pourghasemi, H.R., Dixon, B., 2016. GIS-based groundwater potential mapping using boosted regression tree, classification and regression tree, and random forest machine learning models in Iran. *Environmental Monitoring and Assessment*, 188, 44. <https://doi.org/10.1007/s10661-015-5049-6>
- Naghibi, S.A., Pourghasemi, H.R., Pourtaghi, Z.S., Rezaei, A., 2015. Groundwater qanat potential mapping using frequency ratio and shannon's entropy models in the Moghan watershed, Iran. *Earth Science Informatics*, 8(1), 171-186. <https://doi.org/10.1007/s12145-014-0145-7>
- Naghibi, S.A., and Pourghasemi, H.R., 2015. A comparative assessment between three machine learning models and their performance comparison by bivariate and multivariate statistical methods in groundwater potential mapping. *Water Resources Management*, 29(14), 5217-5236. <https://doi.org/10.1007/s11269-015-1114-8>
- Nakashima, Y., Zhou, H., Sato, M., 2001. Estimation of groundwater level by GPR in an area with multiple ambiguous reflections. *Journal of Applied Geophysics*, 47(3-4), 241-249. [https://doi.org/10.1016/S0926-9851\(01\)00068-4](https://doi.org/10.1016/S0926-9851(01)00068-4)
- Nampak, H., Pradhan, B., Manap, M.A., 2014. Application of GIS based data driven evidential belief function model to predict groundwater potential zonation. *Journal of Hydrology*, 513, 283-300. <https://doi.org/10.1016/j.jhydrol.2014.02.053>

- Neal, A., 2004. Ground-penetrating radar and its use in sedimentology: principles, problems and progress. *Earth-Science Reviews*, 66(3-4), 261-330. <https://doi.org/10.1016/j.earscirev.2004.01.004>
- Oh, H.J., Kim, Y.S., Choi, J.K., Park, E., Lee, S., 2011. GIS mapping of regional probabilistic groundwater potential in the area of Pohang City, Korea. *Journal of Hydrology*, 399(3-4), 158-172. <https://doi.org/10.1016/j.jhydrol.2010.12.027>
- Patra, S., Mishra, P., Mahapatra, S.C., 2018. Delineation of groundwater potential zone for sustainable development: A case study from Ganga Alluvial Plain covering Hooghly district of India using remote sensing, geographic information system and analytic hierarchy process. *Journal of Cleaner Production*, 172, 2485-2502. <https://doi.org/10.1016/j.jclepro.2017.11.161>
- Pourghasemi, H.R., Moradi, H.R., Aghda, S.F., 2013. Landslide susceptibility mapping by binary logistic regression, analytical hierarchy process, and statistical index models and assessment of their performances. *Natural Hazards*, 69(1), 749-779. <https://doi.org/10.1007/s11069-013-0728-5>
- Pourtaghi, Z.S., and Pourghasemi, H.R., 2014. GIS-based groundwater spring potential assessment and mapping in the Birjand Township, southern Khorasan Province, Iran. *Hydrogeology Journal*, 22(3), 643-662. <https://doi.org/10.1007/s10040-013-1089-6>
- Pradhan, B., and Lee, S., 2010. Regional landslide susceptibility analysis using back-propagation neural network model at Cameron Highland, Malaysia. *Landslides*, 7(1), 13-30. <https://doi.org/10.1007/s10346-009-0183-2>
- Prasad, R.K., Mondal, N.C., Banerjee, P., Nandakumar, M.V., Singh, V.S., 2008. Deciphering potential groundwater zone in hard rock through the application of GIS. *Environmental Geology*, 55(3), 467-475. <https://doi.org/10.1007/s00254-007-0992-3>
- Rahmati, O., Pourghasemi, H.R., Melesse, A.M., 2016. Application of GIS-based data driven random forest and maximum entropy models for groundwater potential mapping: a case study at Mehran Region, Iran. *Catena*, 137, 360-372. <https://doi.org/10.1016/j.catena.2015.10.010>
- Rahmati, O., Samani, A.N., Mahdavi, M., Pourghasemi, H.R., Zeinivand, H., 2015. Groundwater potential mapping at Kurdistan region of Iran using analytic hierarchy process and GIS. *Arabian Journal of Geosciences*, 8(9), 7059-7071. <https://doi.org/10.1007/s12517-014-1668-4>
- Rajaveni, S.P., Brindha, K., Elango, L., 2017. Geological and geomorphological controls on groundwater occurrence in a hard rock region. *Applied Water Science*, 7(3), 1377-1389. <https://doi.org/10.1007/s13201-015-0327-6>
- Recknagel, F., 2001. Applications of machine learning to ecological modelling. *Ecological Modelling*, 146(1-3), 303-310. [https://doi.org/10.1016/S0304-3800\(01\)00316-7](https://doi.org/10.1016/S0304-3800(01)00316-7)
- Regmi, N.R., Giardino, J.R., McDonald, E.V., Vitek, J.D., 2015. A review of mass movement processes and risk in the critical zone of Earth. *Developments in Earth Surface Processes*, 19, 319-362. <https://doi.org/10.1016/B978-0-444-63369-9.00011-2>

- Rodriguez-Galiano, V.F., and Chica-Rivas, M., 2014. Evaluation of different machine learning methods for land cover mapping of a Mediterranean area using multi-seasonal Landsat images and Digital Terrain Models. *International Journal of Digital Earth*, 7(6), 492-509. <https://doi.org/10.1080/17538947.2012.748848>
- Samui, P., 2008. Slope stability analysis: a support vector machine approach. *Environmental Geology*, 56(2), 255.
- Schapire, R.E., 2003. The boosting approach to machine learning: An overview. In *Nonlinear estimation and classification*. Springer, New York, 149-171. [https://doi.org/10.1007/978-0-387-21579-2\\_9](https://doi.org/10.1007/978-0-387-21579-2_9)
- Shafizadeh-Moghadam, H., Valavi, R., Shahabi, H., Chapi, K., Shirzadi, A., 2018. Novel forecasting approaches using combination of machine learning and statistical models for flood susceptibility mapping. *Journal of Environmental Management*, 217, 1-11. <https://doi.org/10.1016/j.jenvman.2018.03.089>
- Shih, S.F., Doolittle, J.A., Myhre, D.L., Schellentrager, G.W., 1986. Using radar for groundwater investigation. *Journal of Irrigation and Drainage Engineering*, 112(2), 110-118. <https://doi.org/10.1016/j.jenvman.2018.03.089>
- Tahmassebi, N., Rahmati, O., Noormohamadi, F., Lee, S., 2016. Spatial analysis of groundwater potential using weights-of-evidence and evidential belief function models and remote sensing. *Arabian Journal of Geosciences*, 9, 79. <https://doi.org/10.1007/s12517-015-2166-z>
- Tehrany, M.S., Pradhan, B., Mansor, S., Ahmad, N., 2015. Flood susceptibility assessment using GIS-based support vector machine model with different kernel types. *Catena*, 125, 91-101. <https://doi.org/10.1016/j.catena.2014.10.017>
- Tehrany, M.S., Pradhan, B., Jebur, M.N., 2014. Flood susceptibility mapping using a novel ensemble weights-of-evidence and support vector machine models in GIS. *Journal of Hydrology*, 512, 332-343. <https://doi.org/10.1016/j.jhydrol.2014.03.008>
- Yao, X., Tham, L.G., Dai, F.C., 2008. Landslide susceptibility mapping based on support vector machine: a case study on natural slopes of Hong Kong, China. *Geomorphology*, 101(4), 572-582. <https://doi.org/10.1016/j.geomorph.2008.02.011>
- Youssef, A.M., Pourghasemi, H.R., Pourtaghi, Z.S., Al-Katheeri, M.M., 2016. Landslide susceptibility mapping using random forest, boosted regression tree, classification and regression tree, and general linear models and comparison of their performance at Wadi Tayyah Basin, Asir Region, Saudi Arabia. *Landslides*, 13(5), 839-856. <https://doi.org/10.1007/s10346-015-0614>

### List of caption for Figures

- Figure 1** Location maps of the study area (a), and secondary regions (b1, b2, b3, b4)
- Figure 2** Flowchart of the present study for groundwater potential mapping
- Figure 3** Thematic layers of the groundwater affecting factors (a) elevation, (b) slope, (c) aspect, (d) LS, (e) profile curvature, (f) plan curvature, (g) TWI, (h) distance from rivers, (i) distance from lineaments, (j) lithology, (k) geomorphology, (l) soil, (m) landuse, (n) NDVI, (o) rainfall.
- Figure 4** Groundwater potential maps derived from (a) RF, (b) BRT, and (c) RF-SVM models
- Figure 5** Variables importance in groundwater potential mapping
- Figure 6** ROC curve of the models in study area (a), secondary areas (b1, b2, b3, b4)
- Figure 7** Water table depth from GPR profiles (a) 3.8 m (b) 3.4 m (c) 4.9 m (d) not detectable (e) 4.3 m and (f) 4.1 m

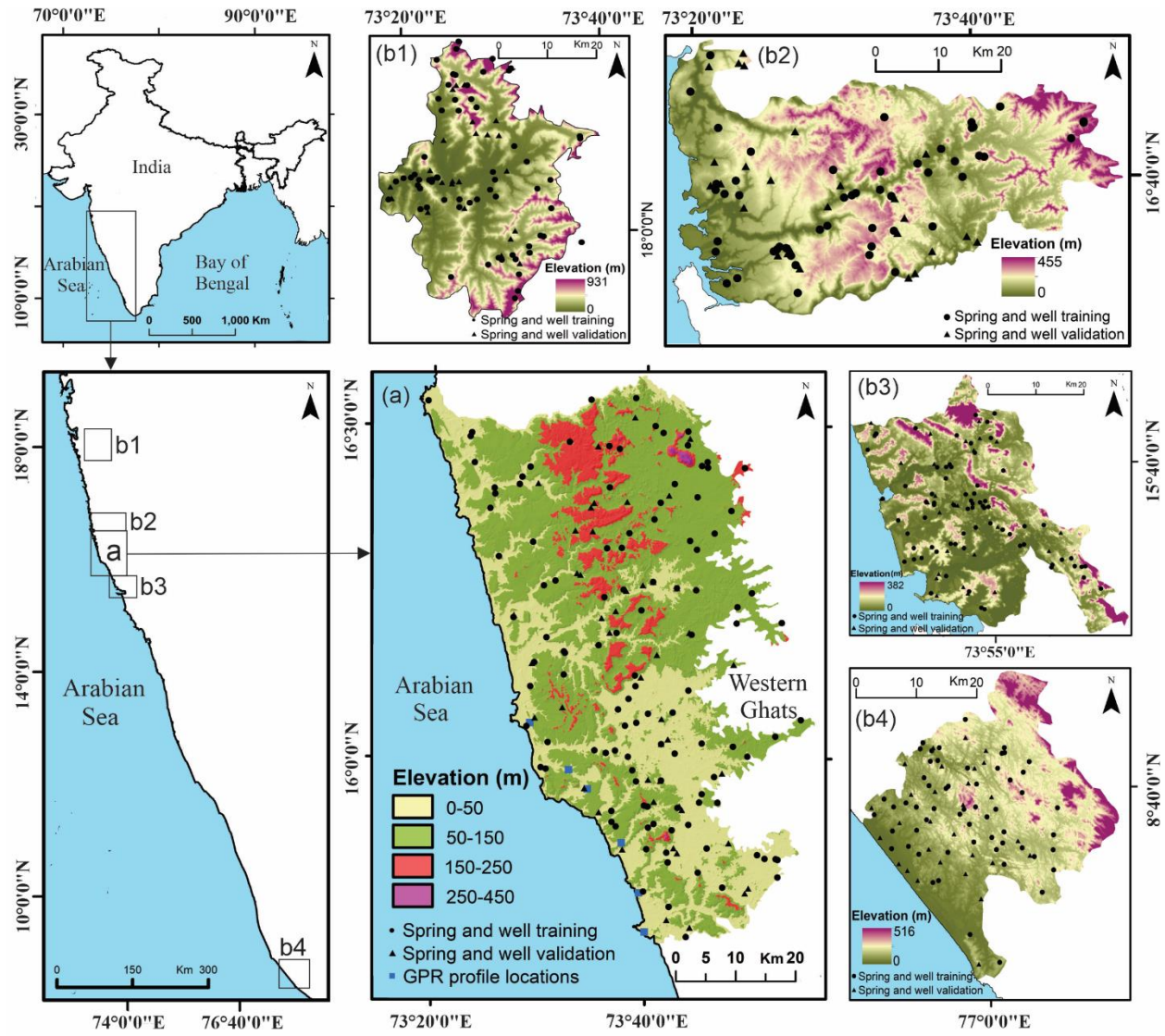


Fig. 1



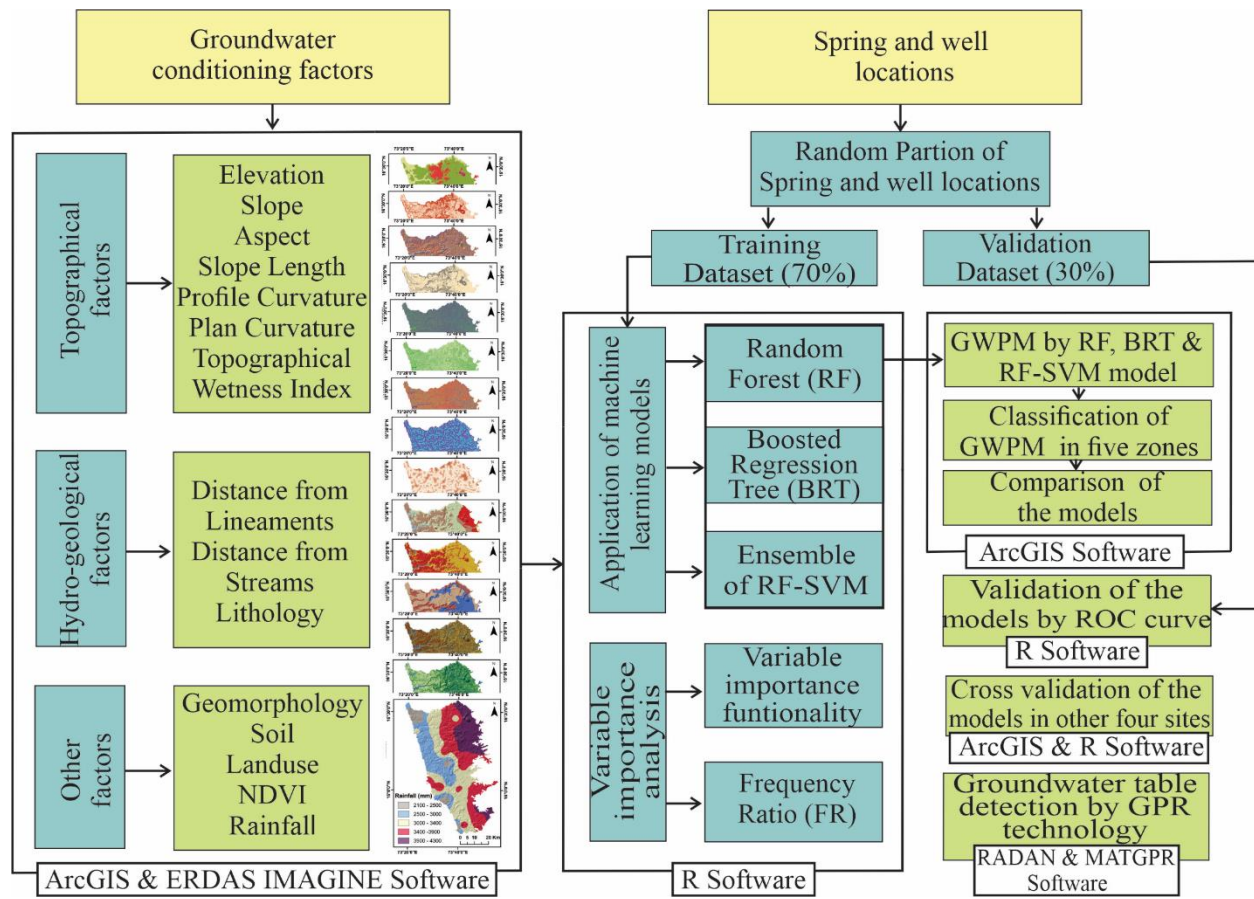
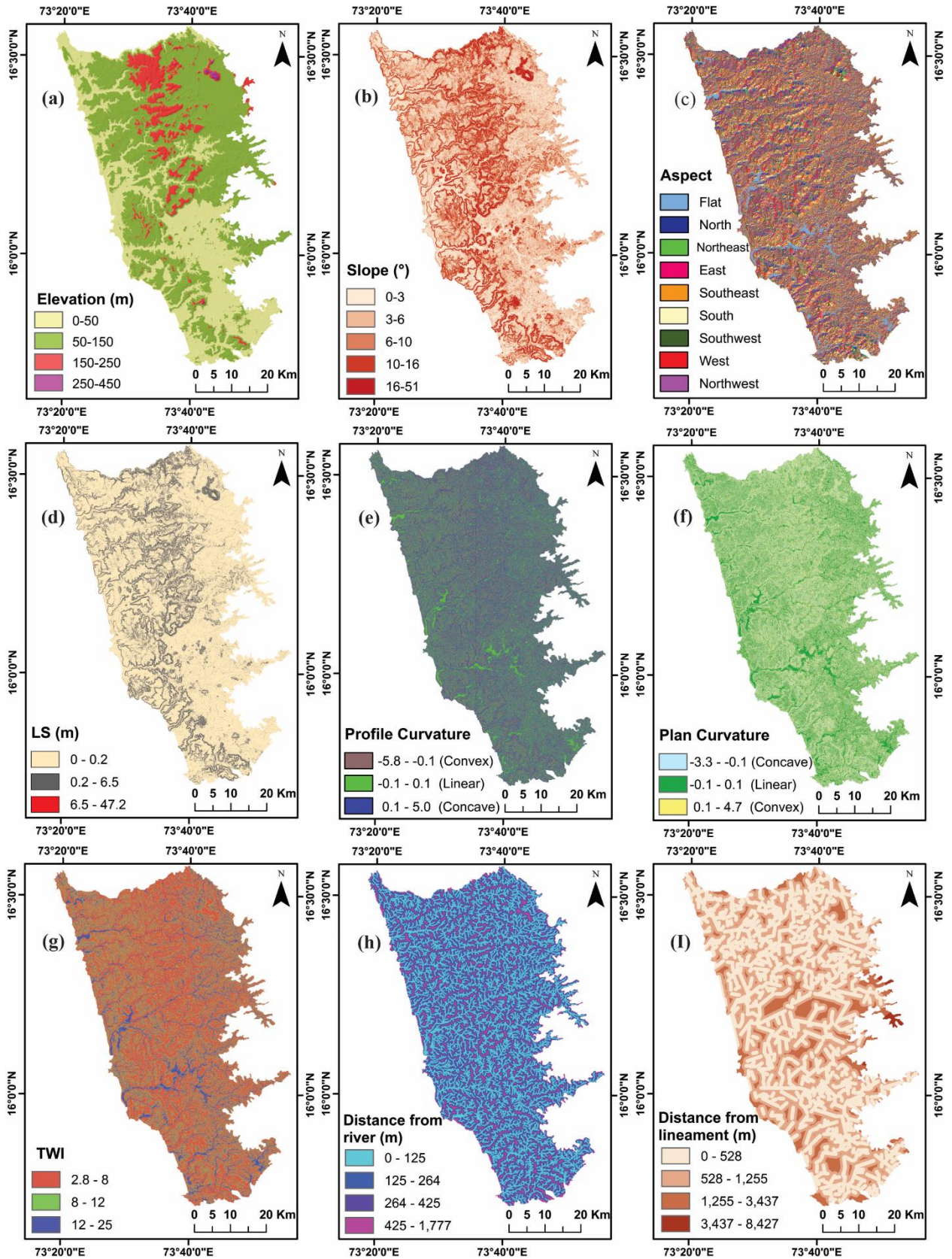


Fig. 2



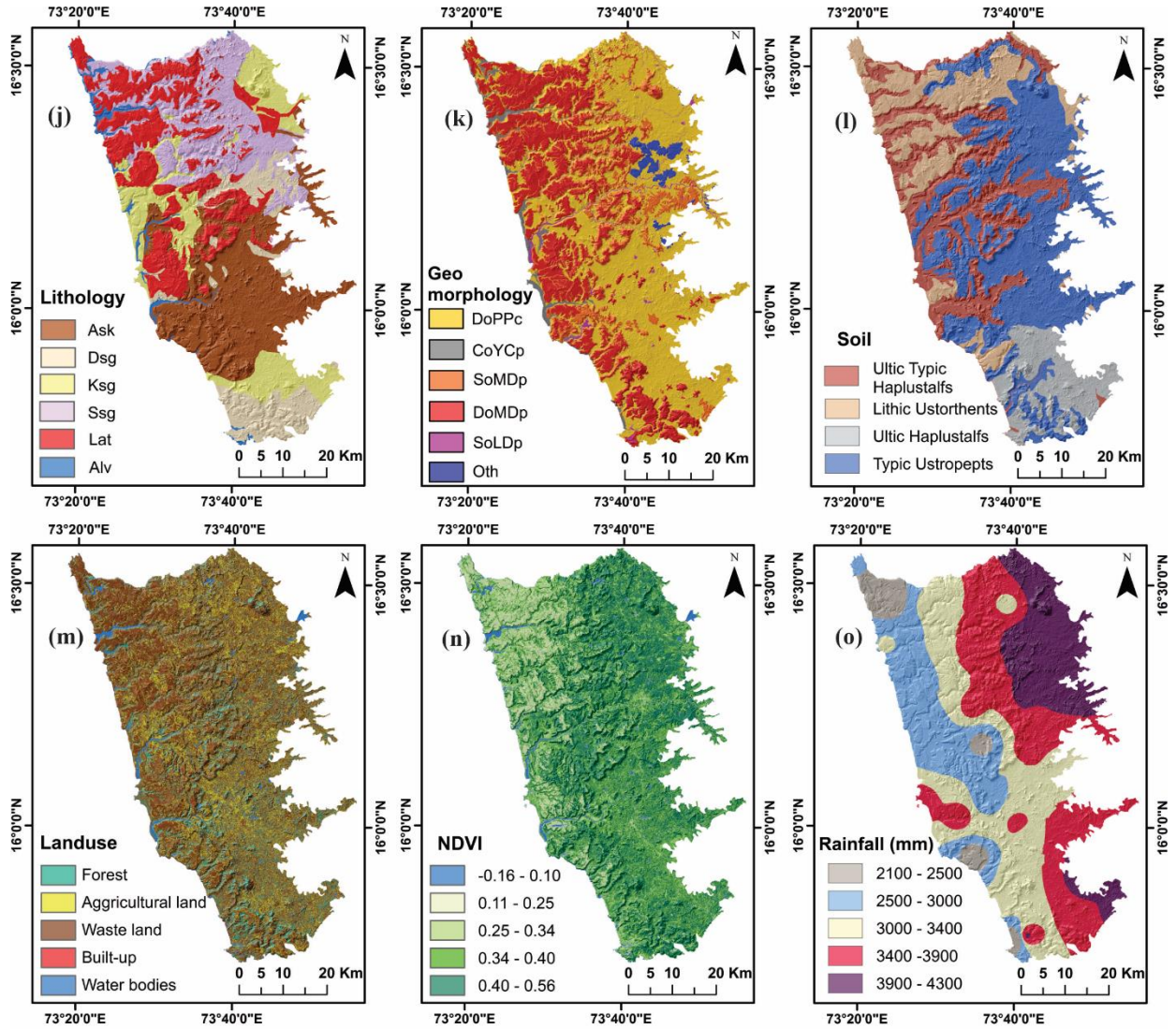


Fig. 3

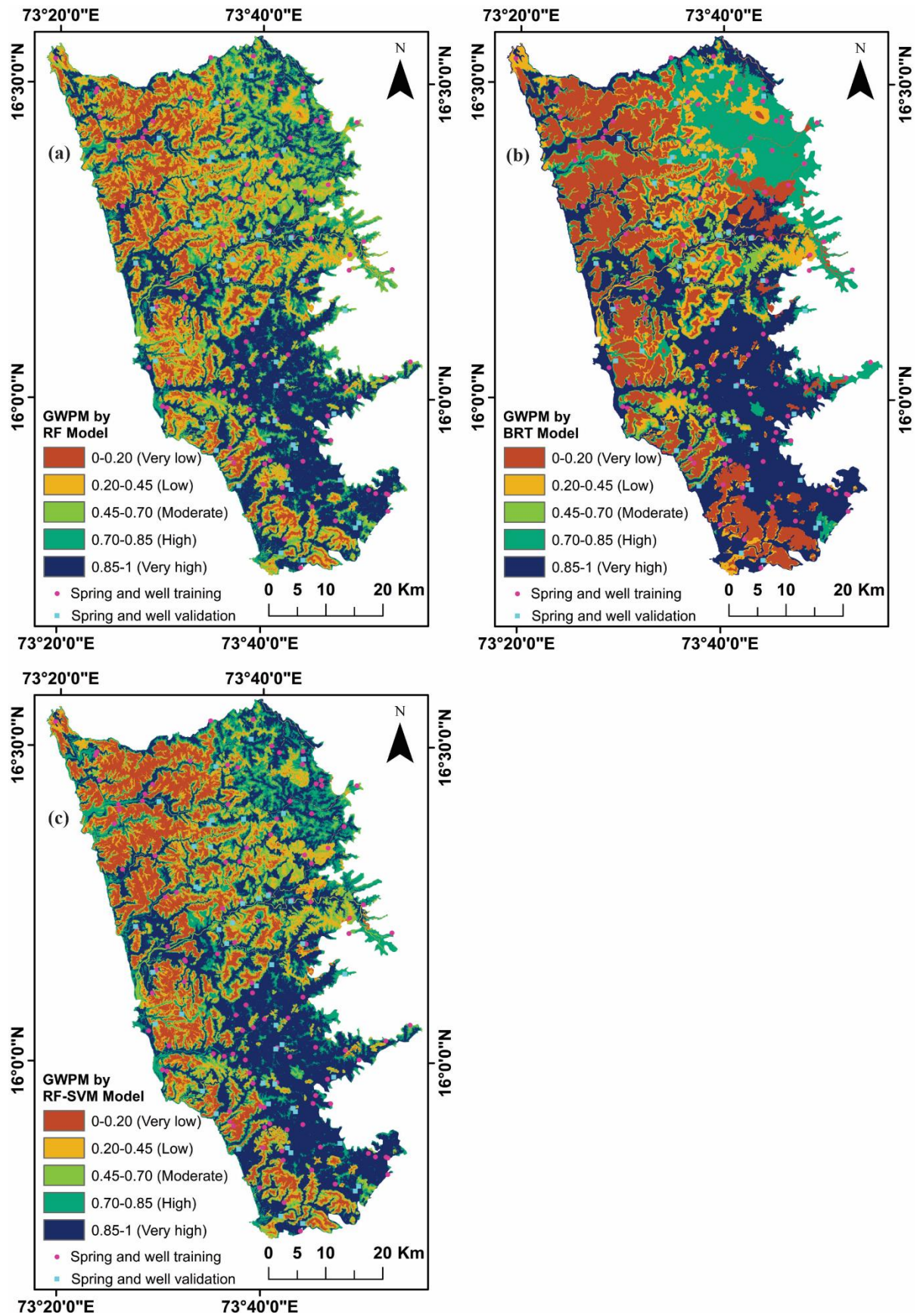
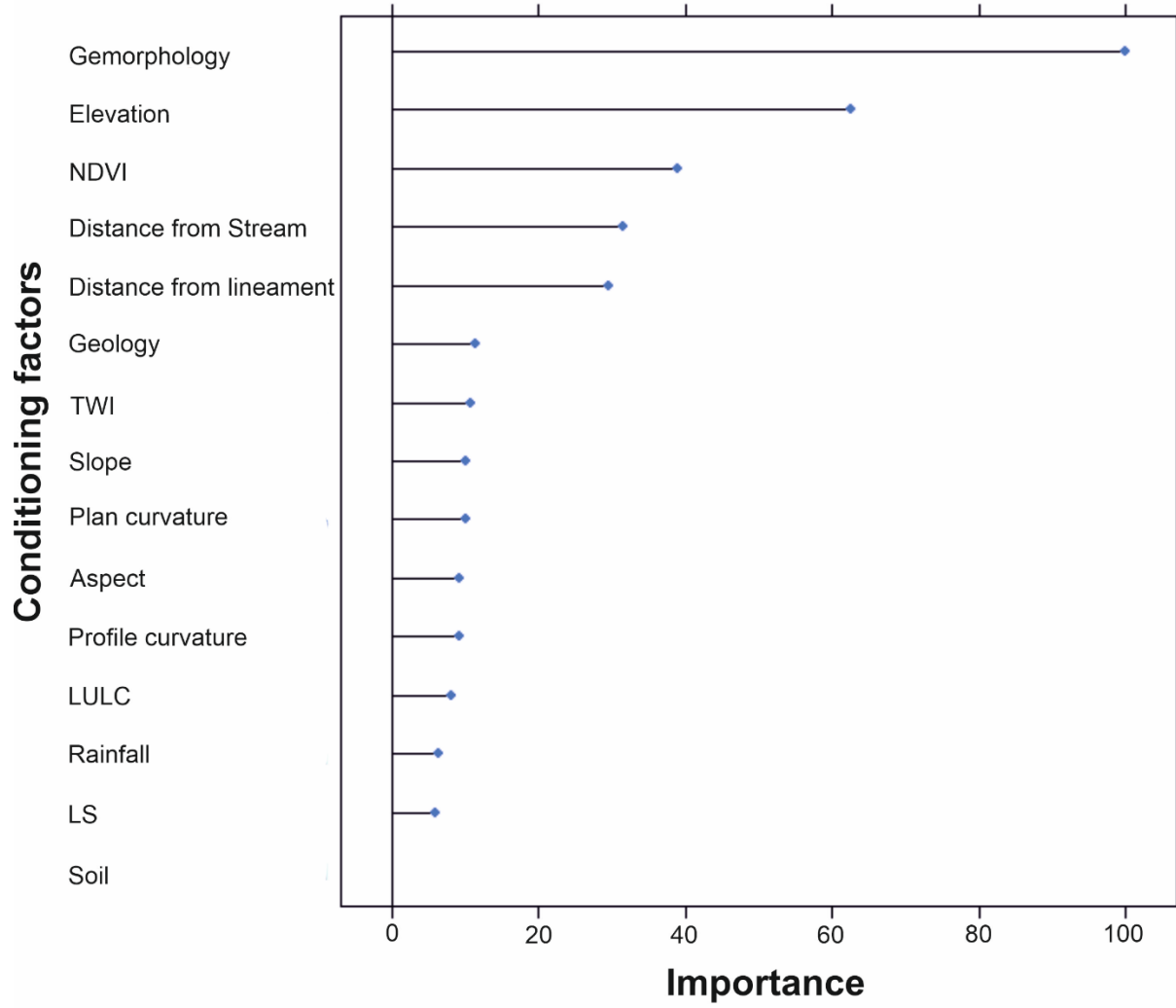
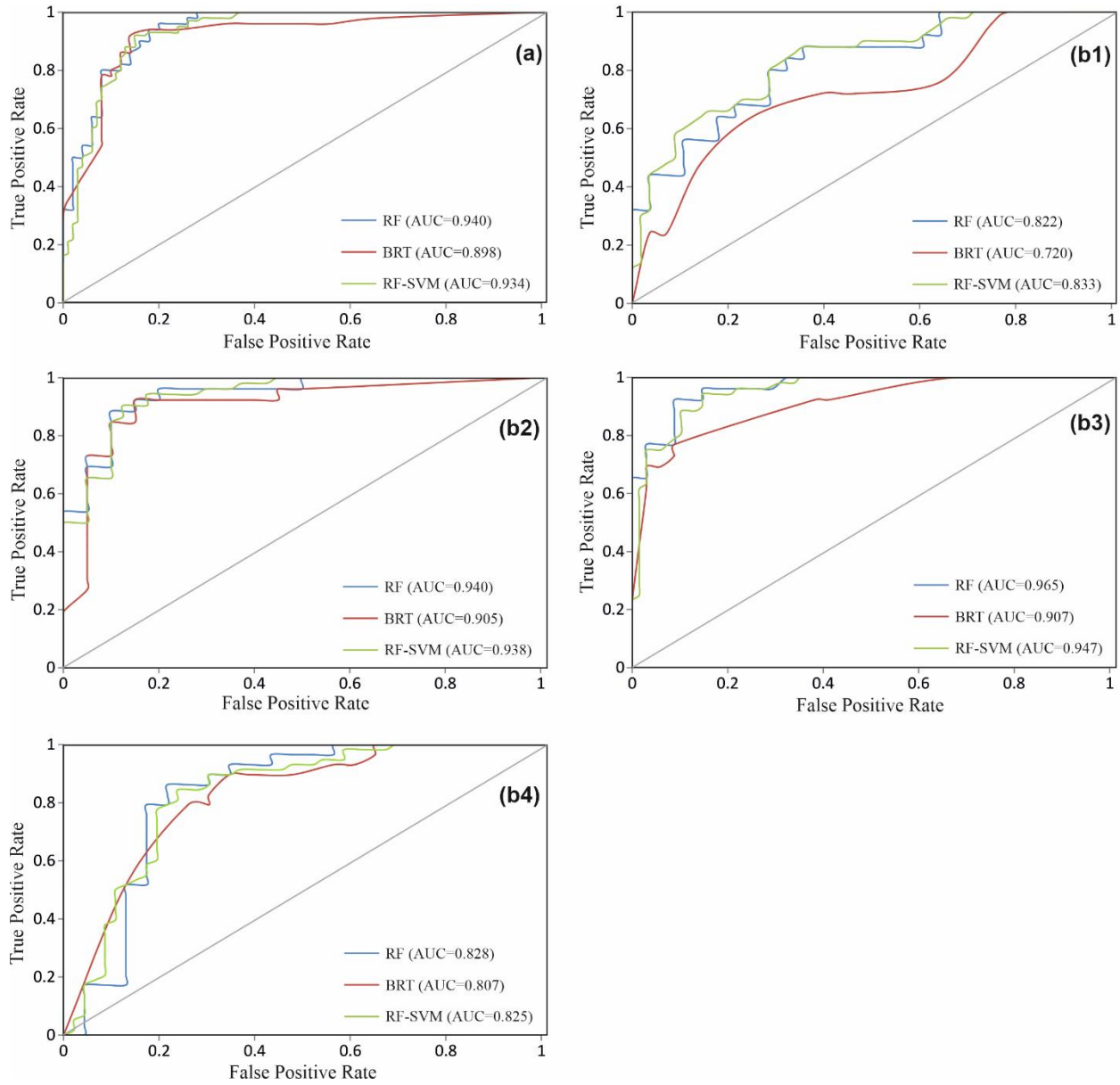


Fig. 4



**Fig. 5**



**Fig. 6**

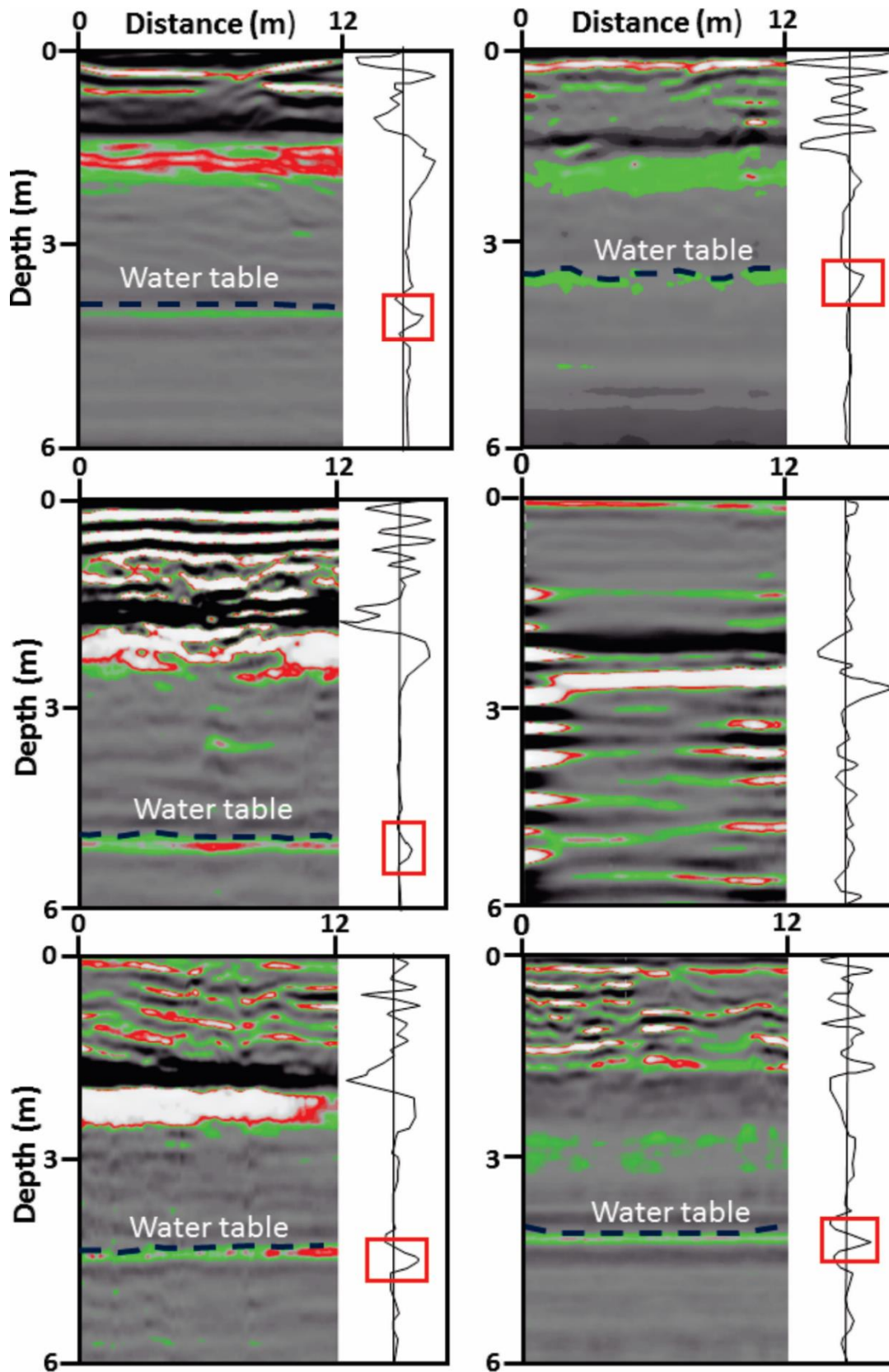
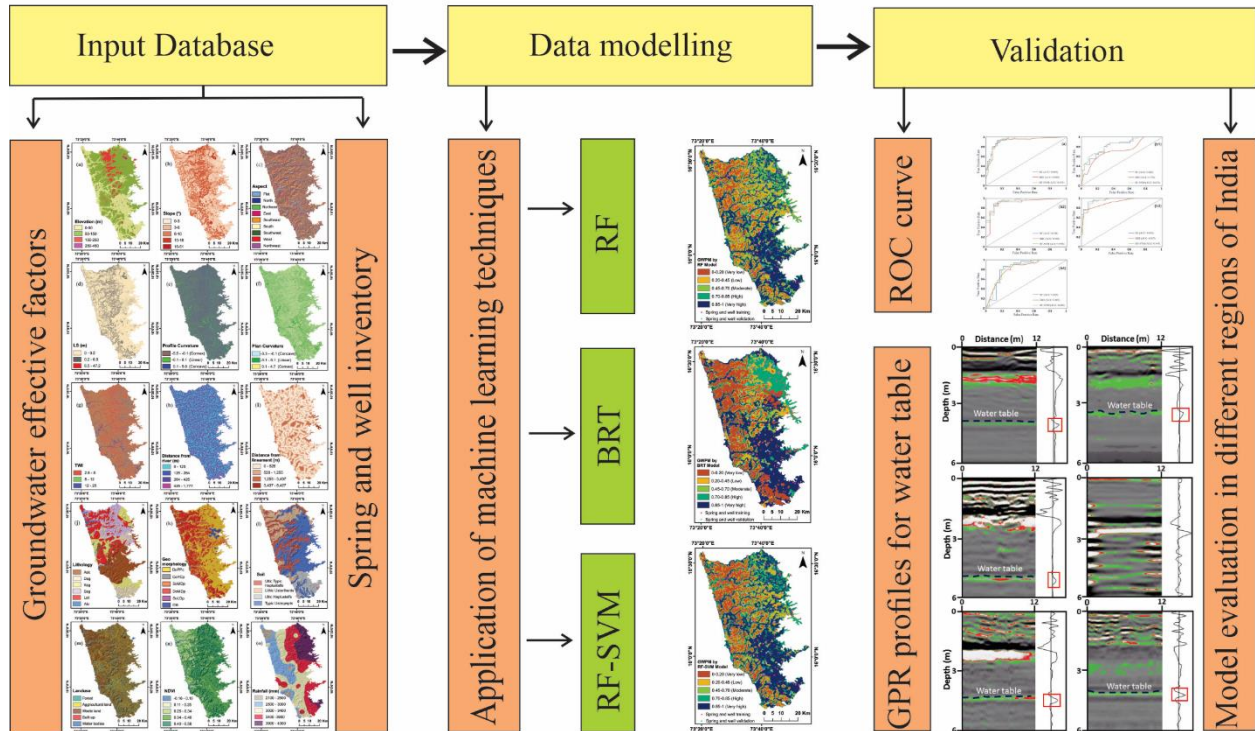


Fig. 7

# Groundwater Potential Mapping



Graphical abstract



**Table 1** Geographical description of selected sites for cross-validation

Regions	Areal extension		No of springs and wells	Area (sq Km)
	Latitude	Longitude		
South-east Raigad (b1)	17°50'57.75"- 18°19'4.39"	73°17'54.76"- 73°40'9.03"	87	1077
South Ratnagiri (b2)	16°29'58.34"- 16°49'44.78"	73°18'20.97"- 73°50'57.34"	76	1221
North Goa (b3)	15°24'22.51"- 15°47'52.15"	73°41'19.26"- 74°7'29.49"	99	933
South Kerala (b4)	8°23'21.62" - 8°50'13.49"	76°48'3.46" - 77°11'53.69"	87	1036

**Table 2** Details of database used in groundwater potential mapping

Data layers	Source of data	Scale	Time period
Spring and well locations	Topographical maps <a href="http://www.surveyofindia.gov.in/">http://www.surveyofindia.gov.in/</a>	1:50000 and 1:25000	1967
	Groundwater surveys and development agency <a href="https://gsda.maharashtra.gov.in">https://gsda.maharashtra.gov.in</a>		1990-2018
	Central Groundwater board <a href="http://cgwb.gov.in/">http://cgwb.gov.in/</a>		1990-2018
	Field survey		
Groundwater level	Groundwater surveys and development agency Central Groundwater board	--	2007-2016
Rainfall	Department of agriculture Maharashtra state <a href="http://maharain.gov.in/">http://maharain.gov.in/</a> India meteorological department <a href="https://mausam.imd.gov.in/">https://mausam.imd.gov.in/</a>	--	2013-2018
Geology	Geological Survey of India <a href="https://www.gsi.gov.in">https://www.gsi.gov.in</a>	1:250,000	2001
Geomorphology	National remote sensing centre <a href="https://bhuvan.nrsc.gov.in">https://bhuvan.nrsc.gov.in</a> SRTM Landsat8 OLI Field survey	1:50000	2005-2006
Soil map	National bureau of soil survey and land use planning <a href="https://www.nbsslup.in/">https://www.nbsslup.in/</a>	1:500,000	1996
Digital elevation model (DEM)	SRTM <a href="https://earthexplorer.usgs.gov/">https://earthexplorer.usgs.gov/</a>	1 arc second,	23rd September,2014
Satellite image	Landsat8 OLI <a href="https://earthexplorer.usgs.gov/">https://earthexplorer.usgs.gov/</a>	30 m spatial resolution	19th October 2016, 17th December, 2016.

**Table 3** Spatial relationship between each groundwater effective factor and inventory of spring and well using frequency ratio (FR) model

Parameters	Classes	% of total area (a)	% of inventory area (b)	Frequency ratio (b/a)
Elevation (m)	0-50	34.29	52.12	1.49
	50-150	56.36	43.64	0.78
	150-250	8.87	4.24	0.48
	250-450	0.48	0.00	0.00
Slope (Degree)	0-3	31.78	30.30	0.95
	3-6	30.86	36.36	1.18
	6-10	18.01	18.79	1.04
	10-16	10.96	7.88	0.72
Aspect	16-51	8.39	7.27	0.87
	Flat	4.82	2.42	0.51
	North	11.04	7.88	0.71
	Northeast	10.71	12.73	1.19
	East	10.70	10.30	0.96
	Southeast	11.59	12.12	1.05
	South	12.37	9.09	0.73
	Southwest	13.22	14.55	1.10
LS (m)	West	13.13	15.76	1.20
	Northwest	12.42	15.15	1.22
	0-0.2	88.80	84.24	1.01
	0.2-6.5	16.17	15.76	0.97
Profile Curvature	6.5-47.2	0.02	0.00	0.00
	Convex	29.31	20.61	0.70
	Linear	35.09	50.30	1.43
Plan Curvature	Concave	35.60	29.09	0.82
	Convex	24.56	16.36	0.67
	Linear	45.71	64.24	1.41
TWI	Convex	29.73	19.39	0.65
	2.8-8	63.26	55.15	0.87
	8-12	22.69	29.09	1.29
Distance from river (m)	12-25	14.05	15.76	1.10
	0-125	38.13	45.45	1.19
	125-264	30.30	28.48	0.94
	264-425	22.80	21.82	0.94
Distance from lineament (m)	425-1777	8.77	4.24	0.51
	0-528	52.52	61.82	1.18
	528-1255	37.79	33.94	0.90
	1255-3437	9.19	4.24	0.46
Lithology	3437-8427	0.49	0.00	0.00
	Ask	27.40	37.58	1.37
	Dsg	11.30	5.45	0.48
	Ksg	14.44	23.64	1.53
	Ssg	20.15	18.79	0.93
	Lat	23.38	13.33	0.57
Geomorphology	Alv	3.33	1.21	0.53
	DoPPc	49.76	76.36	1.55

	CoYCp	2.25	3.03	1.35
	SoMDp	18.65	13.33	0.71
	DoMDp	25.96	5.45	0.21
	SoLDp	2.03	1.82	0.89
	Oth	1.35	0.00	0.00
Soil	Ultic Typic Haplustalfs	25.27	33.94	1.34
	Lithic Ustorthents	18.15	10.91	0.60
	Ultic Haplustalfs	10.92	19.39	1.78
	Typic Ustropepts	45.64	35.76	0.78
Landuse	Forest	15.27	18.79	1.23
	Agricultural Land	26.23	39.39	1.50
	Wasteland	55.00	38.79	0.71
	Built-up	1.98	2.42	1.22
	Water bodies	1.90	0.61	0.32
NDVI	-0.16-0.10	1.84	1.21	0.67
	0.10-0.25	9.06	1.82	0.20
	0.25-0.34	16.10	27.27	1.69
	0.34-0.40	34.23	40.00	1.17
	0.40-0.56	38.77	29.70	0.77
Rainfall(mm)	2100-2500	4.94	45.45	0.99
	2500-3000	19.98	27.88	0.73
	3000-3400	29.14	22.42	1.08
	3400-3900	28.53	4.24	0.98
	3900-4300	17.41	45.45	1.22

---

**Table 4** Area of different classes (%) of groundwater potential maps using RF, BRT and RF-SVM ensemble models

Class	Area in Percentage		
	RF	BRT	RF-SVM
0-0.20 (Very low)	12.48	26.37	16.84
0.20-0.45 (Low)	17.98	12.93	13.99
0.45-0.70 (Moderate)	19.50	5.97	13.14
0.70-0.85 (High)	16.73	19.13	19.18
0.85-1 (Very high)	33.31	35.60	36.85
Total	100	100	100

**Table 5** Water table depths from GPR profiles and measured wells

GPR Profiles	Place name	Latitude (N)	Longitude (E)	Water level depth from the measured wells (m)	Water table depth from GPR profiles (m)	Deviation from well water level to GPR water table (m)
a	Dedoolwad a	16.04812	73.47902	3.6	3.8	0.2
b	Betwa	15.98158	73.54611	3.4	3.4	0
c	Mopar	15.95280	73.57527	4.9	4.9	0
d	Daboli	15.87222	73.62916	5.1	Not detectable	-
e	Vengurla	15.86111	73.63194	4.1	4.3	0.2
f	Redi	15.73811	73.66579	4.4	4.1	0.3

## Highlights

- RF and RF-SVM ensemble models performed very well with AUC>0.9.
- Evaluation of the models in four different regions with high precision results.
- Geomorphology is the most important variable in groundwater potential mapping.
- GPR technique successfully measured the groundwater table.